

TEHNICI DE RECUNOAȘTERE AUTOMATĂ ÎN ANALIZA ȘI DETERMINAREA SIMILARITĂȚII SECVENȚELOR BIOLOGICE

Teză destinată obținerii
titlului științific de doctor
la
Universitatea „Politehnica” din Timișoara
în domeniul AUTOMATICA
de către

Ing. Alina Bogan-Marta

Conducă

of.univ.dr. ing. Nicolae Budișan
prof.univ.dr. ing. Nicolae Robu
prof.univ.dr. ing. Mircea Petrescu
prof.univ.dr. ing. Mircea Ivănescu
prof.univ.dr. ing. Vladimir Crețu

Referenți științifici:

Ziua susținerii tezei: 19 ianuarie 2007

UNIV. "POLITEHNICA"	
TIMIȘOARA	
BIBLIOTECA CENTRALĂ	
Nr. volum	651. 36
Dulap _____	Lit. _____

Seriile Teze de doctorat ale UPT sunt:

- | | |
|------------------------|---|
| 1. Automatică | 7. Inginerie Electronică și Telecomunicații |
| 2. Chimie | 8. Inginerie Industrială |
| 3. Energetică | 9. Inginerie Mecanică |
| 4. Ingineria Chimică | 10. Știința Calculatoarelor |
| 5. Inginerie Civilă | 11. Știința și Ingineria Materialelor |
| 6. Inginerie Electrică | |

Universitatea „Politehnica” din Timișoara a inițiat seriile de mai sus în scopul diseminării expertizei, cunoștințelor și rezultatelor cercetărilor întreprinse în cadrul școlii doctorale a universității. Seriile conțin, potrivit H.B.Ex.S Nr. 14 / 14.07.2006, tezele de doctorat susținute în universitate începând cu 1 octombrie 2006.

Copyright © Editura Politehnica – Timișoara, 2006

Această publicație este supusă prevederilor legii dreptului de autor. Multiplicarea acestei publicații, în mod integral sau în parte, traducerea, tipărirea, reutilizarea ilustrațiilor, expunerea, radiodifuzarea, reproducerea pe microfilme sau în orice altă formă este permisă numai cu respectarea prevederilor Legii române a dreptului de autor în vigoare și permisiunea pentru utilizare obținută în scris din partea Universității „Politehnica” din Timișoara. Toate încălcările acestor drepturi vor fi penalizate potrivit Legii române a drepturilor de autor.

România, 300159 Timișoara, Bd. Republicii 9,
tel. 0256 403823, fax. 0256 403221
e-mail: editura@edipol.upt.ro

Cuvânt înainte

Teza de doctorat a fost elaborată pe parcursul activității mele în cadrul Departamentului de Calculatoare al Universității din Oradea în colaborare cu Departamentul de Automatică al Universității „Politehnica” din Timișoara, grupul de cercetare în „Recunoașterea vorbirii” al Universității Katolice din Leuven, Belgia și laboratorul de „Inteligență artificială și analiza informației” al Universității Aristotel din Tesalonic, Grecia.

Îi mulțumesc d-lui prof. Nicolae Budișan pentru înțelegere și sprijin de-a lungul acestor ani de activitate doctorală.

Adresez cele mai sincere mulțumiri d-lui prof. Nicolae Robu pentru răbdare, susținerea continuă și îndrumările pe parcursul acestui drum, și sub încurajările căruia am reușit să îmi conturez o direcție în activitatea desfășurată.

Le mulțumesc colegilor și supervisorilor din grupurile de cercetare din care am facut parte pentru ajutorul oferit pe timpul perioadelor de activitate desfășurate în cadrul acestora. De asemenea, îndrept mulțumiri și recunoștință tuturor colegilor și persoanelor care și-au oferit sprijinul în realizarea acestei teze.

Nu în ultimul rand le mulțumesc părinților pentru înțelegere, ajutor și dragostea cu care m-au susținut în realizarea acestui obiectiv.

Timișoara, ianuarie 2007

Alina Bogan-Marta

Părinților mei.

Bogan-Marta, Alina

Tehnici de recunoaștere automată în analiza și determinarea similarității secvențelor biologice

Teze de doctorat ale UPT, Seria 1, Nr. 2, Editura Politehnica, 2006/2007, 212 pagini, 27 figuri, 9 tabele.

ISSN: 1842-5208

ISBN: 978-973-625-415-4

Cuvinte cheie:

secvențe biologice, algoritmi de determinare a similarității, tehnici automate de recunoaștere, entropie, modele lingvistice Markov

Rezumat:

Multe cercetări în inteligență bioinformatică sunt focalizate pe algoritmi de învățare și sisteme integrate care să permită transformarea secvențelor biologice, a observațiilor și cunoștințelor în informație structurată și semnificativă pe care biologii o pot interoga, vizualiza și înțelege.

Această teză de doctorat contribuie la îmbogățirea metodelor de determinare automată a similarității secvențelor biologice cu propunerea unei noi abordări folosind o strategie statistică. Ea oferă o cale efectivă de capturare a caracteristicilor comune ale secvențelor comparate evitând dificultățile întâmpinate de utilizatori în practicile actuale. Performanța ridicată și eficiența computațională fac din această nouă abordare o alternativă promițătoare la algoritmii cunoscuți în prezent.

CUPRINS

ABREVIAR.....	9
LISTA DE TABELE.....	9
LISTA DE FIGURI.....	10
INTRODUCERE	11
1. SECVENȚE BIOLOGICE DIN PERSPECTIVA ABORDĂRII BIOINFORMATICE.....	13
1.1. DEFINIREA, ORGANIZAREA ȘI REPREZENTAREA SECVENȚELOR BIOLOGICE.....	13
1.1.1. Sinteza proteinelor	14
1.1.2. Clasificarea structurală a proteinelor	18
1.1.3. Formate de reprezentare ale secvențelor biologice	20
1.1.4. Baze de date de secvențe	23
1.2. ELEMENTE ÎN ANALIZA SECVENȚELOR BIOLOGICE.....	24
1.2.1. Omologia secvențelor în genetică.....	24
1.2.2. Motiv al unei secvențe.....	26
1.2.3. Aliniament.....	26
1.2.4. Similaritatea secvențelor biologice	30
1.3. CONCLUZII.....	30
2. STADIUL ACTUAL AL TEHNICILOR DE SIMILARITATE PENTRU SECVENȚELE BIOLOGICE.....	31
2.1. METODE DE DETERMINARE A SIMILARITĂȚII SECVENȚELOR BIOLOGICE	31
2.1.1. Distanță de editare și similaritate.....	31
2.1.2. Matrice de substituție.....	33
2.2. STAREA ACTUALĂ ÎN DEZVOLTAREA UNOR NOI METODE DE SIMILARITATE	39
2.3. TEHNICI DE REALIZARE A ALINIAMENTULUI DE SECVENȚE BIOLOGICE.....	43
2.3.1. Metode vizuale.....	45
2.3.2. Algoritmi pentru aliniamentul perechilor de secvențe	45
2.3.3. Metode fundamentale de programare dinamică	47
2.3.4. Extensii ale metodelor de bază	49
2.3.5. Algoritmul de programare dinamică a lui Ukkonen.....	50
2.3.6. Aliniamente aproape optimale	51
2.3.7. Metoda regiunilor	52
2.3.8. Localizarea segmentelor lungi de potrivire.....	54
2.3.9. Aliniamentul perechilor folosind modele Markov ascunse(HMMs)....	59
2.3.10. Aliniament multiplu de secvențe	67
2.4. CONCLUZII.....	69
3. APLICAREA TEHNICILOR DE RECUNOAȘTERE ÎN ANALIZA SECVENȚELOR BIOLOGICE	70
3.1. OBIECTIVE DE INTERES BIOLOGIC CE IMPLICĂ TEHNICILE INFORMATICE	70
3.2. IDENTIFICAREA GENELEOR/PROTEINELOR ȘI CLASIFICAREA ÎN CATEGORII	70
3.2.1. Clasificatorul celor mai apropiați k vecini (k-NN)	71
3.2.2. Clasificatori ce folosesc lanțuri sau modele Markov generalizate	71

3.2.3.	Clasificarea secvențelor pe baza caracteristicilor	72
3.2.4.	Rețele neuronale pentru clasificarea secvențelor	73
3.3.	COMPARAREA SECVENȚELOR.....	73
3.3.1.	Căutarea eficientă în baze de date	74
3.3.2.	Un algoritm vectorizat pentru segmente maxime	76
3.3.3.	Metoda regiunilor	76
3.3.4.	Algoritmi de programare dinamică	77
3.3.5.	Un algoritm bazat pe conceptul de modele Markov ascunse	78
3.4.	IDENTIFICAREA REGIUNILOR REGULATORII ȘI TIPARELOR NOI ÎN DATE SECVENȚIALE ...	80
3.4.1.	Căutarea exhaustivă	80
3.4.2.	Metode de învățare automată.....	81
3.4.3.	Metode bazate pe alinierea secvențelor.....	81
3.4.4.	Metode care se bazează pe conceptul de graf	82
3.4.5.	Metode hibride.....	82
3.5.	PRODUSE SOFTWARE DEZVOLTATE PENTRU ANALIZA SECVENȚELOR BIOLOGICE	83
3.6.	CONCLUZII.....	93

4. ANALIZA ȘI EVALUAREA CONTEXTULUI LINGVISTIC CU AJUTORUL MODELELOR STATISTICE PENTRU SISTEMELE DE RECUNOAȘTERE 95

4.1.	FUNDAMENT TEORETIC	95
4.1.1.	Lanțuri Markov lingvistice	95
4.1.2.	Construirea modelelor lingvistice statistice	96
4.1.3.	Măsuri de evaluare a modelelor lingvistice statistice	97
4.2.	IMPLEMENTAREA UNEI APLICAȚII PENTRU ANALIZA ȘI	100
	EVALUAREA MODELELOR LINGVISTICE.....	100
4.2.1.	Descrierea aplicației.....	100
4.2.2.	Procesarea textului.....	100
4.2.3.	Generatorul de modele lingvistice	102
4.2.4.	Estimarea entropiei și analiza rezultatelor	104
4.3.	CONCLUZII.....	111

5. CONTRIBUȚII LA UTILIZAREA TEHNICILOR DE RECUNOAȘTERE ÎN ANALIZA SECVENȚELOR BIOLOGICE 113

5.1.	DESCRIEREA INSTANȚEI DE REZOLVAT	113
5.2.	PROPUNEREA UNEI NOI METODE DE COMPARARE A PROTEINELOR PE BAZA EVALUĂRII MODELELOR LINGVISTICE MARKOV	114
5.2.1.	Noțiuni teoretice utilizate.....	114
5.2.2.	Descrierea metodei.....	115
5.2.3.	Descrierea modul de aplicare al noii măsuri de similaritate	118
5.2.4.	Implementare și detalii de performanță ale noii metode	121
5.3.	IDENTIFICAREA SECVENȚELOR MUTANTE ALE UNEI PROTEINE. EXPERIMENTE	123
5.3.1.	Baza de secvențe	123
5.3.2.	Descrierea experimentelor. Rezultate.	124
5.3.3.	Analiza și interpretarea rezultatelor.....	127
5.4.	DETERMINAREA SIMILARITĂȚII SECVENȚELOR DINTR-O BAZĂ STRUCTURATĂ DE PROTEINE (SCOP). EXPERIMENTE.	128
5.4.1.	Baza de secvențe	128
5.4.2.	Descrierea experimentelor.....	129
5.4.3.	Analiza și interpretarea rezultatelor.....	134
5.5.	CONCLUZII.....	134

6. STUDIU COMPARATIV AL ALGORITMULUI DE SIMILARITATE PROPUȘ.	136
6.1. ALGORITMUL CLUSTAL W	136
6.2. EXPERIMENTE.....	138
6.3. INTERPRETAREA REZULTATELOR.....	139
6.5. CONCLUZII.....	140
7. CONCLUZII FINALE	141
ANEXA 1.....	148
ANEXA 2.....	158
ANEXA 3.....	167
ANEXA 4.....	173
BIBLIOGRAFIE	192
INDEX.....	206

Abreviar

ADN- aciddeoxiribonucleic
NLA-Normalized Local Alignment
ROC- Receiver Operating Characteristic
DM-Data Mining
SSE-Secondary Structure Elements
PSI-Protein Structure Index
GO-Geene Ontology
DAG- Dyrected Acyclic Graph
USM-Universal Similarity Metric
SVM- Support Vector Machine
LSA-Latent Semantic Analysis
SVD-Singular Value Decomposition
kNN- *k* Nearest Neighbour
IMM-Interpolated Markov Models
SMM-Selective Markov Models
NNets-Neural Networks
BNN-Baiesian Neural Networks
LCC-Linear Corelation Coeficient

Lista de tabele

Tabel 4.2.4-1. Rezultatele obținute aplicând procedura de obținere a entropiei simple.....	108
Tabel 4.2.4-2. Evaluarea entropiei pentru cazul în care setul de.....	109
Tabel 5.3.2-1 Scorul ROC obținut pentru ambele metode pentru diferite valori de prag.	127
Tabel 5.4.1-1. Numărul de superfamilii și familii folosite în experimente.....	129
Tabel 5.4.1-2. Familiile SCOP incluse în experimente ^a	129
Tabel 5.4.2-1. Valorile preciziei	132
Tabel 5.4.2-2. Valorile coeficienților de corelație pentru Set1, 2 și 3 folosind statisticile Hubert.	133
Tabel 6.2-1. Valorile preciziei obținute pe baza rezultatelor de similaritate furnizate de programul CLUSTAL W sunt în coloana 'CLUST.W' urmate de coloanele corespunzătoare celor două metode noi de similaritate folosind modele 2,3,4-gram pentru celei trei seturi de date.....	139
Tabel 6.2-2. Valorile de corelație pentru set1, 2 și 3 ale bazei de secvențe obținute din Astral SCOP folosind statisticile Huberts	139

Lista de figuri

Figura 1.1-1. Parte a unui gel de secvențiere marcat radioactiv.....	14
Figura 1.1.1-1. O reprezentare sugestivă a structurii primare.....	15
Figura 1.1.1-2. Reprezentarea simbolică a structurii de alfa-helix.....	16
Figura 1.1.1-3. Reprezentarea simbolică a fâșiiilor plate beta.....	16
Figura 1.1.1- 4. Reprezintă structura 3D a unei proteine (Myoglobin).....	17
Figura 1.1.1- 5. Reprezentarea 3D a unei proteine conținând structuri mai complexe.....	17
Figura 1.1.1- 6. Reprezentare a structurii quaternare.....	18
Figura 1.1.2-1. Reprezentarea ierarhiei folosite în clasificarea proteinelor în SCOP.....	20
Figura 2.3.7-1. Arborele de poziție pentru a=AATAATGC.....	52
Figura 2.3.9-1. Diagrama unui automat cu stări finite pentru un aliniament cu gap.....	58
Figura 2.3.9-2. Modelul probabilistic corespunzător diagramei unui automat cu stări finite pentru un aliniament cu gap	58
Figura 2.3.9-3. Versiunea probabilistică completă a figurii 2.3.9-2	60
Figura 2.3.9.2-1.O pereche HMM pentru aliniament local.....	65
Figura 3.3.5-1. Un model Markov ascuns liniar și exemplu de aliniament.....	79
Figura 3.4.5-1. Reprezentarea relațiilor dintre secvențele biologice și metoda HMM, figură preluată și parțial tradusă din documentația detaliată, disponibilă la adresa de web a institutului Pasteur, Franța (http://bioweb.pasteur.fr/seqanal/motif/hmmer-uk.html).....	88
Figura 3.4.5-2: Rezultatul oferit de Clustal W în urma aplicării metodei dinamice pentru 7 secvențe de proteine din baza de secvențe SwissProt. Figura este reprodusa din lucrarea lui J.D.Thompson ©[THOM'94] care studiază îmbunătățirea metodei de aliniament multiplu de secvențe a lui Clustal W.....	91
Figura 4.2.4-1. Reprezentarea valorilor pentru estimarea entropiei simple.....	108
Figura 4.2.4-2. Reprezentarea evoluției valorilor entropiei când setul.....	109
Figura 4.2.4-3. Reprezentarea evenimentelor n-gram cu frecvență.....	110
Figura 4.2.4-4. Reprezentarea evoluției frecvenței minime și a evenimentelor.....	110
Figura 5.3.2-1. Vizualizarea matricelor conținând disimilaritățile tuturor perechilor determinate de cele 100 de proteine din baza de date utilizată.....	125
Figura 5.3.2-2. Reprezentarea în spațiu dimensional redus a secvențelor de proteine folosind măsura disimilarității pentru setul experimental de 100 secvențe ...	126
Figura 5.3.2-3. Reprezentarea evoluției valorilor raporturilor adevărate și pozitive	127
Figura 5.4.2-1. Vizualizarea matricilor care conțin toate disimilaritățile posibile pentru perechile de proteine ale Set1 format din 497 secvențe pentru modele 2,3,4 gram.....	131
Figura 5.4.2-2 Vizualizarea matricilor care conțin toate disimilaritățile posibile pentru perechile de proteine ale Set2 format din 497 secvențe pentru modele 2,3,4- gram	132
Figura 1.1.1-3. Vizualizarea matricelor care conțin toate disimilaritățile posibile pentru perechile de proteine ale Set3 format din 466 secvențe pentru modele 2,3,4-gram.....	129
Figura 5.4.3-1 Interfața oferită de programul Clustal W	138

INTRODUCERE

Cercetările comunității în inteligență bioinformatică sunt focalizate pe diverși algoritmi de învățare, sisteme de data mining și sisteme integrate care să permită transformarea secvențelor biologice, a observațiilor și cunoștințelor în informație structurată și semnificativă pe care biologii o pot interoga, vizualiza și înțelege. Dintre cele mai importante funcții computaționale identificate în acest domeniu de cercetare sunt: identificarea genelor/proteinelor și clasificarea lor în categorii; metode de comparare a secvențelor biologice la diferite nivele de detaliu; identificarea regiunilor regulatorii și a modelelor noi în secvența de date (căutarea regiunilor care se repetă, similarități și modele rare care au semnificație biologică). Pe baza acestor teme, oamenii de știință încearcă să dezvolte tehnici specializate și un mare volum al cercetării este orientat spre algoritmi destinați prelucrării informației conținute în depozite mari de secvențe. Analiza și compararea automată a secvențelor biologice a beneficiat în ultimul deceniu de un efort considerabil din partea cercetării informatice reușind să substituie cu succes multe din necesitățile impuse de extinderea rapidă a cantității de cunoștințe biologice stocate sub forma bazelor de date de secvențe.

Pentru căutarea în baze de secvențe cele mai frecvent folosite metode folosesc algoritmi de determinare a similarității secvențelor care adesea sunt reprezentați de proceduri anevoioase și costisitoare. Astfel, prin această teză de doctorat se contribuie la îmbogățirea metodelor de determinare automată a similarității secvențelor cu propunerea unei noi abordări. Prin urmare, utilizând o strategie statistică, se oferă o cale efectivă de capturare a caracteristicilor comune ale secvențelor comparate evitând misiunea mai puțin plăcută a utilizatorilor de a fi puși în dificultate în alegerea de parametri, funcții adiționale, matrice de substituție sau diverse metode de evaluare.

Metoda a fost inspirată de folosirea cu succes a conceptului de entropie pentru domeniul *information retrieval* (tradus liber „regăsirea informației”) în modelarea statistică a limbajului pentru sistemele automate de recunoaștere a vorbirii. Performanța ridicată și caracterul facil de implementare, împreună cu eficiența sa computațională fac din această abordare o alternativă promițătoare la algoritmi cunoscuți și sofisticăți de comparare și analiză a secvențelor.

Aparte de contextul biologic, noua metodă poate fi aplicată și pentru studiul general al secvențelor mari de text dar un aspect esențial îl va reprezenta limitarea impusă de mărimea alfabetului considerat. Astfel, ea își poate găsi deocamdată aplicabilitatea în domeniul biologic unde alfabetul este fie de 4 litere cand este vorba de nucleotidele care formează secvențe de ADN/ARN fie de 20 dacă sunt vizate secvențele de proteine, definite printr-un alfabet de 20 amino acizi.

Teza, deservind domeniul biologiei, are un evident caracter interdisciplinar: biologie, automatizări, matematică, informatică și calculatoare. Ea însă poate fi mai puțin atribuită doar domeniului biologiei, matematicii, informaticii sau calculatoarelor, fiind cel mai bine încadrată în sfera mai largă a automatizărilor. Această încadrare este motivată prin urmărirea realizării scopului domeniului căreia îi este dedicată fără intervenția umană („automatos” însemnând „funcționare de la sine” în limba greacă). Așa cum automatica circumscrie automatizarea celor mai diferite procese – tehnice, energetice, militare etc.-, tot așa și preocupările cuprinse

în teză – analiza automată a similarităților secvențelor biologice– urmăresc realizarea scopului în mod automat, ceea ce o poate înscrie în domeniul specializării de “Automatică”, servindu-se de suportul teoretic specific al acesteia.

Teza este structurată în șapte capitole după cum urmează:

În *capitolul 1* se face o introducere a cititorului în domeniul secvențelor biologice subliniind aspectele necesare înțelegerii modului de reprezentare și analiză al acestora. Prin această introducere se vine în sprijinul capitolelor următoare pentru a ilustra și a motiva dezvoltarea metodelor și algoritmilor pentru sisteme de calcul destinate analizei și comparării automate a secvențelor biologice.

Pe parcursul *capitolului 2* este realizată o investigație a metodelor de similaritate a secvențelor biologice deja cunoscute, urmată de o analiză a stării actuale în dezvoltarea de noi metode de similaritate. Deoarece aliniamentul secvențelor este cea mai folosită strategie în analiza și determinarea similarității sau disimilarității dintre secvențe sunt evidențiate cele mai semnificative metode folosite în prezent.

Capitolul 3 conține o multitudine de tehnici de recunoaștere care și-au găsit aplicabilitate directă în analiza secvențelor biologice. Ele au fost întâlnite în procesul de identificare a genelor/proteinelor și clasificarea lor în categorii, în compararea secvențelor, identificarea regiunilor regulatorii și tiparelor noi în date secvențiale, fiind valorificate în produse comerciale specializate.

Ca o continuitate în utilizarea tehnicilor de recunoaștere, în *capitolul 4* este realizată o investigație a potențialului lingvistic din perspectiva modelării statistice pentru sistemele de recunoaștere a vorbirii. Astfel, se descriu premisele teoretice de la care pornește investigația și analiza datelor reprezentate textual. De asemenea, se descrie modul de analiză și rezultatele obținute cu un set de programe experimentale implementate pentru acest scop. Aceste rezultate se vor constitui ca bază de pornire în promovarea ideii de recunoaștere a conținutului unui text pe baza evaluării cu ajutorul măsurilor derivate din conceptul de entropie.

În *capitolul 5* se propune o nouă metodă de evaluare a similarității secvențelor biologice inspirată din rezultatele experimentale descrise în capitolul anterior. Prin urmare, cu scopul identificării secvențelor similare din baze de secvențe biologice este propusă noua metodă de comparare. Ea este aplicată secvențelor de proteine făcând uz de principiul de evaluare al entropiei folosind modele Markov lingvistice. Pentru testarea și validarea acesteia au fost propuse două experimente ale căror rezultate sunt analizate și interpretate.

Capitolul 6 este destinat comparării rezultatelor obținute aplicând noua strategie de determinare a similarității secvențelor cu cele ale unei alte metode larg utilizate. Astfel, a fost ales programul ClustalW și vis-a-vis de rezultatele obținute este făcută o analiză comparativă.

Întreg *capitolul 7* este destinat concluziilor finale subliniind aspectele esențiale obținute cu noua metodă precum și partea de contribuții personale și posibile direcții de dezvoltare.

1. SECVENȚE BIOLOGICE DIN PERSPECTIVA ABORDĂRII BIOINFORMATICE

1.1. Definirea, organizarea și reprezentarea secvențelor biologice

Fața biologiei a fost schimbată de urgențele geneticii moleculare. Printre cele mai apreciate realizări în direcția susținerii progresului în domeniul biologic sunt eforturile secvențierii la scară largă a părților de ADN (cum ar fi proiectul "Human Genome Project") care produce o cantitate imensă de date. Nevoia de înțelegere a datelor a devenit însă și mai presantă. Cererile pentru analize sofisticate ale secvențelor biologice conduc înainte domeniul nou creat și într-o continuă expansiune al cercetării în biologie moleculară, computațională sau bioinformatică.

Pentru a avea o imagine generală asupra a ceea ce reprezintă secvențele biologice și modul lor de reprezentare și gestionare, în acest capitol sunt descrise noțiuni necesare înțelegerii domeniului abordat pe parcursul acestei lucrări.

Secvența de ADN codifică informația necesară pentru viețuitoare de a trăi și a se reproduce. Determinarea secvenței este prin urmare, utilă în cercetare pentru a ști de ce și cum trăiesc organismele. Datorită calității sale de *cheie a vieții*, cunoștințele despre ADN devin practic folositoare în orice domeniul de cercetare biologică.

În genetică și biochimie procesul de secvențiere presupune determinarea structurii primare (sau secvenței primare) a unui fragment de ADN. Mai exact, în terminologia genetică, secvențierea ADN este procesul prin care se determină ordinea nucleotidelor unui fragment de ADN dat. Astfel, se determină structura primară a secvenței respective. În mod curent, aproape toate secvențierile de ADN sunt executate utilizând tehnica dezvoltată de Frederick Sanger ce presupune utilizarea terminației secvenței specifice a unei reacții in vitro a sintezei ADN folosind substraturi modificate de nucleotide¹ [WIKI'05b].

Prin urmare, în urma secvențierii se obține o înșiruire de simboluri de forma :

```
1  ggatccccta cccagctggg acctcccca accccttct tcctcacct ctgcacgaga
61 acgaagtct tcgtcaccag ctttttgga aggatggagg aggggaggcg tacgtgaggg
...
```

unde, fiecare simbol reprezintă una dintre cele patru nucleotide, după cum este ilustrat în figura 1.1-1.

¹ Nucleotidele sunt unități structurale ale ARN, AND și câțiva cofactori. În celulă ei joacă roluri importante în producerea de energie, metabolism și semnalizare [WIKI'05a].

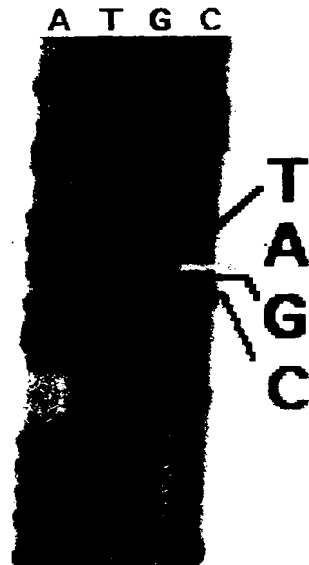


Figura 1.1- 1. Parte a unui gel de secvențiere marcat radioactiv

Fiecare nucleotidă conține o moleculă de fosfat, una de zahar (dezoxiriboza) și una dintre cele patru molecule denumite baze (adenina (A), guanina (G), citozina (C) sau timina(T)). Seria acestor patru molecule determină codul genetic.

Anumite segmente de ADN care conțin instrucțiunile pentru producerea proteinelor specifice organismului se numesc *gene*. Genele sunt situate în anumite segmente de-a lungul ADN-ului uman, împachetate bine în structuri numite cromozomi (oamenii de știință apreciază că ADN-ul uman conține circa 30000 gene). Unele dintre ele sunt responsabile de formarea proteinelor structurale care determină în final trăsăturile fizice, cum ar fi ochii căprui sau parul ondulat al unei persoane. Altele, oferă instrucțiuni pentru ca organismul să producă substanțe chimice importante, numite enzime. În funcție de codurile unei anumite gene, chiar o mică eroare în structura ADN-ului poate însemna probleme grave pentru întregul organism. Spre exemplu, uneori, o eroare într-o singură genă poate duce la scurtarea vieții sau la deficiențe fizice. Anumite probleme genetice sunt cauzate de o singură genă, care este prezentă, dar alterată. Unele modificari ale genelor se numesc *mutații*. Pentru a detecta gena defectă, specialiștii folosesc tehnici de "screening" sofisticate.

1.1.1. Sinteza proteinelor

Prima menționare a cuvântului "proteina", care înseamnă în limba greacă "de prim rang", este atestată de o scrisoare trimisă de Jöns Jakob Berzelius (medic și chimist, 1779-1848) lui Gerhardus Johannes Mulder (biochimist, 1802-1880) pe 10 iulie 1838.

Așa după cum s-a mai menționat, ADN-ul poartă instrucțiunile pentru producerea de proteine. O secvență de trei baze nucleotidice, numite triplet sau

codon, este codul genetic care specifică un anumit amino acid. De exemplu, un triplet GAC (guanina, adenina și citozina) este codul genetic pentru amino acidul *leucina* iar un triplet CAG (citosina, adenina și guanina) este codul genetic pentru amino acidul *valina*. O proteină este compusă din mai multe molecule mici numite amino acizi și structura și funcția proteinei este determinată de secvența acestora. Secvența amino acizilor este la rândul ei determinată de succesiunea bazelor nucleotide în ADN. Proteinele, constituite sub formă de lanțuri de amino acizi se pliază² într-o structură unică 3-dimensională. Forma în care se pliază în mod natural o proteină este cunoscută ca stare nativă și este determinată de secvența de amino acizi. Biochimistii se referă la patru aspecte distincte ale structurii proteinelor. Acestea sunt: structura primară, secundară, terțiară și quaternară.

1. **Structura primară**, reprezintă secvența formată din înșiruirea de amino acizi.

Exemplu de secvență de proteină descrisă prin structura sa primară:

1 madaldlpl laagwshdpe rgvlsktfgf etfvaafgfm trvalwaekw nhhpdwsnys
61 gtvevsltsh dagglterdl klarkidela aa

(*Rhodobacter sphaeroides* 2.4.1, din baza de proteine NCBI, nr de acces: ABA80080.)

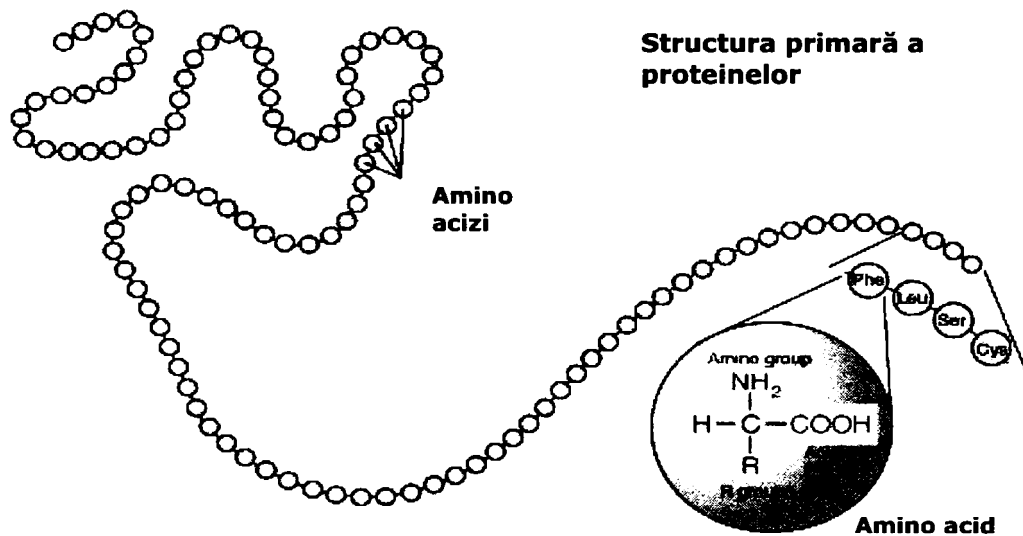


Figura 1.1.1-1. O reprezentare sugestivă a structurii primare.

² Plierea proteinelor este procesul prin care structura unei proteine își asumă forma funcțională sau conformația. Toate moleculele proteinei sunt lanțuri de amino-acizi ramificate heterogen. Prin spiralară și plieri într-o formă 3D specifică ele sunt în stare să își execute funcția biologică [WIKI'05c].

2. **Structura secundară**, descrie forma generală (3D) a regiunilor locale. Ea poate include regiuni ale helixului alfa (alfa helices), fâșii beta (beta sheets), răsuciri (turns) și spire aleatoare (random coil) sau ceva mai puține structuri comune.

Structura de alfa-helix. Peptidele³ sunt răsucite în jurul unui cilindru imaginar și stabilizate de legăturile de hidrogen.

Fâșiile plate beta. Amino acizii adoptă conformația unei foi de hârtie iar structura este stabilizată de legături de hidrogen în diferite fâșii de polipeptide. Se observă că unele dintre fâșii sunt paralele și altele antiparalele.

Alte părți ale structurii nu sunt foarte stabile și adoptă o formă răsucită aleator.

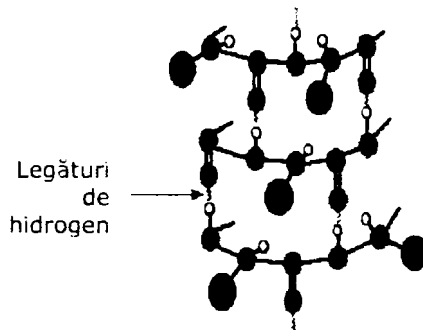


Figura 1.1.1-2. Reprezentarea simbolică a structurii de alfa-helix

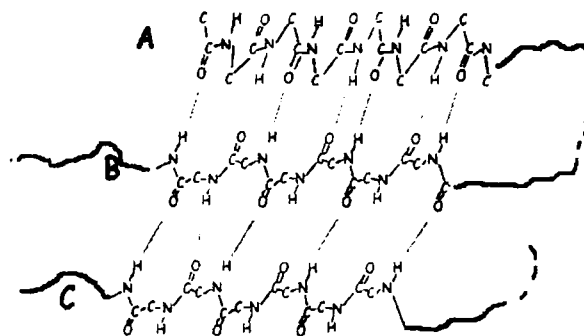


Diagram 1: Beta pleated sheet. The lateral groups (R) are not shown.

Figura 1.1.1-3. Reprezentarea simbolică a fâșiilor plate beta

³ Doi sau mai mulți amino acizi legați printr-o legătură numită "legătură peptidică" (peptide bond). O proteină fiind formată din amino acizi legați în acest mod este uneori referită ca polipeptidă. Unele proteine conțin mai mult decât un singur lanț de polipeptide [WIKI'05d].



Figura 1.1.1-4. Reprezintă structura 3D a unei proteine (Myoglobin). Alfa helix-urile sunt prezentate colorate, răsucirile aleatoare în alb, nu există fâșii beta (beta sheets).

3. **Structura terțiară**, reprezintă forma generală a unei singure molecule de proteine; relația spațială a motivelor⁴ din structura secundară. Conventional, structura terțiară este dedusă prin cristalografie sau rezonanța magnetică nucleară. Domeniul de studiu al structurii terțiare al proteinei este cunoscut ca *biologie structurală*. Ea poate fi considerată drept împachetarea în structura 3D a unui alt aspect al lanțului (side chain packing in the 3D structure).

⁴ Într-o moleculă biologică reprezentată ca un lanț, neramificată, cum ar fi o proteină sau o fâșie de ARN, un motiv structural este un element structural 3D sau o pliere în cadrul lanțului care apare de asemenea într-o varietate de alte molecule. În contextul proteinelor, termenul este adesea folosit pentru domeniul structural, cu toate că un domeniu nu trebuie să fie un motiv și în cazul în care conține un motiv, nu trebuie să fie compus doar dintr-unul singur.

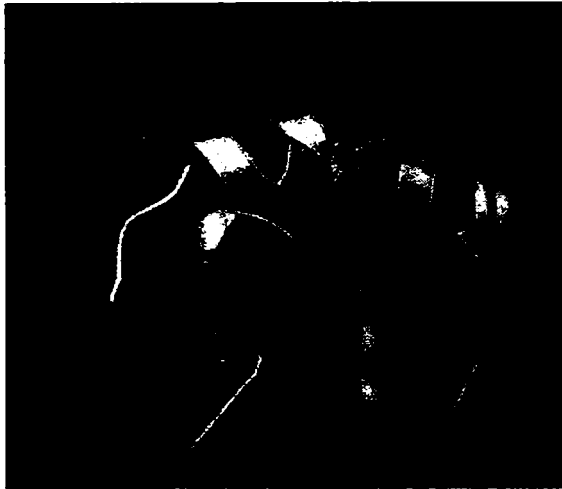


Figura 1.1.1-5. Reprezentarea 3D a unei proteine conținând structuri mai complexe

Amino acizii care sunt foarte distanți unul față de altul în structura primară pot fi apropiați în cea terțiară datorită plierii lanțului și stabilizării lui prin interacțiuni ionice, legături de hidrogen, forțe de dispersie (Van Der Waals), poduri de sulfuri (*sulphur bridges*)

Notă : Unul din scopurile bioinformaticii este de a prezice conformația nativă a unei proteine din secvența sa primară.

4. **Structura quaternară**, reprezintă forma sau structura care rezultă din unificarea mai multor molecule de proteine care în acest context funcționează ca parte a unui ansamblu mai mare de proteine sau complex de proteine. Cu alte cuvinte, această structură este aranjamentul subgrupurilor de polipeptide în cadrul proteinelor complexe compuse din două sau mai multe subunități. Este de menționat faptul că numai proteinele cu mai mult de un singur lanț au o structură quaternară.

Un exemplu de reprezentare în care este evidențiată structura quaternară poate fi observat în figura 1.1.1- 6.



Figura 1.1.1- 6. Reprezentare a structurii quaternare

Proteinele sunt implicate practic în fiecare funcție executată de o celulă. Prin inginerie genetică cercetătorii pot altera secvența și prin urmare structura proteinei. Dintre efectele deficiențelor în proteine sunt simptomele de oboseală, rezistența la insulină, pierderea părului, a pigmentilor din păr, a masei musculare, reducerea temperaturii corporale și neregularități în funcțiile hormonale. De asemenea și excesul poate crea probleme cum ar fi supra reacția sistemului imunitar, disfuncții ale ficatului datorate excesului de reziduuri, scurtarea oaselor datorată gradului crescut de aciditate în sânge.

Notă : Imaginile sunt preluate din materialul biologic documentar realizat pentru proiectul EU Biopattern: Computational Intelligence for biopattern analysis in Support of eHealthcare, Network of Excellence Project No. 508803.

1.1.2. Clasificarea structurală a proteinelor

Conform celor menționate în [MURZ'95], aproape toate proteinele au similarități structurale cu alte proteine și în multe cazuri împart o origine evoluționară comună. Cunoașterea acestor relații aduce contribuții importante biologiei moleculare și altor domenii de știință înrudite. Înțelegerea structurii și evoluției proteinelor este o problemă centrală care joacă un rol important în interpretarea secvențelor produse de proiectele genomice, și prin urmare, înțelegerea evoluției și dezvoltării.

Având în vedere creșterea exponențială în numărul proteinelor ale căror structuri au fost determinate prin cristalografie cu raze X și spectroscopie, se remarcă existența unui corpus larg de informație disponibilă și care crește rapid. Structurile proteinelor pot fi clasificate într-o varietate de moduri înrudite între ele: similaritate funcțională, similaritate evoluționară și similaritate a modului de pliere (fold similarity). Cîteva dintre regulile luate în considerare sunt bazate pe următoarele aspecte:

1. Clasificarea familiilor de proteine la modul general este definită de *funcția lor* care provine din biochimia experimentală.
2. Similaritatea funcțională și clasificarea familiilor poate fi dedusă din similaritatea secvențelor.
3. Dacă există o identitate a secvențelor mai mare de 25% se presupune o oarecare omologie a structurii. Spre exemplu, lanțurile A și B ale hemoglobinei au o omologie de 45%, iar diferit de micile inserări ele au o structură 3D identică.
4. Există o distribuție Gauss a structurilor similare cu un maxim la o valoare în jur de 9% identitate a secvențelor. Aceasta înseamnă că similaritatea secvențelor poate exista și nu poate fi detectată prin metodele standard de comparare a secvențelor. Evoluția divergentă poate merge o cale lungă și încă să păstreze pliarea.
5. Omologia structurii, singură, inferă un tip de aranjamente din punct de vedere energetic favorabile dar nu neaparat o relație funcțională –evoluție convergentă.
6. Este foarte dificil să se distingă evoluția divergentă de cea convergentă.
7. S-au depus eforturi pentru a clasifica proteinele pe baza structurii.

Pentru a facilita înțelegerea și accesul la această informație a fost construită baza de proteine SCOP⁵, care furnizează descrierea detaliată și cuprinzătoare a relațiilor evoluționare și structurale ale proteinelor a căror structură tridimensională a fost determinată. Metoda folosită la clasificare este în esență bazată pe inspectarea vizuală și compararea structurilor cu ajutorul diverselor programe specializate. Clasificarea este făcută manual, pe nivele ierarhice, după cum urmează.

Familia conține un grup de proteine înrudite evoluționar. În primul rând toate proteinele care au o identitate reziduală de 30% și mai mare, apoi proteinele cu o identitate mai mică a secvențelor dar ale căror funcții și structuri sunt foarte simiare.

Superfamilia este constituită din familii ale căror proteine au o identitate a secvențelor scăzută dar ale căror structuri și, în multe cazuri, caracteristici funcționale sugerează o probabilă origine evoluționară comună.

Plierea comună. Superfamiliile și familiile sunt definite ca avînd o pliere comună dacă proteinele lor au aceeași structură secundară majoră în același aranjament cu aceleași conexiuni topologice.

⁵ SCOP: Structural Classification of Proteins. <http://scop.berkeley.edu/>

Clasa. Pentru comoditatea utilizatorilor, diferite plieri au fost grupate în clase. Majoritatea plierilor sunt asociate la una dintre cele cinci clase structurale pe baza structurilor secundare din care sunt compuse.

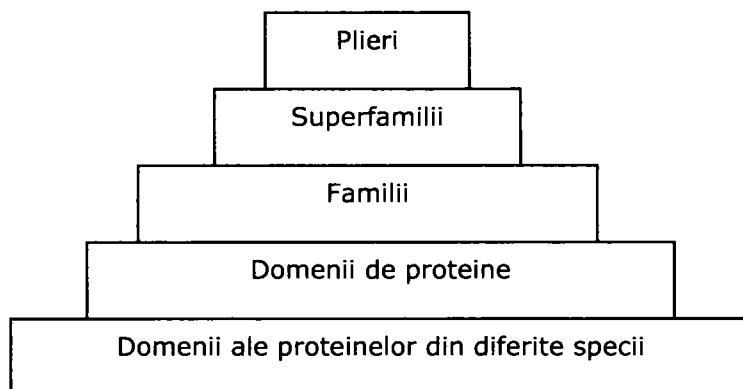


Figura 1.1.2-1 Reprezentarea ierarhiei folosite în clasificarea proteinelor în baza de proteine SCOP. Unitatea de clasificare este de obicei domeniul proteinei.

Exemplu: cytochrome c2 [a.9.3.1]

- **Class:** toate alfa proteinele
- **Fold:** Citocrom c
nucleu: 3 helixuri; fâșie pliată, deschidere
- **Superfamily:** Citocrom c
hema legată covalent completează nucleul
- **Family:** monodomeniul Citocrom c
- **Protein:** Citocrom c2

Pentru studiul intensiv al secvențelor biologice, în prezent se estimează că există în jur de 500 de baze de date publice și comerciale [WIKI'05e]. Acestea conțin informații despre secvențe de nucleotide ale genelor sau secvențe de amino acizi ale proteinelor. Mai mult, ele stochează informații despre funcția, structura, localizarea în cromozomi, efecte clinice ale mutațiilor precum și facilitează găsirea similarității secvențelor biologice. Astfel, pot fi găsite diverse baze de secvențe după cum vor fi menționate în subcapitole următoare.

1.1.3. Formate de reprezentare ale secvențelor biologice

Formatele secvențelor sunt, la modul simplu, felul în care secvențele de amino acizi sau ADN sunt înregistrate într-un fișier pe calculator. Programe diferite așteaptă adesea și formate diferite. Prin urmare, pentru a executa cu succes o funcție este important să se înțeleagă cum arată diverse formate și care este structura lor de bază. Formatele utilizate sunt foarte flexibile dar nu pot face față la prea multă inconsistență.

În general, secvențele sunt descrise sub formă de text ASCII. Ele sunt un aranjament de caractere impus, simboluri și cuvinte cheie care specifică modul de reprezentare pentru: secvență, nume de identificare, comentarii, etc. și indicii despre locul unde ar trebui să caute programul pentru a le găsi. În general nu există caractere ascunse neprintabile în vre-un format cunoscut de secvențe. Toate formatele standard de secvențe pot fi tipărite sau simplu vizualizate prin afișarea simplă a fișierului/filei care le conțin.

Există cel puțin câteva formate de secvențe utilizate în prezent la scară largă, unele mai mult altele mai puțin folosite. Ele au fost proiectate astfel încât să fie capabile să rețină secvența de date și alte informații suplimentare despre secvență. De când au început să se folosească programele software pentru a scrie și citi secvențe, aproape fiecare pachet de analiză a secvențelor și-a inventat propriul ei format. Majoritatea colecțiilor de secvențe care îndrăznesc să se numească baze de date și-au stocat datele în propriul format. O secvență nu presupune vre-un tip de identificare dar în mod sigur aceasta ar ajuta. Cele mai multe formate de secvențe includ cel puțin o formă de nume de identificare (ID), de obicei plasat undeva la începutul formatului de secvență.

Pentru a fi posibilă utilizarea aceluiași set de secvențe în diverse scopuri executate de programe variate s-au implementat convertoare de secvențe dintr-un format în altul. Un asemenea convertor este disponibil la EMBL-EBI (European Bioinformatics Institute <http://www.ebi.ac.uk/cgi-bin/readseq.cgi>), și se numește "READSEQ", tool de conversie al biosecvențelor. EMBOSS, The European Molecular Biology Open Source Software Suite, este un pachet de Open Source software de înaltă calitate pentru analiza secvențelor. Lista completă de formate de secvențe suportată de aplicațiile EMBOSS cu descrierea detaliilor poate fi accesată la adresa oficială sau <http://www.ebi.ac.uk/clustalw/#...formats>, iar unul dintre cele mai cunoscute este formatul Fasta, prezentat în cele ce urmează iar alte câteva exemple sunt menționate în ANEXA 1.

Formatul Fasta (Pearson) este unul dintre cele mai utilizate pentru secvențele biologice și a fost propus de David J. Lipman și William R. Pearson în 1985 în articolul "Rapid and sensitive protein similarity searches" cu numele de FASTA. Programul original a fost proiectat pentru căutarea similarității proteinelor, descris în 1988 ("Improved Programs for Biological Sequence Comparison"). În prezent, în bioinformatică, formatul FASTA este un format de fișier folosit pentru schimbul informației între baze de date cu secvențe genetice. Așa după cum este descris mai jos, o secvență în format FASTA începe cu o descriere biologică într-un singur rând, urmată de rânduri de secvență. Linia de descriere este diferențiată de secvența de date prin simbolul semnului mai mare (">") pe prima poziție. Cuvântul care urmează acestui simbol este identificatorul secvenței și restul liniei este o descriere opțională. De obicei se recomandă să nu fie spațiu între ">" și prima literă a identificatorului. De asemenea, este de preferat ca toate liniile de text să fie mai scurte de 80 de caractere. Secvențele se încheie când o nouă linie începe cu ">".

Un exemplu de format FASTA :

```
>gi|5524211|gb|AAD44166.1| cytochrome b [Elephas maximus maximus]
LCLYTHIGRNIYYGSYLYSETWNTGIMLLLITMATAFMGYVLPWQMSFWGATVITNLFSAIPYI
GTNLV
```

```

EWIWGGFSVDKATLNRFFAFHFILPFTMVALAGVHLTFLHETGSNNPLGLTSDSDKIPFHPYYTI
KDFLG
LLILLLLLLLLALLSPDMLGDPDNHMPADPLNTPLHIKPEWYFLFAYAILRSVPNKLGGVLALFLSI
VIL
GLMPFLHTSKHRSMMLRPLSQALFWTLTMDLLTLTWIGSQPVEYPYTIIGQMASILYFSIILAFPL
IAGX
IENY

```

Sau o altă variantă exprimată la modul general:

```

>SEQUENCE_1
;comment line 1 (optional)
MTEITAAMVKELRESTGAGMMDCKNALSETNGDFDKAVQLLREKGLGKAAKKADRLAAEG
LVSVKVSDDFITIAAMRPSYLSYEDLDMTFVENEYKALVAELEKENEERRRLKDPNKPEHK
IPQFASRKQLSDAILKEAEEKIKEELKAQGKPEKIWDNIIPGKMNSFIADNSQLDSKLT
MGQFYVMDDKKTVEQVIAEKEKEFGGKIKIVEFICFEVGEGLKKTEDFAAEVAAQL
>SEQUENCE_2
;comment line 1 (optional)
;comment line 2 (optional)
SATVSEINSETDFVAKNDQFIALTKDTTAHIQSNLSQVEELHSSTINGVKFEEYLKSQLI
ATTIGENLVRRFATLKAGANGVVNGYIHTNGRVGVVIAAACDSEVASKSRDLLRQICMH

```

Institutul Național de Cercetări în Biotehnologie al SUA a încercat să definească un standard pentru antetul secvențelor stocate. Oricum, nu se oferă o descriere definitivă a acestui format FASTA dar o încercare de a se crea un asemenea format este urmatoarea, unde pentru fiecare bază de date se propune o codificare :

GenBank	<i>gi gi-number gb accession locus</i>
EMBL Data Library	<i>gi gi-number emb accession locus</i>
DDBJ, DNA Database of Japan	<i>gi gi-number dbj accession locus</i>
NBRF PIR	<i>pir entry</i>
Protein Research Foundation	<i>prf name</i>
SWISS-PROT	<i>sp accession name</i>
Brookhaven Protein Data Bank (1)	<i>pdb entry chain</i>
Brookhaven Protein Data Bank (2)	<i>entry:chain PDBID CHAIN SEQUENCE</i>
Patents	<i>pat country number</i>
GenInfo Backbone Id	<i>bbs number</i>
General database identifier	<i>gnl database identifier</i>
NCBI Reference Sequence	<i>ref accession locus</i>
Local Sequence identifier	<i>lcl identifier</i>

Secvențele se așteaptă să fie reprezentate în standard-ul IUB/IUPAC (International Union of Biochemistry (IUB)/ Union of Pure and Applied Chemistry (IUPAC)) pentru amino acizi și coduri pentru acizi nucleici, cu următoarele excepții: literele mici sunt acceptate și transformate în litere mari, o singură linioară poate fi folosită pentru a reprezenta un caracter lipsă (gap character) și în secvențele de amino acizi "U" și "*" sunt caractere acceptabile (a se vedea în cele două tabele din ANEXA 1). Înainte de a începe o interogare, orice cifră din secvența de căutare (query sequence) trebuie fie înlăurată fie înlocuită de codul de literă corespunzător

(e.g. N pentru acid nucleic necunoscut sau X pentru reziduu de amino acid necunoscut).

1.1.4. Baze de date de secvențe

Secvențele biologice rezultate în urma procesului de secvențiere sunt stocate în diverse baze de date respectând cerințele de reprezentare și organizare proprii fiecărei baze. Astfel, există o varietate de „depozite” de secvențe. În general ele sunt organizate în baze de nucleotide, baze de proteine sau baze specializate fie pe o anumită genă fie pe o anumită specie sau proces evolutiv. Câteva dintre cele mai cunoscute vor fi menționate în continuare.

Pentru nucleotide:

a)EMBL Nucleotide DB, European Molecular Biology Laboratory. Baza de nucleotide (cunoscută ca EMBL-Bank) constituie resursa europeană principală de secvențe. Sursele principale de ADN și ARN sunt obținute în urma unor cercetări individuale, proiecte de secvențiere genomică sau aplicații autorizate. (<http://www.ebi.ac.uk/emdl/index.html>.)

b)Genbank (National Center for Biotechnology Information), este o bază de date a Institutului National de Certăări în Biotehnologie din SUA care servește ca arhivă pentru toate secvențele publice de ADN pentru un număr mai mare de 100,000 de organisme diferite. (<http://www.ncbi.nih.gov/Genbank/index.html>.)

c)DDBJ (DNA Database of Japan). Aceasta este singura bancă de ADN din Japonia care este certificată oficial să colecteze secvențe de ADN de la cercetători și de a genera numere de acces internațional recunoscute. Există un schimb de date cu EMBL/EBI și Gen Bank/NCBI. (<http://www.ddbj.nig.ac.jp/Welcome-e.html>.)

d)Entrez este baza de nucleotide care colecționează secvențe din diferite surse incluzând GenBank, RefSeq, și PDB. Numărul de componente continuă creșterea cu o rată exponențială. (<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Nucleotide>)

Pentru proteine :

a)Swiss-Prot, este o bază de secvențe de proteine care încearcă să furnizeze un nivel ridicat de adnotări (cum ar fi descrierea funcției proteinei, structura domeniului ei, modificările post translaționale, variante, etc), un nivel minim de exces și un nivel ridicat de integrare cu al altor baze de date de secvențe. (<http://us.expasy.org/sprot/>)

b)PIR-PSD (Protein Information Resource – Protein Sequence Database), este o bază adnotată ce conține peste 283 000 de secvențe acoperind întreaga arie taxonomică a proteinelor. International Protein Sequence Database-PIR, conține informații privind toate tipurile de proteine care apar natural și ale căror secvențe sunt cunoscute. Un obiectiv major al bazelor de secvențe este de a furniza date cuprinzătoare, non redundante, organizate unic pe baza omologiei și taxonomiei. (<http://pir.georgetown.edu/pirwww/search/textpsd.shtml>.)

c)NCBI (National Center for Biotechnology Information), cuprinde date compilate din diverse surse incluzand: Swiss-Prot, Protein Information Resource (PIR), Protein Data Bank (PDB), Protein Resource Foundation (PRF) in Japan. (<http://www.ncbi.nih.gov/>)

d)Uniprot (Universal Protein Resource), este cel mai cuprinzător catalog de informații despre proteine. Este un depozit central de secvențe de proteine și funcții creat unind informația conținută în alte baze de proteine: Swiss-Prot, TrEMBL, and PIR. (<http://www.uniprot.org>)

e)SCOP (Structural Classification of Proteins) își dorește să ofere o descriere detaliată și cuprinzătoare a relațiilor structurale și evoluționare dintre proteine a căror structură 3D este cunoscută. (<http://scop.berkeley.edu/>)

Baze specializate:

Profile pe genele producătoare de cancer:

-**CGAP** Cancer Genes , <http://cgap.nci.nih.gov/Genes/GeneFinder>,

-**Tumor Gene Database**, <http://condor.bcm.tmc.edu/ermb/tgdb/tgdf.html>.

-**Entrez-Cancer Chromosomes**,

<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=cancerchromosomes>.

-**OMIM** (Online Mendelian Inheritance in Man), este o bază de date de gene umane, tratamente genetice și dezordini ereditare disponibile la National Center for Biotechnology Information (NCBI) și care pot fi accesate pe web.

<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=OMIM>.

- **LSDBs** (Locus Specific Database),

<http://www.genomic.unimelb.edu.au/mdi/dblist/glsdb.html>.

-**P53 Databases**, International Agency for Research on Cancer (Lyon), France (IARC), <http://www.iarc.fr/p53/>, TP53 este o genă responsabilă pentru tumoare.

- **BRCA1&2** database, <http://research.nhgri.nih.gov/bic/>.

Baze de secvențe cu mutații:

- **Human Gene Mutation Database (HGMD)**,

<http://archive.uwcm.ac.uk/uwcm/mg/hgmd0.html>

- **Human Genome Variation database (HGV)**, <http://hgibase.cgb.ki.se/>

Baze de secvențe profilate pe o anumită specie:

- **GDB** (Human Genome Organization), <http://wehih.wehi.edu.au/gdb/>

1.2. Elemente în analiza secvențelor biologice

Pentru înțelegerea modului de relaționare dintre secvențele biologice, în cele ce urmează vor fi descrise pe scurt elementele de interes major.

1.2.1. Omologia secvențelor în genetică

În genetică este folosită omologia cu referire la secvențele de proteine sau ADN, însemnând că secvențele respective au strămoși comuni. Omologia secvențelor poate indica de asemenea o funcție comună iar din punct de vedere al gradelor de comparație ea este o calitate de forma totul-sau-nimic; nu există grade de omologie. Regiunile de secvențe care sunt omoloage pot fi numite regiuni conservate, consensus sau canonice și reprezintă cea mai comună alegere a bazei sau amino acidului la fiecare poziție.

Omologia proteinelor și ADN-ului este adesea concluzionată pe baza similarității secvențelor, și asta în special în bioinformatică. Spre exemplu, în general, dacă două sau mai multe secvențe au o secvența de ADN aproape identică este probabil ca ele să fie omoloage. Oricum este posibil ca similaritatea secvențelor să nu provină din apartenența la un ancestor comun și secvențele scurte pot fi similare accidental sau secvențele pot fi similare deoarece ambele au fost selectate pentru a fi legate la o anumită proteină cum ar fi factorul de transcripție⁶; asemenea secvențe sunt similare dar nu omoloage. Există mulți algoritmi pentru a clasifica secvențele de proteine în familii de secvențe care sunt mulțimi mutual omoloage.

1.2.2. Motiv al unei secvențe

Motiv al unei secvențe este numită o secvență model de nucleotide sau amino acizi care este destul de răspândită și are sau se presupune a avea o semnificație biologică. Motivele au fost descoperite prin studierea genelor similare la diferite specii.

În cadrul unei secvențe sau bază de date de secvențe, cercetătorii caută și găsesc motivele folosind tehnici bazate pe utilizarea calculatoarelor. Asemenea tehnici aparțin disciplinei numite bioinformatică.

⁶ Transcripția (în genetică) este procesul prin care o secvență de ADN este copiată din punct de vedere enzimatic de către un ARN polimerizat pentru a produce un ARN complementar. Sau, în alte cuvinte, transferul informației genetice de la ADN la ARN. În cazul proteinelor, codificarea ADN, transcripția este începutul procesului care în mod ultimativ conduce la translatarea codului genetic în proteine sau peptide funcționale.

1.2.3. Aliniament

Respectând definițiile din [WIKI'05f], aliniamentul secvențelor este un aranjament de două sau mai multe secvențe, subliniind similaritatea lor. Secvențele sunt presărate cu gap-uri (care sunt amino acizi sau secvențe de amino acizi lipsă și sunt de obicei marcate cu linii) astfel că oricând este posibil, coloanele rezultate în urma alinierii conțin caractere identice sau similare din secvențele implicate:

```

          ↙ gap
tcctctgcctctgccaatcat---caaccccaaagt
      ||| ||| |||| |||| | ||||| |||||
tcctgtgcatctgcaatcatgggcaaccccaaagt

```

De obicei se studiază evoluția secvențelor care provin de la un ancestor comun, în special secvențele biologice cum ar fi cele de proteine sau ADN. Nepotrivirile din aliniament corespund mutațiilor iar gap-urile corespund inserărilor sau ștergerilor. Aliniamentul secvențelor poate fi folosit pentru studiul evoluției limbajului sau similarității dintre texte. Termenul de aliniament al secvențelor poate de asemenea să se refere la procesul de construire al unui asemenea aliniament sau de căutare a unor aliniamente semnificative într-o bază de secvențe posibil neînrudite.

Aliniament pentru perechi de secvențe

Metodele de aliniament ale perechilor de secvențe sunt preocupate de găsirea intervalelor de aliniament local sau global ale secvențelor de proteine (amino acizi) sau ADN (acizi nucleici). În mod frecvent, scopul acestora este de a găsi omologii sau înrudiri ale unei gene sau produs al unei gene într-o bază de date de exemple cunoscute. Această informație este utilă în a răspunde unei mari varietăți de întrebări biologice. Cea mai importantă aplicație a intervalului de aliniament este identificarea secvențelor cu structuri sau funcții necunoscute. O altă utilizare importantă o reprezintă studiul evoluției moleculare.

Aliniamentul multiplu

Aliniamentul multiplu este o extensie a aliniamentului perechilor de intervale pentru a incorpora mai mult de două secvențe. Metodele de aliniament multiplu încearcă să alinieze toate secvențele dintr-un set specificat. Aliniamentele ajută la identificarea regiunilor comune dintre secvențe. Există câteva abordări pentru crearea aliniamentelor multiple de secvențe, una dintre cele mai populare fiind strategia de aliniament progresiv utilizată de familia de programe Clustal (a se vedea în capitolul 3-5).

Notă: Aliniamentul multiplu al secvențelor este o problemă dificilă din punct de vedere computațional și este clasificată ca o problemă *NP-hard* (Non-deterministic Polynomial-Time hard).

Aliniamentul global

Aliniamentul global dintre secvențe este un aliniament în care toate caracterele din secvențele comparate participă la aliniament. Aliniamentele globale sunt cel mai utile pentru găsirea secvențelor celor mai înrudite. Deoarece aceste secvențe sunt de asemenea ușor de identificat prin metode de aliniament local, aliniamentul global este oarecum discriminat ca tehnică.

Aliniamentul local

Metodele de aliniament local caută regiuni înrudite în cadrul secvențelor – cu alte cuvinte ele pot să consiste dintr-un subset de caractere aflat în fiecare secvență. Spre exemplu, pozițiile 20-40 ale unei secvențe. A pot fi aliniat cu pozițiile 50-70 ale altei secvențe B.

Aceasta este o tehnică mai flexibilă decât aliniamentul global și are avantajul că regiuni înrudite care apar într-o ordine diferită în cele două proteine pot fi identificate ca înrudite. Aceasta nu este posibil însă cu o metodă de aliniament global.

Aliniamentul structural

Un mod posibil de a valida aliniamentul secvențelor este de a utiliza aliniamentele structurale unde este posibil. Aliniamentul structural este o formă de aliniament care

încearcă să stabilească echivalențe între două sau mai multe structuri de proteine pe baza plierilor lor. Deoarece structura proteinelor este mai conservată în timpul evoluției decât secvența așa după cum afirma biologul Cyrus Chothia și cercetătorul Arthur Lesk în 1986 [LESK'86], aliniamentele structurale pot fi mai de încredere pentru distanțe evoluționare lungi, când secvențele diverg așa de mult încât simpla comparare a secvențelor nu poate detecta similaritatea lor. Oricum, cu toate că aliniamentele structurale sunt utile, ele sunt complet artificiale.

1.2.4. Similaritatea secvențelor biologice

Similaritatea a două sau mai multe secvențe biologice poate fi privită din perspectiva asemănării funcționale care la randul ei, este strâns legată de asemănarea structurală a secvențelor biologice. Una dintre cele mai importante măsuri de similaritate este distanța editată dintre amino acizii a două secvențe.

La prima vedere, decizia că două secvențe sunt similare nu este diferită de decizia conform căreia două șiruri de text sunt similare. O mulțime de metode pentru analiza secvențelor biologice este, prin urmare, înrădăcinată în știința calculatoarelor unde există o literatură extensivă în metode de comparare a șirurilor de caractere. Conceptul de aliniament e crucial. Evoluând, secvențele acumulează inserții și ștergeri precum și substituiri. Prin urmare, înainte de a evalua similaritatea a două secvențe în general se începe prin căutarea unui aliniament plauzibil între ele [DURB'98].

Apropape toate metodele de aliniament găsesc cel mai bun aliniament dintre două șiruri urmând unele scheme/planuri de calcul. Deoarece se dorește un algoritm de calcul care să ofere cel mai probabil aliniament biologic având scorul cel mai mare, se ia în considerare faptul că evoluțiile biologice au istorii evoluționare, structuri 3D pliate, și alte caracteristici care limitează evoluția secvenței lor primare. Prin urmare, ca un adaos la mecanica aliniamentului și a algoritmilor de comparație, sistemul de calcul însuși cere o atenție separată putând fi foarte complex.

Pe baza similarității și/sau omologiei se construiesc familii de secvențe care partajează proprietăți comune. Ele sunt constituite fie prin analiza manuală a specialiștilor fie utilizand metode computaționale.

1.3. Concluzii

Progresul rapid în cercetare legat de secvențe genomice/proteomice ale diverselor organisme, motivează nevoia de conceptualizare, analiză și în final înțelegere a informației stocate de aceste secvențe. Datorită modului lor de reprezentare textuală s-a deschis calea către analiza asistată de calculator, îmbinând cunoștințele biologice cu cele matematice și informatice pentru a eficientiza găsirea de răspunsuri la diverse întrebări biologice.

Dintre multiplele avantaje ale cooperării între aceste domenii ale științei sunt:

- Păstrarea informației biologice sub forma organizată a bazelor de secvențe;
- Posibilitatea de explorare a unei cantități imense de informații în timp scurt;
- Extragerea automată de caracteristici pe baza comparației secvențelor, operațiune de altfel dificilă și mare consumatoare de timp, prin analiza cu metode biologice sau manuale.

2. STADIUL ACTUAL AL TEHNICILOR DE SIMILARITATE PENTRU SECVENȚELE BIOLOGICE

2.1. Metode de determinare a similarității secvențelor biologice

În definirea noțiunii de similaritate, mai întâi depindem de noțiunile de similaritate intuitive și umane. Astfel, similaritatea este exact ceea ce credem noi că este [WATE'84]. Pe parcursul acestui capitol se va încerca ajungerea la o noțiune formală a ceea ce determină ceva să fie mai similar cu un lucru decât cu altul.

Considerând secvențele ACTCCG și ACACCG, ele par a fi similare. De fapt se observă că simpla înlocuire a lui T cu un A este suficientă pentru a obține a doua secvență. Pentru definirea similarității, este utilă mai întâi introducerea noțiunii de „distanță” dintre două șiruri. Astfel, distanța dintre două șiruri este zero dacă ele sunt identice și crește dacă ele devin mai disimilare [WATE'84]. Un mod de a defini distanța dintre două șiruri este de a se observa cantitatea de schimbări care este necesară de a se aplica unei secvențe pentru a o obține pe cealaltă. În cazul în care trebuie să înlocuim T cu A din exemplul anterior, distanța dintre cele două secvențe va fi 1 (dacă înlocuirea poate fi considerată unitatea de măsură). Se poate merge mai departe la prezentarea altor modificări cum sunt „inserarea” și „ștergerea”. Inserarea are loc atunci când se inserează una sau mai multe litere în secvență (la anumite poziții) iar ștergerea atunci când ștergem anumite litere aflate în anumite poziții. Înlocuirea însăși poate fi considerată ca o ștergere urmată de o inserție dar în termeni de complexitate computațională nu creează vreo problemă în a considera înlocuirea drept o unitate de schimb [WATE'84].

2.1.1. Distanță de editare și similaritate

Această *distanță de editare* este definită în general ca numărul minim de schimbări efectuate asupra unei secvențe pentru a o aduce exact în forma celeilalte. Spre exemplu, șirul ACCTGA devine AGCTA prima dată prin înlocuirea celei de-a doua litere, C, cu G și apoi ștergând G din a cincea poziție. Astfel, distanța de editare dintre cele două șiruri este 2.

În [WATE'76] este făcută următoarea formulare: fie S setul de cuvinte finite dintr-un alfabet, incluzând cuvântul vid și fie $\tau = \{ T \mid T : S \rightarrow S \}$ un set de transformări care include transformările identice. Interesul este în transformările T_1, T_2, \dots, T_k din τ astfel încât

$$T_1 \circ T_2 \circ \dots \circ T_k (a) = b. \quad (2.1.1-1)$$

Când ponderi nenegative sunt asociate transformărilor, atunci suma minimă a ponderilor $T_1 \circ T_2 \circ \dots \circ T_k (a) = b$ poate fi văzută ca distanța de la o secvență a la o

alta b . Aceasta este descrisă în mod explicit în [WATE'76] și, cu o simetrizare evidentă, este obținut un spațiu metric. Cazuri interesante apar când τ este restricționat la seturi specifice de transformări și ponderi specifice. Spre exemplu, este important în ce mod există sau nu algoritmi eficienți.

Pentru mutații, inserții și ștergeri, aliniamentul corespunzător sumei minime de ponderi ar trebui să fie afișat atâta timp cât fiecare element al secvenței nu este în mai mult de un singur eveniment evoluționar. Lăsând $d(x,y)$ să fie ponderea substituției y pentru x , chiar având $d(x, x)=0$ pentru toți x , nu face din suma aliniamentului minim de ponderi o metrică. Dacă $d(*.*)$ este o metrică pe mulțimea de litere, atunci Sellers în [SELL'74b] arată că rezultă o metrică pe S . În acest caz poate fi afișat aliniamentul.

O abordare a lui Needelman și Wunch în [NEED'70], prezintă un algoritm pentru maximizarea numărului de potriviri minus numărul de inserări și ștergeri. Acesta este referit ca un criteriu de similaritate maximă în timp ce, cel amintit anterior este un criteriu al distanței minime. Legarea conceptului de distanță cu cel de similaritate este important în psihologie [SHEP'80], iar relația dintre ele este de asemenea importantă în acest context. Uneori, este folosit $d(x, y)<0$ iar minimul rezultat referit ca distanță. În [WATE'84], distanța este rezultatul minimizării ponderilor nenegative, cu $d(x, x)=0$. Situații cu $d(x, y)<0$ pot fi numite *similaritate negativă*.

Fiecare substituție sau pereche (a_i, b_j) într-un aliniament corespunde uneia cu $k = 1, 2, \dots, 16$ perechi (A,A), (A,T), ..,(G,G) și are o similaritate $\alpha_k \geq 0$ sau o distanță $\beta_k \geq 0$.

Observație: Operațiile de inserție și ștergere sunt referite adeseori cu un singur cuvânt: **indels**.

Indels au fiecare o pondere dată w pentru similaritate sau x pentru distanță. Aliniamentul similarității maxime coincide cu aliniamentul distanței minime dacă și numai dacă

$$\beta_i = \max_{1 \leq j \leq 16} (\alpha_j) - \alpha_i \text{ pentru } i = 1, 2, \dots, 16$$

și

(2.1.1-2)

$$x = \left(\max_{1 \leq j \leq 16} \alpha_j \right) / 2 + w.$$

Un rezultat mai general decât acesta poate fi găsit în [SMIT'81]. Problema legării *distanței de similaritate* presupune existența unui criteriu general pentru similaritate care să poată fi complementară unei distanțe metrice. Aceasta însă rămâne o chestiune deschisă, abordată separat de diverși algoritmi.

2.1.1.1. Similaritate versus distanță

Elementele matricelor de scoruri folosite pentru estimarea similarității secvențelor (vor fi tratate mai pe larg în secțiunea 2.1.2) specifică ponderea care trebuie asociată unei comparații date prin:

- Costul înlocuirii unui reziduu cu un altul (distanța); sau o măsură a similarității pentru înlocuire.
- Distanța este în mod natural folosită pentru reconstruirea arborelui filogenetic; similaritatea este folosită pentru căutarea în baze de date.
- Logica algoritmului nu se schimbă: maximizarea similarității este fundamental aceeași cu minimizarea unei distanțe.
- Matricele de distanță și similaritate sunt inter-convertibile prin unele transformări matematice în conformitate cu aplicația dată.

2.1.2. Matrice de substituție

În scopul evaluării gradului de similaritate al secvențelor aliniată sunt necesare planuri de calcul care de cele mai multe ori se bazează pe matrice de ponderi obținute din analiza comportamentelor ce formează secvențele. Alegerea uneia dintre matrice poate influența rezultatul analizei iar ca un alt aspect, ele reprezintă o teorie particulară a evoluției.

Un biolog cu o bună intuiție pentru proteine ar putea inventa un set de 210 termeni de calcul pentru toate perechile posibile de amino acizi dar este extrem de util de a avea o teorie ghid pentru semnificația scorului. Pentru exemplificare, vor fi prezentate derivarea de scoruri de substituție dintr-un model probabilistic. Fie o pereche de secvențe x și y , de lungime n și respectiv m iar x_i este al i -lea simbol în x și y_j al j -lea simbol al lui y . Aceste simboluri vin din dintr-un alfabet A ; în cazul ADN alfabetul va fi format din 4 baze $\{A, G, C, T\}$ și în cazul proteinelor din 20 amino acizi. Simbolurile din acest alfabet A vor fi notate cu litere mici ca a, b . Secvențele care vor fi aliniată simbolic se consideră a nu avea goluri („gap”-uri), astfel vor avea lungime egală.

Fiind dată o pereche de secvențe aliniată se dorește asocierea unui scor acestui aliniament, care să dea o măsură a probabilității relative ca secvențele să fie înrudite sau complet divergente. Aceasta se realizează de regulă folosind modele care asociază o probabilitate unui aliniament pentru fiecare caz; prin urmare, se consideră raportul dintre cele două probabilități.

Modelul independent sau aleator R este cel mai simplu. El presupune că o litera a apare independent cu o frecvență q_a , și prin urmare probabilitatea celor două secvențe este tocmai produsul probabilităților fiecărui amino acid:

$$P(x, y | R) = \prod_i q_{x_i} \prod_j q_{y_j} . \quad (2.1.2-1)$$

În **modelul alternativ** de potrivire M , perechile de reziduuri aliniată apar cu o probabilitate de legătură q_{ab} . Această valoare q_{ab} poate fi considerată drept probabilitatea ca reziduurile a și b sunt fiecare derivate independent dintr-un reziduu original necunoscut c în strămoșul lor comun (c poate fi la fel ca a și/sau b). Aceasta dă o probabilitate întregului aliniament de

$$P(x, y | R) = \prod_i q_{x_i y_i} . \quad (2.1.2-2)$$

Raportul acestor două probabilități este cunoscut ca *grad de asociere* numit în literatura de specialitate și „**odds ratio**”:

$$\frac{P(x, y | M)}{P(x, y | R)} = \frac{\prod_i p_{x_i y_i}}{\prod_i q_{x_i} \prod_i q_{y_i}} = \prod_i \frac{p_{x_i y_i}}{q_{x_i} q_{y_i}}. \quad (2.1.2-3)$$

Pentru a ajunge la un sistem aditiv de calcul, se ia logaritmul acestui raport numit „*log-odds ratio*”:

$$S = \sum_i s(x_i, y_i), \quad (2.1.2-4)$$

unde

$$s(a, b) = \log \left(\frac{p_{ab}}{q_a q_b} \right) \quad (2.1.2-5)$$

este log din raportul de probabilitate al perechii de reziduuri (a, b) care apare ca pereche aliniată, opusă uneia nealiniată.

Așa cum s-a dorit, ecuația 2.1.2-4 este o sumă de scoruri individuale $s(a, b)$ pentru fiecare pereche de reziduuri aliniată. Valorile $s(a, b)$ pot fi aranjate într-o matrice. Pentru proteine, spre exemplu, ele formează o matrice de 20×20 , cu $s(a_i, a_j)$ în poziția i, j în matrice, unde a_i, a_j sunt al i -lea și al j -lea amino acizi (în numerotații). Aceasta este cunoscută ca *matricea de scoruri* sau *matricea de substituție*. Un exemplu de matrice obținute în acest fel sunt matricele PAM și BLOSSUM.

Observații:

Când se consideră matrice de scoruri se folosește convenția conform căreia matricele au indici numerici ce corespund liniilor și coloanelor matricei. Astfel, M_{11} se referă la intrarea din prima linie și prima coloană. În general, M_{ij} se va referi la conținutul corespunzător liniei i și coloanei j . Pentru aliniamentul secvențelor, se asociază fiecărei litere din alfabet valori numerice.

Urmând abordarea lui Wheeler în [WHEE'96], dacă alfabetul este $A = \{A; C; G; T\}$, atunci $A=1, C=2$, etc. Astfel la M_{12} se va afla scorul dintre A și C. Deoarece vor fi considerate diverse matrice de scoruri, ele vor fi distinse folosind diferite litere ca R_{ij} pentru matricea de înlocuiri, S_{ij} pentru matricea log odds⁷, și așa mai departe.

Distanța dintre două secvențe este calculată ca suma de diferențe dintre secvențe. Diferențele se pot datora inserărilor și ștergerilor de nucleotide sau amino acizi (în funcție de tipul secvențelor comparate). În acest scop au fost dezvoltate tabele de frecvență numite Percent Accepted Mutations sau matrice PAM care măsoară ratele de mutații ale amino acizilor în cadrul familiilor de proteine.

⁷ Matricele **log-odds** sumarizează înlocuirile observate care au avut loc în timp ce s-au conservat proprietățile esențiale a multor familii de proteine.

Există mai multe tipuri de matrice PAM (PAM 250⁸, PAM 45) [DAYH'78]. Un alt tip de matrice des utilizate sunt BLOSUM (Blocks Substitution Matrix) care sunt matrice de substituție derivate din observarea frecvențelor înlocuirilor de amino acizi în regiunile foarte conservate ale aliniamentelor locale fără discontinuități. Datele pentru scorurile de substituție în aceste matrice provin de la aproximativ 2000 blocuri de segmente de secvențe aliniată caracterizând mai mult de 500 de grupuri de proteine înrudite [HENI'92].

2.1.2.1. Matrice de substituție pentru nucleotide

1. Matricea identitate (similaritate)

	A	T	C	G
A	1	0	0	0
T	0	1	0	0
C	0	0	1	0
G	0	0	0	1

Pentru elementele în linia i pe coloana j : $S_{ij} = 1, i=j$; $S_{ij} = 0, i \neq j$

2. Matricea BLAST (similaritate)

	A	T	C	G
A	5	-4	-4	-4
T	-4	5	-4	-4
C	-4	-4	5	-4
G	-4	-4	-4	5

3. Matricea de Tranziție/Transversie

	A	T	C	G
A	0	5	5	1

⁸ Numărul 250 în PAM250 corespunde unei medii de 250 înlocuiri de amino acizi per 100 reziduuri dintr-un set de date de 71 de secvențe aliniată [DAYH'78]. Cu cât este mai mare numărul matricii cu atât este mai mare distanța evoluțională dintre secvențele comparate.

T 5 0 1 5

C 5 1 0 5

G 1 5 5 0

Folosirea unei matrice de Tranziție/Transversie reduce "zgomotul" în comparații de secvențe distante ca înrudire.

2.1.2.2. Matrice de substituție pentru proteine

a) Matricea identitate

$$R_{ij} = 1, i=j$$

$$S_{ij} = 0, i \neq j$$

- b) Matricea Codului Genetic, ce presupune un scor bazat pe numărul minim de schimbări de baze necesar pentru convertirea unui amino acid în altul (matricea distanță conține costurile mutațiilor de amino acizi [WHEE'96]).
- c) Caracteristici Fizico/chimice, care încearcă să cuantifice unele atribute fizice sau chimice ale reziduurilor și asignează arbitrar ponderi pe baza similarității reziduurilor în proprietatea aleasă (ex. Matricea hidrofobicității [WHEE'96] construită din date hidofile (M. Levitt, J. Mol. Biol. 104, 59 [1976]) derivate de George et al. 1990, și prezentate în lucrarea: Mutation Data Matrix and Its Uses, Methods in Enzymology 183, 333.)

2.1.2.3. Matricea "Log odds"

$$S_{ij} = \log \frac{q_{ij}}{p_i p_j} \quad (2.1.2.3-1)$$

S este raportul a două probabilități: probabilitatea ca două reziduuri, i și j sunt aliniate prin descendență evoluționară și probabilitatea ca ele să fie aliniate din întâmplare.

- q_{ij} sunt frecvențele cu care se observă alinierea reziduurilor i și j în secvențe cunoscute ca înrudite. Ele sunt derivate dintr-o "matrice a probabilităților de tranziție."
- p_i și p_j sunt frecvențele de apariție a reziduurilor i și j în setul de secvențe.
- ex., PAM250, BLOSUM62 și altele.

2.1.2.4. Matricea PAM(250)

Pașii parcurși:

1. Se aliniază secvențe care sunt cel puțin 85% identice. În plus, se minimizează ambiguitățile în aliniamente și numărul mutațiilor identice.

2. Se reconstruiește arborele filogenetic și se inferă secvențele originale (cu rol de ancestor). Au fost folosiți 71 arbori conținând 1 572 schimbări.
3. Punctarea înlocuirilor acceptate prin selecția naturală, în toate comparațiile de perechi de secvențe (fiecare A_{ij} este frecvența cu care amino acidul j a fost înlocuit de amino acidul i în toate comparațiile).
4. Se calculează mutabilitatea amino acizilor, m_{ij} , adică tendința unui aminoacid j de a fi înlocuit;
5. Se combină datele de la 3. și 4. pentru a produce o matrice a probabilității mutațiilor pentru o distanță evoluționară PAM, în conformitate cu formula

$$M_{ij} = \frac{m_j A_{ij}}{\sum_i A_{ij}} \quad (2.1.2.4-1)$$

$$M_{jj} = 1 - m_j$$

6. Se calculează Matricea *log odds* pentru scorurile de similaritate:
Se împarte fiecare element al matricei de mutații M , la frecvența aparițiilor fiecărui reziduu:

$$R_{ij} = \frac{M_{ij}}{f_i} \quad (2.1.2.4-1)$$

R este o matrice a înrudirilor matricei de asociere (a *Relatedness Odds Matrix*), f_i este frecvența reziduuului i . Matricea *log odds*, S_{ij} este calculată din matricea înrudirilor de asociere, R_{ij} simplu, luând *log* din fiecare R_{ij} . Familii diferite de proteine manifestă rate PAM diferite.

2.1.2.5. Proprietăți ale Matricei Probabilităților de Mutații

1. Suma lui m_{ij} pentru orice coloană j , este unu (trivial). A se nota faptul că probabilitatea cu care un amino acid se va modifica este de ordinul 1% pentru fiecare amino acid. Probabilitatea că el va rămâne același este de ordinul 99% pentru fiecare amino acid.
2. Matricea Probabilităților de Mutații (The Mutation Probability Matrix), M_1 , definește un element de schimbare evoluționară: anume, 1 PAM (Accepted Point Mutation per 100 residues). Matricea poate fi folosită pentru a simula evoluția prin folosirea unui generator de numere aleatoare pentru a selecta soarta fiecărui reziduu în secvență în conformitate cu probabilitatea dată în tabel.
3. Aplicarea succesivă a M_1 pe o secvență rezultă în 2, 3, 4... PAM de modificări evoluționare.
4. Matricea conține informații compoziționale, atâta timp cât ea depinde de frecvențele relative ale amino acizilor din baza de secvențe din care sunt deduse potrivirile. Pentru cazurile extreme se obțin:

- a. Elementele matricei 0 PAM sunt 1 pentru M_{ij} și 0 pentru M_{ij} ;
- b. Elementele ∞ ale matricei PAM abordează o compoziție a amino acizilor asimptotică.

Presupuneri în modelul PAM:

1. Înlocuirile în fiecare poziție depind numai de amino acidul din acea poziție și probabilitatea dată de tabel (model Markov).
2. Secvențele comparate au compoziția medie a amino acizilor.

Surse de eroare în modelul PAM

1. Multe secvențe se abat de la compoziția medie.
2. Înlocuirile rare nu se observă prea frecvent pentru a rezolva acuratețea probabilităților relative (pentru 36 de perechi nu s-a observat nici o înlocuire!).
3. Erorile în 1PAM sunt intensificate în explorarea la 250 PAM.
4. Procesul Markov este o reprezentare imperfectă a evoluției: Secvențele înrudite ca distanțe au de obicei insule (blocuri) de reziduuri conservate. Aceasta implică faptul că înlocuirea nu este egal probabilă pe întreaga secvență.

2.1.2.6. Matricea BLOSUM (Blocks Substitution Matrix)

Înțiatorii acestui tip de matrice sunt Steven Henikoff și Jorja G. Henikoff (1992) în lucrarea "Amino acid substitution matrices from protein blocks" apărută în Proceedings of Natural Academic Sciences în 1989.

1. Datele de start sunt blocuri conservate extrase dintr-o bază specializată de blocuri și conțin:
 - o Secvențe aliniate, fără goluri (ungapped sequences);
 - o Prezintă un grad de similaritate foarte variată dar măsurile sunt luate astfel încât să evite limitarea secvenței de test cu secvențele foarte înrudite care apar frecvent.
2. Înregistrările înlocuirilor sunt făcute direct prin consemnarea tuturor perechilor de reziduuri aliniat f_{ij} .

Frecvența observată a fiecărei perechi este: $q_{ij} = f_{ij} / (\text{numărul total al perechilor de reziduuri})$

- o Aceasta include cazul în care $i = j$ (i.e. nici o înlocuire observată);
 - o Frecvența așteptată a fiecărei perechi este în esență produsul frecvențelor fiecărui reziduu din setul de date.
3. Secvențele similare dintr-un bloc, deasupra unui prag procentual de similaritate, sunt grupate iar membrii ai grupului contribuie fracțional la înregistrarea finală. Această etapă:
 - o Reduce numărul perechilor identice (AA, SS, TT, etc., potriviri) în înregistrările finale;
 - o Este într-o oarecare măsură analogă cu creșterea distanței PAM.
 - o Dacă pragul de grupare este 80%, matricea finală este BLOSUM 80.

- Gruparea la 62% reduce numărul de blocuri care contribuie la tabel cu 25% - încă 1.25×10^6 perechi contributante!
- Cea mai puțin frecventă înlocuire a perechilor de amino acizi a fost observată de 2369 ori!

2.1.2.7. Noi matrice de scoruri

O actualizare a matricei PAM folosind metoda lui Dayhoff este propusă de autori în [DAVI'92] unde sunt marcate 59 190 mutații acceptate în 16 130 secvențe.

O altă matrice de scoruri rezultată din aliniamentul întregii baze de date SWISS-PROT pentru care s-au identificat 1.7×10^6 potriviri din secvențe prin deducție de 6.4 la 100 PAM este propusă de Gonnet et al. în [GONN'92].

Alte matrice specializate de scoruri

- În [LEE'94] se propune crearea unei matrice "step-matrix" bazată numai pe blocuri aliniate de receptori cuplați de G-Proteine în care elementele matricei sunt proporționale cu raritatea substituției. Constituie un avantaj excelent în construirea unei filogenii coerente a larg divergenților receptori ai G-proteinei.
- Matrice pentru detectarea mutațiilor "frame shift" care duc la noi secvențe codificate (sau care provin din erori de secvențiere) [CLAV'93].
- Matrice de scoruri create din substituții observate de reziduuri găsite în medii structurale similare 3D [BORD'91].

2.2. Starea actuală în dezvoltarea unor noi metode de similaritate

În această secțiune sunt prezentate o varietate de metode oferite publicului pentru exprimarea similarității dintre secvențele biologice. Datorită versatilității acestor tehnici se va evita impunerea vreunui model de organizare a prezentării lor. În schimb ele vor fi descrise succesiv în funcție de ordinea lor cronologică, incluzând ideile principale, datele folosite, rezultatele și unele concluzii.

O abordare interesantă este întâlnită în [SJOL'96], unde este descrisă o metodă pentru detectarea omologiei slabe dintre proteine utilizând modelarea mixturilor Dirichlet. Densitățile incorporate sunt proiectate pentru a fi combinate cu frecvențele observate ale amino acizilor, în scopul formării de estimări ale probabilităților așteptate ale amino acizilor pentru fiecare poziție într-un profil HMM sau orice alt model statistic. Background-ul teoretic al acestei lucrări a fost stabilit în publicații anterioare [BERG'85], [SANT'89].

În [KATT'00], au fost selectate un set de 80 000 de proteine din baza de date SWISS-PROT și analizate pentru repetițiile în tandem folosind tehnica alunecării ferestrelor ("sliding window technique"). Modelele observate că se repetă au fost analizate în privința impactului pe care îl au asupra structurii generale și funcției

proteinelor. Rezultatele obținute au fost organizate sub forma unei baze de date⁹ public accesibilă.

Autorii în [ARSL'01] afirmă că "algoritmul Smith-Waterman găsește aliniamentul local cu scor maximal dar este incapabil să găsească aliniamentul local cu gradul maxim de similaritate" precum și procentul maxim de potriviri. De asemenea, este menționat faptul că încă nu există un algoritm eficient care să poată oferi un răspuns la întrebarea dacă două secvențe partajează un fragment suficient de lung având un grad de similaritate mai mare de 70%. Ca rezultat, uneori aliniamentele locale oferă un "mozaic" de fragmente bine conservate, conectate artificial de fragmente slab conservate sau neînrudite. O remarcă similară a fost făcută de Zhang în [ZHAN'99]. El spune că asemenea cazuri pot conduce la probleme în compararea secvențelor genomice lungi și predicția comparativă a genelor. Pornind de la această motivație, autorii din [ARSL'01] propun un nou algoritm de comparație al secvențelor numit aliniament local normalizat (NLA) care este capabil să raporteze regiuni cu grad maxim de similaritate. Problema este formulată sub următoarea formă: scorul $s(I, J)$ al aliniamentului local ce implică subșirurile I și J poate fi ajustat prin împărțirea lui $s(I, J)$ cu lungimea totală a regiunilor aliniate, $s(I, J)/(|I| + |J|)$. Problema NLA este astfel de a găsi subșirurile I și J care maximizează $s(I, J)/(|I| + |J|)$ dintre toate subșirurile I și J cu $|I| + |J| \geq T$, unde T este un prag pentru lungimea generală minimă a lui I și J . Cu scopul de a fi important din punct de vedere biologic pentru această problemă, autorii au introdus aici un obiectiv puțin diferit sub forma $s(I, J)/L(|I| + |J| + L)$, pentru un parametru dat L . Ideea a fost de a controla normalizarea variind pe L și au implementat algoritmul lui Dinkelbach [DINK'67] cunoscut ca programare fracțională, ce folosește o metodă parametrică drept tehnică de optimizare. Noua metodă s-a dovedit a fi de 3-5 ori mai înceată decât algoritmul standard Smith-Waterman și sunt alți parametri care influențează rezultatele. Se concluzionează că nu o singură alegere a lui L elimină toate efectele nedorite și relevă toate aliniamentele importante în același timp. Totuși, autorii argumentează că există o valoare a lui L cu care poate fi detectat, folosind algoritmul NLA, un aliniament important dacă are un scor normalizat suficient de mare.

În [ESKI'01], este găsită o motivație biologică pentru folosirea unui model mixt de strămoși comuni pentru a estima probabilitatea de distribuție din familiile de proteine. Autorii consideră metoda a fi mai bună decât mixturile Dirichlet deoarece ea este "simplu de calculat pentru alfabet mari". Rezultatele experimentale sunt evaluate folosind metoda ROC (folosită în secțiunea 5.5) și ele sunt comparabile cu cele obținute cu metodele anterioare. Munca fundamentală care descrie utilizarea mixturilor de "ancestori" comuni în estimarea probabilităților evenimentelor discrete în secvențele biologice este oferită în [JAAK'00].

O metodă mai puțin relevantă dar care utilizează conținutul secvențelor de proteine este dezvoltată în [SARK'02]. Ideea principală este de a extrage relațiile de interacțiune dintre proteine din texte științifice (ex. abstractele din MEDLINE) înrudite cu descrierea proteinelor (ca și cum ar fi selectate de către un expert).

⁹ <http://www.ncl-india.org/trips>

Prin modificarea unui tool destinat procesării limbajului natural, autorii construiesc o analiză bazată pe reguli pentru regăsirea datelor text necesare găsirii proteinelor similare, cu similaritatea exprimând noțiunea de atribute funcționale comune. Sunt selectate chiar și proteinele care nu sunt similare structural dar sunt înrudite prin faptul că provoacă tipuri similare de dezordini neuronale. Rezultatul principal al acestui studiu a constat în faptul că, rulând peste 70 de propoziții conținând 40 *predicații* marcate, au fost identificate 27 relații de interacțiune a proteinelor dintre care 18 corecte. Prin urmare, procentajul nefavorabil (recall-ul) a fost 45% și precizia de 67%. Concluzia finală a acestui studiu este că toolurile lingvistice nu vor putea niciodată substitui cercetarea medicală.

În [WU'03] autorii propun o nouă tehnică pentru compararea perechilor de secvențe biologice pe baza unor modele mici asociate cu regiuni bine conservate în unele secvențe de proteine. Aceste modele sunt folosite ca termeni index în regăsirea informației (information retrieval) fără a considera gap-urile. În acest sens, ei crează grupuri de amino acizi similari dând fiecărui grup un cod. Pentru fiecare secvență de patru amino acizi sunt create toate modelele posibile iar similaritatea bazată pe aceste modele este măsurată folosind oricare din următoarele scoruri:

a) $S_1(p_1, p_2) = c \times Match(p_1, p_2) / (|Pattern(p_1)| + |Pattern(p_2)|)$,
unde $Match(p_1, p_2)$ este setul de modele partajat de secvențele p_1 și p_2 ; c este un factor normalizator care este folosit când se ignoră lungimea proteinelor;

b) $S_2(p_1, p_2) = c \times Match(p_1, p_2) / (|Pattern(p_1)| + |Pattern(p_2)| - |Length(p_1) + Length(p_2)|)$, când cele două secvențe de proteine se cer a avea aceeași lungime.

Această metodă a fost aplicată pentru a construi arbori filogenetici, pentru gruparea proteinelor (proteins clustering) și predicția structurii secundare. Concluzia este că predicția structurii secundare folosind modele se pare să depășească alte metode existente; este ușor de implementat și are o senzitivitate și specificitate relativ mare. În plus, mecanismul de codificare reduce numărul fragmentelor candidate pentru a fi verificate.

În [CAMO'03] scopul este de a găsi similaritățile în structura bazelor de date de proteine extrăgând vectori de caracteristici de triplete de elemente din structura secundară (SSE) și acești vectori de caracteristici sunt indexați folosind o structură multidimensională de indecși. Această structură de indecși este folosită pentru a înlătura automat din domeniul de căutare toate protenele potențial nepromițătoare. Pentru restul de proteine este folosit un tool de căutare de aliniamente (numit VAST) pentru efectuarea aliniamentului perechilor [MADE'95]. Tehnica este numită de autori PSI (Protein Structure Index) și include un pas de modelare statistică pentru luarea deciziei de a considera sau nu proteinele drept candidați în căutarea subsecvență. Experimentele au fost efectuate comparând performanțele lui PSI și VAST pe aceleași date. Acestea arată că noua tehnică este mai rapidă în timp ce se menține aceeași senzitivitate.

Un tip diferit de măsuri de similaritate, ca tooluri pentru explorarea ontologică a genelor, este folosit în [LORD'03]. Pe baza ontologiei genelor (GO), descrieri ca funcția moleculară, procesul biologic și componenta celulară sunt

comparate via trei măsuri de similaritate semantică. Tipul de date este adnotat și provine din setul de date GO¹⁰ disponibil accesului public. GO reprezintă secvențele adnotate din cadrul unui graf aciclic (DAG - Dyrected Acyclic Graph) ce constă dintr-un număr de termeni reprezentați ca noduri în graf și conectați de muchii reprezentând relațiile dintre ei. Termenii pot avea "părinți" multipli precum și mai mulți "copii". Acest mod de manipulare a secvențelor poate fi folosit pentru a reprezenta toate bazele de date [LORD'03]. În această lucrare experimentele au fost realizate numai pentru acele asocieri care au fost identificate între termenii GO și proteinele din SWISS-PROT-Human. Toate măsurile testate se bazează pe conținutul de informație al fiecărui termen, care este definit ca numărul de apariții în corpus al fiecărui termen sau termen copil și este exprimat ca probabilitate. Fiecare măsură folosește conținutul de informație al "părinților" partajați ai celor doi termeni comparați. Remarcile concluzive sugerează că toate cele trei măsuri dovedesc o corelare puternică dintre similaritatea secvențelor și similaritatea semantică a funcției moleculare și nici una nu se dovedește a avea o superioritate clară față de celelalte.

O abordare diferită a similarității este utilizată în [KRAS'04]. Autorii propun folosirea metricii universale de similaritate (USM) pentru structurile proteinelor. Este menționat faptul că această metodă a fost introdusă și îmbunătățită de [LI'01] și [LI'03] respectiv, și se bazează pe conceptul de complexitate Kolmogorov. Pe scurt, citez "complexitatea Kolmogorov este o măsură obiectivă a cantității de informație conținute întru obiect dat" și se menționează că o măsură apropiată este complexitatea Kolmogorov condiționată a lui o_1 fiind dat o_2 , descrisă de relația:

$$K(o_1|o_2) = \min\{P \mid P \text{ a program and } U(P, o_2) = o_1\}, \quad (2.2-1)$$

unde U este o Mașină Turing Universală, necesară pentru a produce un anume obiect [LI'97]. Această relație este folosită pentru definirea distanței informației (ID-Information Distance) dintre două obiecte după cum urmează [KRAS'04]:

$$ID(o_1|o_2) = \max\{K(o_1|o_2), K(o_2|o_1)\}. \quad (2.2-2)$$

Similaritatea structurală dintre perechile de proteine a fost aplicată pe patru seturi de date diferite și toate distanțele dintre perechi sunt stocate într-o matrice de similarități care este apoi alimentată de un tool de grupare deja implementat și public accesibil¹¹. După experimente extensive, concluzia a fost că USM este robust, poate diferenția între familii de proteine și sub-familii dar nu oferă vre-un indiciu despre locația exactă a (di)similarității de-a lungul celor două secvențe.

Folosind abordarea SVM (Support Vector Machines), metodele discriminative sunt luate în considerare de către autorii [SAIG'04], care le dovedesc a fi cele mai eficiente metode pentru problema recunoașterii superfamiliiilor în clasificarea structurală a proteinelor în baza de date SCOP. Ei explică performanța ridicată a SVM ca un rezultat al funcțiilor nucleu (kernel-functions) folosite pentru a cuantifica

¹⁰ <http://www.godatabase.org/dev>

¹¹ <http://www2.biology.ualberta.ca/jbrzusto/cluster.php>.

similaritatea dintre secvențe (ei propun noțiunea de nucleu pentru șirul care este adaptat secvențelor biologice și este numit nucleu de aliniament local). Aceste funcții nucleu măsoară similaritatea dintre două secvențe prin însumarea scorurilor obținute din aliniamente locale ale secvențelor ținând cont de gap-uri. Experimentele au fost efectuate folosind 4352 proteine din baza de date Astral¹², grupate în familii și super-familii. Urmând procedura de marcaj din [LIAO'03] pentru fiecare familie, domeniile proteinelor din cadrul familiei au fost considerate exemple pozitive de **testare** și domeniile proteinelor din cadrul superfamiliei dar dinafara familiei au fost considerate exemple pozitive de **antrenare**. Exemple negative au fost luate dinafara clasei secvențelor pozitive și au fost împărțite aleator în seturile de testare și antrenare. Evaluarea metodei a fost făcută folosind scorurile furnizate de caracteristicile ROC și rezultatele de concluzie ale autorilor au fost că această metodă depășește performanțele altor metode considerate „state of the art” și care folosesc SVM.

Analiza semantică latentă (LSA) este o metodă bazată pe tehnica SVD (Singular Value Decomposition). Ea este o abordare extrem de utilă în procesarea limbajului natural pentru generarea sumarelor (rezumatelor), compararea documentelor, generarea de tezaure lingvistice și mai departe pentru regăsirea informației (information retrieval) [BELL'00], [LAND'98], [BOGA'03a]. Modul în care LSA capturează relațiile conceptuale din text, pe baza distribuției cuvintelor în documente, autorii în [GANA'04] îl folosesc pentru capturarea propensităților (tendențelor) structurii secundare în secvențele de proteine, folosind diferite vocabulare. Ei consideră documentele d_1, d_2, \dots, d_{N_1} a fi segmente nesuprapuse de proteine pentru care sunt cunoscute categoriile structurale C_1, C_2, \dots, C_{N_1} și t_1, t_2, \dots, t_{N_2} segmente de test nesuprapuse cu lungime cunoscută pentru care trebuie prezisă structura secundară. Structura secundară a datelor de test este prezisă folosind ca metodă un model de clasificare kNN (k Nearest Neighbour). Pentru fiecare segment de test t_i , similaritatea cosinus a lui t_i este calculată față de toate segmentele de antrenare d_1, d_2, \dots, d_{N_1} și sunt identificate cele k segmente având similaritatea cu

¹² www.cs.columbia.edu/compbio/svmpairwise

similaritatea lui t_i maximă. Aceste k segmente sunt k NN al lui t_i . Categoria prezisă a lui t_i este categoria structurală căreia îi aparțin cei mai mulți k NN. Acest proces este repetat pentru fiecare dintre segmentele de test. În experimente, autorii consideră ca vocabular cei 20 de amino acizi, grupurile chimice și tipurile de amino acizi. Ca documente ei folosesc structurile de: helix, fâșii și răsuciri (helix, strand and coil structures) identificate în structura secundară a fiecărei proteine. După metoda LSA a fost folosită metoda VSM (Vector Space Model) pentru aceleași modele și date. Subsecvențele de proteine au fost tratate ca documente în matricea vocabular/documente. Rezultatele sunt comparate folosind măsurile: precizie și recall preluate din teoria prelucrării/regăsirii informației (information retrieval theory). Autorii concluzionează că alfabetele distincte diferă prin cantitatea de informație pe care o poartă. Astfel, helixurile și fâșiile (helices and sheets) au fost cel mai bine clasificate folosind LSA cu tipurile de amino acizi ca vocabulare, în timp ce răsucirile (coils) sunt caracterizate cu o mai mare precizie când amino acizii sunt folosiți ca vocabulare și pentru analiză se folosește metoda SVM.

2.3. Tehnici de realizare a aliniamentului de secvențe biologice

În cele ce urmează vor fi considerate metode clasice de aliniament și anume, cei mai larg aplicați algoritmi în cercetarea secvențelor de gene/proteine. Pe baza numărului de secvențe implicate în aliniament, aceste metode sunt împărțite în algoritmi pentru compararea de perechi de secvențe sau algoritmi pentru compararea de secvențe multiple. O mare parte din temele acestui capitol urmează descrierea făcută de unul dintre fondatorii metodeleor de aliniament al secvențelor biologice, și anume Michael S. Waterman [WATE'84].

Combinatorica aliniamentelor

Fie secvențele $a = a_1, a_2, \dots, a_n$, și $b = b_1, b_2, \dots, b_m$. Un aliniament poate fi produs prin creșterea lungimii fiecărei secvențe cu inserarea lui Δ . Dacă lungimea unui asemenea aliniament este L , atunci acesta ar putea fi scris

$$\begin{array}{ccccccc} & * & & * & & & * \\ a & 1 & a & 2 & & & a & L \\ & * & & * & & & * \\ b & 1 & b & 2 & \dots & & b & L \end{array}$$

unde subsecvența $a^*(b^*)$ de elemente neegale cu Δ este $a(b)$.

Interesul este în calcularea numărului de aliniamente. Dacă acesta nu este prea mare, atunci o căutare directă este posibilă pentru căutarea aliniamentelor optimele. După cum se va putea observa însă, numărul crește foarte repede.

Un aliniament a lui $a = a_1, a_2, \dots, a_n$, cu $b = b_1, b_2, \dots, b_m$ poate sfârși în unul din următoarele trei moduri:

$$\begin{array}{ccc} \dots a_n & \dots a_n & \dots \Delta \\ \dots \Delta & \dots b_m & \dots b_m \end{array}$$

(alinieră a lui Δ peste Δ este eliminată deoarece nu contribuie cu vreo informație). Dacă $f(n,m)$ calculează configurații generate de aliniamente terminate recursiv ca mai sus, atunci

$$f(n,m) = f(n-1,m) + f(n-1, m-1) + f(n,m-1) .$$

Numerele generate de această ecuație recursivă sunt cunoscute ca numere Stanton-Cowan [STAN'70], unde ele provin din calcularea volumului unei sfere de rază m în n dimensiuni folosind metrica Lee. Asimptotele sunt date pentru o generalizare a acestor numere de către H.T. Laquer [LAQU'81]. El arată că

$$f(n,n) \approx (1 + \sqrt{2})^{2n+1} \sqrt{n}. \quad (2.3.-1)$$

După o examinare atentă, s-a observat că $f(n,m)$ supracalculează aliniamentele sau cel puțin o definiție rezonabilă a aliniamentelor. Spre exemplu, cele două aliniamente

$$\begin{array}{ccc} A\Delta & & \Delta A \\ & \text{și} & \\ \Delta T & & T\Delta \end{array}$$

pot să nu fie distincte într-un sens biologic. Pentru a proiecta o recursivitate care nu calculează dublu aceste ștergeri „tandem”, fie $g(n,m)$ numărul unor asemenea aliniamente. Dacă un aliniament se termină într-un Δ există trei posibilități:

$$\begin{array}{ccc} \dots a_{n-1} a_n & \dots a_{n-1} a_n & \dots \Delta a_n \\ \dots b_m \Delta & \dots \Delta \Delta & \dots b_m \Delta; \end{array}$$

și dacă un aliniament se termină în b_m , există trei posibilități:

$$\begin{array}{ccc} \dots a_n \Delta & \dots \Delta \Delta & \dots a_n \Delta. \\ \dots b_{m-1} b_m & \dots b_{m-1} b_m & \dots \Delta b_m \end{array}$$

Prin urmare,

$$g(n, m) = g(n-1, m) + g(n, m-1) + g(n-1, m-1) - g(n-1, m-1)$$

sau

$$g(n,m) = g(n-1, m) + g(n, m-1) ,$$

o revenire a caracterului mai simplu decât recursivitatea Stanton-Cowan. Condițiile de limită pentru recursivitatea adoptată de [WATE'84] trebuie considerate în aceasta fază.

Evident,

$$g(0, 0) = g(1,0) = g(0, 1) = 1,$$

astfel soluția explicită poate fi scrisă

$$g(n,m) = \binom{n+m}{n}. \quad (2.3-2)$$

Cu formula lui Stirling,

$$g(n, n) \approx 2^{2n} / (4\sqrt{nn}) \quad (2.3-3)$$

pentru $n = 1000$, $g(n, n) > 10^{600}$ și examinarea directă a tuturor aliniamentelor este imposibilă. Acesta este unul din motivele care a dus la dezvoltarea de algoritmi eficienți.

2.3.1 Metode vizuale

Există o metodă importantă de analiză a secvențelor care este cel mai bine descrisă de expresia „doar uită-te la...”. Un analist poate observa o succesiune ca GCGCGC... și să găsească o caracteristică de interes fără vreo metodă matematică sau computațională sofisticată. Este posibil să se formuleze o problemă analitică dintr-o asemenea observație: „Cât de probabilă este observația unui asemenea eveniment în secvențe aleatoare?” Un biolog este puțin probabil surprins, informat sau impresionat dacă i se spune că apariția succesiunii GCGC... a fost neobișnuită. În schimb, problema funcției regiunii GC captează atenția biologului. Doar privirea secvențelor poate fi utilă.

Metoda „dot matrix” tradusă ca „matricea punct” este o metodă vizuală larg răspândită care în general folosește calculatorul. În acest caz, este formată o matrice $M = (m_{ij})$, unde $m_{ij} = 0$ dacă al i -lea element al secvenței a este diferit de al j -lea element al secvenței b altfel $m_{ij} = 1$. Execuția potrivirilor exacte arată ca diagonale de elemente 1. Matricea punct a fost descoperită independent de câteva ori, spre exemplu Maizel și Lenk (1981), Novotny (1982), Harr et al. (1982), Jagadeeswaran și McGuire (1982), Gibbs și McIntyre (1970).

Multe dintre aceste implementări filtrează potrivirile pentru a afișa doar execuția unui element sau mai multe. Este posibilă combinarea unor metode mai sofisticate pentru filtrarea potrivirilor.

Metoda *matricei punct* este utilă în localizarea regiunilor de mare potrivire dintre două secvențe. Este normal astfel de a fi interesați de probabilitatea distribuției celei mai lungi potriviri dintre două secvențe. Această distribuție va permite analistului să localizeze potrivirile semnificative.

2.3.2. Algoritmi pentru aliniamentul perechilor de secvențe

Problema obținerii aliniamentului global optimal dintre două secvențe permițând gap-uri a fost rezolvată la început folosind un algoritm de programare dinamică. Metoda produsă este cunoscută în analiza secvențelor biologice ca algoritmul Needleman-Wunch [NEED'70], după numele celor doi autori. Ideea principală este de a construi un aliniament optimal folosind soluții anterioare pentru aliniamente optimale ale unor subsecvențe mai mici. Algoritmul a fost îmbunătățit mai departe de Sellers [SELL'74] și mai târziu optimizat de O. Gotoh [GOTO'82].

Cu intenția de a căuta cel mai bun aliniament local util pentru aflarea dacă două proteine împart un domeniu comun, pentru a compara secțiuni extinse ale secvențelor genomice de ADN sau pentru a compara două secvențe foarte divergente, a fost dezvoltat algoritmul Smith-Waterman [SMIT'81a].

O metodă alternativă aliniamentului, bazată pe algoritmi de programare dinamică, folosește modelul stărilor automatului finit (Finite State Automata) și utilizează tooluri interactive construite automat [SEAR'95].

Metodele de aliniament amintite sunt considerate corecte în privința unui scor optimal. Ele se bazează pe crearea dinamică a unor matrice și ca rezultat complexitatea timp este de ordinul $O(nm)$, (i.e. produsul lungimii secvențelor) [DURB'98]. Aceasta înseamnă că pentru o bază de date mare de proteine/secvențe, timpul de execuție poate deveni considerabil de mare. Pentru ajustarea acestei probleme, au fost introduse tehnicile euristice. În aceste tehnici există un compromis între acuratețe și eficiență [LIAO'03]. Dintre cele mai cunoscute metode sunt BLAST [ALTS'90], [ALTS'97] și FASTA [PEAR'88]-[PEAR'98], care au fost dezvoltate pentru detectarea omologiilor în baze de date mari cu ajutorul scorurilor obținute din aliniamente locale. Ambele au fost incorporate în pachete de programe accesibile public^{13,14}.

În toate tehnicile menționate mai sus se cere o atenție suplimentară în alegerea modelului de calcul al scorului de similaritate atâta timp cât scorul total asociat unui aliniament va fi o sumă de termeni corespunzând diferențelor dintre reziduurile perechilor aliniate plus termenii corespunzători execuției operațiilor de inserare/ștergere. Schema de calcul incorporată este bazată pe una dintre matricele de substituție disponibile (fiecare având propriile caracteristici). Cele mai populare dintre ele sunt PAM, BLOSSUM, GONNET și matricea de identitate ADN. Semnificația scorurilor obținute este calculată făcând comparații cu un model Baiesian sau folosind abordarea clasică a comparării statistice a valorilor extreme ale scorurilor de potrivire [DURB'98].

Pentru o analiză mai detaliată vor fi prezentate în continuare metodele cele mai importante de aliniere a perechilor de secvențe.

Problematika acestui subcapitol ar putea fi descrisă ca cea a „găsirii numărului cel mai mic” de pași care să schimbe o secvență în cealaltă. Dacă cele două secvențe ar fi formate din acizi nucleici: $a = a_1, a_2, \dots, a_n$, și $b = b_1, b_2, \dots, b_m$, unde a_i și b_j sunt unul dintre nucleotidele adenina (A), timina (T), guanina (G) sau citozina (C). Aceste secvențe sunt cuvinte finite formate din cele patru litere ale alfabetului folosit iar pașii menționați mai sus corespund evenimentelor evoluționare ce pot altera secvența. Cel mai simplu set de evenimente evoluționiste constă din mutații, unde o literă (sau mai multe) este substituită de o alta și inserări sau ștergeri unde o literă sau mai multe sunt inserate respectiv șterse.

Dacă $a = AATAG$, atunci o substituție a lui T pentru $a_2 = A$ transformă a în b :

$$a = AATAG \rightarrow b = ATTAG.$$

¹³ <http://fasta.bioch.virginia.edu/>

¹⁴ http://www.ncbi.nlm.nih.gov/Class/BLAST/blast_course.short.html

Această corespondență este arătată de obicei prin reprezentarea într-un aliniament:

a : AATAG;
 b : ATTAG.

Dacă $a_4=A$ este șters, transformarea în $c = AATG$ este reprezentată de

a : AATAG;
 c : ATTΔG,

unde ștergerea lui $a_4=A$ este indicată de un „Δ” inserat în c . Prin corespondență, inserarea lui C între a_3 și a_4 este reprezentată prin

a : AATΔAG
 b : AATCAG.

Kruskal și Sankoff [KRUS'83] plasează aceste probleme într-un context general. O altă formulare generală este în [WATE'76] care a mai fost amintită și presupune următoarele: Fie S setul de cuvinte finite dintr-un alfabet, incluzând cuvântul vid și fie $\tau = \{ T \mid T : S \rightarrow S \}$ un set de transformări care include transformările identice. Interesul este în transformările T_1, T_2, \dots, T_k din τ astfel încât

$$T_1 \circ T_2 \circ \dots \circ T_k (a) = b. \quad (2.3.2-1)$$

Pe baza acestei exprimări pot fi formulate probleme variate. Pentru problema celui mai mic număr de pași trebuie căutat k minim unde $\tau = \{\text{mutații de tipul unei singure litere, inserții și ștergeri}\}$. Funcția obiectiv poate fi schimbată iar setul de transformări mărit. Există un echilibru între realitatea biologică și algoritmi computaționali.

Unul dintre algoritmi timpurii de comparare a șirurilor s-a datorat lui Levenstein [LEVE'66] cu toate că munca sa nu a influențat dezvoltările menționate în continuare. Lucrările lui Fitch și Margoliash [FITC'67] și [FITC'69] au adus problema comparării secvențelor în atenția unui mare număr de persoane.

2.3.3. Metode fundamentale de programare dinamică

Un algoritm bazat pe programare dinamică furnizează o soluție elegantă $O(n^2)$, n fiind lungimea celei mai mari dintre cele două secvențe. Programarea dinamică este o tehnică algoritmică ce folosește cunoștințe învățate din subprobleme anterioare în calcularea soluției unei probleme mai mari.

La început se aplează la intuiție în spatele programării dinamice pentru calcularea aliniamentului cel mai bun (în termenii distanței minime).

Până acum, secvențele $a = a_1, a_2, \dots, a_n$ și $b = b_1, b_2, \dots, b_m$ au fost folosite pentru a genera aliniamente,

$$\begin{matrix} a_1^* & a_2^* & \dots & a_L^* \\ b_1^* & b_2^* & \dots & b_L^* \end{matrix},$$

unde subsecvențele de elemente a^* și b^* neegale cu Δ sunt secvențele originale. Pentru cazul similarității și distanței, ecuația de revenire pentru $f(n,m)$ (din introducerea secțiunii 2.3) este modificată pentru a oferi un mod eficient pentru calcularea similarității și distanței.

Din punct de vedere istoric, aceste metode au început în biologie cu Needleman și Wunsch [NEED'70] care au prezentat algoritmul descris în cele ce urmează. După aceea, Sankoff [SANK'72] și Sankoff și Sellers [SANK'73] găsesc metode de programare dinamică pentru aliniamente optimale cu un număr dat de indels. Sellers [SELL'74a,b] vine cu algoritmi pentru calcularea distanței $D(a,b)$, problemă ridicată de Ulam [ULAM'72]. Gordon [GORD'73] și Delcoigne și Hansen [DELC'75] au contribuții utile în compararea secvențelor. Toate aceste metode sunt incluse într-o clasă de tehnici cunoscute ca programare dinamică ce a fost introdusă de Richard Bellman. În Byers și Waterman [BYER'84] poate fi găsită o discuție generală pe această temă.

Considerând alfabetul $\{A, C, G, T, \Delta\}$, lărgit pentru a include Δ , fie o funcție de pondere $s(a, b)$ definită pe perechi de litere din alfabet. O funcție de tipul folosit de Needleman-Wunsch este

$$s(a,b) = \begin{cases} 1 & \text{daca } a = b, \\ 0 & \text{daca } a \neq b, a \neq \Delta, b \neq \Delta, \\ -1.5 & \text{daca } a \neq b \text{ si unul dintre } a, b = \Delta. \end{cases} \quad (2.3.3-1)$$

Aici potrivirile primesc pondere pozitivă iar indels primesc ponderi negative. Următoarele declarații apar în [SMIT'81b] unde dovezile merg de-a lungul ecuației pentru $f(n,m)$, așa cum a fost menționat în secțiunea 2.3. La început se definește unde se ia maximum peste toate aliniamentele.

$$S(a,b) = \max \sum_{k=1}^L s(a_k^*, b_k^*) \quad (2.3.3-2)$$

Algoritmul Needleman-Wunsch

Fie

$$S_{0j} = \sum_{k=1}^j s(\Delta, b_k), \quad S_{00} = 0 \quad \text{si} \quad S_{i0} = \sum_{k=0}^i s(a_k, \Delta)$$

Dacă $S_{ij} = (a_1 a_2 \dots a_i, b_1 b_2 \dots b_j)$, atunci

$$S_{ij} = \max \{ S_{i-1,j} + s(a_i, \Delta), S_{i-1,j-1} + s(a_i, b_j), S_{i,j-1} + s(\Delta, b_j) \}$$

Dacă o funcție distanță inițială d este specificată pe $\{A, C, G, T, \Delta\}$, atunci este obținut un rezultat similar. O asemenea funcție este

$$d(a,b) = \begin{cases} 1 & \text{daca } a = b \\ 0 & \text{daca } a \neq b. \end{cases}$$

Se definește, ca mai sus,

$$D(a,b) = \min \sum_{k=1}^L d(a_k^*, b_k^*) \quad (2.3.3-3)$$

Algoritmul Sellers

Fie $D_{0j} = \sum_{k=1}^j d(\Delta, b_k)$, $D_{00} = 0$ și $D_{i0} = \sum_{k=0}^i d(a_k, \Delta)$

Dacă $D_{ij} = S(a_1 a_2 \dots a_i, b_1 b_2 \dots b_j)$, atunci

$$D_{ij} = \min \{ D_{i-1,j} + d(a_i, \Delta), D_{i-1,j-1} + d(a_i, b_j), D_{i,j-1} + d(\Delta, b_j) \}$$

Fiecare dintre cei doi algoritmi au timpul computațional proporțional cu

$$\sum_{i=1}^n \sum_{j=1}^m 1 = nm. \quad (2.3.3-4)$$

Spațiul cerut pentru calcularea $D(a,b)$ este $\min\{n,m\}$ dar $D(a,b)$ este câteodată cu valoarea în afara setului asociat de aliniamente optime.

Există două tehnici de a produce aliniamente. Prima este de a salva pointeri la fiecare (i,j) pentru a arăta care dintre $D_{i-1,j}$; $D_{i-1,j-1}$; $D_{i,j-1}$ sunt folosiți în calcularea lui $D_{i,j}$ optim. Pointerii sunt salvați în timpul calculului înainte (forward) astfel, în timpul parcurgerilor înapoi (tracebacks), aceștia pot fi urmați pentru a produce un aliniament optimal. În cazul în care există *optimi* multipli, pointerii nefolosiți pot fi păstrați și în această manieră (breath-first search with stacking), pot fi produse toate aliniamentele optimale. În al doilea caz, dacă pointerii nu sunt salvați, recalculând care dintre $D_{i-1,j}$; $D_{i-1,j-1}$; $D_{i,j-1}$ rezultă în $D_{i,j}$ va fi simplu de realizat dacă matricea D este salvată. În oricare din cele două cazuri, spațiul de stocare necesar este $O(nm)$.

2.3.4. Extensii ale metodelor de bază

Cele mai importante transformări ale evoluției tratate în compararea secvențelor sunt mutațiile unei singure baze și cele determinate de *indels*. Până acum au fost tratate doar situațiile indels.

Metodele clasice de comparare a secvențelor de AND au presupus că secvențele pot fi numai "mutate" prin operații care acționează individual asupra nucleotidelor i.e. substituții, inserări și ștergeri. Mai recent însă, studii adiționale au considerat rearanjarea evenimentelor genomice la scară largă cum ar fi inversiunile, transpozițiile, și translocările. În continuare se va extinde cazul pentru indels

mai lungi și se va descrie o problemă a inversiunii lungi de segmente ale unei secvențe.

În timp ce indels a multe baze, fie 100, ar putea fi suma a 100 de baze indels, explicația probabilă este că a fost un singur eveniment. Într-un studiu de aliniament de parametri, Fitch și Smith [FITC'83] arată că pentru anume secvențe, sunt necesare indels mai lungi pentru a obține aliniamentul corect. Ponderile indels-urilor mai lungi nu ar trebui să fie socotite ca sumă de indels singulare.

Fie x_k ponderea aleasă pentru un indel de k litere, $k \geq 1$. Următoarele rezultate apar în [WATE'76]. Dacă $x_1 \leq x_2 \leq \dots$ și d este o metrică în mulțimea elementelor secvențelor, atunci D este o metrică pe mulțimea secvențelor.

Algoritmul Waterman-Smith-Byers

$$\text{Fie } D_{0j} = x_j, \quad D_{0j} = x_j \quad D_{00} = 0, \quad \text{si } D_{ij} = D(a_1 a_2 \dots a_j, b_1 b_2 \dots b_j)$$

Atunci,

$$D_{ij} = \min \left\{ D_{i-1, j-1} + d(a_i, b_j), \min_{k \geq 1} \{ D_{i, j-k} + x_k \}, \min_{k \geq 1} \{ D_{i-k, j} + x_k \} \right\}$$

Un algoritm corespunzător va calcula un S generalizat.

Este important a se nota faptul că timpul de calcul este mărit cu

$$\sum_{i=1}^n \sum_{j=1}^m (ij) = O(nm^2 + n^2m) \quad (2.3.4-1)$$

Un algoritm $O(n^3)$, pentru două secvențe de lungime $n = 1000$, este un preț semnificativ de plătit pentru multiple indels. O abordare pentru a evita aceasta este de a presupune x_k liniar. Acesta este,

$$x_k = a + bk..$$

Pentru calculele structurii secundare această presupunere este exploatată în [WATE'76] și [KANE'82]. Gotoh în [GOTO'82] (1982) derivă un algoritm înrudit pentru x_k liniar cu timp de execuție $O(nm)$. Taylor [TAYL'84] vine cu mai multe rezultate pe linia deschisă de Gotoh. Desigur, după x_1 și x_2 , un asemenea x_k se comportă mai mult ca un indel singur. Este prin urmare de dorit să se extindă algoritmul, în mod special să prelucreze funcțiile indel astfel încât

$$x_k = a + b \log(k).$$

Aceasta a fost realizată în [WATE'84] și are aplicație în problemele de structură secundară unde indels sunt analoge salturilor interioare și salturilor ramificate. Algoritmul pentru x_k *concav* are timp de execuție $O(nm)$.

Problema includerii inversiunilor este foarte interesantă. Interschimbarea a două litere adiacente este o transformare propusă cu ajutorul computerelor. Wagner [WAGN'83] a arătat că această transformare poate fi inclusă cu timpul de calcul $O(n m 4^0)$, unde

$$a \leq \min \{4 \max d(a, b), 2x_1\} / \gamma + 1,$$

unde y este costul transpunerii. Includerea inversiunilor lungi pare să fie o activitate dificilă. În mod cert, aceste inversiuni au loc în secvențele de ADN și ar trebui incluse. O problemă relevantă pentru biologie este de a include inversiuni într-un algoritm unde există un singur cost al inversiunii plus distanța dintre segmente. Richard et al. [REIC'73], Wong et al. [WONG'74] și Cohen et al. [COHE'75] produc o serie de algoritmi de programare dinamică motivați de mașina Turing și teoria informației. Acești algoritmi pot fi văzuți ca și cazuri speciale ale algoritmilor generali de programare dinamică [WATE'84].

Notă: O exemplificare a modului de aplicare al acestor algoritmi dinamici este ilustrată în ANEXA 2.

2.3.5. Algoritm de programare dinamică a lui Ukkonen

În timp ce algoritmul de programare dinamică $O(n^2)$ este mult mai rapid decât forța brută a algoritmului $O(2^{2n})$, el devine productiv pentru un n foarte mare. E. Ukkonen [UKKO'83], [UKKO'84] vine cu algoritmi semnificativ mai rapizi pentru calcularea distanței dintre două secvențe. Algoritmii, evidențiați în cele ce urmează, calculează distanța s dintre secvențe de lungime n și m , de-a lungul aliniamentului, în timp și spațiu de stocare de $O(s \cdot \min\{n, m\})$ (cu spațiu de stocare de $O(s^2)$ în unele cazuri). Dacă nu este dorit nici un aliniament, stocarea este de $O(s)$. Cel mai rău caz de comportament este echivalent cu algoritmul standard, când pentru s cu valoare mică îmbunătățirea este dramatică. Dacă nu este dorită nici o distanță mai mare de un prag t atunci timpul nu este mai mare de $O(t \cdot \min\{n, m\})$.

Un algoritm similar în esență cu cel al lui Ukkonen este propus de J.W. Ficket [FICK'84], dar acesta din urmă nu este așa de eficient în termeni de timp. (Ambii autori tratează cazul pentru indels singulare.)

Ukkonen [UKKO'83] prezintă algoritmul său pentru nepotiviri, indels singulare, și transpunerii de 2 litere. Mai jos sunt aceleași rezultate pentru nepotiviri și indels multiple. Fie ca mai sus:

$$D_{ij} = \min \left\{ D_{i-1, j-1} + d(a_i, b_j), \min_{k \geq 1} \left\{ D_{i, j-k} + x_k \right\}, \min_{k \geq 1} \left\{ D_{i-k, j} + x_k \right\} \right\},$$

și se presupune $d(a, b) = 1$ doar dacă $a = b$. Ukkonen demonstrează o lemă cheie care se arată a fi adevărată pentru cazul general de indels multiple.

Lema U.: Pentru toți (i, j) , $D_{ij} - 1 \leq D_{i-1, j-1} \leq D_{ij}$.

Demonstrație: Demonstrația este prin inducție pe $i+j$. Partea stângă reiese imediat din recursivitate. Dacă $D_{ij} = D_{i-1, j-1} + d(a_i, b_j)$ atunci urmează $D_{ij} \leq D_{i-1, j-1}$. Altfel, fără a pierde caracterul de generalitate, se presupune $D_{ij} = D_{i-k, j} + x_k$. Ipoteza inducției implică $D_{i-k, j} \geq D_{i-(k+1), j-1}$ astfel încât $D_{i-k, j} \geq D_{i-(k+1), j-1} + x_k \geq D_{i-1, j-1}$.

Lema U care este elementară, este cheia metodei lui Ukkonen. Ea declară că $D_{i, i+c}$ este o funcție a lui i nedescrescătoare. Aceasta implică o structură pentru matricea $D_{i, j}$: ea este de forma unei văi cu ridicături de-a lungul liniilor constante $i-j$. Cea mai joasă ridicătură este $D_{00} = 0$. Mai jos, atenția este focalizată pe limitele schimbărilor când $D_{i, i+c} = k$ se schimbă cu $D_{i+1, i+1+c} = k+1$.

Se presupune că toate indels au costul de 1, i.e. $x_k = k$. Ideea de bază a algoritmului lui Ukkonen [UKKO'83] este de a începe la $D_{00} = 0$ și a se extinde de-a lungul $j - i = 0$ până când $D_{ij} = 1$. În general, vor exista $2k + 1$ limite ale regiunii $D_{ij} \leq k$. Fiecare limită $j - i = c$ este extinsă până când $D_{ij} = k + 1$ (pentru $i - j = c$). Extinderea limitei la $k + 1$ poate fi determinată din limitele pentru $k, k - 1, \dots$ și teste ale $a_i = b_j$. Această procedură este urmată până când este atins $D_{n,m}$. Dacă $D_{n,m} = s$, este clar că nu au fost calculate mai mult de $(2s + 1)\min\{n, m\}$ intrări. Este suficient doar să se stocheze aceste limite astfel spațiul necesar de stocare este $O(s^2)$.

2.3.6. Aliniamente aproape optimale

Cu toate că algoritmi localizează aliniamentele optimale, ponderile sunt determinate de utilizator. Calibrarea ponderilor prin aliniamente deja cunoscute a fi corecte poate fi făcută dar nici un set de ponderi nu poate fi considerat a fi absolut corect. Chiar dacă ponderile sunt corect alese, constrângeri biologice necunoscute pot determina ca aliniamentul adevărat să fie diferit de cel optimal, generat de computer.

Prin urmare, o problemă naturală este de a găsi toate aliniamentele din cadrul distanței optimale. Motivația vine din faptul că aliniamentul corect ar trebui să fie lângă cel optimal și că biologul sau cel care analizează secvențele să îl poată recunoaște. Această problemă a fost soluționată în [WATE'83] și primește un tratament mai general în Byers și Waterman [BYER'84].

În mod formal, problema este de a găsi toate aliniamentele cu distanța sau scorul de aliniament pâna la ϵ al distanței $D_{n,m}$ dintre cele două secvențe a și b . Se presupune că toți $D_{n,m}$ sunt calculați și stocați.

La poziția (i, j) se presupune o parcurgere înapoi (traceback) de la (n, m) la $(0, 0)$ ce poate rezulta într-un aliniament cu scor mai mic sau egal cu $D_{n,m} + \epsilon$. Scorul de la (n, m) la (i, j) dar fără a include (i, j) , este de $T_{i,j}$. $T_{i,j}$ este suma posibilului aliniament non-optimal de a atinge (i, j) . De la (i, j) , ca de obicei, sunt posibili trei pași: $(i-1, j)$, $(i-1, j-1)$ și $(i, j-1)$. Fiecare pas este într-un aliniament dorit dacă și numai dacă

$$\begin{aligned} T_{i,j} + d(a_i, \Delta) + D_{i-1,j} &\leq D_{n,m} + \epsilon \\ T_{i,j} + d(a_i, b_j) + D_{i-1,j-1} &\leq D_{n,m} + \epsilon \\ T_{i,j} + d(\Delta, b_i) + D_{i-1,j} &\leq D_{n,m} + \epsilon. \end{aligned}$$

Aliniamentele multiple aproape optimale pot fi produse prin gruparea direcțiilor neexplorate.

Un studiu al sensibilității aliniamentului de secvențe la ponderi și multiple indels a fost rerealizat de Fitch și Smith [FITC'83].

2.3.7. Metoda regiunilor

O examinare a matricelor punct (abordate în secțiunea 2.3.1) poate sugera construirea unei liste cu regiuni de potrivire. Atunci când un aliniament este doar o submulțime ordonată a unei asemenea liste, algoritmi pot fi divizați pentru găsirea aliniamentelor optimale. Această abordare trebuie să fi fost continuată de câteva grupuri independente deoarece este importantă găsirea de regiuni lungi de potrivire. Pentru predicția structurii secundare, Studnicka et al. [STUD'78] urmează o asemenea direcție. Martinez [MART'80], [MART'83] vine cu un algoritm mai matematic, care de fapt adaptează algoritmul său pentru structura secundară în aliniament de secvențe. O încercare de descriere a acestui algoritm va fi dată în cele ce urmează.

Pentru început este necesară o listă L de regiuni de potrivire. O regiune R este definită a fi un triplet de forma $(w; i, j)$, care înseamnă o potrivire a cuvântului w care începe la $a_i = b_j$. În alte cuvinte, dacă $l = |w|$ este lungimea cuvântului,

$$a_i = b_j, a_{i+1} = b_{j+1}, \dots, a_{i+l-1} = b_{j+l-1}.$$

Pentru a obține o asemenea listă, Martinez [MART'83] mai întâi concatenează secvențele într-una singură S . El folosește tehnica sortării rapide a lui S folosind ordinea lexicografică. O primă sortare grupează toate elementele egale ale lui S împreună. A doua sortare operează asupra fiecărui asemenea grup de elemente egale și grupează împreună elemente care sunt urmate de elemente egale în secvența originală S . La sfârșitul celei de-a k sortări, două elemente ale lui S permutat vor aparține la același grup dacă și numai dacă locațiile lor i și j în șirul original S sunt astfel încât elementele la locațiile $i+l$ și $j+l$ sunt egale pentru $l = 0, 1, \dots, k-1$.

După cum este arătat în Martinez [MART'83], viteza acestei proceduri de sortare pentru a genera regiuni este, în cazul așteptat, de ordinul $N \log N$, unde N este lungimea totală a secvențelor concatenate. Procedura este, prin urmare, comparabilă în viteză cu metoda standard utilizată de calculatoare pentru construirea „arborilor poziționali” (position trees) pentru identificarea substringurilor comune a două sau mai multe secvențe, așa după cum o descrie Aho, Hopcroft și Ullman [AHO'74], dar are avantajul implementării simple.

Pentru a ilustra conceptul arborelui de poziție fie $a = \text{AATAATGCS}$, unde S semnalează finalul secvenței. Pentru fiecare i , $i = 1$ la 8 , fie subșirul S cel mai scurt subșir care începe la i care nu apare altundeva în a . Acest subșir este spus că identifică pe i . Spre exemplu, poziția $i = 4$ este identificată de AATG. Aceste subșiruri de identificare sunt organizate într-un arbore de poziții care reprezintă informația:

Poziția	Identificând subșirul
1	AATA
2	ATA
3	TA
4	AATG
5	ATG
6	TG
7	G

8	CS
9	S
.	

Cele n noduri terminale ale arborelui de poziție pentru $a = a_1, a_2, \dots, a_n$ constă din $1, 2, \dots, n$. Secvența de etichete pe muchii de la rădăcină spre nodul terminal i este subșirul identificator pentru poziția i . Arborele de poziție pentru lungimea secvenței 8 de mai sus este reprezentat în Figura 2.3.7-1. Două secvențe (sau mai multe) pot fi procesate simultan pentru a da un arbore de poziție unde cea mai lungă potrivire poate fi găsită simplu.

O altă metodă pentru găsirea rapidă de regiuni poate să se bazeze pe conceptul de „hashing” așa cum este folosit în problemele clasice de căutare lexicografică. Cea mai timpurie referință este Dumey [DUME'56]. Descrisă mai detaliat în Dumas și Nino [DUMA'82], idea de bază a acestui concept este de a asocia cu fiecare poziție a unei secvențe echivalentul numeric al lui k -mer începând cu acea poziție. Echivalentul numeric este obținut cu privire la alfabetul secvenței ce definește baza unui sistem de numere. Astfel, un alfabet format din patru numere dă un sistem de numere bazei patru, și pentru k fix există 4^k numere posibile. Aceste numere pot fi folosite pentru a identifica poziții ale unui vector de mărime 4^k de liste ale pozițiilor locațiilor din secvență la care apare k -merul corespunzător. Această metodă este folosită de Wilbur și Lipman [WILB'83] în realizarea căutarilor rapide în baze de date, și aparent la fel de către Karlin et al. [KARL'83] pentru găsirea repetițiilor exacte. Vectorul poate fi construit într-un timp de ordinul N . Cele mai lungi repetiții, și astfel regiuni, sunt găsite prin simpla îmbinare a tuturor repetițiilor de lungime k , iar viteza pare să fie de ordinul $N \log N$ (sau ordinul N ?).

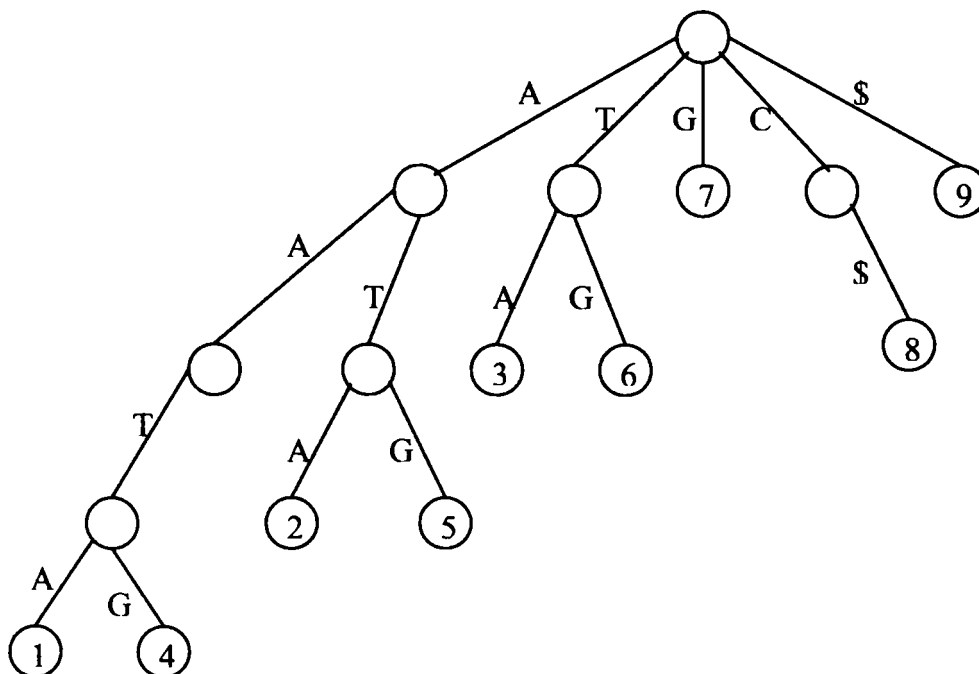


Figura 2.3.7-1. Arborele de poziție pentru $a=ATAATGC$.

Fiind dată lista regiunilor găsite de oricare dintre aceste metode, ele vor fi utilizate pentru găsirea aliniamentelor optimale. Două regiuni $R_1 = (w_1; i_1, j_1)$ și $R_2 = (w_2; i_2, j_2)$ se spune că satisfac $R_1 < R_2$ dacă $i_1 + |w_1| - 1 < i_2$ și $j_1 + |w_2| - 1 < j_2$. $R_1 \leq R_2$ înseamnă că există $i_2 - i_1 - |w_1|$ baze ale lui a și $j_2 - j_1 - |w_1|$ baze ale lui b între regiuni, și acele R_1 este stânga lui R_2 . Trebuie dată o formulă pentru ponderea acestor baze care nu se potrivesc. Dacă sunt x baze din a și y baze din b , fie $z(x,y)$ această pondere. Dacă nepotrivirile costă 1 și indels costă δ cu $1 < 2\delta$, atunci o alegere rezonabilă pentru $z(x,y)$ este

$$z(x,y) = |x-y|\delta + \min\{x,y\}.$$

Formula poate fi divizată pentru alte situații și $z(x,y) = x + y$ poate fi folosită.

Fiecare regiune R poate fi considerată a fi la sfârșit stânga al unui aliniament optimal $A(R)$, începând cu R și continuând cu a_n și b_m . Fie $D(A(R))$ un scor al unui asemenea aliniament. Optimizarea urmează din

$$D(A(R)) = \min\{w(k - i - |w|, l - j - |w|) + D(A(R^*))\} : R=(w: i, j) < R^*=(w*: k, l).$$

Algoritmii generali sunt cunoscuți cu timpul de execuție $O(|L|^2)$.

O implementare a acestui algoritm a fost realizată de către Martinez și Sobel și descrisă în Martinez [MART'83]. Ei produc de asemenea un aliniament aproape optimal prin adaptarea ideilor din 2.3.6.

2.3.8. Localizarea segmentelor lungi de potrivire

Problemele din această secțiune se centreză pe căutarea de segmente cu grad ridicat de similaritate care fac parte din două secvențe biologice. Există două abordări de bază: căutarea pentru potriviri *lungi exacte* sau pentru potriviri *lungi inexacte*. Noi dezvoltări în teoria probabilităților asistă aceste căutări.

2.3.8.1. Potriviri lungi exacte

Prima metodă eficientă pentru localizarea potrivirilor exacte a fost dată de Korn et al. [KORN'77]. Abordarea lor utilizează conceptul de arbore de poziții discutat deja. Aceeași metodă poate fi folosită pentru a găsi repetiții lungi într-o secvență fixată. Repetițiile vor partaja ramuri în arborele de poziții. Spre exemplu, în Figura 2.3.7-1, poziția 1 este identificată de AATA în timp ce poziția 4 este identificată de AATG, care implică repetiția „lungă” AAT. Korn et al. aranjează subșirurile de identificare într-o ordine lexicografică pentru a găsi aceste repetiții printre celelalte lucruri. Aceasta este o aplicație directă și utilă a științelor computaționale moderne pentru analiza secvențelor biologice și cu precădere a celor de ADN.

Algoritmii dat în Aho et al. [AHO'74] pentru construirea arborilor de poziții are cel mai dezavantajos timp de execuție de $O(n^2)$ pentru o secvență de lungime n . Timpul de execuție este proporțional cu numărul de vârfuri (vertices=noduri) ale arborelui. Oricum dacă literele cuvântului sunt independente și identic distribuite atunci timpul de execuție așteptat este $O(n)$. De asemenea ei puntează existența unui algoritm care rulează în $O(n)$ pentru toate intrările.

În lucrările lui Karlin et al. [KARL'83] și [KARL'84] se realizează analiza secvențelor prin localizarea repetițiilor directe lungi. Tehnica lor folosește o metodă de *hashing* și localizează toate repetițiile directe exacte.

O determinare teoretică a distribuției statistice a similarității a fost dată. Dacă $M(n, m)$ este lungimea celei mai lungi regiuni de potrivire exactă dintre două secvențe, atunci

$$E(M(n, m)) = \log((1-p)nm + \gamma/\lambda - 1/2 + r_1(n, m) + o(1))$$

$$\text{și}$$

$$\sigma^2(n, m) = n^2/6\lambda^2 + 1/12 + r_2(n, m) + o(1),$$

unde $p = P$ (două potriviri aleatoare de nucleotide), $\log = \log_{1/p}$, $\gamma = 0.577\dots$ este constanta Euler-Mascheroni, $\lambda = \log(1/p)$ și $r_1(n, m)$ și $r_2(n, m)$ sunt mici. Se observă că $\sigma(n, m)$ este esențial independentă de n și m . Acest rezultat este din Arratia et al. [ARRA'84a]. Karlin et al. [KARL'84] declară un rezultat ce diferă doar printr-o constantă de acesta. De asemenea se poate consulta Arratia și Waterman [ARRA'84b] pentru legi înrudite ale numerelor mari. Potrivirile dintre secvențe sunt considerate semnificative dacă depășesc $E(M(n, m)) + 2\sigma$.

Collins și Coulson [COLL'84] vin cu un algoritm de procesare paralelă pentru a produce matricea punct a tuturor potrivirilor de lungime mai mare sau egală cu un prag fixat. Implementarea acestora acceptă secvențe de până la 49152 baze și este un indicator al rezultatelor aplicării noii tehnologii acestor probleme. Rezultatele probabilităților anterioare pot fi folosite pentru fixarea pragului.

2.3.8.2. Potriviri lungi inexacte cu ajutorul programării dinamice

Se consideră problema localizării segmentelor similare dintre două secvențe fără a cere segmentelor să fie identice. Primul care a considerat această problemă a fost Sellers [SELL'79], [SELL'80]. El definește un interval I pentru a care seamănă cel mai bine cu b la modul general dacă $D(I, b) \leq D(J, b)$ pentru toate segmentele J ale lui a . Sunt necesare ambele matrice: „forward” și „backward”. Una dintre probleme este găsirea segmentelor de potrivire dorite din multele produse. Sankoff și Kruskal [SANK'83] estimează că n^4 potriviri rezultă din secvențe de lungime n . Goad și Kanehisa [GOAD'82] modifică tehnica lui Sellers. Erickson și Sellers [ERIC'83], discută pe larg această metodă și dau două aplicații non-triviale. În acest articol, Sellers rafinează analiza și vine cu un alt algoritm pentru găsirea „celor mai bune segmente”. El folosește conceptul de densitate de potrivire a lui Goad și Kanehisa [GOAD'82] pentru găsirea celor mai lungi segmente ale unei densități de potrivire recomandate. Similaritatea calculează

AAAA	versus	A
AAAA		A

ca patru potriviri versus una, atata timp cât distanța este zero în fiecare caz. Pentru a filtra porțiuni ale secvenței cu potrivire negativă, se definește H_{ij} a fi similaritatea maximă a două segmente care se termină în a_i și b_j , sau zero, indiferent care este mai mare.

Se definește:

$$H_{ij} = \max \{0, S(a_x a_{x+1} \dots a_j, b_y a_{y+1} \dots b_j): 1 \leq x \leq i \text{ și } 1 \leq y \leq j \}.$$

În acest articol, Seller sugerează faptul că cele mai bune segmente sunt în esență segmente de maximă similaritate care:

- (1) au similaritate non-negativă;
- (2) au scoruri cel puțin la fel de largi ca orice alt segment cu căi care se intersectează și
- (3) au scoruri cel puțin la fel de mari ca niște valori izolate.

Smith și Waterman [SMIT'81a,b] au recomandat o procesare secvențială a matricei H , găsind aliniament cu cele mai mari, apropiat de cele mai mari, etc. valori de similaritate cu direcții neintersectate. După stabilirea algoritmului pentru H , [WATE'84] vine cu un nou și mai complet algoritm pentru găsirea de aliniamente satisfăcând recomandările (1), (2) și (3) ale lui Seller.

Algoritmul Smith – Waterman

Fie $H_{i0} = H_{0j} = 0$ pentru $1 \leq i \leq n$ și $1 \leq j \leq m$. Atunci

$$H_{ij} = \max \left\{ H_{i-1, j-1} + s(a_i, b_j), \max_{1 \leq k \leq i} \{ H_{i-k, j} - x_k \}, \max_{1 \leq k \leq j} \{ H_{i, j-k} - x_k \}, 0 \right\}.$$

Valorile care pot fi folosite pentru ponderi sunt

$$s(x, y) = \begin{cases} 1 & \text{daca } x = y \\ -1/3 & \text{daca } x \neq y \end{cases}$$

și

$$x_k = 1 + k/3.$$

O reducere în timpul de calculare de la $O(n^3)$ la $O(n^2)$ pentru funcții liniare sau concave de ștergere pot fi obținute ca în secțiunea 2.3.4.

Când se construiește matricea H , se pun în stivă toate (i, j, Y) cu $Y = H_{ij}$ și $H_{ij} \geq C =$ valoare extremă (cut-off value). Stiva este ordonată de „>”, unde

$(i, j) > (k, l)$ dacă:

- (1) $H_{ij} > H_{kl}$ sau
- (2) $H_{ij} = H_{kl}$ și $i + j < k + l$ sau
- (3) $H_{ij} = H_{kl}$, $i + j = k + l$ și $i < k$.

În timpul parcurgerii înapoi, pentru unele intrări în stivă aliniamentele multiple sunt rezolvate în modul următor: dacă două aliniamente multiple se termină în (i, j) , și (k, l) o produce pe aceea care se termină în (i, j) , dacă: (1) $i + j < k + l$, sau dacă (2) $i + j = k + l$ atunci și $i > k$. Odată ce o parcurgere înapoi este completată cu succes, intrările aliniamentelor în matrice sunt înmulțite cu -1, (i.e. aliniamentul produs), și intrarea corespunzătoare în stivă este înlăturată din stivă. Elementele negative ale lui H nu sunt folosite în vre-unul dintre aliniamentele viitoare.

Dacă o intrare în stivă are un element negativ corespunzător al matricei, se înlătură și se merge mai departe. Dacă o parcurgere înapoi întâlnește un element negativ al matricei ea nu se va putea continua. Dacă cel mai bun aliniament generat de departe are scorul $= Y \geq C$, atunci (i, j, Y) trebuie înlocuit în stiva ordonată.

Boswell and McLachlan [BOSW'84] sugerează de asemenea folosirea valorilor similarității pentru localizarea segmentelor similare. Matricea lor directă (forward) este calculată de către formula

$$F(i, j) = s(a_i, b_j) + \lambda \max \{F_{i-2, j-1} - w_1, F_{i-1, j-1}, F_{i-1, j-2} - w_1\}. \quad (2.3.8.2-1)$$

Matricea inversă R este găsită prin inversarea secvențelor. Atunci

$$M(i, j) = F(i, j) + R(i, j) - s(a_i, b_j).$$

Idea este că $M(i, j)$ este suma a $s(a_i, b_j)$ plus căile ponderate cele mai bune care se extind în orice direcție. Parametrul $\lambda \in (0, 1)$ este un factor geometric de amortizare (damping).

Problema distingerii valorilor statistice semnificative a

$$H^* = \max_{i,j} H_{ij} \quad (2.3.8.2-2)$$

este în mod evident importantă.

Arratia et al. [ARRA'84a] a arătat că lungimea $M(n, m)$ a celei mai lungi potriviri întreruptă de k nepotriviri satisface relația:

$$E(M(n, m)) = \log(nm) + k \log(nm) + (k+1)\log(1-p) - \log(k!) + k + \gamma/\lambda - 1/2 \\ + r_1(n, m) + o(1)$$

și

$$\sigma^2(n, m) = \pi^2/6\lambda^2 + 1/12 + r_2(n, m) + o(1),$$

unde $p = P$ (două potriviri aleatoare de nucleotide), $\log = \log_{1/p}$, $\gamma = 0.577\dots$ este constanta Euler-Mascheroni, $\lambda = \log(1/p)$ și r_1, r_2 sunt mici. Se poate folosi spre exemplu $H_{ij} \geq E(M(n, m)) + 2\sigma(n, m)$ pentru a decide care H_{ij} sunt de intreres pentru a fi produse. Se observă că $\sigma(n, m)$ este încă odată esențial independentă de n și m .

Într-un studiu empiric, Smith et al. [SMIT'84] arată că, pentru

$$s(x, y) = \begin{cases} 1 & \text{daca } x = y \\ -0.9 & \text{daca } x \neq y \end{cases}$$

și

$$x_k = \begin{cases} 2 & \text{daca } k = 1 \\ \infty & \text{daca } k > 1 \end{cases}$$

valorile lui $E(H^*)$ și σ sunt

$$E(H^*) = 2.5 \log(nm) - 9$$

$$\sigma = 1.78$$

unde $\log = \log_{1/p}$ ca mai sus.

Un rezultat asimptotic cum că legea lui $\log(n)$ este de acord cu indels precum și cu nepotrivirile este dată în Arratia și Waterman [ARRA'84a] (1984). Studiul empiric raportat mai sus este o evidență a robusteții acestei distribuții.

2.3.8.3. Potriviri lungi inexacte folosind regiuni

Un algoritm pentru găsierea potrivirilor lungi inexacte și care nu folosește programarea dinamică este propus de Korn et al. [KORN'77]. El are câteva dezavantaje serioase care sunt punctate mai jos, dar este o metodă folositoare și a fost larg răspândită printre analiștii secvențelor biologice, a se vedea Queen et al. [QUEE'82].

Algoritmul începe la poziția (i, j) unde $a_i = b_j$ și $a_{i+1} = b_{j+1}$. Această potrivire de lungime doi este extinsă într-o manieră recursivă, unde regulile pentru extindere sunt:

- (1) următoarele baze se potrivesc (i.e. $a_{i+2} = b_{j+2}$),
- (2) prin ștergerea a 1, 2 sau 3 baze din secvența a există o execuție a 3 potriviri,
- (3) prin ștergerea a 1, 2 sau 3 baze din secvența b există o execuție a 3 potriviri sau
- (4) prin nepotrivirea lui a_{i+2} și b_{j+2} , două dintre următoarele 3 perechi se potrivesc. Programul lor nu caută perechi (i, j) unde (a_i, b_j) sunt deja într-o regiune de identificată.

Sankoff și Kruskal [SANK'83] remarcă faptul că această metodă, când compară AACAAA și AAAAA, va găsi AACAAA și AAΔAAA, dar cu secvențele inversate nu va găsi această regiune. Aceasta ar putea fi o proprietate nedorită. Sankoff și Kruskal au remarcat de asemenea că AACCGT și AACGT vor produce

AAC		CGT
	și	
AAC		CGT

Care au o bază în comun, în loc de

AACCGT		AACCGT
	și	
AACΔGT		AACΔGT

Timul de execuție al acestui algoritm pentru două secvențe de lungime n este proporțional cu n^2 . Constanta este mai mare decât așteptarea unei variabile geometrice aleatoare. Pentru baze egal probabile aceasta înseamnă: constanta va

depăși 4. De când algoritmi de programare dinamică au o constantă de 3, aceasta sugerează puternic folosirea algoritmilor matematici mai riguroși. Metoda Wilbur-Lipman descrisă în secțiunea destinată căutării eficiente în baze de secvențe (3.3.3) este de asemenea o metodă a regiunilor ce poate fi aplicată acestei probleme.

2.3.9. Aliniamentul perechilor folosind modele Markov ascunde (HMMs)

Folosind teoria HMM, a fost dezvoltat un anumit HMM care modelează perechi aliniată de secvențe. Această abordare oferă unele moduri de estimare a acurateții unui aliniament și calculează scorul care exprimă similaritatea fără a referi un aliniament particular. HMM-urile specifice se bazează pe conceptul de FSA (Finite State Automata) care sunt folosite pentru aliniamentul perechilor de secvențe ca descriptori convenabili ai unor algoritmi de programare dinamică mai complecși. Este executată o transformare în care toate alternativele sunt ponderate probabilistic. Acest punct de vedere oferă posibilitatea de a calcula similaritatea a două secvențe independent de orice aliniament specific. Algoritmii incluși în această categorie cuprind: algoritmul Viterbi pentru perechi de HMM, calcularea „înainte” (forward calculation) pentru perechi de HMM și calcularea „înapoi” (backward calculation) pentru perechi de HMM. O descriere detaliată a acestora poate fi găsită în [DURB'98].

Utilitatea profilurilor HMM este extinsă la căutarea familiilor de secvențe, permițând celor preocupați de biologia computațională să infere aproape de trei ori mai multe omologii decât oricare alt algoritm simplu pentru perechi de secvențe [PARK'98].

Dintre toolurile disponibile care folosesc profile de HMM sunt SAM¹⁵ (Sequence Alignment and Modelling System) și cele accesibile pe serverul HMMER¹⁶, care este o implementare de software specializat pe profile HMM pentru analiza secvențelor de proteine distribuită gratuit.

În aliniamentul de perechi folosind modele Markov ascunde (HMMs) se presupun 3 stări, în care M corespunde stării de potrivire și două stări ce corespund inserărilor și care vor fi notate cu X și Y în următoarea figură.

¹⁵ <http://www.cse.ucsc.edu/research/compbio/HMM-apps/HMM-applications.html>

¹⁶ <http://bioweb.pasteur.fr/seqanal/motif/hmmer-uk.html>

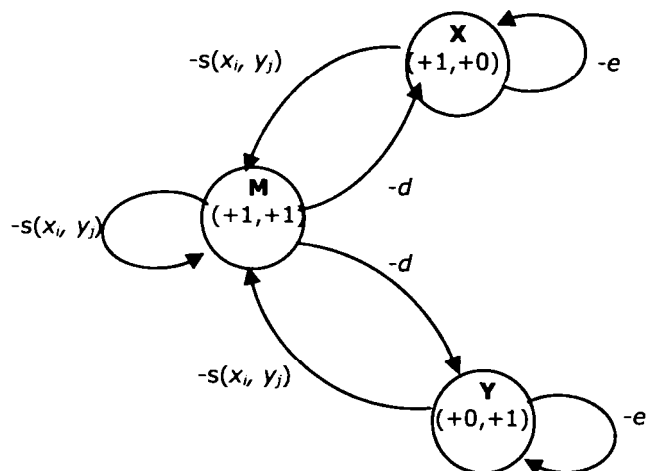


Figura 2.3.9-1. Diagrama unui automat cu stări finite pentru un aliniament cu gap

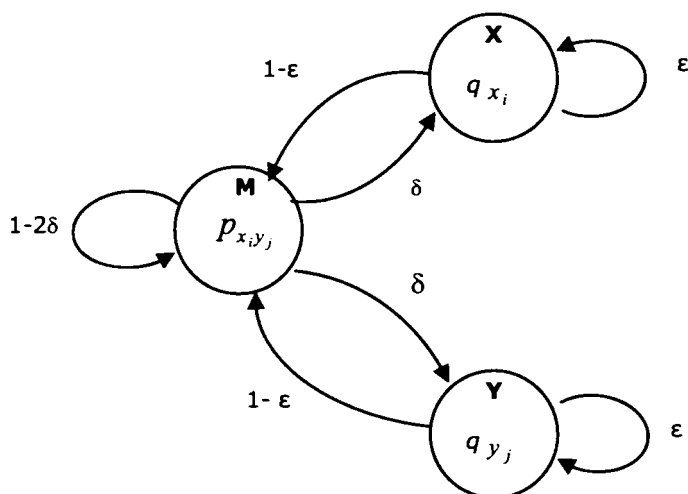


Figura 2.3.9-2 . Modelul probabilistic corespunzător diagramei unui automat cu stări finite pentru un aliniament cu gap din figura anterioară.

Relațiile de recurență pentru actualizarea valorilor acestor stări în matricea de programare dinamică sunt:

$$v^M(i, j) = s(x_i, y_j) + \max \begin{cases} v^M(i-1, j-1), \\ v^X(i-1, j-1), \\ v^Y(i-1, j-1); \end{cases} \quad (2.3.9-1)$$

$$v^X(i, j) = \max \begin{cases} v^M(i-1, j) - d, \\ v^X(i-1, j) - e; \end{cases}$$

$$v^Y(i, j) = \max \begin{cases} v^M(i, j-1) - d, \\ v^Y(i, j-1) - e. \end{cases}$$

Aceste ecuații sunt corespunzătoare aliniamentului global. Pentru aliniamentul local se intervine cu observațiile adecvate privind modificările necesare.

Așa după cum se poate observa din figurile anterioare, sunt necesare două seturi de schimbări asupra unui FSA pentru a-l transforma într-un HMM. Mai întâi trebuie date probabilități pentru emiterea de simboluri din stări precum și pentru tranziții între stări. Spre exemplu, starea M are distribuția probabilității de emisie p_{ab} pentru emiterea unei perechi aliniată $a:b$, și stările X și Y vor avea distribuții q_a pentru emiterea simbolului a opus unui gap. Deoarece starea X emite simboluri x_i

din secvența x , se va scrie q_{x_i} în interiorul cercului ce reprezintă starea X. De asemenea, probabilitățile de tranziție dintre stări, trebuie să satisfacă cerința ca suma tuturor probabilităților pentru toate tranzițiile care pleacă dintr-o stare să fie unu. Permițând simetria, există doi parametri liberi pentru probabilitățile de tranziție între cele trei stări principale. Se notează forma de tranziție de la M la o stare de inserție (X sau Y) cu δ , iar probabilitatea rămânerii în starea de inserare cu ϵ .

Oricum, modelul rezultat din figura a doua nu generează un model complet care să ofere distribuția de probabilități peste toate secvențele posibile. Pentru aceasta este necesară definirea unei stări de Begin și End așa cum o ilustrează Figura 2.3.9-3. Ca efect, acestea formalizează condițiile de inițializare și terminare necesare pentru algoritmi de programare dinamică [DURB'98]. După cum se explică în [DURB'98] aranjamentele mai complexe a stărilor de Begin și End pot corespunde aliniamentelor locale sau altor tipuri de aliniament. Adăugarea unei stări End explicite introduce nevoia unui alt parametru, probabilitatea unei tranziții în starea de End, care acum se presupune a fi aceeași pentru fiecare din M, X și Y și este notată cu τ . Ca efect, aceasta va determina lungimea medie a unui aliniament dintr-un model. Pentru moment, se vor stabili tranzițiile de la starea Begin ca fiind aceleași cu cele de la starea M (se putea spune și că se pornește în M, dar se recurge la această formă pentru a clarifica faptul că atât inițializării cât și terminării i se pot da considerente independente.)

Această metodă ne oferă un model probabilistic foarte apropiat de HMM. Diferența constă în faptul că în loc să se emită o singură secvență se emite un aliniament de perechi. Modelul astfel construit este referit ca *pereche* HMM pentru a fi distins de alte timpuri standardizate care emit secvențe singulare.

Așa cum un HMM standard poate genera o secvență, această *pereche* HMM poate genera o pereche aliniată de secvențe. Aceasta se realizează pornind din starea de Begin și ciclând asupra următorilor doi pași: (1) ia următoarea stare în conformitate cu distribuirea probabilităților de tranziție părăsind starea curentă; (2) ia simbolul unei perechi ce va fi adăugată aliniamentului în conformitate cu

distribuția de emisie în noua stare. Procesul se oprește când o tranziție este făcută în starea de End.

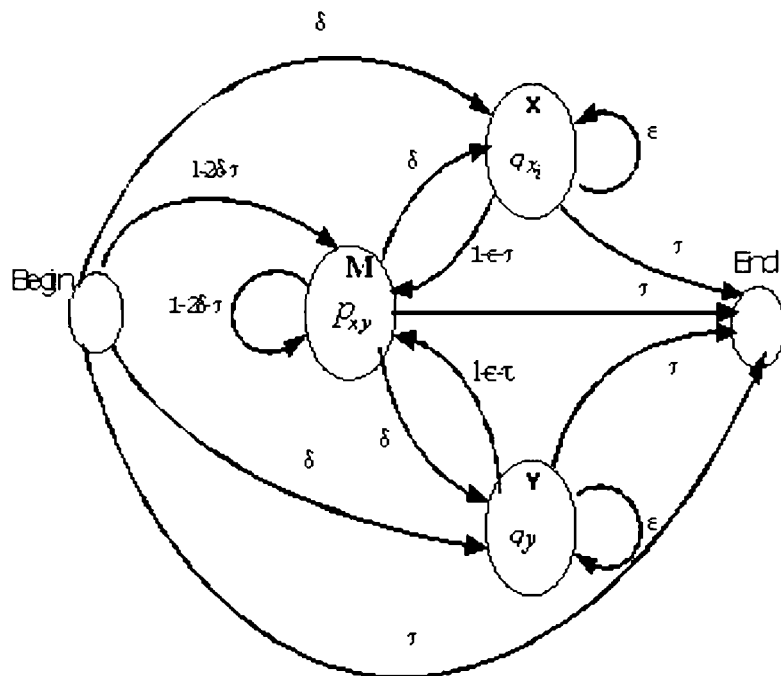


Figura 2.3.9-3. Versiunea probabilistică integrală a figurii anterioare

Deoarece fiecare pas are probabilități se poate, de asemenea, păstra probabilitatea totală de generare a unui aliniament particular pe care l-am produs. Acesta este doar produsul probabilităților fiecărui pas individual.

2.3.9.1. Cea mai probabilă cale - aliniamentul FSA optimal

Algoritmul Viterbi permite găsirea celei mai probabile căi într-o pereche HMM fiind date secvențele x și y . Forma corectă pentru perechea globală HMM din Figura 2.3.9-3 este după cum urmează. Pentru a face ecuația mai simplă, se definește starea de Begin cu M . De asemenea, se folosesc simboluri cu litere mici $v^*(i,j)$ pentru valorile de probabilitate și cu litere mari $V^*(i,j)$ scorurile log-odds (care simbolizează valorile probabilistice asociate, calculate de regulă ca logaritmul natural sau în baza 2 din $q_{ij}/p_i p_j$ (a se vedea detaliile în secțiunea 2.1.2.3.) Algoritmul Viterbi în termeni probabilistici este următorul.

Algoritmul Viterbi pentru perechi HMMs

Inițializare:

$V^*(0,0)=1$. Toți ceilalți $v^*(i,0)$, $v^*(0,j)$ sunt setați la 0.

Recurența: $i=1,\dots,n$, $j=1,\dots,m$;

$$v^M(i, j) = p_{x_i y_j} \max \begin{cases} (1 - 2\delta - \tau) v^M(i-1, j-1), \\ (1 - \varepsilon - \tau) v^X(i-1, j-1), \\ (1 - \varepsilon - \tau) v^Y(i-1, j-1); \end{cases}$$

(2.3.9.1-1)

$$v^X(i, j) = q_{x_i} \max \begin{cases} \delta v^M(i-1, j), \\ \varepsilon v^X(i-1, j); \end{cases}$$

$$v^Y(i, j) = q_{y_j} \max \begin{cases} \delta v^M(i, j-1), \\ \varepsilon v^Y(i, j-1); \end{cases}$$

Terminare:

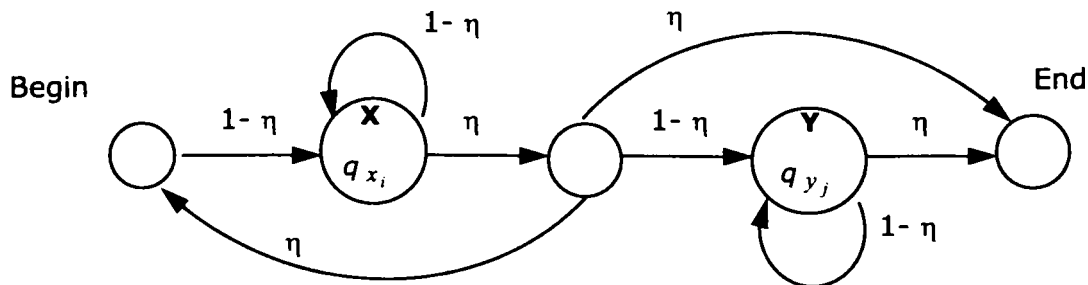
$$v^E = \tau \max(v^M(n, m), v^X(n, m), v^Y(n, m))$$

Pentru găsirea celui mai bun aliniament se pastrează pointerii și se parcurge drumul înapoi. Desigur, pentru obținerea aliniamentului însuși se păstrează reziduurile emise la fiecare pas din traiectorie pe timpul parcurgerii înapoi precum și secvența de stări.

Cu toate că este evident faptul că ecuațiile de recurență ale algoritmului Viterbi pentru perechi HMM au aceeași formă ca și cele pentru versiunea automatului cu stări pentru aliniament de perechi, este utilă urmărirea formei exacte de corespondență.

Mai întâi trebuie transformate în rapoarte log-odds cu raportare la modelul aleator. De fapt acum se are un model probabilistic complet pentru aliniamentul în cauză și trebuie de asemenea unul pentru modelul aleator cu o condiție de terminare corespunzătoare. Până acum s-a ignorat faptul că modelul aleator nu poate produce secvențe de lungime variabilă într-o formă probabilistică adecvată.

În figura de mai jos este un nou model aleator care de asemenea este o pereche HMM.



Stările principale sunt X și Y, care emit cele două secvențe, independente una de cealaltă. Fiecare are un salt înapoi la ea însăși cu probabilitatea $(1 - \eta)$. La fel ca stările Begin și End există, de asemenea, o stare de silent/tăcere între X și Y, indicată de un cerc mai mic. Aceasta nu emite vre-un simbol dar este folosită pentru a acumula intrări din ambele stări X și Y. Definit în acest fel, modelul permite secvențe x sau y de lungime zero, la fel cum o permite și modelul de perechi HMM în Figura 2.3.9-3 și generează o formă simplă pentru distribuția modelului aleator asupra secvențelor. Probabilitatea unei perechi de secvențe x și y conform acestui model este

$$\begin{aligned} P(x, y|R) &= \eta(1 - \eta)^n \prod_{i=1}^n q_{x_i} \eta(1 - \eta)^m \prod_{j=1}^m q_{x_j} \\ &= \eta^2 (1 - \eta)^{n+m} \prod_{i=1}^n q_{x_i} \prod_{j=1}^m q_{x_j}. \end{aligned} \quad (2.3.9.1-2)$$

Acum se dorește alocarea termenilor în această expresie celor care constituie probabilitatea aliniamentului Viterbi, astfel încât raportul de probabilitate asociat pentru întregul aliniament poate fi exprimat ca un produs al rapoartelor de probabilitate a termenilor individuali (și corespunzători, astfel încât raportul de probabilitate al aliniamentului este o sumă a termenilor probabilistici). Aceasta se face prin alocarea unui factor $(1 - \eta)$ și a factorului corespunzător q_a la fiecare reziduu care este emis într-un pas Viterbi. Astfel, tranzițiile de potrivire vor fi alocate $(1 - \eta)^2 q_a q_b$ unde a și b sunt două reziduuri care se potrivesc, iar stările de inserare $(1 - \eta) q_a$, unde a este reziduu inserat. Deoarece calea Viterbi trebuie să țină cont de toate reziduurile, vor fi folosiți exact $(n+m)$ termeni, și toate din eq. 2.3.9.1-2 cu excepția factorului inițial de η^2 .

În termeni de probabilități asociate, acum se poate calcula un model aditiv cu scorul emis de probabilitățile log-odds și scorurile de tranziție log-odds. În practică aceasta este cea mai normală cale de a implementa HMMs. Din aceasta, este posibilă combinarea scorurilor de emisie cu cele de tranziție după exemplul următor:

$$\begin{aligned} s(a, b) &= \log \frac{p_{ab}}{q_a q_b} + \log \frac{(1 - 2\delta - \tau)}{(1 - \eta)^2}, \\ d &= -\log \frac{\delta(1 - \varepsilon - \tau)}{(1 - \eta)(1 - 2\delta - \tau)}, \\ e &= \log \frac{\varepsilon}{1 - \tau}, \end{aligned} \quad (2.3.9.1-3)$$

pentru a produce scoruri care corespund termenilor standard folosiți în aliniamentul de secvențe prin programare dinamică. A se nota: contribuția lui q_a la d și e este inutilă datorită factorilor anulați din modelele Viterbi și aleatoare. De asemenea, pentru anihilarea diferențelor în tranzițiile ce provin din stările de potrivire și gap a fost necesară ceva îndemânare pentru exprimarea lui s și d . Se intenționează folosirea lui $s(a, b)$ ca scor pentru fiecare potrivire, indiferent dacă este urmat de o

altă potrivire sau o inserare. Cu scopul de a face această strategie să funcționeze corect, s-a construit în d o ajustare care să corecteze diferențele în scorul de potrivire când revine dintr-o inserare. Aceasta presupune că termenii matricei de programare dinamică nu mai corespund exact rapoartelor probabilistice de asociere (log-odds ratios) de a fi în aceleași stări, cu toate că rezultatul final va fi corect. Astfel, acum se poate oferi versiunea log-odds a algoritmului de aliniament Viterbi într-o formă care arată ca algoritmul dinamic clasic de aliniere pentru perechi de secvențe.

Aliniamentul log-odds optimal

Inițializare:

$$V^M(0, 0) = 2 \log \eta, \quad V^X(0, 0) = V^Y(0, 0) = -\infty.$$

Toți $V^*(i, +1)$, $V^*(-1, j)$ sunt setați la $-\infty$.

Recurența:

$i = 0, \dots, n, j = 0, \dots, m$ mai puțin $(0, 0)$;

$$v^M(i, j) = s(x_i, y_j) + \max \begin{cases} v^M(i-1, j-1), \\ v^X(i-1, j-1), \\ v^Y(i-1, j-1); \end{cases} \quad (2.3.9.1-4)$$

$$v^X(i, j) = \max \begin{cases} v^M(i-1, j) - d, \\ v^X(i-1, j) - e; \end{cases}$$

$$v^Y(i, j) = \max \begin{cases} v^M(i, j-1) - d, \\ v^Y(i, j-1) - e; \end{cases}$$

Terminare:

$$v = \max(v^M(n, m), v^X(n, m) + c, v^Y(n, m) + c)$$

Acestea sunt identice cu relațiile 2.3.9-1 cu excepția constantei $2 \log \eta$ la inițializare, și constantei $c = \log(1-2\delta - \tau) - \log(1 - \epsilon - \tau)$ la terminare, care este necesară pentru corectarea ajustării descrise mai sus în d .

Procedura, după cum este descrisă, arată cum pentru orice pereche HMM de tipul din Figura 2.3.9-3 se poate deriva un FSA echivalent pentru a obține cel mai probabil aliniament. Aceasta ne permite observarea unei interpretări probabilistice riguroase pentru termenii folosiți în aliniamentul secvențelor. Pentru a parcurge invers acest proces, adică de la un algoritm de programare dinamică exprimat via FSA la perechi HMM este mai complicat. În acest caz ar fi nevoie în general de un nou parametru λ care să acționeze ca un factor global de scalare pentru scoruri, și pentru orice set de scoruri date pot fi constrângeri în alegerea lui η și τ .

2.3.9.2. Perechi de HMM pentru aliniament local

Modelul prezentat în Figura 2.3.9-3 este potrivit pentru găsirea unei potriviri globale între secvențe. După cum s-a constatat în [DURB'98], multe dintre căutările cele mai sensibile pentru perechi de secvențe sunt locale. De obicei, când se introduce un algoritm de aliniament local sau alte variante cum ar fi algoritmi de repetare sau suprapunere se explică în termeni de schimbări în ecuațiile de actualizare și condiții de limită. Toate acestea sunt făcute explicit în formalismul perechilor HMM prin adăugarea de stări și tranziții. Astfel se poate defini un model de perechi HMM pentru fiecare variantă. În Figura 2.3.9.2-1 este prezentat un model pentru aliniament local. Acesta arată mai complicat decât modelul global din Figura 2.3.9-3, dar este construit din bucăți simple și într-un mod direct.

Un model probabilistic complet trebuie să țină seama de toate secvențele x și y : nu numai de aliniamentul local dintre x și y ci și de secvențele nealiniat care le flanchează. Prin urmare pot fi adăugate secțiuni de modele suplimentare înainte și după segmentul de potrivire cu trei stări din Figura 2.3.9-3. Fiecare segment alăturat este o copie a modelului complet fundamental, deoarece secvențele în regiunile laterale sunt nealiniat. Cei mai mulți termeni în contribuțiile probabilistice a acestor secțiuni vor fi anulate cu termenii echivalenți în modelul aleator când se calculează scorurile log-odds ale unei potriviri în comparație cu modelul aleator, lăsând numai scorul local de potrivire din partea centrală a modelului, și unii termeni suplimentari. Modele similare mixte pot fi construite pentru modele de suprapunere și repetiție și hibride (a se vedea cap 2 din [DURB'98]).

Având perechi HMM se poate face mai mult decât oferirea unei alternative raționale pentru aliniament standard prin programare dinamică. Și anume, se poate merge la discutarea semnificației potrivirilor. Spre exemplu, atunci când similaritatea este scăzută este dificilă identificarea aliniamentului corect și testarea semnificației. Astfel, există diverși algoritmi și metode care se dezvoltă pe principiile de bază ale perechilor HMM. A se vedea [DURB'98].

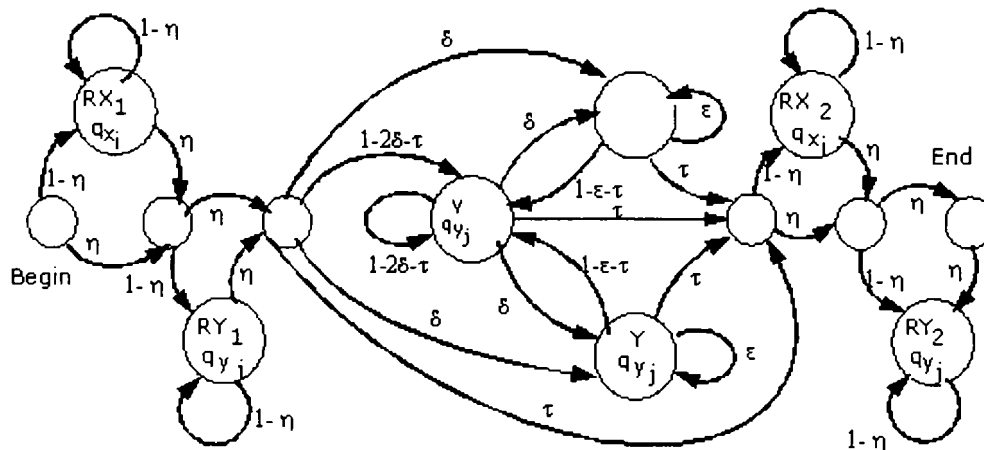


Figura 2.3.9.2-1: O pereche HMM pentru aliniament local. Aceasta este compusă din modelul global (stările M , X și Y) mărginite de două copii ale modelului aleator (stările RX_1 , RY_1 și RX_2 , RY_2).

2.3.10. Aliniament multiplu de secvențe

Tehnicile pentru aliniamentul multiplu al secvențelor sunt folosite pentru a recunoaște/caracteriza familiile de proteine și a identifica regiuni partajate de omologie de-a lungul a multiple secvențe (aceasta se întâmplă în general când căutarea unei secvențe relevă omologii la unele secvențe). În plus, aceste tehnici sunt importante pentru determinarea secvenței de consens (**consensus sequence**), adică secvența ideală pentru interacțiunea cu o proteină reguloare a câtorva secvențe aliniate și predicția structurilor secundare și terțiare a secvențelor noi. În concluzie, ele sunt considerate ca un pas de bază preliminar în evoluția analizei moleculare cu scopul construirii arborelui filogenetic pe baza metodelor filogenetice.

Metodele de grupare tradiționale pot executa automat aliniamente dar ele sunt subordonate unor restricții. Adesea ele necesită prelucrări suplimentare și interacțiune cu utilizatorul. În unele cazuri mai dificile este necesară chiar și validarea vizuală. În [BART'96] au fost identificate diferite categorii de tehnici de aliniament multiplu:

- Extensii ale algoritmilor de programare dinamică pentru perechi de secvențe;
- Extensii ierarhice a aliniamentului de perechi de secvențe;
- Metode pentru segmente;
- Metode pentru consens sau 'regiuni'.

Prima categorie își are originea în sugestia autorilor algoritmului original de programare dinamică pentru perechi de secvențe și anume, că metoda lor poate fi utilizată pentru a alinia mai mult de două secvențe [NEED'70], [WATE'76]. În practică, aceasta înseamnă extinderea spațiului de stocare de N ori (unde N este numărul de secvențe disponibile) și datorită acestei extensii a fost aplicat la pâna la trei secvențe. Una dintre aplicațiile timpurii ale programării dinamice este realizată de Sankoff [SANK'76] care aliniază 9 secvențe. Algoritmul său construiește o structură de arbore iar problema alinierii simultane a 9 secvențe este descompusă în definirea a 7 aliniamente de trei triplete. Aliniamentele sunt executate secvențial parcurgând arborele de jos în sus. Nu există însă tooluri larg cunoscute care să utilizeze această metodă.

Dintre abordările destinate aliniamentelor multiple, cea mai practică și mai populară este extensia ierarhică a aliniamentului de perechi. Principiul esențial constă în faptul că aliniamentele multiple pot fi efectuate prin aplicarea succesivă a metodelor pentru perechi de secvențe. Pașii obișnuiți, așa cum sunt descriși în [BART'96], sunt:

1. Se compară toate secvențele după metoda aplicată perechilor. Vor fi $N(N-1)/2$ perechi pentru un număr N de secvențe date.
2. Se generează o ierarhie reprezentând toate datele aflate sub formă de perechi. Aceasta poate fi prezentată sub forma unui arbore binar sau a unei simple liste ordonate.
3. Se construiesc aliniamentele multiple aliniind la început cea mai « populară » pereche de secvențe apoi perechea cea mai similară și așa mai

departe. Odată ce a fost aliniată o pereche de secvențe aceasta rămâne fixă. Astfel, pentru un set de secvențe A,B,C,D, după ce a fost aliniat A cu C și B cu D, aliniamentul lui A, B, C, D este obținut prin compararea aliniamentelor lui A și C cu acela al lui B și D.

Cel mai popular tool care se încadrează în această categorie este CLUSTAL W și care este accesibil publicului.

Metodele segment sunt identificate în [BACO'86] unde problema aliniamentului a N secvențe este rezolvată considerând aliniamentul într-o ordine specifică. Mai întâi, secvența unu este comparată cu secvența doi și primele M perechi de segmente calculate ca cele mai similare sunt stocate. Următoarea secvență este comparată cu segmentele selectate iar cele mai ridicate scoruri dintre cele trei secvențe sunt reținute. Continuând acest proces se obține o listă de M aliniamente (i.e. perechi de segmente) cu cele mai ridicate scoruri dintre cele N secvențe. Barton în [BART'96] consideră că acest algoritm este foarte promițător pentru identificarea similarităților importante dintre segmente scurte de secvențe. Metoda însă nu furnizează un aliniament general al secvențelor și nu consideră gap-urile în mod explicit. În [JOHN'86], autorii modifică algoritmul reducând numărul comparațiilor segment-cu-segment care trebuie executat progresiv. Ei restricționează evaluarea aliniamentului dintre segmente într-o "fereastră" specificată anterior. Această abordare poate genera un aliniament complet al secvențelor considerând gap-uri. Din păcate timpul de execuție devine nerezonabil de mare pentru secvențe cu mai mult de 50 reziduuri.

Ultima categorie a aliniamentelor multiple, așa numita *consensus* sau *metoda regiunilor*, a fost introdusă în [WATE'86]. Aici este prezentat un algoritm care execută aliniament multiplu de secvențe prin potrivire de cuvinte de lungime dată. Mărimea cuvintelor și gradul de nepotrivire permis sunt selectate de utilizator. Aliniamentul maximizează funcția de calcul corespunzătoare. Metoda este bazată pe o extindere nouă a metodelor consensus anterioare și permit estimarea semnificației statistice.

2.4. Concluzii

Indiferent de obiectivul biologic urmărit prin analiza și compararea secvențelor biologice este necesară folosirea de metode și tehnici algoritmice adecvate. Din diversitatea de metode existente larg cunoscute se conturează două tendințe:

- Metode de aliniament bazate pe programarea dinamică;
- Metode orientate spre folosirea modelelor Markov ascunse.

Deoarece fiecare dintre aceste direcții implică o serie de limitări, cum ar fi timpii de execuție relativ mari, spațiul de memorie semnificativ pentru prelucrarea unei cantități mari de date sau căutarea unor valori adecvate pentru diverși parametri necesari estimărilor probabilistice, a existat și încă există o preocupare permanentă de a simplifica aceste metode și de a găsi metode alternative care să servească aceluiași scopuri.

Prin urmare s-a realizat și o investigație a tehnicilor de similaritate propuse relativ recent. Ele vin oarecum ca alternativă la metodele de aliniament fundamentale și chiar dacă nu se pretinde o acoperire totală a acestora ele pot să ofere o imagine de ansamblu asupra noilor metode propuse/experimentate pentru determinarea similarității secvențelor (reprezentate ca structuri primare în cazul proteinelor).

Dată fiind complexitatea problemei de comparare a secvențelor atât cu scopul determinării similarității cât și al identificării de segmente cu rol biologic important, s-a ajuns la următoarele concluzii:

- 1) Pentru determinarea similarității sunt necesare fie funcții pentru evaluarea aliniamentelor în vederea stabilirii distanței dintre secvențe, care implică și folosirea matricelor de substituție, fie tehnici algoritmice independente de aliniamentul secvențelor;
- 2) Aliniamentele multiple folosesc de obicei tehnicile pentru compararea de perechi de secvențe pe care le încorporează în algoritmi de parcurgere și evaluare a tuturor secvențelor implicate în procesul de comparare.
- 3) Din diversitatea de metode nu poate fi stabilit cu precizie cel mai performant algoritm de determinare a similarității întrucât el este dependent de scopul urmărit în compararea secvențelor. Totuși, ceea ce se urmărește în general este ca timpul și spațiul computațional necesar să fie cât mai redus;
- 4) Metodele de aliniament și determinare a similarității cunosc o continuă îmbunătățire având ca scop obținerea de performanțe ridicate la costuri cât mai mici.
- 5) Este de remarcat implicarea estimărilor statistice atât în evaluarea rapoartelor de asociere pentru ponderea elementelor în matricele de substituție cât și funcțiile de evaluare ale aliniamentelor cu precădere pentru modelele Markov ascunse.

Aceste metode de similaritate se constituie în nuclee care stau la baza diverselor programe comerciale destinate comparării, analizei și extragerii de informații relevante din baze mari de secvențe.

3. APLICAREA TEHNICILOR DE RECUNOAȘTERE ÎN ANALIZA SECVENȚELOR BIOLOGICE

3.1. Obiective de interes biologic ce implică tehnicile informatice

Cercetările comunității în inteligență bioinformatică sunt focalizate pe diverși algoritmi de învățare, sisteme de data mining și sisteme integrate care să permită transformarea secvențelor biologice, a observațiilor și cunoștințelor în informație structurată și semnificativă pe care biologii o pot interoga, vizualiza și înțelege [KASI'99]. În conformitate cu descrierea făcută de Simon Kasiff, cele mai importante funcții computaționale identificate în acest domeniu de cercetare pot fi grupate în următoarele categorii:

- Identificarea genelor/proteinelor și clasificarea lor în categorii;
- Metode de comparare a secvențelor biologice la diferite nivele de detaliu;
- Identificarea regiunilor regulatorii și a modelelor noi în secvența de date (căutarea regiunilor care se repetă, similarități și modele rare care au semnificație biologică).

Pe baza acestor teme, oamenii de știință încearcă să dezvolte tehnici specializate și cel mai mare volum al cercetării corespunzătoare este orientat spre:

- algoritmi pentru învățare automată destinați prelucrării informației conținute în depozite mari de secvențe;
- noi metode de declanșare automată a algoritmilor pornind de la date (bootstrapping algorithms);
- integrarea surselor multiple de informare într-un singur sistem de învățare și descoperire;
- îmbunătățirea vitezei și acurateții sistemelor probabilistice de raționament.

În secțiunile următoare sunt prezentate cu precădere tehnici bine cunoscute în domeniul sistemelor de recunoaștere automată, puse în serviciul atingerii obiectivelor biologice mai sus menționate așa după cum au fost identificate în lucrări de specialitate relativ recente. Scopul este de a furniza o privire cuprinzătoare.

3.2. Identificarea genelor/proteinelor și clasificarea în categorii

Acest subcapitol include o prezentare sumară a algoritmilor aparținând marii familii a clasificatorilor. Descrierea detaliată a acestor algoritmi poate fi găsită în lucrările citate.

3.2.1. Clasificatorul celor mai apropiați k vecini (k -NN)

Urmând procedura standard a clasificării, pentru a clasifica o nouă secvență (test) fiind dată o mulțime de secvențe (set de antrenare) cu clasificare cunoscută (și prin urmare caracterizare funcțională cunoscută) mai întâi trebuie identificați primii k cei mai apropiați vecini de secvențe de test în setul de antrenare. Acest pas folosește o măsură de similaritate așa încât să ordoneze secvențele în setul de antrenare și în final să le localizeze pe k cele mai similare. Mai departe, și bazat pe etichetele claselor celor mai apropiați k vecini (k -NNs), eticheta cea mai frecventă este asociată secvenței de testare și cu această simplă schemă (vot majoritar), clasificarea ei este îndeplinită. Cea mai critică parte a practicii acestei clasificări este cea a estimării similarității perechilor de secvențe. În diferite versiuni ale clasificatorului k -NN, sunt utilizate funcții alternative de similaritate și diferențele dintre ele constă în modul în care sunt aliniate secvențele. Pot fi folosite fie scorurile aliniamentului global (care consideră secvențele aliniate de-a lungul întregii lungimi) sau scoruri ale aliniamentului local (unde numai porțiuni ale celor două secvențe sunt aliniate). Clasificarea k -NN poate fi găsită adesea în lucrări legate de compararea secvențelor sau predicția funcțiilor [GANA'04], [MARK'03].

3.2.2. Clasificatori ce folosesc lanțuri sau modele Markov generalizate

Tehnicile bazate pe lanțuri Markov sunt foarte populare în modelarea secvențelor datorită abilității lor de a captura constrângeri secvențiale existente în date. În [DESP'02] sunt descrise câteva abordări de clasificări înrudite. Ele pot folosi lanțuri Markov simple (ordinul 1), lanțuri Markov de ordin superior, modele Markov interpolate și modele Markov selective.

În primul caz, secvențele de antrenare sunt partiționate în clase diferite pe baza etichetelor asociate. Apoi, un singur lanț Markov M este construit pentru fiecare din aceste grupuri de secvențe. O secvență de test dată S_r este clasificată calculând la început probabilitatea condiționată $P(S_r|M)$ a acelei secvențe fiind generată de fiecare din acele lanțuri Markov. Lanțul Markov ce corespunde probabilității maxime este selectat pentru a împrumuta eticheta clasei corespunzătoare secvenței de test. Pentru problema a două clase aceasta este făcută prin calcularea rației \log a probabilității furnizată de funcția de clasificare:

$$L(S_r) = \log \frac{P(S_r|M_+)}{P(S_r|M_-)}, \text{ unde } M_+ \text{ și } M_- \text{ sunt lanțurile Markov corespunzătoare clasei}$$

de exemple pozitive respectiv negative.

În cazul lanțurilor Markov de grad mai mare, probabilitatea de tranziție pentru un simbol este calculată pe baza a k simboluri precedente, obținându-se astfel un lanț Markov de ordinul k . În general, aceste modele de ordin superior au o acuratețe a clasificării ridicată comparativ cu modelele de ordin mai mic datorită faptului că ele sunt capabile să captureze secvențe mai lungi din setul de date cu rol constrictiv dar, de asemenea, pot întâmpina probleme. În principiu, problema algerii modelelor de ordin ridicat este dependentă de mărimea setului de

antrenare. O variantă a acestor metode sunt modelele Markov interpolate (IMM) unde o serie de lanțuri Markov sunt construite începând de la ordinul 0 până la ordinul k . k este o variabilă definită de utilizator și fixată în avans. Motivația acestei abordări este că, în ciuda faptului că stările Markov de ordin ridicat capturează un context mai lung, ele au suport redus. Tehnicile IMM au fost dezvoltate inițial pentru problema prezicerii simbolului următor în secvență și nu pentru clasificarea secvențelor. Dar, din moment ce ele constituie o metodă alternativă pentru estimarea diverselor probabilități condiționate pot fi folosite, de asemenea, la clasificare.

Modelele Markov selective (SMM) sunt un alt set de tehnici care adresează problema estimării lanțurilor Markov de ordin ridicat și, ca IMM-urile, ele au fost dezvoltate pentru predicția simbolului următor în secvențe. Aceste tehnici construiesc lanțuri Markov de ordin variat după care retează stările non-discriminative din lanțurile Markov de ordin ridicat. Atribuția principală este de a decide care stări sunt non-discriminative pentru a putea fi înlăturate. Autorii din [DESP'02] au folosit metoda bazată pe pragul simplu de frecvență cu scopul de a înlătura stările care nu au loc frecvent.

3.2.3. Clasificarea secvențelor pe baza caracteristicilor

Clasificarea tradițională folosind algoritmi de învățare automată cum sunt arborii de decizie, clasificatorii bazați pe reguli și mașini cu suport vectorial (Support Vector Machines notate presuratat SVM), pot fi folosiți și pentru clasificarea mulțimilor de secvențe dacă se dovedește că secvențele au fost mai întâi modelate într-o formă potrivită respectivului algoritm dar fără a ignora natura lor secvențială.

Autorii [DESP'02] au dezvoltat un cadru pentru modelarea secvențelor astfel încât algoritmi tradiționali de învățare automată cum ar fi SVM pot fi aplicați ușor. Ei consideră N numărul simbolurilor distincte din diferite secvențe (ex. 4 pentru nucleotide și 20 pentru amino-acizi) după care clasificarea este făcută pe baza funcției de clasificare a lanțurilor Markov (deja menționate). Prin urmare $L(S_i) = u'w$, unde u și w sunt vectori de lungime N^2 . Fiecare dintre dimensiunile acestor vectori corespunde unei perechi unice de simboluri. Dacă ab este perechea de simboluri corespunzătoare dimensiunii j a acestor vectori, atunci $u(j)$ este egal cu numărul care reprezintă frecvența apariției lui ab în secvența S_i , iar $w(j)$ este egală cu logaritmul corespunzător care apare în funcția de clasificare. Aspectele importante ale acestei transformări constau în faptul că permit abordarea fiecărei secvențe ca și un vector într-un spațiu nou ale cărui dimensiuni sunt constituite din toate perechile, tripletele, etc. de simboluri posibile și că fiecare din tehnicile bazate pe lanțurile Markov nu sunt altceva decât un clasificator liniar [MITC'97], unde planul de decizie este definit de vectorul w care corespunde logaritmului din raportul de incertitudine (log-odds¹⁷) al estimării probabilităților maxime pentru diverse probabilități de tranziție. Acuratețea rezultatelor obținute sugerează că modelul liniar de clasificare învățat de SVM este mai bun decât modelele liniare învățate de abordarea bazată pe lanțurile Markov

¹⁷ Un sistem de notare în care valorile sunt logaritmul probabilității relative a unei comparații datorate omologiei sau simplei șanse [BIG'05].

3.2.4. Rețele neuronale pentru clasificarea secvențelor

O aplicație interesantă a rețelelor neuronale (NNets) poate fi întâlnită în [WANG'00a], unde sunt extrase caracteristici din date aflate sub forma proteinelor și folosite în combinație cu rețelele neuronale Baesiene (BNN) pentru a clasifica secvențele de proteine. Într-un mod formal, problema studiată consideră secvența de proteină nemarcată/neclasificată S și o superfamilie F . Scopul este de a verifica în ce măsură S aparține sau nu lui F . Aici, atât similaritatea locală cât și cea globală sunt analizate în mod simultan. Metoda lor include o etapă de extragere a caracteristicilor. În acest pas și pentru a considera similaritatea globală este folosită o metodă de codificare de forma 2-gram [Wu'00] (cu și fără substituție). Aceasta rezultă într-un număr total de 436 modele posibile de 2 litere care ar fi variabilele de intrare în rețeaua neuronală. Pentru a evita „capcana dimensionalității”, este urmat un pas de selectare a caracteristicilor pe baza unei măsuri probabilistice de distanță. Cele mai relevante N_g caracteristici sunt folosite ca intrări pentru BNN. Separat, pentru a compensa pierderea de caracteristici este inclus un coeficient de corelare liniar (LCC). Pentru a incorpora în schema clasificării informația din similaritatea locală este folosită o descriere a lungimii (LS) pentru fiecare secvență. Această măsură LS este utilizată ca intrare separată pentru rețeaua neuronală. Astfel, în final, o secvență de proteine este reprezentată folosind N_g+2 numere reale. Comparând rezultatele cu cele obținute folosind BLAST [ALTS'97], ce se bazează pe aliniamentul secvențelor și SAM [KARC'98], ce folosește modele Markov ascunse, autorii concluzionează că cele trei metode se completează și “combinându-le se obțin rezultate mai bune decât dacă s-ar folosi fiecare clasificator separat”.

O abordare diferită a clasificării pe baza învățării supervizate folosind rețelele neuronale a fost găsită în [ANAS'03], unde intenția este de a demonstra faptul că dacă unii algoritmi deja cunoscuți sunt corect utilizați ei pot duce la soluții mai bune pentru unele probleme cerute în bioinformatică. Îmbunătățirile soluțiilor acestor probleme s-au obținut cu ajutorul rețelelor neuronale feed-forward aplicând scheme mai avansate pentru învățarea supervizată în rețea. Rezultatele sunt evaluate comparativ cu cele ale altor clasificatori cunoscuți.

3.3. Compararea secvențelor

Compararea secvențelor biologice conduce, în mod indirect, la măsuri de similaritate. Există o varietate de tehnici pe această temă așa cum au fost menționate în *capitolul 2*. În acest subcapitol însă va fi făcută o mai mult o prezentare a modului în care acestea se aplică pentru căutari eficiente în bazele de secvențe.

Cele mai frecvent folosite metode pentru aflarea similarității secvențelor biologice sunt cele care se bazează pe aliniamentul perechilor de secvențe (cum ar fi algoritmul Needleman-Wunsch [NEED'70] și algoritmul Smith-Waterman [SMIT'81b]) sau pe principiul modelelor Markov ascunse (HMM) (ex. [KROG'94], [BALD'94]). Conform [LIAO'03], dezvoltarea metodelor puternice pentru detectarea similarității proteinelor (și a secvențelor biologice în general) poate fi împărțită în patru generații.

Metodele timpurii sunt reprezentate de similaritățile perechilor (pairwise similarity). Algoritmul dinamic al lui Smith-Waterman (1981) a rămas referință datorită acurateții rezultatelor, în timp ce algoritmi euristici cum ar fi cei dezvoltati pentru BLAST (Basic Local Alignment Search Tool)¹⁸, FASTA (după FAST –All, sau comparație rapidă a proteinelor/nucleotidelor)¹⁹ sau CLUSTAL²⁰ (pentru aliniamentul mai multor secvențe) negociază o acuratețe redusă în schimbul îmbunătățirii eficienței.

A doua generație ar fi reprezentată de profile și lanțuri Markov ascunse și conțin metode bazate pe ideea de familie. Ei permit biologului să deducă/obțină apropape de trei ori mai multe omologii decât un simplu algoritm bazat pe compararea perechilor [LIAO'03].

În a treia generație, algoritmi ca PSI-BLAST [ALTS'97] și SAM (Sequence Alignment and Modeling System)²¹ folosesc informație stocată în baze de date mari și îmbunătățesc rezultatele peste metodele bazate pe profile, colectând secvențe omoloage și incorporând statisticile rezultate într-un model central.

În generația a patra creșterea acurateții s-a câștigat modelând diferența dintre exemplele pozitive și negative de alinament [LIAO'03].

Toate aceste metode sunt construite pe scoruri de similaritate dintre secvențele alinate, revelând evoluția măsurilor de similaritate [KRAS'04], [LIAO'03], dar încă mai există interes în căutarea unor metode noi deoarece costurile calculului interogărilor cresc liniar cu volumul bazelor de date de secvențe biologice [CAMO'03]. În plus, cele mai multe dintre metodele existente implicit consideră că este importantă alegerea unui algoritm de căutare adecvat, matrice de scoruri sau funcții care pot fi definite și un set de parametri opționali pentru care valorile optime corespund celor mai bune scoruri posibile de similaritate dintre două secvențe.

3.3.1. Căutarea eficientă în baze de date

Mărimea secvenței de acizi nucleici este într-o creștere continuă iar interesul în compararea secvențelor biologice noi cu cele existente este de asemenea în creștere.

Problema de interes matematic este cum să se caute rapid în baze de date mari. Spre exemplu, pentru căutarea a 2000 secvențe a 500 perechi de baze(bp)²², fiecare comparată cu 500 noi secvențe bp ar lua, cu tehnici de programare dinamică, timp proporțional cu $2000(500)^2 = 7.5 \times 10^8$. Acesta este un număr

¹⁸ <http://www.ncbi.nlm.nih.gov/BLAST/>

¹⁹ <http://www.ebi.ac.uk/fasta33/>

²⁰ <http://www.ebi.ac.uk/clustalw/>

²¹ <http://www.cse.ucsc.edu/research-compbio-sam.html>

²² În genetică, două nucleotide în straturi complementare opuse de ADN sau ARN care sunt conectate via legături de hidrogen sunt o pereche de baze (adesea prescurtat bp). Deoarece ADN este de obicei strat dublu, numărul de perechi de baze în "dsDNA" este egal cu numărul de nucleotide într-unul dintre straturi. În ADN, adenina (A) și timina (T), precum și guanina și citozina, pot fi câte o pereche de baze. În ARN, timina este înlocuită de uracil wikipedia.org/wiki/Base_pairs.

inacceptabil de mare și, în cele ce urmează, vor fi prezentate câteva abordări la problema căutării rapide.

O metodă deja cunoscută pentru aceste căutări este cea a lui Wilbur și Lipman [WILB'83], [WILB'84]. Aceasta ține cont de unele descoperiri biologice și în plus metoda de programare dinamică a segmentelor maxime este folosită în acest scop. Waterman [WATE'84] sugerează conexiuni între metoda Wilbur-Lipman și cea a lui Martinez.

Metoda Wilbur-Lipman

În două lucrări importante, Wilbur și Lipman (1983, 1984) dezvoltă ceea ce ei numesc comparare dependentă de context. Metoda lor este conturată de următoarele aspecte (1) produce o listă L de regiuni de potrivire, toate de lungime fixă, (2) ordonează regiunile ca în secțiunea 2.3.7, (3) obțin un aliniament optimal prin procesarea listei L . De fapt ei dezvoltă o teorie pentru un context mult mai general al dependenței și în plus prezintă condiții pentru măsurile lor de similaritate de a avea o distanță asociată care este metrică. Pe mai departe discuția va fi limitată la cazul regiunilor.

În secțiunea 2.3.7 s-a discutat cazul regiunilor. Aici, lista L a regiunilor potrivite exact poate fi restricționată spre exemplu la acele regiuni de lungime exact 4. Wilbur și Lipman [WILB'83] descriu un algoritm de hashing în timp liniar pentru a produce L cu regiuni de lungime fixă.

Se notează o regiune r prin $(w; i, j)$ unde w este un cuvânt de lungime $|w|$ care începe la poziția i în secvența a și poziția j în secvența b . Ca și în secțiunea 2.3.7, $r_1 < r_2$ dacă $i_1 + |w_1| - 1 < i_2$ și $j_1 + |w_2| - 1 < j_2$. De asemenea fie regiunea $r_0 = (\phi; 0, 0)$ un element minimal și $r_* = (\phi; n, m)$ un element maximal. $\Gamma = (r_1, r_2, \dots, r_l)$ este o cale dacă $p < q$ implică $r_p < r_q$. Scorul unei căi Γ este dat de următoarea relație

$$\text{score}(\Gamma) = \sum_{k=1}^l s(r_k) - \sum_{k=1}^{l-1} g(i_{k+1} - |w_k| - i_k - 1, j_{k+1} - |w_k| - j_k - 1), \quad (3.3.1-1)$$

unde $s(\cdot)$ este un scor de similaritate pentru o regiune r_k , ca $|w_k|$, și $g(\cdot, \cdot)$ este o penalitate pentru gap. Atunci,
 $S(a, b) = \max\{\text{score}(\Gamma) : \Gamma \text{ este o cale de la } r_0 \text{ la } r_*\}$.

Algoritm Wilbur - Lipman.

Algoritmul este după cum urmează: se fac două liste a două regiuni L^- , ordonat de „<”, și L^+ ordonat după ordinea obișnuită „□” a celor mai bune scoruri din r_0 la regiunea listată în L^+ .

(0) Fie $L^- = L$ și $L^+ = \phi$

(1) $r_q = \{\min r : r \in L^-\}$
 $\text{score}(r_q) = S(r_q)$

(2) Începe la cel mai mare element a lui L^+

(A) Mută în jos L^+ până când $r_u < r_q$ astfel încât

$$\gamma = \text{score}(r_u) - g(i_q - |w_{r_u}| - i_u - 1, j_q - |w_{r_u}| - j_u - 1) + s(r_q) > \text{score}(r_q).$$

Dacă nu există vreo regiune r_u în L^+ mai jos de r_q sub γ cu un scor mai mare decât scorul $(r_q) - s(r_q)$, sau dacă inegalitatea nu poate fi satisfăcută, se merge la (C)

(B) Se fixează score $r_q = \gamma$ și salt la (A)

(C) Se înlătură r_q din L^- și se inserează în L^+ după regula impusă de „□”.

Dacă $L^- \neq \phi$, salt la (A).

Este posibil să se modifice acest algoritm, pentru a obține segmente de similaritate maximă în loc de aliniamente de similaritate maximă.

3.3.2. Un algoritm vectorizat pentru segmente maxime

O altă abordare de calcul rapid a fost propusă de T.F.Smith la Los Alamos national Laboratory. Smith a modificat algoritmul Smith-Waterman astfel, utilizând arhitectura de vector este posibil să se realizeze comparații între un foarte mare număr de secvențe de acizi nucleici în timp rezonabil. Spre exemplu, un studiu al lui Smith et al. [SMIT'84] raportează toate comparațiile de perchi de secvențe dintre 204 secvențe de la vertebrate care au fost realizate în aproximativ 170 min, la o rată de peste 240 secvențe per minut cu o medie a lungimii secvenței de 800 nucleotide.

Vectorizarea algoritmilor este un subiect relativ nou. Ideea esențială este că mașina execută un număr de operații simultan și câștigă acest factor peste operațiile obișnuite liniare ale secvențelor. Dacă atenția este focalizată pe calcularea lui H , nici o calculare a lui H_{ij} executată nu poate depinde de rezultatele unor alte calcule executate simultan. Aceasta elimină linie cu linie sau coloană după coloană construirea matricei. Ce rămâne este calculare de blocuri ale lui H_{ij} pe diagonale negative, i.e. $i + j = \text{constant}$.

Collins și Coulson [COLL'84] discută algoritmi de procesare paralelă pentru metoda lui Sellers. Abordarea lui Smith ar trebui să permită implementări mult mai eficiente pentru procesarea paralelă.

3.3.3. Metoda regiunilor

Metoda regiunilor lui Martinez (descrisă în secțiunea 2.3.7) și metoda lui Wilbur-Lipman sunt evident asemănătoare, cu toate că Wilbur-Lipman vine cu o tehnică diferită de optimizare. Desigur că cele două tehnici au fost dezvoltate independent dar subiectul este evident. Pentru calcularea rapidă se limitează logic lista regiunilor și aliniamentele vor fi produse mult mai rapid.

Martinez [MART'83] a sugerat abordarea acestor probleme via regiuni. Localizarea tuturor repetițiilor a R secvențe de lungime n , poate fi făcută în timp proporțional cu nR . „Repet” în acest context se referă la un cuvânt w care apare în toate secvențele R . O asemenea repetare se numește regiune și este plasată într-o listă L . O ordine parțială „<” în R dimensiuni (în loc de două dimensiuni) este plasată în L .

Exact aceiași algoritmi din secțiunile 2.3.7 și cei menționați în această secțiune pot fi folosiți pentru a produce un aliniament multiplu de secvențe. Limitarea aici este că cerința de potriviri exacte în toate cele R secvențe este imperios necesară. De asemenea, porțiunile aliniate ar fi foarte convingătoare și un asemenea aliniament ar putea fi produs într-un timp rezonabil.

Karlin et al. [KARL'84] nu declară un algoritm pentru alinierea regiunilor lungi care se potrivesc, cu toate că din aceste potriviri rezultă aliniamente. Este evident că folosirea regiunilor ca intrări într-o metodă bazată pe regiuni ar produce aliniamente semnificative. În plus, dacă lista cu regiuni nu este total ordonată, o asemenea metodă de aliniament ar fi foarte utilă în rezolvarea posibilelor aliniamente.

3.3.4. Algoritmi de programare dinamică

Într-un articol elaborat de Sankoff și Cedergren [SANK'83] sunt revizuite metodele care realizează aliniamente a R secvențe fiind dat un arbore care le conectează. Fiecare nod interior al arborelui dat T are un grad de cel puțin trei și cele R secvențe sunt atașate de cele R noduri terminale. Algoritmul construiește secvențe pentru fiecare nod interior și dă un aliniament legând secvențele originale și cele reconstruite. Dacă arborele are N noduri interioare acesta este un aliniament de R+N secvențe.

Waterman [WATE'84] încearcă să dea aspecte ale programării dinamice acestei probleme. Se presupune ca cele R secvențe sunt a, b, ..., r. Costul unui aliniament general al lui R original plus N secvențe construite este suma costului aliniamentului de perechi în arborele T. Ideea este de a gândi la alinierea $a_1a_2, \dots, a_i, b_1b_2 \dots b_j \dots, r_1 r_2 \dots r_s$. Ultima coloană a aliniamentului poate fi înlăturată de lângă coloanele inițiale. Pentru a, b, ..., r ultima coloană va apare ca

$$\begin{array}{c} \varepsilon_1 a_i \\ \varepsilon_2 b_j \\ \dots \\ \varepsilon_R r_s \end{array}$$

unde $\varepsilon_i \in \{0, 1\}$, $0 \cdot a_i = \Delta$, și $\varepsilon = (\varepsilon_1, \dots, \varepsilon_R) \neq 0$. Această ultimă restricție ține doar ultima coloană de a nu fi toată Δ . Pasul de programare dinamică este

$$D_{i,j,\dots,s} = \min_{\varepsilon \neq 0} \left\{ D_{i-\varepsilon_1, j-\varepsilon_2, \dots, s-\varepsilon_R} + \min_{\varepsilon \neq 0} \text{length} \left(\begin{array}{c} \varepsilon_1 a_i \\ \varepsilon_2 b_i \\ \cdot \\ \cdot \\ \varepsilon_R r_s \\ x_1 \\ \cdot \\ \cdot \\ x_N \end{array} \right) \right\} .$$

Ultimul termen este pentru a indica faptul că literele $x_1 x_2 \dots x_N$ pentru N noduri interioare a fost determinat într-un mod optimal. Pentru acest scop este folosită metoda economică a lui Fitch [FITC'71], generalizată adecvat. Timpul de calcul al acestui algoritm pentru R secvențe de lungime n este $O(2^R n^R N)$ unde 2^R provine din $\varepsilon_1, \dots, \varepsilon_R$ la fiecare pas, n^R vine de la i, j, \dots, s , și N vine de la metoda economică a lui Fitch. Pentru secvențe de lungime 100, aceasta este aproximativ $O(10^{2 \cdot 3^R} N)$ astfel $R=3$, este atît de larg cît este posibil.

Independent de munca lui Sankoff [SANK'76], Waterman et al [WATE'76] prezintă un algoritm similar care nu presupune în mod explicit un arbore. În cadrul muncii lui Sankoff aceasta ar corespunde unui arbore, un nod interior și R noduri terminale. A fost folosită o funcție $d(x, y, \dots, z)$ de R variabile pentru a generaliza metrica pe litere astfel, este posibil de a avea o ponderare diferită, dar ideea esențială este în lucrarea lui Sankoff.

3.3.5. Un algoritm bazat pe conceptul de modele Markov ascunse

Secvențele biologice funcționale de obicei apar în familii și multe dintre cele mai puternice metode de analiză sunt bazate pe identificarea relațiilor dintre o secvență individuală cu o familie de secvențe. Prin urmare, identificarea faptului că o secvență aparține unei familii și aliniind-o cu ceilalți membri adesea permite inferențe despre funcția sa. Dacă spre exemplu avem un set de secvențe care aparțin unei familii, se poate face o căutare într-o bază de date pentru mai mulți membri folosind aliniament de perechi cu unul dintre membrii familiei cunoscute ca secvență de query. Pentru a fi mai aproape de realitate se poate face căutare după fiecare dintre membrii familiei unul câte unul. Oricum căutarea pe baza comparației perechilor pentru oricare dintre secvențe poate să nu găsească secvențele apropiate ca distanță cu cele avute deja. O abordare alternativă este de a folosi în căutare caracteristici statistice ale întregului set de secvențe. Similaritatea, chiar și atunci când apartenența la o familie este clară, datorită acurateții aliniamentului, poate fi adesea îmbunătățită prin concentrarea pe caracteristici care sunt conservate în întreaga familie.

În acest caz intervine aliniamentul de secvențe multiple în care scopul este de a construi un model probabilistic. În particular, este dezvoltat un anumit tip de HMM care să se potrivească bine modelării multiple de aliniamente. Acestea sunt denumite profile HMM după profilele standard, care sunt structuri apropiate nonprobabilistice introduse anterior de Gribskov, McLachlan și Eisenberg [GRIBS'87]. Profilele HMMs sunt probabil cea mai populară aplicație a HMM în biologia moleculară până în jurul anului 2000 [DURB'98].

Un model liniar Markov ascuns este o secvență de noduri fiecare corespunzând unei coloane într-un aliniament multiplu. În cazul HMM folosit în aplicația SAM²³ fiecare nod are, de asemenea, o stare de potrivire/matching, de inserare/insertion și o stare de ștergere/deletion. Fiecare secvență folosește o serie a acestor stări pentru a traversa modelul de la început până la sfârșit. Folosirea unei stări de *potrivire (matching)* indică faptul că secvența are un caracter în acea coloană în timp ce folosirea unei stări de *ștergere (deletion)* semnalează contrariul. Stările de *inserare* permit secvențelor de a avea caractere adiționale între coloanele de aminoacizi. Stările de poziții vecine sunt conectate prin linii unde, pentru fiecare dintre aceste linii, există asociată o probabilitate de tranziție care este probabilitatea de trecere dintr-o stare în alta. Topologia de bază a modelului este ilustrată în Figura 3.3.5-1.

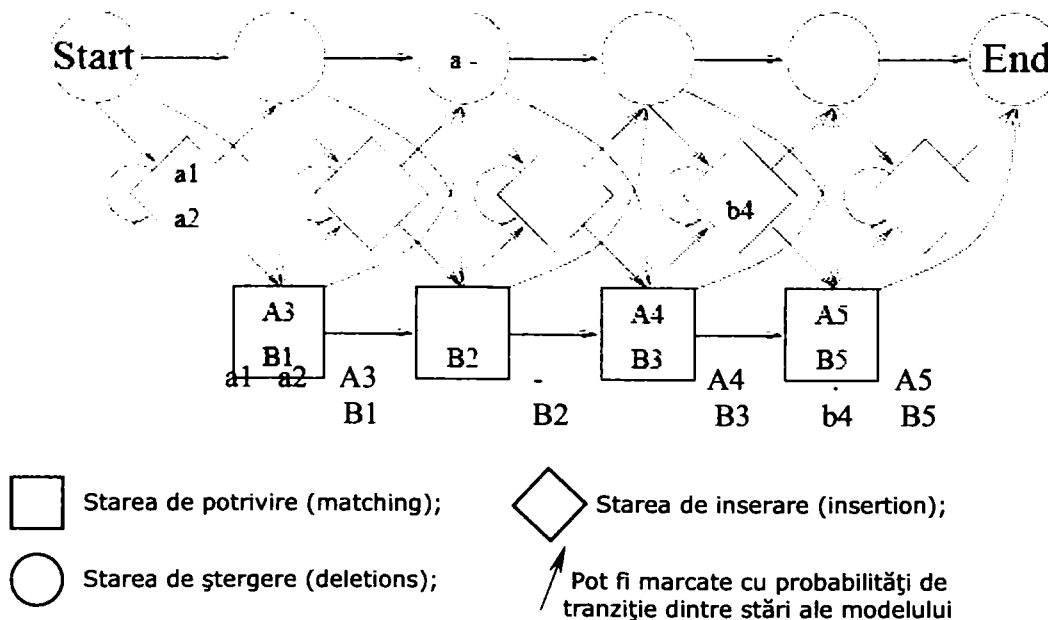


Figura 3.3.5-1. Un model Markov ascuns liniar și exemplu de aliniament.

²³ Sequence Alignment and Modeling System,
<http://www.cse.ucsc.edu/research/compbio/HMM-apps/HMM-applications.html>

Topologia acestui model a fost aleasă pentru a caracteriza inserările și ștergerile în mod similar cu tehnicile de aliniament biologic al secvențelor pe baza penalităților caracterelor lipsă (gap penalties). Modelele sunt antrenate pe o familie de secvențe de proteine utilizând un algoritm de maximizare al așteptărilor și o varietate de algoritmi euristici. Un model antrenat poate să fie folosit atât pentru a genera aliniament multiplu cât și pentru căutare în baze de date a noi membri ai familiei.

Avantajul principal al acestor modele peste metodele standard de căutare îl reprezintă abilitatea lor de a caracteriza o întreagă familie de secvențe. Astfel, fiecare poziție are o distribuție a bazelor, așa cum o au tranzițiile între stări. Aceasta înseamnă că HMMs liniare au distribuții ale caracterelor dependente de poziție, precum și inserții și penalități de ștergeri (gap penalties) dependente de poziție. Aliniamentul fiecărei familii la un model antrenat va genera automat aliniamentul multiplu al secvențelor componente. În multe feluri aceste modele corespund profilelor. Metoda implementată în SAM-T2K, de detectare a omologiilor distante, este o metodă HMM iterativă de căutare pentru crearea unui HMM dintr-o singură secvență de proteine sau aliniament sursă folosind căutarea iterativă a secvenței într-o bază de date. Metoda este considerată de autori [HUGH'96] cel mai sensibil algoritm bazat pe determinarea de la distanță a omologiei unei secvențe. Detalii care descriu pe larg strategia de funcționare a SAM pot fi consultate la adresa: http://www.cse.ucsc.edu/research/compbio/html_format_papers/hughkrogh96

În concluzie, această trecere în revistă a principalelor metode nu epuizează tipurile de probleme care prezintă interes. Algoritmii existenți pot fi modificați pentru a folosi la soluționarea unei mari varietăți de probleme.

3.4. Identificarea regiunilor regulatorii și tiparelor noi în date secvențiale

Algoritmii și tehnicile folosite în acest scop acoperă o mare varietate de abordări. O listă detaliată poate fi găsită în [BREJ'00]. În continuare vor fi identificate și descrise pe scurt principalele categorii.

3.4.1. Căutarea exhaustivă

Cel mai bun algoritm pentru identificarea tiparelor semnificative în secvențe biologice lungi nu este ușor de proiectat și prin urmare, în multe aplicații, căutarea aleasă este cea exhaustivă. Una dintre cele mai simple metode în descoperirea tiparelor este de a enumera toate tiparele posibile care satisfac constrângerile impuse de utilizator. Pentru fiecare tipar este estimată mai întâi frecvența de apariții în secvența de intrare și pe baza acesteia este atribuit un scor sau o semnificație statistică. Apoi, sunt identificate tiparele cu cele mai mari scoruri sau tipare cu scoruri peste un anumit prag. Această metodă este adecvată doar când este vorba de tipare scurte deoarece timpul de execuție crește exponențial cu lungimea tiparelor ce trebuie identificate. Pe de altă parte, această metodă poate întotdeauna să găsească cele mai interesante tipare [VANH'98], [TOMP'99], [SMIT'90].

Un mod mai eficient este de a aplica metoda de înlăturare a enumerării tiparelor. În această metodă este identificat un tipar lung și fără goluri (gap) care este posibil să apară cu unele nepotriviri în cel puțin K secvențe. El poate începe de la tipare scurte care apar în cele mai multe K secvențe și să le extindă atâta timp cât suportul nu scade sub K . La fiecare pas este nevoie de extensia tiparului (identificat curent) în toate felurile posibile. Această extensie este urmată de căutarea aparițiilor tiparelor extinse încă în cel puțin K secvențe. Odată ce obținem un tipar care nu poate fi extins fără pierderea de suport, acest tipar este considerat cel maximal și poate fi înregistrat ca final. Acest algoritm de căutare realizează de fapt o primă căutare în adâncime în arborele tuturor secvențelor posibile iar ramurile care nu realizează suport pentru vreun tipar sunt înlăturate. Teoretic, această metodă poate funcționa bine dar creșterea exponențială a timpului de căutare rămâne o problemă [RIGO'98].

Căutarea exhaustivă într-un graf ar fi o altă metodă din acest grup și sunt multe posibilități de determinare a aparițiilor unui tipar. În [PEVZ'00], este descrisă problema căutării tiparului de lungime dată L care apare în toate secvențele de intrare n cu cel mult d nepotriviri. Pentru tipare mai complicate este necesară folosirea abordărilor euristice care nu găsesc neaparat cel mai bun tipar dar pot converge la o soluție suboptimă. O descriere a unei asemenea tehnici, bazată pe partiționarea Gibbs poate fi găsită în [BREJ'00] sau [LAWR'93].

3.4.2. Metode de învățare automată

Când un tipar este exprimat sub forma unui model stohastic cum este modelul Markov ascuns sau matricea cu ponderi a pozițiilor (care este tot o versiune mai simplă a HMM), pot fi folosite tehnicile iterative de maximizare a așteptărilor care nu converg în mod necesar la maximul global. În [LAWR'90], este folosit un algoritm de învățare pe baza maximizării așteptărilor (Expectation Maximization notat și cu EM) al cărui scop este de estimare a parametrilor unui model probabilistic/stohastic al tiparului, care apare o singură dată la o poziție necunoscută în fiecare secvență din familia de secvențe date. Această abordare poate fi extinsă la modele mult mai complicate ca cea cu goluri (gaps) flexibile sau modelul mixt finit [BAIL'94]. O descriere detaliată este făcută în [BREJ'00].

Modelele Markov ascunse pot fi folosite ca modele pentru familii de secvențe. Există trei chestiuni importante legate de utilizarea lor: topologia HMM, procesul de antrenare și căutarea secvențelor. După cum este descris în lucrări ca [KROG'94], [HUGH'96], modelele sunt formate din trei tipuri de stări: starea de potrivire (match state) care specifică probabilitatea de distribuție a caracterelor (amino acid sau nucleotid) în fiecare poziție conservată, stări de inserare care modelează goluri posibile dintre stări și stări de ștergere (deletion states) care determină devierea unora din stările de potrivire. Odată ce au fost setați toți parametrii modelului, poate fi calculată cea mai probabilă cale de-a lungul modelului pentru o secvență dată. În partea de antrenare, topologia HMM este dată în paralel și în concordanță cu o familie de secvențe pe care dorim să o caracterizăm cu HMM. Parametrii modelului sunt estimați într-un mod în care el poate să genereze secvențe similare cu cele din familia specificată. În [DURB'98] pot fi găsite explicații detaliate despre această abordare. Fiind dată o secvență, putem calcula (folosind algoritmul Viterbi) cel mai probabil aliniament al acestei secvențe folosind tiparul reprezentat de HMM. Este menționat în [HUGH'96] că în timp ce această metodă funcționează, în practică trebuie aplicată căutarea în mod simultan în mai multe secvențe datorită necestății unei cantități mari de date.

HMM au fost îmbunătățite prin reducerea numărului de parametri implicați și prin descoperirea de subfamilii. În primul caz înseamnă că topologia poate fi ajustată în timpul antrenării [HUGH'96] sau prin detectarea în secvențe de tipare scurte care sunt imediat transformate în stări de potrivire în HMMs. Al doilea înseamnă că uneori familiile pot fi organizate în subfamilii așa cum a fost făcut în [KROG'94] combinând câteva HMMs mici, cu topologie standard, într-unul mai mare. Un avantaj al acestui algoritm este că odată ce o parte a modelului este limitată la o subfamilie numai secvențele acestei subfamilii vor fi folosite pentru antrenarea acestei părți de model.

3.4.3. Metode bazate pe alinierea secvențelor

Găsirea tiparelor semnificative și executarea aliniamentului local sunt funcții strâns apropiate [BREJ'00]. Odată ce s-a executat aliniamentul a multiple secvențe pot fi identificate mai ușor tiparele.

Propunerea autorilor din [GORO'97] este de a găsi cele mai semnificative secvențe comune și structuri dominante într-o mulțime de secvențe ARN sau ADN. Ideea este de a construi o schemă computațională care aliniază local o colecție de secvențe ARN folosind constrngerii pentru secvențe și structuri. Ei identifică $M \leq N$ secvențe conținând cel mai semnificativ motiv comun care va apare în aliniamentul a M secvențe. Subseturile slab aliniat sunt eliminate. În acest scop este folosit algoritmul lui Sankoff [SANK'85], pentru a alinia optim multiple secvențe de ARN, dar el este simplificat datorită complexității de $O(L^{3N})$ a timpului, pentru N secvențe de lungime L . Simplificarea este făcută în două moduri. Un mod este adaptarea algoritmului local Smith-Waterman pentru aliniamentul local al secvențelor, maximizând un scor bazat pe o combinație a similarității secvențelor și structură. Această schimbare permite aflarea scorului cel mai ridicat al aliniamentului local al secvențelor ARN. A doua simplificare este de a reduce complexitatea timpului la $O(L^4)$ prin nepermiterea structurilor ramificate. Autorii revendică faptul că nu se așteaptă identificarea completă a motivelor structurale într-o singură trecere. De aceea, este aplicat un aliniament progresiv pentru perechi de secvențe ca și în tehnica folosită de Clustal 5.5 și este urmată strategia algoritmului „greedy” în CONSENSUS [HERT'90] pentru a menține soluții intermediare cu scopul minimizării probabilității lipsei soluției optimale.

Datele folosite pentru investigare sunt din experimentele SELEX [TUER'90] (Systematic Evolution of Ligands by Exponential Enrichment) pentru care s-a propus structura CONSENSUS. Constând în programare dinamica 4D și un algoritm greedy pentru compararea perechilor de secvențe, metoda a fost capabilă să recunoască complet aliniamentele publicate ale motivelor conservate. Nu întotdeauna a fost obținută întreaga structură dar aliniamentul central este considerat a fi o dată de intrare utilă la metodele existente, în scopul completării identificării celor mai importante secvențe și motive structurale.

3.4.4. Metode care se bazează pe conceptul de graf

O presupunere despre aflarea tiparelor în secvențele biologice este aceea că regiunile conservate în timpul evoluției sunt importante funcțional. Interesul constă în căutarea unui element regulator diferit de regiunile regulatorii de la multe gene ale aceleiași specii. În acest sens, pot fi folosite regiunile regulatorii ale aceleiași gene luate de la mai multe specii. Aceasta se poate realiza folosind reprezentarea arborelui filogenetic al datelor biologice. În consecință, pot fi găsite tiparele scurte cel mai bine conservate în timpul evoluției [BREJ'00], [BLAN'00].

Folosirea conceptului de graf este găsită în [ESTE'04] unde a fost implementată o metodă top-down pentru minarea tiparelor celor mai frecvente în datele biologice. Fiind dată o bază de date de secvențe D , un graf conceptual Γ și o valoare suport minimă δ se poate descoperi mulțimea celor mai specifice tipare frecvente. Metoda descrisă este bazată pe următoarele observații cheie: tiparele frecvente trebuie să fie subsecvențe ale secvențelor de date; tiparele mai frecvente ar trebui generate înaintea celor mai puțin specifice pentru a evita generarea inutilă a tiparelor frecvente care nu sunt cele mai specifice. Pentru a expluata aceste observații, autorii din [ESTE'04] pornesc cu toate subsecvențele identificate în datele de intrare (specificând o lungime maximă) și continuă să le scurteze cu câte un element până când devin frecvente. Astfel, tiparele cu frecvență redusă sunt generalizate folosind conceptul graf până când ele devin frecvente.

3.4.5. Metode hibride

Există o varietate de alte metode care sunt realizate ca și combinații ale metodelor deja prezentate. Pentru găsirea repetițiilor tandem în secvențele biologice, autorii [COWA'98] descriu o metodă iterativă pentru identificarea repetițiilor pe baza propriului aliniament al unei secvențe genomice date, în timp ce în [BENS'99] aceeași problemă este abordată folosind principiul din teoria probabilității.

În această categorie sunt incluse metodele de minare bazate pe motoare web pentru căutarea de text care încorporează tehnici pentru interogări ale bazelor de date, potriviri și pruning²⁴. Lucrări în această direcție sunt [WANG'00a] sau [REBH'98].

În [CHIA'04] sistemul de căutare în text este focalizat pe 4 tipuri de informații referitoare la gene: funcții biologice, boli asociate, gene înrudite și relații gene-gene. Sistemul de informații dezvoltat pentru gene (GIS) are două module: unul pentru filtrare (screening) și unul pentru extragerea relațiilor gene-gene. Primul furnizează informații despre funcțiile biologice, boli asociate și gene înrudite pentru gena interogată și funcția acestui modul este îndeplinită prin trei agenți ai căutării de informație în documente (colectează documentele medicale din colecția PubMed), selectarea propozițiilor (selectează părțile importante ale abstractelor pentru o procesare ulterioară) și analiza lexiconului (caută, numără și indexează cuvintele cheie ale funcției biologice). Al doilea modul (activat după extragerea relațiilor genă-genă descrise în abstracte) estimează în ce măsură relația dintre o pereche de gene este pozitivă, cooperativă sau negativă.

3.5. Produse software dezvoltate pentru analiza secvențelor biologice

Datorită cantității mari de date stocate automat este imperios necesară dezvoltarea de aplicații complexe în vederea explorării și extragerii de informații cât mai utile. Prin urmare, există o multitudine de produse software destinate analizei secvențelor biologice și câteva dintre ele vor fi menționate pe baza funcțiilor pe care le servesc.

Căutare în baze de date

Pentru simpla căutare în baze de date de secvențe sunt menționate câteva programe specializate. Majoritatea acestora folosesc tehnici bazate pe potrivire sau detalii de identificare existente în bazele de secvențe. În general programele prezentate oferă mai mult decât simpla căutare sau identificare a unei secvențe, aceasta putând fi doar o funcție din ceea ce oferă toolul respectiv.

SRS-EMBL, **SRS6-EBI** (<http://srs.embl-heidelberg.de:8000/srs5/>, își derivă numele din Sequence Retrieval Systems, și sunt programe utilizate pentru

²⁴ metodă care presupune înlăturarea ramificațiilor inutile în pașurarea unui algoritm

interogare generală a secvențelor din baze de date aparținând celor două instituții EMBL și EBI. <http://srs6.ebi.ac.uk/srs6bin/cgi-bin/wgetz?-page+top>)

Pfam –este atât o bază de secvențe organizată pe grupuri de familii, ce conține aliniamente multiple ale secvențelor de ADN și proteine cât și un complex de programe care oferă posibilitatea de căutare de informații în aceasta. Profilele modelelor Markov ascunse (profile HMM) construite din aliniamentele Pfam pot fi foarte utile pentru a recunoaște dacă o nouă proteină aparține sau nu unei familii existente chiar dacă omologia este slabă. Spre deosebire de metodele de aliniere a perechilor de secvențe (e.g. BLAST, FASTA), HMMs Pfam lucrează sensibil cu proteinele care aparțin mai multor domenii. Ca detaliu, versiunea Pfam cea mai recentă, **18.0** (August 2005) conține aliniamente și modele pentru **7973** familii de proteine, furnizate de bazele de date de proteine **Swissprot 47.0** și **SP-TrEMBL 30.0**. <http://pfam.wustl.edu/>. Din punct de vedere structural, Pfam este format în două moduri separate. Pfam-A sunt aliniamente multiple realizate de experți umani în timp ce Pfam-B este o grupare generată automat a restului unei baze de date de proteine nonredundante derivată din baza de date PRODOM [BATE'04].

BLAST (**B**asic **L**ocal **A**lignment **S**earch **T**ool), este un program informatic complex care caută regiuni de similaritate locală între secvențe. Programul compară secvențele de nucleotide sau proteine cu secvențe obținute din diverse baze de date și calculează semnificația statistică a potrivirilor. BLAST poate fi folosit și pentru a deduce relații evoluționare și funcționale între secvențe precum și pentru identificarea membrilor familiilor de gene. (<http://www.ncbi.nlm.nih.gov/blast/Blast.cgi>)

Prosite. Acest tool permite parcurgerea secvențelor de proteine (fie din baza de cunoștințe, Swiss-Prot/TrEMBL, fie furnizată de user) pentru detectarea aparițiilor de tipare, profile și motive stocate în baza de date PROSITE sau caută alte baze de date de proteine pentru potriviri după motive specifice. (<http://au.expasy.org/tools/scanprosite/>)

PRINTS, reprezintă un compediu pentru ampretele proteinelor. O amprentă este un grup de motive conservate utilizate pentru a caracteriza o familie de proteine; puterea sa de diagnostic este rafinată de scanarea iterativă a amestecului de date oferite de SWISS-PROT/TrEMBL. De obicei motivele nu se suprapun, dar sunt separate de-a lungul unei secvențe, deși ele pot fi contigue (apropiate) în spațiul 3D. Ampretele pot codifica plierile și funcționalitățile proteinelor mai flexibil și mai puternic decât o pot face motivele singure, potențialul de diagnostic complet derivând din contextul mutual oferit de vecinii motivelor. (<http://umber.sbs.man.ac.uk/dbbrowser/PRINTS/>)

InterPro este atât interfața web a unui motor de căutare cât și o bază de date ce conține familii de proteine, domenii și locații funcționale în care caracteristicile identificabile găsite în proteinele cunoscute pot fi aplicate secvențelor de proteine necunoscute. (<http://www.ebi.ac.uk/interpro/>)

SMART (a Simple Modular Architecture Research Tool). Acest program permite identificarea și adnotarea domeniilor genetice mobile și analiza arhitecturii lor. Aceste domenii sunt extensiv adnotate ținând cont de distribuția pe linie de evoluție, clasa funcțională, structuri terțiare și reziduuri importante din punct de vedere funcțional. Fiecare domeniu găsit într-o bază de date nonredundante de proteine precum și parametrii de căutare și informațiile taxonomice sunt stocate într-un sistem de baze de date relaționale. Interfețele utilizator la această bază de date permit căutări de proteine care conțin combinații specifice de domenii în gruparea definită. Toate detaliile sunt disponibile în publicațiile asociate prezentării acestui tool. (<http://smart.embl-heidelberg.de/>)

Notă: În general, fiecare bază de date de secvențe biologice oferă pe lângă simpla stocare a datelor și interfețe cu mecanisme de căutare și explorare a acestora.

Compararea de secvențe biologice

Compararea perechilor de secvențe biologice

Pentru compararea perechilor de secvențe biologice există programe bazate pe diverse tehnici de similaritate care, conform afirmațiilor existente pe site-ul oficial al Institutului European de Bioinformatică (EBI), recunosc omologiile apropiate cu o identitate mai mare de 30%.

În principal, algoritmi implicați în căutarea similarității și aliniamentului de secvențe folosesc două metode diferite de calculare (cantitativă) a similarității pentru o pereche de secvențe. Ele sunt măsura distanței și măsura similarității (de obicei reprezentate sub forma unei matrice). În general, pentru orice pereche dată de secvențe acești doi factori sunt invers proporționali unul față de altul, astfel cu cât este mai mare distanța dintre secvențe cu atât este mai mică similaritatea și vice versa.

În cele ce urmează atenția va fi acordată la trei dintre cele mai utilizate moduri pentru parcurgerea bazelor de date de ADN și de proteine în scopul determinării similarității unei secvențe interogatoare cu o alta sau cu cele conținute într-o bază de secvențe. Primele două, FASTA și BLAST sunt metode și tooluri în același timp bazate pe algoritmi euristici care folosesc asumții și aproximări. Ambele programme, cu toate că folosesc abordări diferite, mai întâi identifică potrivirile exacte și foarte scurte dintre secvența de test și cea/cele conținute în baza de date. După care, cele mai bune potriviri scurte din prima etapă sunt extinse pentru a se face căutarea pentru secțiuni mai lungi de similaritate. În final, cele mai bune potriviri sunt optimizate cu ajutorul unor forme de programare dinamică.

BLAST (Basic Local Alignment Search Tool, <http://www.ebi.ac.uk/blast/>), amintit și în categoria de programe cu ajutorul cărora se caută în baze de date de secvențe biologice este folosit și pentru a compara o secvență nouă cu cele conținute în bazele de date de nucleotide sau proteine accesate prin alinierea noii secvențe cu genele caracterizate anterior. Importanța acestui tool este de a găsi regiuni de secvențe similare care vor conține chei funcționale și evoluționare despre structura și funcția noii secvențe. Regiunile de similaritate detectate via acest tip de

aliniament pot fi ori la nivel local unde regiunea de similaritate este bazată pe o singură locație sau la nivel global, unde regiunile de similaritate pot fi detectate de-a lungul codului genetic de altfel neînrudit.

FASTA (pronunțat FAST-Aye, <http://www.ebi.ac.uk/fasta33/>) provine de la FAST-All, reflectând faptul că programul poate fi folosit pentru o comparare rapidă a proteinelor sau nucleotidelor. El atinge un nivel înalt de sensibilitate în căutarea similarității la o viteză ridicată. Aceasta performanță este realizată prin executarea căutărilor optimizate pentru aliniamente locale folosind matricea de substituție. Viteza ridicată a acestui program este obținută folosind formele observate de potriviri de cuvinte (unde prin cuvânt se înțelege o secvență formată din 5-6 nucleotide sau amino acizi, în funcție de dorința utilizatorului) pentru a identifica potențialele potriviri înainte de începerea optimă a căutării, care este mare consumatoare de timp. Compromisul dintre viteză și sensibilitate este controlat de un parametru care specifică mărimea cuvântului. Crescând valoarea acestui parametru descrește numărul de potriviri de fond. Nu este investigată fiecare potrivire de cuvânt dar inițial se caută segmente care să conțină ceva potriviri. Astfel, acest program demonstrează similaritatea secvențelor și căutarea omologiilor în baze de date de nucleotide și proteine. El poate fi foarte specific când identifică regiuni lungi de similaritate scăzută în special pentru secvențe puternic divergente. De asemenea se poate conduce căutarea similarității și omologiei în baze de date complete, proteomice sau genomice.

Ssearch (<http://ori.nibb.ac.jp/SIT/SSEARCH.html>) este cea de-a treia abordare, numită Smith-Waterman și care de asemenea este implementată sub forma unui tool. Este în totalitate bazată pe programare dinamică ce realizează efectiv toate comparațiile posibile de perechi de secvențe conținute de baza de date. Prin urmare, este o tehnică mult mai sensibilă în comparație cu FASTA și BLAST dar este mult mai costisitoare din punct de vedere al efortului computațional și al timpului. În general, sensibilitatea și selectivitatea căutărilor făcute de FASTA și BLAST sunt comparabile cu cele ale căutărilor în cazul Smith-Waterman datorită incorporării unor parametri statistici în cazul primelor două programe. Oricum, Smith-Waterman rămâne "*standardul de aur*" pentru aliniamentul perechilor de secvențe proteină-proteină sau nucleotidă-nucleotidă, viteza sa putând fi îmbunătățită prin utilizarea unui hardware specializat sau o platformă Linux extinsă [SSEA'05].

Software specializat pentru compararea de secvențe multiple:

Compararea de secvențe biologice multiple este strâns legată cu generarea de aliniamente multiple a secvențelor. În acest scop, majoritatea programelor dezvoltate au la bază algoritmi destinați aliniamentului multiplu după care trecerea la utilizarea unei metode de determinare a similarității dintre secvențele comparate. În general, toolurile existente sunt capabile să recunoască omologii distanțe cu o identitate mai mică decât 30%

PSI-BLAST. Denumirea provine de la Position Specific Iterative BLAST (PSI Blast) și se referă la o caracteristică a lui BLAST 2.0 în care un profil (o tabelă care afișează frecvențele²⁵ fiecărui amino acid în fiecare poziție a secvenței de proteine) este construit automat dintr-un aliniament multiplu al potrivirilor cu cele mai ridicate scoruri într-o căutare BLAST inițială. Această strategie iterativă de căutare conduce la o creștere a sensibilității. (<http://www.ncbi.nlm.nih.gov/BLAST/>)

PDB BLAST. Acesta este o variantă îmbunătățită a lui PSI BLAST pentru găsirea omologilor distanți ale secvenței interogatoare într-o bază de date specifică. Metodele de căutare a profilului pot găsi secvențe omoloage mai distanți decât în cazul căutării simple. PDB BLAST caută în baza de date PDB cu PSI BLAST folosind un profil generat din baza de date NR (Natural Resource, care este baza primară de date din NCBI). O asemenea îngustare a scopului căutării a fost ilustrat pentru a îmbunătăți eficiența căutării în baza de date NR. Prin urmare, căutarea PDB BLAST este compusă din două execuții ale lui PSI BLAST: prima caută în baza de date NR sau alt stoc de secvențe apoi, cu acest profil, execuția finală caută în baza de date PDB (Protein Data Base). (http://bioinformatics.burnham-inst.org/pdb_blast/)

HMMer Server. Profilele modelelor Markov ascunse pot fi folosite pentru a face căutare senzitivă în bazele de date folosind descrieri statistice ale consensului²⁶ unei secvențe de familii. În prezent, HMM-ul are o implementare software a profilului pentru analiza secvențelor de proteine distribuită gratuit. (<http://bioweb.pasteur.fr/seqanal/motif/hmmer-uk.html>)
O reprezentare sugestivă a relațiilor dintre secvențele biologice și această metodă HMM este sugerată în figura următoare. Procedurile reprezentate și care relaționează cu bazele de secvențe sunt implementate la adresa oficială indicată.

Sistemul software **SAM** (Sequence Alignment and Modelling System), <http://www.soe.ucsc.edu/research/compbio/sam.html> este o colecție tooluri software flexibile pentru crearea, rafinarea și utilizarea modelelor Markov ascunse liniare în vederea analizei secvențelor biologice.

SAM Include programe și scripturi pentru metoda SAM-T2K de detectare a omologilor distanți. Acesta este o metodă HMM iterativă de căutare pentru crearea unui HMM dintr-o singură secvență de proteine sau aliniament sursă folosind căutarea iterativă a secvenței într-o bază de date. Metoda este considerată de autori cel mai sensibil algoritmul bazat pe determinarea de la distanță a omologiei unei secvențe.

²⁵ Frecvențele sunt calculate din aliniamente multiple de secvențe care conțin un domeniu de interes.

²⁶ În biologia moleculară și bioinformatică, o secvență consensus este o cale de reprezentare a rezultatelor unui aliniament multiplu, unde sunt comparate una cu cealaltă secvențe înrudite și sunt găsite motive funcționale similare. Secvența consensus arată ce reziduuri sunt conservate (sunt întotdeauna aceleași), și care sunt variabile [WIKI'05h].

Detalii care descriu pe larg strategia de funcționare a **SAM** pot fi consultate la adresa: http://www.cse.ucsc.edu/research/compbio/html_format_papers/hughkrogh96.

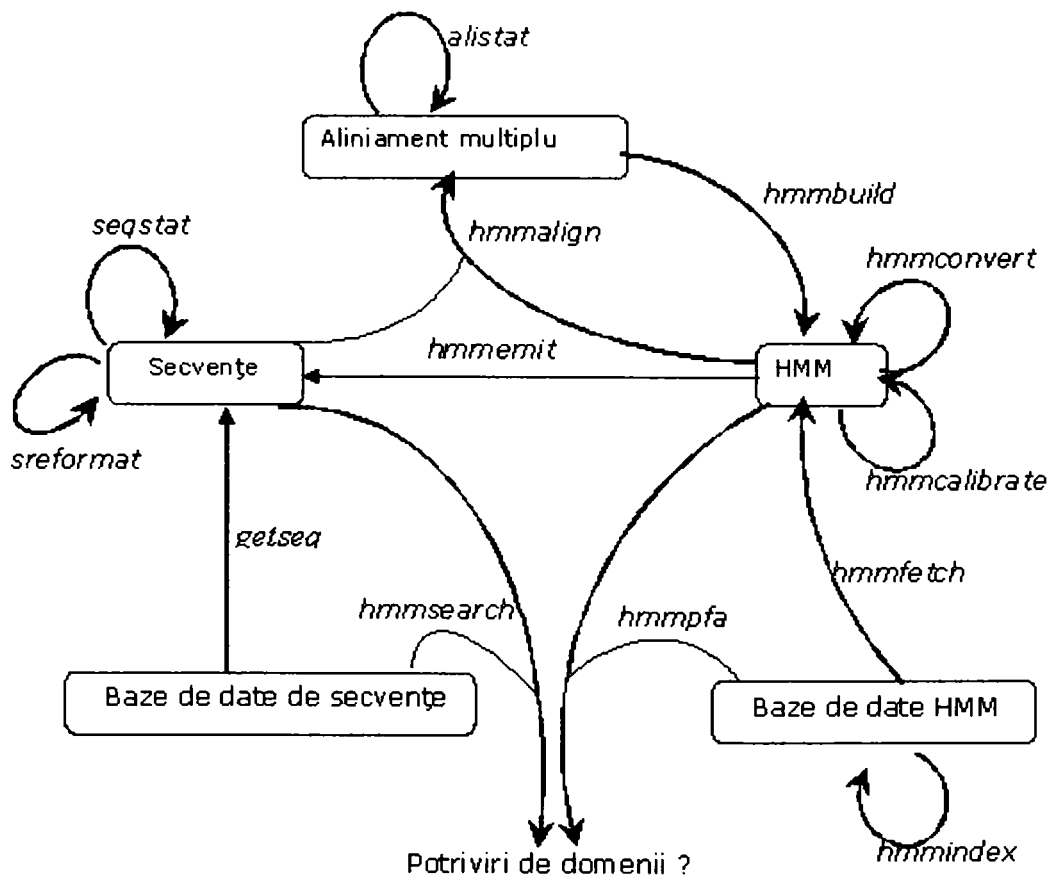


Figura 3.4.5-1. Reprezentarea relațiilor dintre secvențele biologice și metoda HMM, figură preluată și parțial tradusă din documentația detaliată, disponibilă la adresa de web a institutului Pasteur, Franța (<http://bioweb.pasteur.fr/seqanal/motif/hmmer-uk.html>).

Tcoffee este un software care realizează aliniament multiplu de secvențe, evaluarea aliniamentului multiplu, combinarea aliniamentelor multiple, alinierea de structuri (3DCoffee). După cum o descriu autorii metodei [NOTR'00], T-Coffee este o metodă pentru aliniamentul multiplu al secvențelor care furnizează o îmbunătățire

categorică a acurateții în comparație cu cele mai comune alternative utilizate. Metoda este bazată în general pe abordarea uzuală a aliniamentului multiplu progresiv dar evită capcanele datorate naturii "greedy" a acestui algoritm. Cu T-coffee se preprocesează un set de date al tuturor perechilor de aliniamente dintre secvențe. Acesta va furniza o bibliotecă de informație de aliniere care poate fi folosită pentru ghidarea aliniamentului progresiv. Prin urmare, aliniamentele intermediare nu se bazează doar pe următoarea secvență care trebuie aliniată dar și pe felul în care toate secvențele sunt aliniate una cu alta. Aceste informații despre aliniament pot fi derivate din surse heterogene cum ar fi o mixtură de programe de aliniament și/sau o superpoziție a structurii. Biblioteca poate fi construită folosind o combinație de aliniament global sau local. Ca detalii, unele teste confirmă faptul că acest program complex are o acuratețe mai mare decât Clustal W (care urmează să fie descris) pentru secvențe cu o identitate sub 30% dar, în schimb, este mai încet. Secvențele de intrare sunt în format FASTA (Pearson) iar numărul maxim de secvențe este 30 cu o lungime maximă de maximum 10 000 caractere.

http://igs-server.cnrs-mrs.fr/Tcoffee/tcoffee_cgi/index.cgi),

Clustal W este un program general pentru aliniament multiplu de secvențe de proteine sau ADN. Aliniamentul se aplică pentru secvențe cu diverse grade de divergență și are o semnificație biologică, calculând cele mai bune potriviri pentru secvențele selectate. El le aliniază astfel încât identitățile, similaritățile și diferențele pot fi văzute. Relațiile evoluționare pot fi văzute via cladograme²⁷ sau filograme²⁸ ce pot fi generate.

Metoda de bază a acestui aliniament constă din trei etape principale:

- (i) toate perechile de secvențe sunt aliniate separat cu scopul calculării matricii distanță rezultând divergența fiecărei perechi de secvențe;
- (ii) din matricea de distanță este calculat un arbore ghid;
- (iii) secvențele sunt aliniate progresiv în conformitate cu ordinul de ramificare din arborele ghid.

Programul oferă posibilitatea alegerii modului de aliniament dintre o metodă rapidă de aproximare [BASH'87] ce permite să fie aliniate un număr foarte mare de secvențe și o alta mai încetă dar cu un grad mai mare de precizie. În cazul primei metode, scorurile sunt calculate folosind numărul de k-tupluri de potriviri (înseamnă reziduuri identice, de obicei de lungime 1 sau 2 pentru proteine sau de lungime 2-4 pentru secvențe de nucleotide) în cel mai bun aliniament dintre două secvențe minus o penalitate fixă pentru fiecare gap. A doua metodă, după cum o afirmă autorii, oferă o mai mare acuratețe a scorurilor, fiind complet bazată pe programarea dinamică și folosește două penalități pentru gap-uri (pentru deschidere sau extindere) și o matrice de ponderi a amino acizilor. Scorurile în acest caz sunt calculate ca număr de identități în cel mai bun aliniament împărțit la numărul de reziduuri comparate (pozițiile de gap sunt excluse).

Scorurile în ambele metode sunt calculate inițial ca procent de scor identitate și sunt convertite în distanțe prin împărțirea cu 100 și apoi scăzând din

²⁷ Cladogramele sunt un arbori de diagrame folosiți în clasificarea biologică a membrilor din punct de vedere al evoluției pentru a ilustra relațiile filogenetice [WIKI'05g]

²⁸ Filogramele sunt arbori filogenetici care indică relațiile dintre diverse specii de plante sau animale grupate în funcție de relațiile dintre ele și de asemenea exprimă un sens al timpului sau ratei de evoluție. Aspectul temporal al filogramei lipsește dintr-o cladogramă [WIKI'05g]

1.0 pentru a furniza numărul de diferențe per site. În aceste distanțe inițiale nu se aplică nici o corecție pentru substituții multiple.

Pentru exemplificare, în Figura 3.4.5-2, preluată din materialul original [HIGG'94] sunt date 7 secvențe biologice din baza de date SwissProt, al căror aliniament este calculat folosind metoda dinamică a programului Clustal W.

Arborele neorientat afișează toate ramurile cu lungimile proiectate pentru scalare. În arborele orientat, toate lungimile ramificațiilor (semnifică numărul de diferențe per reziduu de-a lungul fiecărei ramuri) sunt date la fel ca și ponderile pentru fiecare secvență. În aliniamentul multiplu s-au folosit parametrii implicați din Clustal W și seria PAM(3) de matrice de ponderi.

Spre deosebire de alte tooluri de aliniament multiplu, sensibilitatea metodelor de aliniament multiplu obținute a fost semnificativ îmbunătățită pentru aliniamentul secvențelor de proteine divergente. După cum afirmă autorii lucrării [HIGG'94], Clustal W conține unele caracteristici cum ar fi :

- Ponderile individuale sunt asociate fiecărei secvențe într-un aliniament parțial cu scopul de a reduce ponderea secvențelor aproape duplicate și a o crește pe aceea aparținând celor mai divergente.
- Matricele de substituție ale amino acizilor sunt variate la diferite stadii de aliniament în conformitate cu divergența secvențelor care trebuie aliniate.
- Penalitățile specifice lipsei de reziduuri (residue-specific gap penalties) și penalitățile reduse locale ale gap-urilor în regiunile hidroflice încurajează noi gap-uri în potențialele regiuni de nod (loop) mai degrabă decât structura secundară regulată.
- Pozițiile în primele aliniamente unde gap-urile au fost deschise primesc penalități locale reduse pentru a încuraja deschiderea de noi gap-uri la aceste poziții. (<http://www.ebi.ac.uk/clustalw/>)

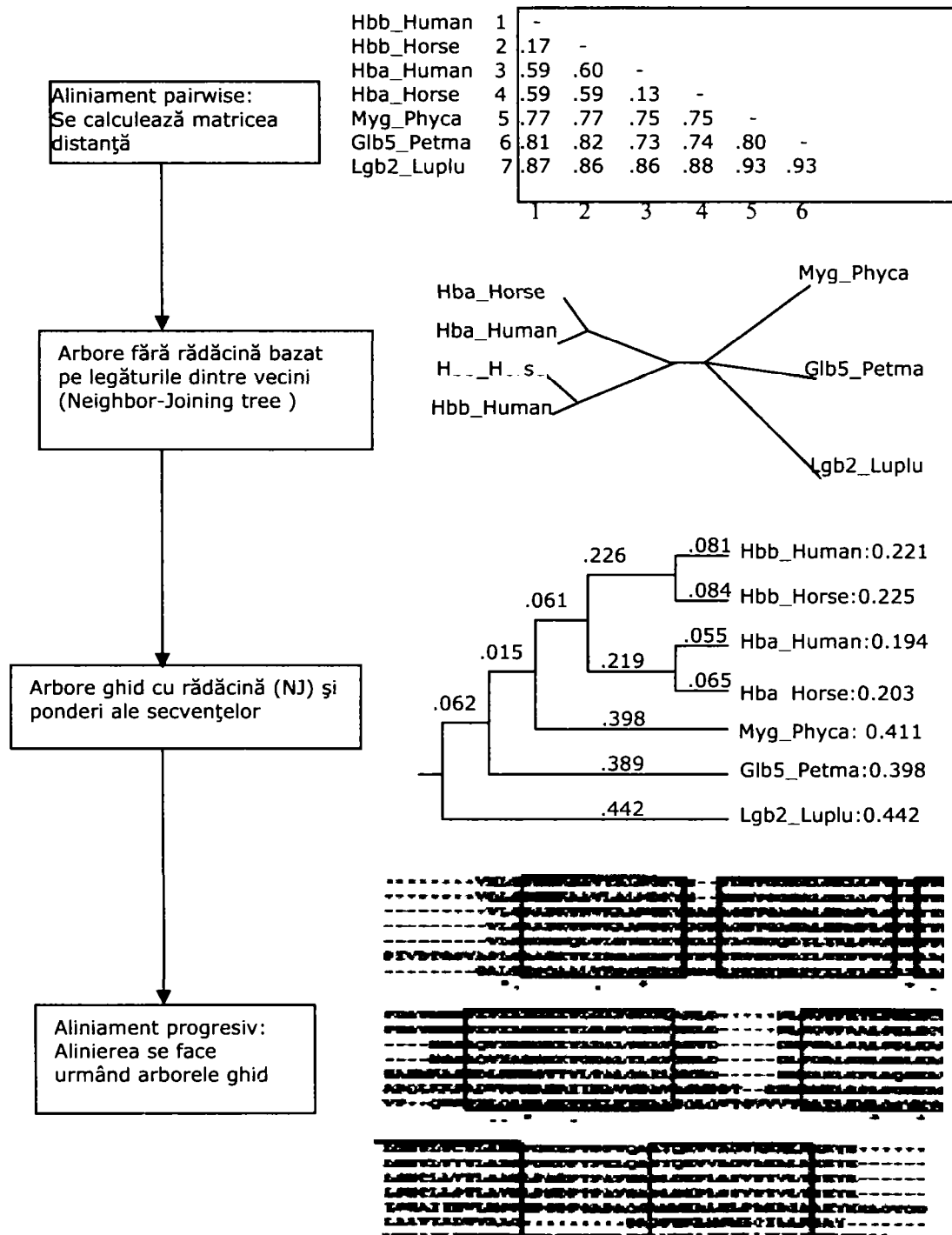


Figura 3.4.5-2: Rezultatul oferit de Clustal W în urma aplicării metodei dinamice pentru 7 secvențe de proteine din baza de secvențe SwissProt. Figura este reprodusa din lucrarea lui J.D.Thompson ©[THOM'94] care studiază îmbunătățirea metodei de aliniament multiplu de secvențe a lui Clustal W.

Software pentru recunoașterea plierilor și descompunerii secvențelor biologice în structurile componente:

Proteinele pot fi clasificate fie după omologia secvențelor (homology) fie după structura lor (folds=plieri). Ambele metode sunt pur informative. Schemele de clasificare în funcție de plieri încearcă să "descompună" esența formelor structurale complicate ale protenei în chei, determinanți comuni ai structurii lor (o codificare a structurii 3D). Aceste relații structurale abstracte pot revela în schimb relații evoluționare și/sau funcționale între proteine cu un grad mare de divergență. Prin urmare, descompunerea secvențelor (sequence **threading**) este o metodă de predicție a structurii care este înrudită cu analiza profilului. Fiecare pliere cunoscută a unei proteine este reprezentată într-o matrice abstractă, bidimensională, de distanțe între reziduuri. Câteva dintre produsele software existente sunt amintite în cele ce urmează.

123D. 123D-este un program care combină profile de secvențe, predicția structurii secundare și potențialele capacități de contact pentru a descrie o secvență de proteine prin setul de structuri. Hidrofobicitatea este forța majoră conducătoare pentru plierea proteinelor. Fără perechile de contact potențiale se poate folosi un algoritm rapid de programare dinamică pentru alinierea unei secvențe la o structură. Acest program apelează și la serviciile alor subprograme specializate cum ar fi Fold library, Compare protein structures, Domain assignment. (<http://123d.ncifcrf.gov/123D+.html>) Un exemplu este reprezentat în ANEXA 3.

FFAS (Fold & Function Assignment System) un program utilizat pentru analiza secvențelor biologice din punct de vedere structural. Metoda de asociere a plierilor se bazează pe un algoritm de potrivire profil-la-profil. Profilul secvenței interogatoare este generat cu PSI-BLAST folosind secvențe din baza de date NR. El este potrivit cu profilele ale secvențelor din baza de date grupată PDB. Partea crucială a algoritmului este noua, bidimensională schemă de ponderare care consideră topologia arborelui evoluționar al tuturor proteinelor în familia omoloagă [RYCH'00]. (<http://ffas.ljcrf.edu/ffas-cgi/cgi/ffas.pl>). Detalii suplimentare în ANEXA 3.

3D PSSM Fold Recognition Server, este o metodă web rapidă pentru recunoașterea plierilor proteinelor folosind profilele secvențelor reprezentate 1D sau 3D cuplate cu structura secundară și Informația Potențială de Solvatare (Solvation Potential Information). (<http://www.sbg.bio.ic.ac.uk/~3dpssm/index2.html>) O exemplificare este prezentată în ANEXA 3.

Cn3D, este o aplicație ajutătoare pentru căutarea web care permite vizualizarea structurilor 3D din baze de secvențe gestionate de NCBI Entrez. Această aplicație afișează simultan structura, secvența și aliniamentul iar în prezent are o adnotare performantă și caracteristici de editare ale aliniamentului. (<http://www.ncbi.nlm.nih.gov/Structure/CN3D/cn3d.shtml>). Se poate vedea un exemplu în ANEXA 3.

Metodele și exemplele de analiză și comparare a secvențelor biologice prezentate în acest ultim capitol nu și-au propus să epuizeze toată cantitatea de produse și metode existente la ora actuală ci să definească într-un mod cât mai sugestiv potențialul oferit de metodele computaționale aplicate în biologie. Astfel, toolurile prezentate fac parte din diverse categorii, în funcție de aplicabilitatea lor și tehnicile de prelucrare și extragere a informației din secvențele biologice.

3.6. Concluzii

Deoarece tehnicile de recunoaștere au o largă aplicabilitate în domeniul bioinformaticii, iar obiectivul principal al tezei este de a veni cu o nouă soluție la problemele de analiză și căutare în baze de secvențe biologice mari, pe parcursul acestui capitol s-a realizat o identificare a acestor tehnici în funcție de scopul la care servesc, precum și o investigare a produselor informatice existente. Astfel au fost identificate:

- metode de clasificare precum k-NN, clasificatori ce folosesc lanțuri sau modele Markov generalizate, clasificări pe baza identificării de caracteristici sau cu ajutorul rețelelor neuronale;
- algoritmi pentru căutare eficientă în baze de secvențe având ca nuclee diverse metode de similaritate și comparație a secvențelor, care au fost deja prezentate în capitolul anterior;
- metode folosite în căutarea și identificarea anumitor regiuni și tipare folosind căutarea exhaustivă, metode de învățare automată, metode bazate pe alinierea secvențelor, metode care se bazează pe conceptul de graf și metode hibride;
- produse informatice specializate care utilizează metodele de recunoaștere și algoritmi de analiză și comparație dezvoltate.

Până în momentul de față nu se cunoaște vreo evaluare a performanțelor comparative a mai multor metode care să decidă asupra superiorității totale a uneia sau alteia. Aceasta se explică fie prin:

1. tipul de format al datelor prelucrate, deoarece nu este același lucru a se obține secvențe similare prin realizarea de aliniamente de secvențe ca și prin interogarea lingvistică a unor descrieri textuale despre aceste secvențe;
2. scopul urmărit de utilizator. Scopul urmărit prin analiza și/sau căutarea de secvențe poate să fie de la unul simplu de a identifica obiectul de interes până la cel de a obține detalii suplimentare despre respectivul „obiect”.

Totuși, dacă se presupune că se pornește strict de la aceleași premise și se urmăresc aceleași obiective se poate concluziona că până la urmă vor fi alese, ca și în cazul metodelor de similaritate, metodele cu timpi computaționali cât mai reduși, consum minim de resurse și gradul de acuratețe al rezultatelor ridicat. Astfel, metodele de comparare bazate pe profile și lanțuri Markov ascunse și care folosesc ideea de familie, permit biologului să deducă/obțină apropape de trei ori mai multe omologii decât un simplu algoritm bazat pe compararea perechilor [LIAO'03].

Multe dintre cele mai puternice metode de analiză a secvențelor sunt acum bazate pe principiile modelării probabilistice. Exemple de asemenea metode includ:

- folosirea matricelor cu scoruri derivate probabilistic pentru a determina semnificația aliniamentelor secvențelor,

- folosirea modelelor Markov ascunse ca baze pentru căutări de profile în scopul identificării membrilor la distanță ai familiilor de secvențe,
- și inferența arborilor filogenetici folosind abordările probabilității maxime.

Caracteristic tuturor metodelor de comparație și analiză este însă gradul ridicat de dependență de estimări de diverși parametri, sau estimări biologice colectate din diverse surse. Astfel, dacă vorbim de metodele de aliniament este evidentă necesitatea implicării matricelor de substituție și a evaluărilor de penalități de gap sau inserții și ștergeri. În cazul în care explorăm folosind modelele Markov trebuie avută în vedere fie calcularea de diverse ponderi de tranziție de la o stare la alta fie extragerea de „profile” din cadrul familiilor de secvențe pentru a putea evalua gradul de apartenență al unei secvențe la o familie sau alta.

Scopul principal al acestui capitol a fost de a include cele mai importante aspecte în practica actuală privind analiza secvențelor biologice. Lista publicațiilor menționate poate fi completată fără îndoială. Totuși, cele mai specifice au fost întâmpinate. Datorită progresului continuu în cercetarea biologică, datele disponibile continuă să se extindă nu numai în termeni cantitativi dar de asemenea și în ce privește formatele structurilor de stocare. Este de presupus astfel faptul că tehnicile de căutare și analiză sunt adaptate acestui progres rapid rezultând în noi metodologii și funcții computaționale.

Desigur că dezvoltarea aplicațiilor software nu se limitează la utilizarea ADN/proteine sub formă de secvențe biologice ci ele pot fi produse și pentru o altă formă de reprezentare a genelor cum ar fi microvectorii (microarrays). Pentru acestea există de asemenea baze de date specializate cum ar fi ArrayExpress (locată la European Bioinformatics Institute <http://www.ebi.ac.uk/arrayexpress/>) și analiza lor presupune tehnici specializate pentru procesarea imaginilor.

Atât acest capitol cât și cel referitor la metodele de comparație a secvențelor biologice pot fi considerate ca un cadru informațional pentru a motiva și a atrage atenția asupra contribuției personale la dezvoltarea unei noi metode de comparație, analiză și căutare în baze de secvențe, ce va fi descrisă și validată în capitolele următoare.

4. ANALIZA ȘI EVALUAREA CONTEXTULUI LINGVISTIC CU AJUTORUL MODELELOR STATISTICE PENTRU SISTEMELE DE RECUNOAȘTERE

Acest capitol își dorește să prezinte la un nivel de percepție cât mai facil, câteva fundamente matematice și teoretice și aplicarea lor în modelarea lingvistică din punct de vedere statistic, ce a devenit așa după cum o arată și multiple exemple, un instrument puternic și eficient în procesarea și extragerea informației dintr-o cantitate mare de date .

O importanță majoră în multe direcții de cercetare o au lanțurile Markov astfel, ele își demonstrează utilitatea inclusiv în domeniul procesării lingvistice. În capitolul de față este exemplificată utilizarea lanțurilor Markov în modelarea limbajului natural, constituind și elementul de bază al abordării experimentale conținute de această teză. Evaluarea acestor modele se realizează cu măsuri ca entropia și perplexitatea. Conform "Teoriei Informației", există mai multe tipuri de entropie cunoscute și utilizate la ora actuală dar pentru scopul propus, atenția s-a limitat la formele de entropie și cross-entropie (sau entropie încrucișată).

Valorificarea noțiunilor teoretice este regăsită în aplicația experimentală, motivând utilizarea lor precum și deducerea de noi raționamente.

Deoarece partea statistică și a teoriei decizionale este cea care contribuie în mod esențial la obținerea de rezultate eficiente în rezolvarea problemelor atât în recunoaștere cât și în procesarea și explorarea informației utilizând modelele lingvistice, lucrarea de față i-a acorat acest spațiu.

Aplicația practică privind estimarea și evaluarea modelelor lingvistice permite observarea și analiza diverselor aspecte ale rezultatelor obținute în vederea unor posibile direcții de investigare concretizate ulterior.

4.1. Fundament teoretic

4.1.1. Lanțuri Markov lingvistice

În matematică, un lanț Markov (în timp discret), numit după matematicianul rus Andrei Andreievici Markov, este un proces stohastic în timp discret cu proprietatea Markov. Într-un asemenea proces trecutul este irelevant pentru previziunile viitorului cunoscând prezentul. Există de asemenea lanțuri Markov în timp continuu dar care nu fac parte din sfera de interes în acest moment.

Un lanț Markov este o secvență X_1, X_2, X_3, \dots de variabile aleatoare. Acoperirea acestor variabile, adică mulțimea valorilor posibile este numită spațiul stărilor, valoarea lui X_n fiind starea procesului la momentul n . Dacă distribuția probabilității condiționate a lui X_{n+1} în stările anterioare depinde doar de X_n , singur, atunci:

$$P(X_{n+1} = x | X_0 X_1 X_2, \dots, X_n) = P(X_{n+1} = x | X_n) \quad (4.1.1-1)$$

Unde x este una dintre stările procesului. Identitatea menționată este proprietatea Markov [MARK'71]. Andrei Markov a venit cu primele rezultate pentru asemenea procese în 1906.

În lucrarea de importanță majoră "A *Mathematical Theory of Communication*" [SHAN'48], Claude Shannon propune utilizarea lanțului Markov pentru a crea un model statistic din secvențe de litere dintr-un fragment de text în limba engleză. Shannon a aproximat structura statistică a unui fragment de text folosind un model matematic Markov. Astfel, un model Markov de ordinul 0 prezice că fiecare literă din alfabet apare cu o probabilitate fixă. Prin urmare, pentru o secvență "agggcagggcg", modelul Markov de ordin 0 prezice că litera "a" apare cu probabilitatea 2/13, "c" cu probabilitatea 3/13 și "g" cu probabilitatea 8/13. Secvența următoare este un exemplu tipic pentru acest model:

a g g c g a g g g a g c g g c a g g g . . .

Un model de ordinul 0 presupune că fiecare literă este aleasă independent. Aceasta nu coincide cu proprietățile statistice ale limbajului (limbii engleze în acest caz) deoarece există o mare corelație între literele successive dintr-un cuvânt sau propoziție. Spre exemplu, litera "h" și "r" cel mai adesea urmează literei "t" decât "c" sau "x". Se obține un model mult mai rafinat dacă se permite exprimarea probabilității fiecărei litere successive care să depindă de cea anterioară sau cele anterioare. Prin urmare, un model Markov de ordinal k prezice că fiecare literă apare cu o probabilitate fixă dar acea probabilitate poate depinde de k litere consecutive anterioare (k -gram). De exemplu, dacă textul are 100 de apariții ale lui "th", cu 60 de apariții ale lui "the", 25 apariții cu "thi", 10 apariții cu "tha" și 5 apariții cu "tho", atunci modelul Markov de ordinal 2 prezice litera următoare urmând modelul 2-gram "th" și este: "e" cu probabilitatea 3/5, "i" cu probabilitatea 1/4, "a" cu probabilitatea 1/10 și "o" cu probabilitatea 1/20. Reprezentarea stărilor și o simulare a generării de text pseudo-aleator poate fi consultată în ANEXA 4.

4.1.2. Construirea modelelor lingvistice statistice

În toate abordările modelelor lingvistice din diverse domenii de aplicabilitate definirea lor este aceeași. Prin urmare, un model lingvistic statistic este considerat o probabilitate de distribuție peste toate propozițiile sau alte unități lingvistice într-un limbaj [ROSE'00]. Acesta poate fi văzut de asemenea ca un model statistic pentru

generarea de text. În general, funcția modelării lingvistice răspunde la întrebarea: cât de probabil al i -lea cuvânt dintr-o propoziție s-ar produce cunoscând cele $i-1$ cuvinte precedente? În cele mai multe aplicații ale modelării lingvistice, cum ar fi recunoașterea vorbirii și explorarea informației, probabilitatea unei propoziții este descompusă într-un produs de probabilități n -gram.

Dacă S ar fi o secvență compusă din k cuvinte,

$$S = w_1, w_2, \dots, w_k.$$

Un model lingvistic n -gram consideră secvența de cuvinte S a fi un proces Markov cu probabilitatea

$$P_n(S) = \prod_{i=1}^k P(w_i | w_{i-1}, w_{i-2}, w_{i-3}, \dots, w_{i-n+1}), \quad (4.1.2-1)$$

unde n referă ordinul procesului Markov. Când $n=2$ spunem că avem un model lingvistic numit bigram ce este estimat folosind informații despre coexistența perechilor de cuvinte. În cazul în care $n=1$, spunem că avem unigram ce folosește numai estimări ale probabilităților cuvintelor individuale. Pentru aplicații ca recunoașterea vorbirii sau traducătoare automate ordinea cuvintelor este importantă și sunt folosite modelele de ordin mai înalt (de obicei trigram). În extragerea de informații (information retrieval) rolul ordinii cuvintelor este mai puțin clar și modelele unigram au fost folosite extensiv.

Pentru stabilirea modelului lingvistic al cuvintelor, estimatorii probabilistici sunt de obicei derivați din frecvența modelelor n -gram din setul de date de antrenare. De obicei multe elemente de tipul n -gram este posibil să nu apară în setul de date curent utilizat pentru estimare, chiar dacă mărimea setului este imensă și valoarea lui n este mică. Ca o consecință, pentru evenimentele rare sau neașteptate, estimările probabilităților care sunt direct bazate pe numărarea frecvențelor de apariție devin o problemă. Această situație este adesea calificată drept „inconsistența setului de date”. Ajustarea (smoothing) este folosită pentru rezolvarea problemei și constituie o parte importantă în orice model lingvistic. Dintre cei mai cunoscuți algoritmi sunt: „Add-one”, „Witten-Bell Discounting”, „Good-Turing Discounting”, „Backoff”, „Deleted Interpolation” majoritatea propuși pentru recunoașterea vorbirii ce sunt descriși în detaliu în [JURA'00]. Există însă și alte metode de smoothing utilizate în regăsirea informației (IR) și care sunt: *ajustarea parametrilor* cunoscută și ca „parameter smooting” sau *ajustare semantică* („semantic smoothing”) descriși în [LIU'04].

4.1.3. Măsurile de evaluare a modelelor lingvistice statistice

În practică, *entropia* și *perplexitatea* sunt cele mai cunoscute metrice folosite în evaluarea sistemelor bazate pe lanțuri Markov, cunoscute și ca n -grame [JURA'00] și pot fi exprimate în modele matematice într-o formă exactă.

În rezolvarea reală a problemelor este necesară însă estimarea probabilităților din observări ale proceselor aleatoare.

Vizând în mod direct domeniul lingvistic, din punct de vedere al eficienței, un alfabet sursă întâlnit în practică ar trebui să aibă o distribuție a probabilității mai mică decât cea optimă. Dacă alfabetul sursă este compus din n simboluri, atunci el poate fi comparat cu un "alfabet optimizat" cu n simboluri a căror distribuție e uniformă. Rata entropiei alfabetului sursă cu entropia versiunii sale optimizate reprezintă eficiența alfabetului sursă ce poate fi exprimată ca procent. Aceasta implică faptul că eficiența unui alfabet sursă cu n simboluri poate fi simplu definită ca fiind egală cu entropia sa n -ară [WIKI'05a]. Uneori evaluarea modelelor lingvistice este realizată și folosind o măsură numită „perplexitate”.

4.1.3.1. Entropia

Originea conceptului de entropie este legată de Ludyig Boltzmamm (1877) și i-a fost dată o interpretare probabilistică de Claude Shannon.

Adoptând descrierea făcută de Richard O. Duda în cartea sa *Pattern classification* [DUDA'01], din punct de vedere matematic, pentru o mulțime discretă de simboluri $X = \{v_1, v_2, \dots, v_m\}$ cu probabilitățile asociate P_i , entropia distribuției discrete – o măsură a caracterului aleator și impredictibil al unei secvențe de simboluri obținută din aceasta – este:

$$H(x) = - \sum_{i=1}^m P_i \log_2 P_i , \quad (4.1.3.1-1)$$

unde entropia este măsurată în *biți* atunci când folosim logaritmul în baza 2. Pentru distribuțiile continue se folosește adesea baza e sau logaritmul natural notat \ln , caz în care entropia este exprimată în "*nats*". În cazul în care o probabilitate dispăre se apelează la faptul că $\lim_{p \rightarrow 0} p \log p = 0$ pentru a defini $0 \log 0$.

Făcând referire la aplicabilitatea entropiei la modul general, cartea lui Dorian Pyle, *"Data Preparation for Data Mining"* [PYLE'99] vine cu exemple de rapoarte ale analizei entropice alături de discuții despre cum pot fi folosite aceste informații într-o prima fază de parcurgere a investigării datelor.

Utilizând calitatea de estimator al distanței informației, s-au dezvoltat algoritmi cu aplicabilitate în bioinformatică și domeniul lingvistic [KALT'04]. Cele mai cunoscute și aplicate forme ale entropiei sunt: entropia încrucișată, entropia de legatura, entropia condiționată și informația mutuală. În acest moment însă interesul va fi orientat către entropia încrucișată și doar la nivel de definiție spre perplexitate.

4.1.3.2. Entropia încrucișată

Entropia încrucișată ("Cross-Entropy") tradusă și "cross-entropie", este utilă când nu cunoaștem probabilitatea de distribuție care a generat datele actuale. Prin urmare, cross-entropia dintre o variabilă aleatoare X cu distribuția reală a probabilității $p(x)$ și o altă funcție model cu probabilitate $q(x)$ (de obicei un model al lui p) este dată de expresia:

$$PH(X, q) = -\sum_x p(x) \log q(x), \text{ unde } x \text{ parcurge } X \text{ [MANN'00]}. \quad (4.1.3.2.-1)$$

Întâlnită în aplicațiile ce vizează modelarea lingvistică, această entropie "încrucișată" presupune că termenul pentru care se aplică \log este dedus dintr-un set diferit de cel folosit pentru estimarea ponderilor pentru modelul curent, numit și set de antrenare și, conform [VANC'03] se definește în modul următor:

$$C(x) = -\sum_i p_{\text{test}}(x_i) \log_2 p_{\text{train}}(x_i) \\ C(x) \geq H(x) \quad (4.1.3.2-2)$$

↓

$C(x)$ este considerată o *limită superioară* pentru estimarea unui proces aleator așa după cum rezultă și din [JURA'00].

În practică, această măsură a devenit larg utilizată. În sistemele de recunoașterea vorbirii și în modelarea lingvistică ea este aplicată pentru explorări ale conținutului textual util în diverse domenii [JURA'00], [MANN'00], [YOUN'97], [DEMA'95].

Mai mult decât atât, după cum este descris în [YAO'04], entropia încrucișată este utilizată ca măsură a performanței unei noi metode de identificare a sesiunilor²⁹ în explorarea bazelor de date. Alegerea acestei măsuri este motivată de relația dintre predicție și compresie (a se vedea în [YAO'04]) iar avantajul utilizării ei în cazul modelelor n -gram pentru detectarea sesiunilor se datorează faptului că pentru a calcula cross-entropia nu este necesară cunoașterea limitelor reale ale sesiunii de căutare în datele de test. Această abordare are un impact semnificativ în domeniul bazelor de date și data mining. Metoda de identificare îmbunătățită a sesiunilor poate fi de asemenea aplicată în web logs³⁰ pentru a determina date mai bune și de o mai mare acuratețe pentru *web log mining*.

²⁹ o sesiune db este o secvență de cereri prezentată sistemului db în scopul obținerii unei anume operațiuni

³⁰ Un weblog (în prezent cunoscut și sub numele de blog) este o publicație web ce conține în mod obișnuit articole periodice. Cu toate că majoritatea weblog-urilor erau manual actualizate, în prezent există aplicații care să automatizeze întreținerea unor asemenea site-uri. Blog-urile pot fi de diverse forma, pornind de la jurnale individuale pînă la campanii politice, programe media și corporații. Ele pot fi scrise de un autor ocazional sau ca urmare a colaborarilor unei comunități de autori. Multe weblog-uri permit vizitatorilor să lase comentarii publice. Totalitatea weblog-urilor sau blog- site-uri înrudite este adesea numită blogsferă. Când o mare cantitate de activități, informații, opinii erup în jurul unui subiect particular sau controversat în blogsferă se consideră că are loc o furtuna_blog sau roi_blog (conform dicționarului wikipedia)

Pentru a întregi viziunea supra aplicabilității cross-entropiei, în [KROE'04] ea este propusă și ca alternativă la metodele clasice de grupare cum ar fi *k-medii*, *fuzzy k-medii* și *cuantificarea vectorială liniară*. Autorii susțin că algoritmul propus este rapid și sigur prezentând și avantajele simplității eliminând condițiile de start și erorile de eșantionare comparativ cu alte euristici care presupun simularea unor normalizatori sau căutări locale ghidate.

4.1.3.3. Perplexitatea

Perplexitatea este și ea o metrică dedusă din valoarea entropiei și este definită în modul următor:

$$PP(x) = 2^{H(x)}. \quad (4.1.3.3-1)$$

Pentru o variabilă aleatoare uniform distribuită cu N valori posibile, entropia și perplexitatea în acest caz sunt N .

Așa cum există mai multe moduri de exprimare a entropiei, perplexitatea este și ea denumită în funcție de entropia care o determină. Se poate întâlni astfel perplexitate condiționată [SHAN'51] care pentru un corpus suficient de mare este de obicei un indicator al cantității de informație furnizate de model. Cu cât este mai mică perplexitatea condiționată cu atât modelul conține mai multă informație; prin urmare un model mai bun. Aceasta se datorează faptului că modelul captează cât mai multă informație și orice incertitudine rămasă este reflectată în perplexitatea condiționată [ROSE'94]. Cea mai mare aplicabilitate o are în modelarea lingvistică pentru sistemele de recunoaștere a vorbirii, în măsurarea confuzibilității modelelor lingvistice. Ea depinde de mărimea vocabularului (care este indicator dacă nu există un model lingvistic relevant) și factorul de ramificare (în medie, numărul de cuvinte care pot urma după un cuvânt dat).

4.2. Implementarea unei aplicații pentru analiza și evaluarea modelelor lingvistice³¹.

Urmând etapele standard necesare în estimarea și evaluarea modelelor lingvistice, a fost dezvoltat un set de programe utilizate pentru explorarea potențialului lingvistic al modelelor statistice și evaluarea lor. Ca urmare, în acest capitol vor fi prezentate și analizate rezultatele obținute în fiecare stadiu experimental. Este necesar de menționat faptul că aceste rezultate exprimă gradul de predictibilitate al unui limbaj (limba engleză în acest caz) la scară redusă, pe motivul limitării corpusului experimental dar și vine cu observații noi privind posibilitatea investigării mult mai profunde a unor aspecte particulare cu aplicabilitate în explorarea informației (Information Retrieval).

³¹ Implementarea programelor și rezultatele experimentale au fost realizate în timpul unei burse de studiu și cercetare oferită de Universitatea Katolică din Leuven, Belgia, departamentul ESAT-PSI, Speech group.

4.2.1. Descrierea aplicației

Implementarea acestei aplicații presupune un set de programe independente din punct de vedere al implementării dar dependente din punct de vedere funcțional. Ele sunt scrise în limbajul Perl³² și oferă posibilitatea de a explora într-un mod flexibil comportamentul măsurilor destinate informației aflate sub formă de text. În continuare va fi dată o descriere a fiecăruia dintre aceste programe alături de detaliile necesare privind funcționalitatea lor. Ordinea de execuție este: procesarea textului, generarea modelului lingvistic, evaluarea entropiei simple și evaluarea entropiei după aplicarea unei metode de ajustare. Programele pot să fie executate atât pe platforma oferită de sistemele de operare Unix cât și Windows fiind rulate din linia de comandă.

4.2.2. Procesarea textului

Textele de obicei conțin nu numai cuvinte ci și semne de punctuație sau semne adiționale pe care dorim sau nu să le considerăm în timpul procesului de evaluare, în funcție de scopul dorit. În termeni tehnici, aceasta este considerată curățarea textului ("text cleaning").

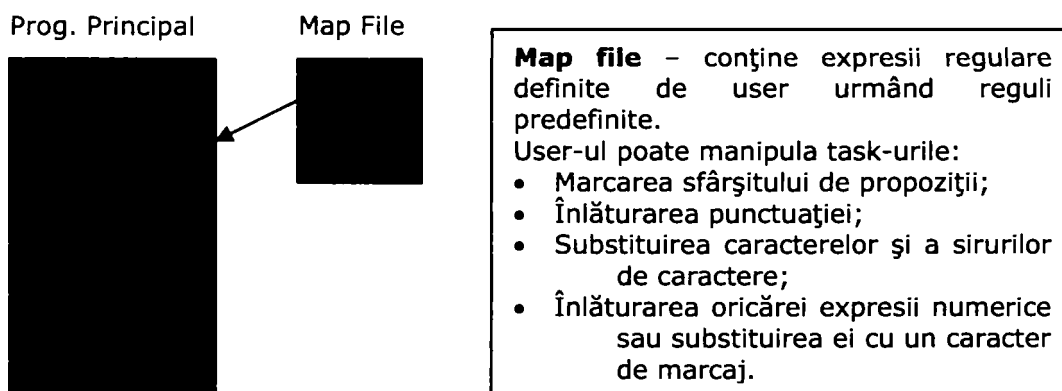
Normalizarea este un aspect semnificativ în dezvoltarea multor aplicații de procesare a textelor. Una din activitățile importante implică rezolvarea problemelor de ambiguitate privind literele scrise cu majuscule. Spre exemplu, în cazurile textelor mixte, cuvintele scrise cu majuscule denotă de obicei nume proprii dar există poziții speciale în text unde acestea sunt așteptate. Rezolvarea ambiguității acestor cuvinte conduce în general la identificarea numelor proprii cu toate că în multe cazuri numele proprii pot coincide cu simple cuvinte existente în vocabularul normal. Ex: "Pasare", "Negru"... În general, această problemă presupune o analiză mai amănunțită. O altă funcție a normalizării textului este determinarea limitelor propozițiilor care de asemenea poate conține situații ambigue presupunând la rândul lor o analiză mai aprofundată. În cele mai multe cazuri, împărțirea unui text în propoziții este o chestiune simplă: un punct, un semn al exclamării sau întrebării semnaleză sfârșitul de propoziție. Mai sunt însă și cazuri în care un punct semnifică un punct zecimal sau este o parte a unui abreviații. Prin urmare nu va semnala sfârșitul de propoziție dar pot exista situații în care o abreviere însăși poate fi ultima componentă a unei propoziții. Această problemă este de asemenea exacerbată de faptul că abrevierile nu formează un set închis, adică nu pot fi listate toate abrevierile posibile. Poate fi chiar problematic faptul că unele abrevieri coincid cu cuvinte obișnuite: "in", poate fi abreviere pentru "inch", "no" poate fi abreviere pentru "number", "bus" pentru "business" etc. Prin urmare, în timpul procesării textului utilizatorul trebuie să se gândească foarte bine ce fel de modificări generale dorește să execute. În acest sens, aplicația dezvoltată ajută utilizatorul să "curețe textul" după propriile criterii pe care le menționează într-un fișier standard. Programul folosește un fișier adițional care conține expresii regulate definite de utilizator sau asupra cărora acesta poate să intervină. Pentru experimentele curente, userul are posibilitatea de a manipula următoarele funcții:

³² Practical Extraction and Reporting Language, este un limbaj de programare folosit adesea pentru crearea de programe CGI (*Common Gateway Interface*)

- Marcarea limitelor propozițiilor considerând că fiecare propoziție începe la un rând nou și este finalizată la întâlnirea unuia sau mai multor semne de punctuație urmate de caracterul ce desemnează linie nouă;
- Înlăturarea semnelor de punctuație;
- Substituirea unor caractere sau șiruri de caractere;
- Înlăturarea oricărei expresii numerice sau substituirea ei cu caracterul dorit.

În cazul absenței acestui fișier nu va fi făcută nici o modificare asupra textului. Ieșirea generată de execuția acestui prim program o constituie un nou fișier text care se va constitui fișier de intrare pentru următorul program. Pentru detalii despre modul în care acest program funcționează poate fi utilizată opțiunea de help.

Reprezentarea procesului 1:



Programul principal – execută modificările stabilite de user în map file => un nou fișier text.

Dacă Map file lipsește – nu se execută vreo modificare asupra textului de intrare.

Experimente:

În această primă etapă este utilizată o colecție de 536 fișiere text din Wall Street Journal (WSJ) corpus, conținând 1527 linii, 304694 cuvinte și 1845026 caractere. Acest corpus este deja considerat normalizat prin urmare este corectă presupunerea conform căreia fiecare propoziție începe la o linie nouă.

Cantitatea totală de text a fost împărțită în două. Un set de 183 fișiere (5167 linii, 103970 cuvinte, 629966 caractere) și un set de antrenare de 352 fișiere (10109 linii, 200724 cuvinte, 1215060 caractere). Asupra fiecărui set a fost aplicată aceeași procesare:

- Punctuația considerată nesemnificativă în estimarea entropiei la nivel de literă a fost înlăturată. Prin urmare, caractere ca "- () [] : ; ' | " au fost înlăturate păstrând în schimb punctele din interiorul propozițiilor deoarece există ambiguități în decizia funcției lor sintactice. De asemenea, sunt păstrate marcasele apostrof întâlnite în forma condensată a cuvintelor cu toate că există situații în care sunt utilizate cu rol de marcaj de citat;
- Sunt păstrate toate cifrele însoțite sau nu de £ , \$, # ce pot simboliza tipuri de monede;

- Spațiile multiple din interiorul propozițiilor sunt substituite cu spații simple și pentru ușurința interpretării au fost înlocuite cu caracterul ^ ;
- Pentru a urmări limitele propozițiilor, ele sunt marcate.

4.2.3. Generatorul de modele lingvistice

Funcția clasică a modelării lingvistice este de a prezice cuvântul următor cunoscându-le pe cele anterioare. Urmând experimentul lui Shannon din 1951 de a ghici următorul caracter într-un text au fost dezvoltate multe aplicații. În mod asemănător cu explicațiile pentru modelele n -gram la nivel de cuvânt [MANN'00] și la nivel de litere avem aceeași abordare. Funcția de predicție a următorului caracter/literă poate fi considerată ca tentativă de estimare a funcției de probabilitate P :

$$P(X_n | X_1, \dots, X_n). \quad (4.2.3-1)$$

Într-o asemenea abordare probabilistică este folosită o clasificare pe baza caracterelor anterioare, *istoria*, pentru prezicerea caracterului următor. Având la bază examinarea unei mari cantități de text, putem prezice care literă tinde să urmeze alte litere. Prin urmare, este necesară o metodă de a grupa *istoriile* care sunt oarecum similare, cu scopul de a furniza predicții rezonabile referitoare la litera așteptată să urmeze. O cale posibilă de a le grupa o constituie presupunerea Markov conform căreia numai contextul local anterior –ultimele puține litere– afectează următoarea literă. Prin urmare, dacă este construit un model în care toate istoriile care au aceleași $(n-1)$ litere plasate în aceeași clasă de echivalență atunci avem un model Markov de ordinul $(n-1)$ sau un n -gram model de litere (ultima literă a n -gram-aticii fiind litera pe care o prezicem).

În aplicația dezvoltată, utilizatorul poate genera modele lingvistice pentru orice valoare dorită a lui n . Valoarea implicită este 1. Astfel, se poate obține mărimea curentă a vocabularului, valoare necesară în evaluarea entropiei. O funcție suplimentară o constituie și construirea unei histograme pentru a urmări frecvența aparițiilor unor anumite n -gram-atici.

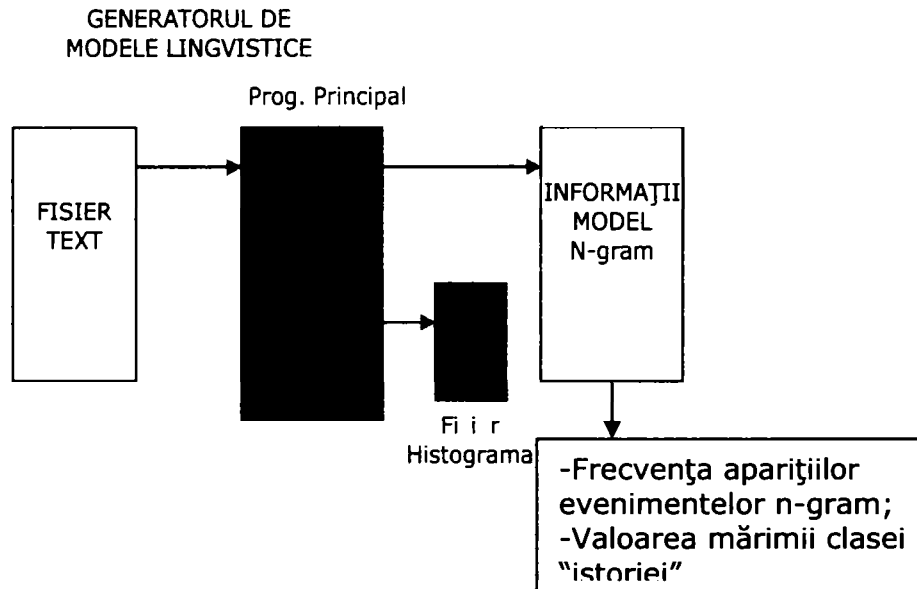
Exemplu:

Numărul evenimentelor n -gram ce au apărut de 5 ori = 14 (40.94 %)

Aceasta înseamnă că există 14 evenimente n -gram distincte care se repetă de 5 ori fiecare și ele reprezintă în jur de 41% din numărul total de unități n -gram identificate. Acest tip de informație poate fi utilizat pentru o analiză detaliată a conținutului textual.

Experimente:

Utilizând fișierele rezultate în etapele de procesare a textului au fost generate modele pentru 2,3,4,...20 –gram-e. Dimensiunea vocabularului utilizat este obținută din generarea 1-gram și are valoarea de 74 (poate fi văzut în tabelul din ANEXA 4).

Reprezentarea procesului 2:**4.2.4. Estimarea entropiei și analiza rezultatelor**

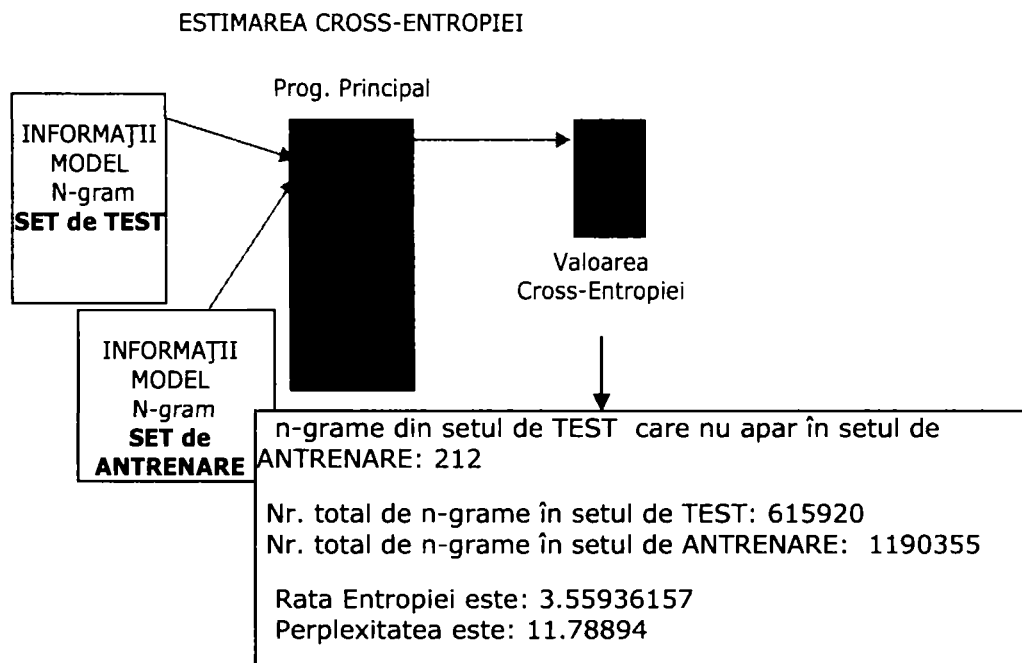
Pentru estimarea entropiei în cazul fiecărui model sunt utilizate fișiere rezultate din etapa anterioară astfel: pentru fiecare unic eveniment n -gram din setul de testare este evaluată ecuația 4.1.3-1, echivalentă în acest caz cu entropia încrucișată exprimată în (4.1.3.2-1). Prima probabilitate este obținută din frecvența relativă a evenimentelor n -gram din setul de testare și estimarea logaritmică este \log_2 din probabilitatea condiționată a aceluiași eveniment dar estimată din setul de antrenare. Și această etapă furnizează informații suplimentare privind cantitatea de evenimente implicate în estimarea entropiei, informație ce poate fi de asemenea utilizată pentru o analiză suplimentară și aprofundată a textului utilizat.

Experimente:

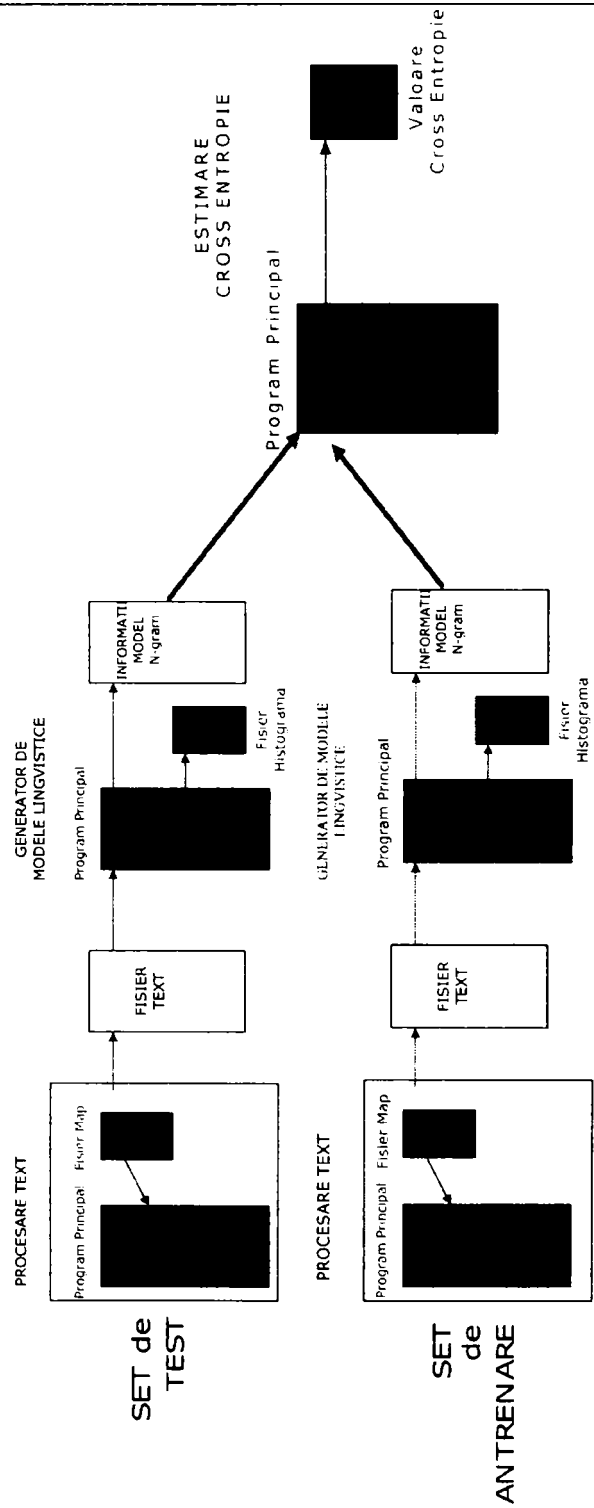
Pentru evaluarea simplă a entropiei ale cărei valori sunt reprezentate în Figura 4.2.4-1 este evident faptul: cu cât este mai mare numărul de caractere consecutive cunoscute cu atât este mai ușor de prezis următorul caracter. În cazul de față, simpla entropie atinge valoarea 0 după aproximativ 11 caractere consecutive. Întrebarea care se ridică acum este în ce măsură această estimare este bună sau nu. Aceasta implică o altă întrebare și anume: "când entropia atinge valori mici?". Răspunsul imediat ar fi: atunci când corpusul de antrenare nu acoperă sau acoperă doar parțial setul de testare și acest lucru poate fi observat analizând numărul de evenimente care nu se găsesc în setul de antrenare. Valorile pot fi observate în coloana a 4-a a

Tabel 4.2.4-1. Aceste evenimente n -gram participă la entropia finală cu valoarea 0. Un alt răspuns ar veni din situații când textele au subiecte tematice comune -setul de testare și de antrenare aparțin unui domeniu foarte specific- caz în care sunt mult mai puține variații în propoziții decât în propoziții la modul general (Ex: rapoartele radiologice).

Reprezentarea procesului 3



Descrierea procesului de execuție



Căutând răspuns la prima întrebare și anume: în ce măsură și dacă estimarea este bună sau nu, sunt realizate câteva investigații. Un indicator poate fi valoarea lui n din modelul n -gram. Acesta ar însemna că dacă n este un număr mare și valoarea entropiei tinde să fie foarte mică este posibil ca evaluarea să fie acceptată. Această afirmație este susținută de interpretarea generală a valorii entropiei/perplexității. Dar ca și măsură de comparație se propune valoarea entropiei obținută când setul de testare este acoperit în mod strict. Cu alte cuvinte, aceasta înseamnă valoarea entropiei când setul de testare și cel de antrenare coincid. Valorile obținute sunt prezentate în Tabel 4.2.4-2 iar reprezentarea grafică a acestor valori este realizată în Figura 4.2.4-2. Intenția este de a urmări cum se comportă valorile entropiei obținute utilizând un set de antrenare care să coincidă cu cel de testare în comparație cu situația în care cele două seturi sunt diferite. Ceea ce se așteaptă este faptul că valorile entropiei obținute în primul caz ar trebui să fie un reper pentru evaluarea entropiei unui model. După cum o arată graficul din Figura 4.2.4-2, simpla entropie se comportă mai bine decât în cazul strictei acoperiri. Diferența dintre cele două evoluții provine de la conținutul setului de antrenare și evenimentele n -gram purtătoare de valoare 0 în evaluarea entropiei. Până la nivelul 6 al modelului, cele două estimări evoluează aproape similar. După acesta însă diferența este vizibilă. O presupunere pentru secțiunile cu evoluție similară este faptul că evenimentele non-participante nu influențează într-un mod semnificativ evaluarea finală sau că distribuția evenimentelor n -gram în setul de antrenare este foarte similară cu acela din setul de testare. Pentru valorile de la 6 la 15 pentru model, diferența provine din valoarea rației diferite a probabilităților condiționate din cele două seturi. După nivelul 15 însă diferența este din nou nesemnificativă ceea ce sugerează faptul că majoritatea evenimentelor n -gram reprezintă expresii regulate care au aceeași structură în ambele seturi.

O altă întrebare ar fi în ce măsură influențează valoarea finală a entropiei evenimentele n -gram care nu participă la evaluarea ei. Prin urmare, atenția este acum direcționată spre frecvența apariției acestor evenimente. Spre exemplu, în Figura 4.2.4-3 sunt reprezentate evenimentele n -gram cu frecvența minimă în setul de testare.

După procesare, setul de testare conține 626254 unități (tokens) și setul de antrenare 1210573 (aproape dublu). Analizând reprezentarea grafică a frecvențelor minime se poate observa spre exemplu, că în setul de testare sunt 283 evenimente 2-gram care au loc o singură dată și reprezintă 0.05% din numărul total de evenimente 2-gram (mențiune: total nu înseamnă evenimente 2-gram unice.) În, pe coloana a patra, sunt prezentate numărul de evenimente care sunt în setul de testare și nu sunt găsite în cel de antrenare. Pentru 2-gram numărul de evenimente este de 212. Cei doi indicatori, frecvența minimă și evenimentele absente, sunt reprezentate în

Figura 4.2.4-4 și se poate observa evoluția lor similară. Oferă acestastă evoluție a indicatorilor o măsură a termenilor particulari sau expresiilor unice utilizate în text?

Într-adevar, ar putea oferi o referință asupra conținutului setului de testare atâta timp cât indică numărul de evenimente n -gram mai puțin frecvente în setul de testare alături de cel al evenimentelor n -gram inexistente în setul de antrenare. Să fie o simplă coincidență este prea puțin posibil dar experimente viitoare vor ajuta la clarificarea acestei observații.

n-grams	Nr. of n-grams in TEST set	Nr. of n-grams in TRAIN set	Nr. of events which do not occur in TRAIN set	Entropy	Perplexity
2	615920	1190355	212	3.55936	11.78894
3	610753	1180246	2153	2.91602	7.54766
4	605586	1170138	10397	2.21442	4.64096
5	600419	1160031	31379	1.63220	3.09987
6	595252	1149933	69916	1.19669	2.29214
7	590086	1139854	123185	0.85135	1.80419
8	584921	1129788	184062	0.58090	1.49579
9	579763	1119730	245694	0.38028	1.30160
12	564326	1089635	388032	0.09469	1.06784
15	548928	1059705	454837	0.02116	1.01478
18	533583	1029945	475142	0.00574	1.00399
20	523397	1010234	476943	0.00306	1.00212

Tabel 4.2.4-1. Rezultatele obținute aplicând procedura de obținere a entropiei simple

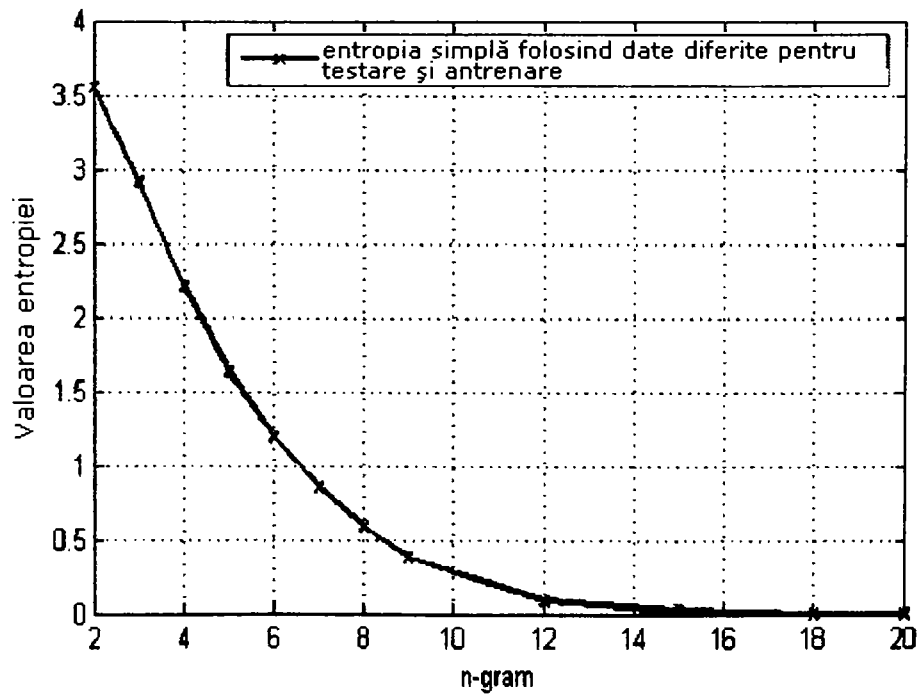


Figura 4.2.4-1. Reprezentarea valorilor pentru estimarea entropiei simple

n-grams	Entropy	Perplexity
2	3.54644	11.68382
3	2.87977	7.36034
4	2.16972	4.49936
5	1.61828	3.07010
6	1.22700	2.34081
7	0.92361	1.89686
8	0.67934	1.60142
9	0.48557	1.40014
12	0.16516	1.12129
15	0.05409	1.03821
18	0.02068	1.01444
20	0.01175	1.00818

Tabel 4.2.4-2. Evaluarea entropiei pentru cazul în care setul de testare și cel de antrenare coincid (overlapping data)

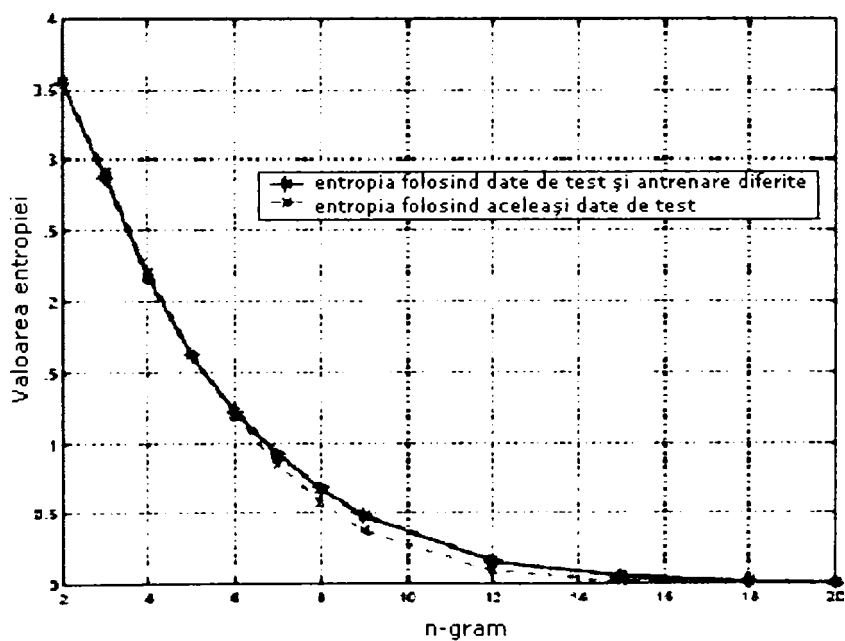


Figura 4.2.4-2. Reprezentarea evoluției valorilor entropiei când setul de testare coincide cu setul de antrenare.

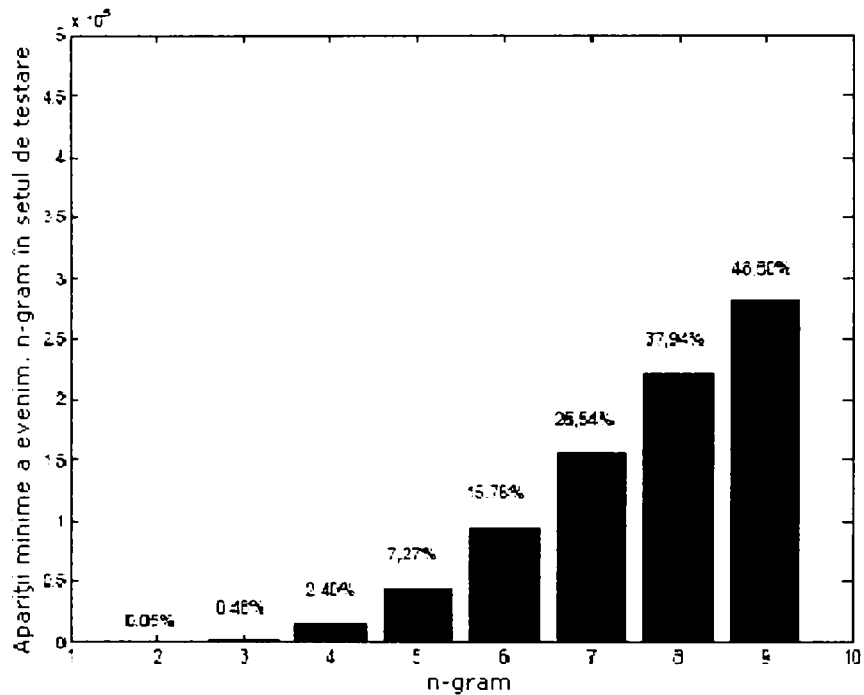


Figura 4.2.4-3. Reprezentarea evenimentelor n-gram cu frecvență minimă în setul de testare

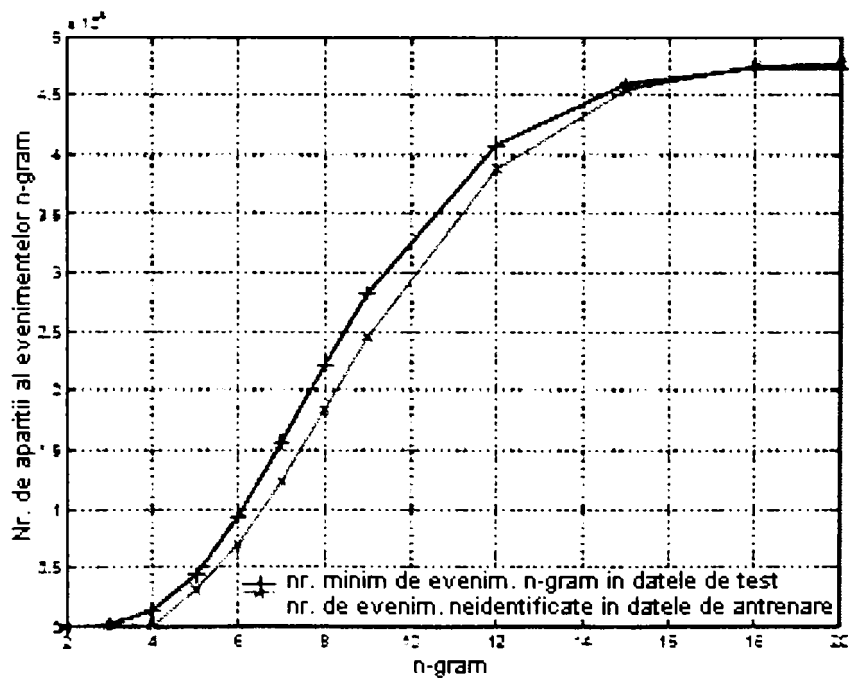


Figura 4.2.4-4. Reprezentarea evoluției frecvenței minime și a evenimentelor n-gram unice identificate în setul de testare

Experimentele de investigare au continuat și cu aplicarea unei metode de smoothing în estimarea entropiei iar rezultatele au fost prezentate în [BOGA'03] dar pentru scopul propus în această teză ele sunt mai puțin relevante. Totuși, efectul aplicării metodei va fi menționat în secțiunea de concluzii a acestui capitol.

4.3. Concluzii

La finalul acestui capitol se poate concluziona cu faptul că modelele lingvistice statistice se pot identifica drept un domeniu de importanță majoră în analiza și exploatarea unei cantități mari de informație fiind totodată ușor de manipulat și cu rezultate eficiente în aplicarea lor [BOGA'03c], [BOGA'04a].

Aplicația descrisă a fost dezvoltată cu scopul de a oferi posibilitatea investigării într-un mod flexibil a măsurilor aplicate în evaluarea textelor pentru modelarea lingvistică. În ceea ce privește rezultatele obținute vor fi punctate aspectele cele mai importante:

- După cum se menționează și în [GIVE'96], acoperirea este un aspect foarte important în recunoașterea bazată pe modele lingvistice.
- Există întotdeauna o negociere între fiabilitate/robustețe și detalii [BOGA'04b]. Aceasta este total dependentă de mărimea setului de antrenare. Cu cât este mai mare setul de antrenare, cu atât se vor obține rezultate acceptabile pentru valorile mai mari ale lui n . Un corpus de antrenare redus va impune o alegere a valorii lui n mai mică pentru a fi siguri că un asemenea model poate să existe. În cazul în care gradul de fiabilitate devine o problemă se poate folosi o tehnică de smoothing.
- Metoda de smoothing care nu a mai fost descrisă în acest capitol dar a fost testată în [BOGA'03b] și a dus la concluzia că "add-one" nu este o alegere bună în evaluarea entropiei modelelor lingvistice statistice pentru sisteme de recunoaștere. Această afirmație este susținută și în [JURA'00], [GOOD'98]. Alegerea metodei de smoothing optime este dificilă întrucât există o mare diversitate de metode al căror studiu presupune o analiză amplă. Așa după cum o afirmă și autorii în [CHEN'98], există o lipsă evidentă a studiilor care să compare în mod sistematic performanța a mai mult de câteva dintre aceste tehnici pe seturi multiple de date. Prin urmare, din literatura de specialitate, este imposibil de a măsura performanța relativă a algoritmilor existenți decât în puține cazuri particulare. Se poate spune că fiecare algoritm poate funcționa bine în anumite situații. În acest sens este relevantă lucrarea [CHEN'98].

Evaluarea este în mod argumentat cea mai importantă parte a oricărui proiect de cercetare [GIVE'96]. Fără metode adecvate și măsuri/metrici general acceptate este dificil de evidențiat vre-un progres. O măsură a acurateții, care nu este una directă, este perplexitatea. În această evaluare însă au fost utilizate valorile entropiei ca măsură a indeciziei și predictibilității deoarece perplexitatea este

derivată automat din aceasta (a se vedea explicațiile din definiția dată în (4.1.3.3-1). Pentru evaluarea performanței unei model lingvistic ea trebuie realizată asupra acelorași date de antrenare deoarece sunt mulți factori care o pot influența. Cu alte cuvinte, entropia/perplexitatea luate în afara contextului nu mai au nici o semnificație [GIVE'96]. Investigațiile făcute asupra acestui set de antrenare au avut ca scop observarea modului diferit în care se manifestă modelele lingvistice și de a căuta estimatori care să poată oferi repere în evaluarea modelelor lingvistice.

Ca și subiect principal pentru o viitoare investigație este considerat căutarea modului în care poate fi acoperit un text particular, ce fel de parametri sunt implicați în evaluarea textelor particulare, cum se comportă modelele lingvistice în textele particulare și care este acuratețea măsurilor estimate?

5. CONTRIBUȚII LA UTILIZAREA TEHNICILOR DE RECUNOAȘTERE ÎN ANALIZA SECVENȚELOR BIOLOGICE

5.1. Descrierea instanței de rezolvat

Proteomica se referă la studiul unei colecții celulare de proteine (în același mod în care genomica se referă la studiul genelor) și găsește o largă aplicație în bioinformatica actuală. Întrebările tipice care se aplică aproape tuturor genelor sunt următoarele: ce proteină produce fiecare genă, când această proteină este produsă, și care este rolul său funcțional? În timp ce secvența genomică ne poate informa despre ce proteine are potențial celula de a le produce, și analiza expresiilor de microvectori poate oferi un răspuns aproximativ despre ce proteine sunt produse, doar abordarea proteomică poate oferi o imagine concretă a biochimiei fundamentale a unei celule. Una dintre abordările proteomice cele mai importante o reprezintă comparațiile dintre secvențele de proteine. Necesitatea pentru asemenea comparații provine din interesul în detectarea omologiilor dintre proteine care pot, la rândul lor, să implice similarități structurale și funcționale. Așa după cum a mai fost menționat în capitolul introductiv, proteinele sunt molecule largi, complexe, compuse din amino acizi și compararea și gruparea lor în conformitate cu gradul de similaritate presupune algoritmi specializați.

Cele mai frecvent folosite metode au fost deja descrise în capitolele 2 și 3 și în general sunt bazate pe proceduri algoritmice complexe pentru aliniament de secvențe. Desigur că dezvoltarea metodelor de similaritate a cunoscut o dinamică remarcabilă dar, în ciuda maturității metodologiilor dezvoltate în această direcție, derivarea de noi măsuri de similaritate a proteinelor este încă un domeniu de cercetare activ. Interesul este reînnoit, datorită creșterii continue a volumului bazelor de secvențe disponibile care reclamă proceduri algoritmice alternative eficiente din punct de vedere al costurilor și care pot cuantifica în mod corect similaritatea secvențelor fără a depinde de un anumit tip de aliniament. Separat de eficiență, o a doua specificare de importanță egală pentru stabilirea măsurilor de similaritate este evitarea parametrilor care trebuie setați de utilizator (o caracteristică inerentă în majoritatea metodologiilor descrise până acum). Acesta este de obicei cazul abordărilor similarității în mod clasic, în care utilizatorul întâmpină o mulțime de dificultăți în alegerea unui algoritm de căutare adecvat, matrice de calcul sau funcții, precum și setul de parametri opționali ai cărori valori optime să corespundă celei mai bune similarități. O varietate de noi metode alternative pentru exprimarea similarității dintre secvențele biologice au devenit disponibile în diverse aplicații așa după cum au fost evidențiate în secțiunea 2.2.

În capitolul de față este introdusă o nouă abordare pentru măsura similarității dintre două secvențe de proteine. Ea a fost inspirată de aplicarea cu succes a conceptului de entropie pentru analiza informației în domeniul modelării statistice a limbajului (Young and Bloothoof [YOUN'97], Manning and Schütze

[MANN'00], Jurafsky and Martin[JURA'00]). Mai concret, modelarea n -gram este aplicată fiecărei secvențe de proteine iar o nouă măsură derivată din conceptul de cross-entropie este apoi folosită pentru compararea perechilor de secvențe. Cu toate că acest concept de n -gram a fost regăsit în lucrări anterioare, e.g. [GANA'04], [ERHA'80], [KARL'91], [KARL'96], munca actuală este de fapt rezultatul primelor încercări de a adopta acest pas dual pentru compararea secvențelor biologice.

Pornind de la fundamentele teoretice necesare evaluării textului, descrise și experimentate în capitolul anterior, care au relevat potențialul măsurilor probabilistice de evaluare a contextului, se va trece spre o nouă propunere de abordare a similarității secvențelor biologice. Metoda se aplică pentru secvențe de proteine descrise de structura lor primară.

5.2. Propunerea unei noi metode de comparare a proteinelor pe baza evaluării modelelor lingvistice Markov

5.2.1. Noțiuni teoretice utilizate

Există diverse tipuri de modele lingvistice care pot fi folosite pentru a captura diferite aspecte ale regularităților limbajului natural [WANG'03]. Lanțurile Markov sunt în general considerate printre cele mai fundamentale concepte pentru construirea modelelor lingvistice. Folosind explicațiile din capitolul 4 cu privire la definirea și evaluarea modelelor lingvistice Markov aici se va face o particularizare a lor pentru a susține modul în care sunt folosite.

După cum este descris în [MANN'00] și [BROW'92], entropia unei variabile aleatoare X care

ia valori într-un domeniu κ , și are o probabilitate a funcției densitate, $P(X)$ este definită ca:

$$H(X) = - \sum_{X \in \kappa} P(X) \log P(X). \quad (5.2.1-1)$$

Relativ recent, în lucrarea lui Van Uytsel și Comparnolle [VANU'98], a fost adoptată ideea generală a entropiei în cazul particular în care o secvență scrisă $W = \{w_1, w_2, \dots, w_{k-1}, w_k, w_{k+1}, \dots\}$ este tratată ca bazată pe compoziția modelului lingvistic L , rezultând următoarea formulă de estimare:

$$\hat{H}_L(X) = - \frac{1}{N} \sum_{W^*} \text{Count}(w_i^n) \log_2 p_L(w_{i+n} | w_i^{n-1}), \quad (5.2.1-2)$$

unde variabila X este de forma unei n -gram $X = w_i^n \Leftrightarrow \{w_i, w_{i+1}, \dots, w_{i+n-1}\}$ și $\text{Count}(w_i^n)$ este numărul aparițiilor lui w_i^n . Suma se execută asupra tuturor combinațiilor consecutive posibile w_i de lungime n (i.e

$W^* = \{\{w_1, w_2, \dots, w_n\}, \{w_2, w_3, \dots, w_{n+1}\}, \dots\}$ și N este numărul total de n -grame din secvența investigată. Al doilea termen din suma ecuației 5.2-2, adică $\log_2 p_L(w_{i+n} | w_i^{n-1})$, este \log din probabilitatea condiționată care conectează al n -lea element dintr-o n -gram cu cele $n-1$ elemente precedente. Urmând principiul estimării probabilității maxime (MLE), el poate fi estimat folosind frecvențele relative:

$$\hat{p}(w_{i+n} | w_i^{n-1}) = \frac{\text{Count}(w_{i+n})}{\text{Count}(w_i^{n-1})} \quad (5.2.1-3)$$

Echivalent cu ecuația 5.2-2 se poate folosi expresia

$$\hat{H}_L(X) = - \sum_{W^*} p(w_{i+n}) \log_2 p_L(w_{i+n} | w_i^{n-1}) \quad (5.2.1-4)$$

unde

$$\hat{p}(w_{i+n}) = \frac{\text{Count}(w_{i+n})}{N} \quad (5.2.1-5)$$

considerată ca frecvența relativă a evenimentului n -gram w_i^n .

Cross-entropia dintre probabilitatea de distribuție actuală $P(X)$ (care parcurge o variabilă aleatoare X) și probabilitatea distribuției $Q(X)$ estimată dintr-un model este definită astfel:

$$H(X, Q) = - \sum_{X \in \mathcal{X}} P(X) \log Q(X). \quad (5.2.1-5)$$

Aici trebuie menționate două remarci importante. Prima, cross-entropia unui proces stocastic, măsurată prin folosirea unui model, este o limită superioară pentru entropia procesului (i.e. $H(X) \leq H(X, Q)$) [MANN'00], [BROW'92]). A doua, așa după cum este menționată în [JURA'00], între două modele, cel mai bun este cel cu valoarea cross-entropiei minimă.

Estimarea entropică de mai sus (luată împreună cu forma generală din ecuația 5.2.1-1 și 5.2.1-2) sugerând o cale directă de trecere de la entropie la formularea cross-entropiei) a fost baza în construirea noii măsuri de similaritate a proteinelor, ce va fi descrisă în continuare.

5.2.2. Descrierea metodei

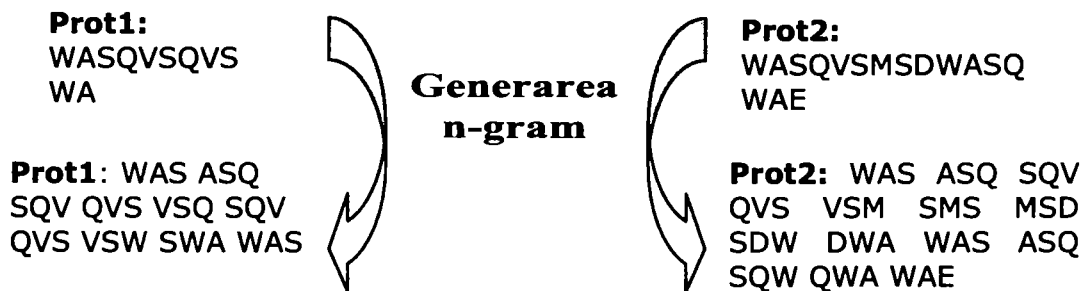
Așa după cum s-a mai menționat la începutul acestui capitol, secvențele de proteine provenind de la diverse organisme pot fi abordate ca text scris într-un limbaj universal în care alfabetul constă din 20 simboluri distincte, amino acizii. Asocierea unei secvențe de proteine structurii sale, dinamica funcțională și rolul său biologic devin astfel analoge cu asocierea dintre cuvinte și înțelesul lor semantic din limbajul natural. S-a sugerat că această analogie poate fi exploatată prin aplicarea modelării lingvistice statistice și a tehnicilor de clasificare a textelor pentru progresul « înțelegerii » secvențelor biologice. Oamenii de știință din acest domeniu hibrid de

cercetare consideră că identificarea regulilor de gramatică/sintaxă pot releva sistematici de o importanță ridicată pentru științele biologice și medicale. Pentru exemplificarea noii abordări alegem o secvență de proteine ipotetică WASQVSENR. În modelarea 2-gram rezultată avem următoarele elemente (tokens/cuvinte) {WA AS SQ QV VS SE EN NR}, în timp ce în reprezentarea 3-gram token-ii identificați sunt {WAS ASQ SQV QVS VSE SEN ENR}. Pe baza frecvenței de apariție a acestor « cuvinte » (determinată prin numărare) și prin formarea unor rapoarte de frecvență corespunzătoare, poate fi estimată cu ușurință entropia unui model n -gram folosind (5.2-2). Această măsură este indicativă pentru a determina cât de bine o secvență specifică de proteine este reprezentată de către modelul n -gram corespunzător. Dacă această măsură ar putea fi aplicată pe două proteine distincte (și ne ajută să decidem care proteină este cel mai bine reprezentată de model), ieșirile nu pot fi folosite pentru compararea directă a celor două proteine. Acest neajuns a direcționat spre derivarea cross-entropiei corespunzătoare, în care modelul n -gram este mai întâi construit pe baza calculării « cuvintelor » unei secvențe de proteine Y și apoi folosit pentru secvența X , contrastând cele două secvențe. Astfel, conținutul de informație comună dintre două proteine X, Y este exprimat via formula :

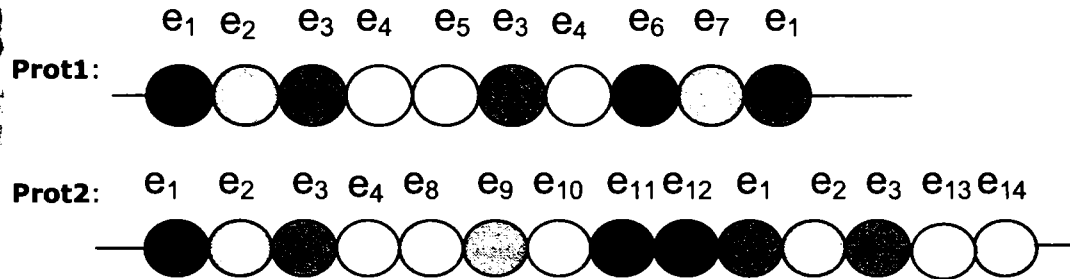
$$E(X, Y) = - \sum_{\text{all } w_i^n} P_X(w_i^n) \log P_Y(w_{i+n} | w_i^{n-1}). \quad (5.2.2-1)$$

Primul termen $P_X(w_i^n)$ din (5.3-1) corespunde secvenței de proteine X (și rezultă din calcularea frecvenței cuvintelor acestei secvențe specifice). Al doilea termen se referă la secvența Y pe baza căreia modelul trebuie estimat (de asemenea rezultatele sale vor veni din calcularea frecvenței de apariție acuvintelor componente). Variabila w_i^n va parcurge toate cuvintele (reprezentate de n -game) ale secvenței de proteine considerată referință. Pentru o înțelegere mai facilă se va exemplifica folosind două secvențe aleatoare de proteine.

Procesarea secvențelor de proteine.



Pentru fiecare secvență descompusă în "cuvinte" n -gram, numite și evenimente, se reprezintă figurativ succesiunea acestora. Astfel, fiecărui eveniment unic îi corespunde o culoare.



Folosind informațiile organizate sub formă de evenimente, pentru fiecare secvență este păstrat un tablou de date care conține informațiile privitoare la frecvența de apariție a fiecărui eveniment sub următoarea formă:

Prot1:	Eveniment	Frecvențe relative din (1)	Probabilități condiționate din (2)
WAS	e1	2/10	2/2
ASQ	e2	1/10	1/1
SQV	e3	2/10	2/2
QVS	e4	2/10	2/2
VSQ	e5	1/10	1/2
VSW	e6	1/10	1/2
SWA	e7	1/10	1/1

Prot2:	Eveniment	Frecvențe relative din (1)	Probabilități condiționate din (2)
WAS	e1	2/14	2/3
ASQ	e2	2/14	2/2
SQV	e3	2/14	2/2
QVS	e4	1/14	2/2
VSM	e8	1/14	1/1
SMS	e9	1/14	1/1
MSD	e10	1/14	1/1
SDW	e11	1/14	1/1
DWA	e12	1/14	1/1
QWA	e13	1/14	1/2
WAE	e14	1/14	1/3

Odată introdusă noua măsură de similaritate, se va trece la descrierea modului său de folosire pentru a efectua căutări în baza de date de proteine. Punctul esențial al acestei abordări constă în faptul că atât proteina necunoscută (ex. o nouă proteină descoperită) aflată în postura de query cât și fiecare proteină din baza de secvențe sunt reprezentate via codificare n -gram și similaritatea introdusă anterior este folosită pentru a le compara sub această formă de reprezentare (descompuse în evenimente n -gram).

5.2.3. Descrierea modul de aplicare al noii măsuri de similaritate

Au fost identificate două moduri în care similaritatea bazată pe n -grame este angajată în căutarea eficientă într-o bază de date de secvențe [BOGA'05a]. Cea mai directă implementare este denumită în mod justificat: "**metoda directă**". Un al doilea algoritm, "**metoda alternantă**", a fost identificat în scopul medierii cazurilor în care proteinele care se compară pot fi de lungimi foarte diferite. Este ușor de observat implicarea acestui aspect în valoarea raportului dintre numărul de cuvinte din secvența de referință implicată în calcularea valorii similarității și numărul total de cuvinte din respectiva secvență (implicat în prima probabilitate din ecuația 5.3-1). Modul de experimentare cu ambele metode și de a compara performanța lor oferă posibilitatea de a verifica sensibilitatea măsurii propuse privind lungimea secvențelor.

Metoda directă: Fie S_q secvența de proteină query și $\{S\}=\{S_1, S_2, \dots, S_N\}$ baza de date de proteine. Primul pas este calcularea scorului 'perfect' (PS) sau scor 'referință' pentru proteina query. Aceasta este realizată prin calcularea $E(S_q, S_q)$ folosind secvența de query atât ca referință cât și ca model în ecuația 5.2.2-1. În al doilea pas, fiecare proteină S_i , $i=1\dots N$, din baza de date servește ca secvență model în calcularea unui scor de similaritate $E(S_q, S_i)$, folosind aceeași ecuație 5.3-1, cu proteina query servind ca secvență de referință. În acest mod, sunt calculate N similarități $E(S_q, S_i)$, $i=1, \dots, N$. În final, aceste similarități sunt comparate cu scorul perfect PS. Calculând diferențele absolute $D(S_q, S_i)=|E(S_q, S_i)-PS|$, sunt exprimate 'discrepanțele' în sensul conținutului informației dintre proteina query și proteinele din baza de date. Ordonând aceste N diferențe putem ușor identifica cele mai similare proteine cu proteina query ca fiind acelea cu cea mai mică valoare a diferenței $D(S_q, S_i)$.

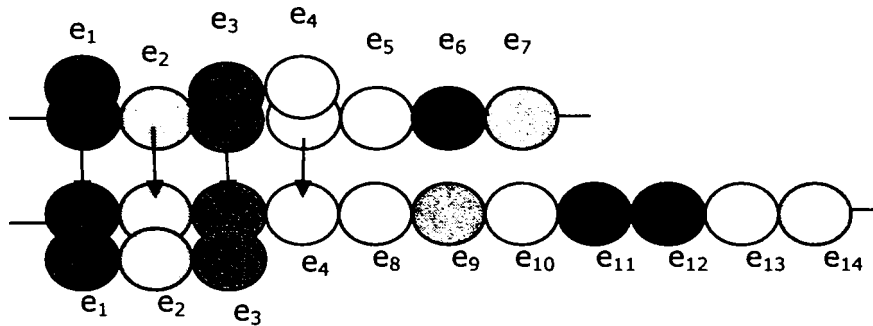
Reprezentarea metodologiei pentru metoda directă:

Pas1. Calcularea scorului de referință: Calculează pentru fiecare secvență scorul său "perfect" folosind (5.2-1) unde P este exprimat din prima și a două coloană a aceluiași tabel

$$S_{\text{Prot1/Prot1}} = -\left(\frac{2}{10} \log_2 \frac{2}{2} + \frac{1}{10} \log_2 \frac{1}{1} + \frac{2}{10} \log_2 \frac{2}{2} + \frac{2}{10} \log_2 \frac{2}{2} + \frac{1}{10} \log_2 \frac{1}{2} + \frac{1}{10} \log_2 \frac{1}{2} + \frac{1}{10} \log_2 \frac{1}{1}\right) = \\ = -(0+0+0+0+0-0.1-0.1+0) = -0.2$$

Pas2. Calculează scorul Prot1/Prot2 folosind (5.2.2-1) cu P_x exprimat în prima coloană a Prot1 și P_y din a doua coloană a Prot2.

$$S_{\text{Prot1/Prot2}} = -\left(\frac{2}{10} \log_2 \frac{2}{3} + \frac{1}{10} \log_2 \frac{2}{2} + \frac{2}{10} \log_2 \frac{2}{2} + \frac{2}{10} \log_2 \frac{1}{1}\right) = -[0.2(1-1.58)] = 0.116$$



Pas3. Calcularea similarității pentru metoda directă :

$$\text{Dist}(\text{Prot1}, \text{Prot2}) = |S_{\text{Prot1}/\text{Prot1}} - S_{\text{Prot1}/\text{Prot2}}| = |0.2 - 0.116| = 0.084$$

Metoda alternantă: Singura diferență față de metoda directă este la pasul doi, atunci când se compară secvența query cu fiecare secvență din baza de date. Rolul de secvență referință și model se poate interschimba în funcție de care are lungimea cea mai mică (secvența cea mai scurtă joacă rolul de referință în (5.3-1)). Restul pașilor, estimarea scorului perfect, ordonarea și selectarea urmează aceeași procedură ca și în cazul metodei directe.

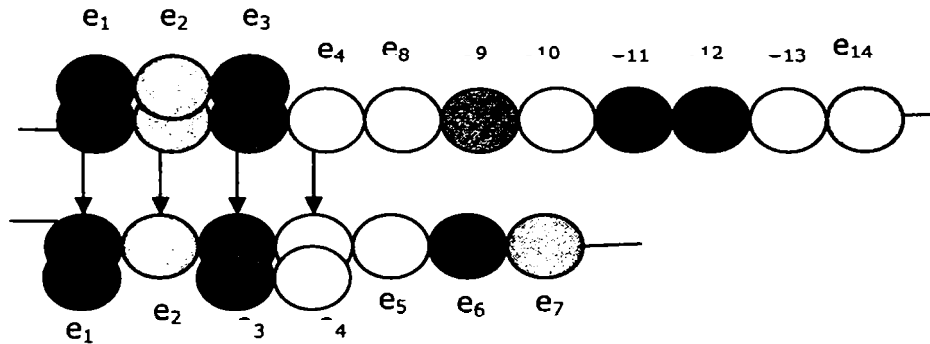
Reprezentarea metodologiei pentru metoda alternantă:

Pas1. Calcularea scorului de referință: Calculează pentru fiecare secvență scorul său "perfect" folosind (5.2-1) unde P este exprimat din prima și a doua coloană a aceluiași tabel

$$\begin{aligned} S_{\text{Prot1}/\text{Prot1}} &= -\left(\frac{2}{10} \log_2 \frac{2}{2} + \frac{1}{10} \log_2 \frac{1}{1} + \frac{2}{10} \log_2 \frac{2}{2} + \frac{2}{10} \log_2 \frac{2}{2} + \frac{1}{10} \log_2 \frac{1}{2} + \frac{1}{10} \log_2 \frac{1}{2} + \frac{1}{10} \log_2 \frac{1}{1}\right) = \\ &= -(0+0+0+0+0-0.1-0.1+0) = -0.2 \end{aligned}$$

Pas2. Calculează scorul Prot1/Prot2 folosind (5.3-1) cu P_x exprimat în prima coloană a Prot2 și P_y din a doua coloană a Prot1.

$$S_{\text{Prot2}/\text{Prot1}} = -\left(\frac{2}{14} \log_2 \frac{2}{2} + \frac{2}{14} \log_2 \frac{1}{1} + \frac{2}{14} \log_2 \frac{2}{2} + \frac{1}{14} \log_2 \frac{2}{2}\right) = -(0+0+0+0) = 0$$



Pas3. Calcularea similarității pentru metoda alternantă :

$$\text{Dist}(\text{Prot1}, \text{Prot2}) = |S_{\text{Prot1}/\text{Prot1}} - S_{\text{Prot2}/\text{Prot1}}| = |0.2 - 0| = 0.2 .$$

Descrierea algoritmică a celor două metode

Metoda directă:

Fie setul S de proteine, P_X este întotdeauna estimat ca frecvența relativă urmând ecuația 5.2.1-5 și P_Y este estimat folosind ecuația 5.2.1-3.

Repetă

*Pentru fiecare proteină din set, obține scorul referință(perfect) aplicand ecuația 5.2.2-1 în modul următor:

X = secvența de proteină curenta (query);

w_i^n parcurge toate evenimentele n-gram din X ;

P_X este estimat din X ;

P_M este estimat de asemenea din același X ;

RS = rezultat;

Repetă

*Pentru fiecare proteină din mulțimea $S-X$, aplică ecuația 5.3-1

unde:

X = secvența de proteină curenta (query);

Y = secvența de proteină curentă din mulțimea $S-X$;

X_n parcurge toate elementele n-gram din X ;

P_X este estimat din X ;

P_Y este estimat din Y ;

Res = rezultat;

Score = $\text{abs}(\text{RS} - \text{Res})$;

*colectează toate scorurile obținute pentru fiecare secvență

X considerată query și restul din setul de date;

până când mulțimea $S-X$ este vidă;

până când mulțimea S este vidă.

Metoda alternantă:

Fie setul S de proteine. P_X este întotdeauna estimat urmând ecuația 5.2.1-5 și P_Y este estimat folosind ecuația 5.2.1-3.

Repetă

*Pentru fiecare proteină din set, obține scorul referință(perfect) aplicand ecuația 5.2.2-1 în modul următor:

X secvența de proteină curenta (query);

w_i^n parcurge toate n-gram din X;

P_X este estimat din X;

P_Y este estimat de asemenea din același X ;

RS :=abs(rezultat);

Repetă

*Pentru fiecare proteină din mulțimea S-X, aplică ecuația 5.2.2-1 unde:

X secvența de proteină curenta (query);

Y secvența de proteină curentă din mulțimea S-X;

*Dacă lungimea proteinei X este mai mare decât lungimea

lui Y

atunci

X_n parcurge toate elementele n-gram din Y;

P_X este estimat din Y;

P_Y este estimat din X;

Res := abs(rezultat);

Score := abs(RS-Res);

*colectează toate scorurile obținute pentru fiecare secvență X considerată query și restul din setul de date;

până când mulțimea S-X este vidă;

până când mulțimea S este vidă.

În ambele cazuri, variabila Score reprezintă distanța dintre două numere reale cunoscută ca $\text{dist}(a, b) = |a-b|$, unde $a, b \in \mathbf{R}$ și $|x|$ este valoarea absolută a lui x.

5.2.4. Implementare și detalii de performanță ale noii metode

Urmând metodologia deja descrisă, implementarea propriu-zisă respectă parcurgerea acestor etape. Astfel, utilizând limbajul de programare Perl s-au construit programe specializate care să execute operațiile de: construire a modelelor n-gram, stocare a informației conținute de fiecare secvență precum și de evaluare și estimare a disimilarității/similarității dintre secvențele comparate. Ca și în cazul programelor care constituie setul destinat pentru analiza lingvistică descris în capitolul anterior, și aceste programe sunt independente din punct de vedere al implementării dar dependente funcțional. Ele se constituie ca module de scripturi

lansate în execuție din linia de comandă, care pot fi rulate atât pe platforma oferită de sistemul de operare Windows cât și Unix/Linux, beneficiind de un grad înalt de portabilitate.

Conform studiilor din ingineria proceselor software, în general produsele soft pot fi evaluate fie în mod direct fie indirect. Prin evaluarea directă se înțelege determinarea costurilor și a eforturilor asociate. Ea presupune calculul numărului liniilor de cod (LOC - lines of code) scrise, determinarea vitezei de execuție, a dimensiunii memoriei, precum și a numărului de defecte raportat într-un anumit interval de timp [FLOR'96]. Evaluarea indirectă a produsului reprezintă în fapt o analiză a funcționalității, calității, complexității, eficienței, fiabilității, întreținerii și multor altor caracteristici. Pentru implementarea algoritmilor introduși în acest capitol însă, atenția nu a fost orientată în mod special pe aspectele de implementare cât pe testarea performanței în ce privește acuratețea rezultatelor generate. Aceasta, deoarece este evidentă simplitatea și avantajul utilizării noii metode și a resurselor informaționale folosite. În plus, implementarea metodei s-a dorit a fi realizată cât mai modular și execuția este secvențială, pentru a putea accesa și modifica ușor orice porțiune care ar necesita intervenții pe parcursul execuției sau în vederea unei alternative la ideea de bază. Toate acestea fac din setul de programe implementat un instrument experimental ușor de manevrat.

Prin urmare, o analiză a complexității în această situație se poate realiza folosind noțiunea de complexitate algoritmică. Este foarte convenabilă clasificarea algoritmilor pe baza cantității relative de spațiu necesar și specificarea creșterii cerințelor de timp/spațiu ca o funcție dependentă de mărimea datelor de intrare. Astfel există noțiunea de complexitate a timpului, care reprezintă timpul de execuție al programului ca funcție dependentă de mărimea datelor de intrare și noțiunea de complexitate a spațiului, semnificând cantitatea de memorie necesară în timpul execuției programului, tot în funcție de mărimea datelor de intrare.

Cea mai comună metrică pentru calcularea complexității timpului și spațiului este notația „O” (o cu majusculă). Aceasta înlătură toții factorii constanți astfel încât timpul de execuție poate fi estimat în relație cu N, pentru N reprezentând o valoare care poate tinde spre infinit. Deoarece $O(\text{expresie})$ reprezintă toate funcțiile care cresc mai încet sau egal cu *expresie*, în general el este folosit pentru a exprima limitele superioare.

În cazul de față, s-a adoptat ca măsură a complexității numărul de pași pe care îi execută algoritmul astfel: dacă expresia numărului de pași este polinomială, se va considera termenul dominant neglijându-se coeficientul numeric. Analog și la expresie exponentială. Adică: dacă expresia este numărul de pași $N=2*n^3*e^n+7*n+1$ (unde n este numărul de elemente și e este baza logaritmilor), atunci complexitatea $C=n^3*e^n$ sau $O(n^3*e^n)$ se dorește a fi o complexitate cât mai mică [AHO'85].

Având în vedere faptul că pentru obținerea similarității este necesară o procesare prealabilă a secvențelor și stocarea de informație relevantă, vor fi approximate complexitățile separat pentru fiecare etapă. Astfel, procesarea unei secvențe presupune o complexitate de timp și spațiu de $O(N)$, unde N este lungimea secvenței de intrare minus $(n-1)$, n fiind ordinul modelului. Pentru evaluarea propriu-zisă a similarității însă se aproximează o complexitate de timp $O(N \times M)$, pentru ambele abordări (metoda directă și alternantă), unde N și M sunt lungimile datelor de intrare (organizate ca succesiuni de evenimente n -gram). În ce privește

complexitatea spațiului stocat pentru *metoda directă* se estimează o complexitate $O(N)$, cu N lungimea secvenței procesate cu rol de query, iar pentru *metoda alternantă* complexitatea spațiului poate fi aproximată de $O(\min\{N,M\})$. Aceste expresii, care exprimă o complexitate relativ redusă se datorează facilității de lucru cu vectori de asociere, oferită de limbajul Perl (Practical Extraction and Report Language).

5.3. Identificarea secvențelor mutante ale unei proteine. Experimente.

Progresul rapid în cercetare legat de bazele de date biologice care conțin secvențe genomice ale diferitelor organisme deschid provocarea de a decifera "misterul vieții". Cu scopul deservirii acestui obiectiv unele dintre cele mai importante interese se pot folosi de tehnicile de data mining (DM). *Data Mining* este definită în general ca un ansamblu de diferite tehnici utilizate pentru a extrage informație ascunsă, predictivă din baze de date multiple și de dimensiuni mari stocate sub diverse forme [BOGA'06a]. Aceste tehnici sunt orientate spre descoperirea de cunoștințe și merg dincolo de simpla caracterizare statistică a datelor disponibile. Secvențele biologice pot fi un domeniu foarte important al DM, fiind capabile să deservească unele dintre cele mai importante arii de interes ale analizei secvențelor biologice. Noua strategie propusă pentru determinarea similarității dintre secvențele biologice este testată în aceste experimente ca metodă de DM în scopul identificării secvențelor mutante ale secvenței de căutare [BOGA'06a].

5.3.1. Baza de secvențe

Strategia propusă pentru măsurarea similarității proteinelor a fost demonstrată și validată folosind o bază de date ce conține în total 100 de secvențe de proteine [BOGA'05b]. Au fost formate două grupuri distincte de proteine după cum urmează. Primele 50 de intrări din baza de date corespund proteinelor selectate aleator din baza de proteine publică NCBI [NCBI'05]. Ultimele 50 de intrări corespund proteinelor rezultate ca urmare a unor mutații diverse ale genei responsabile pentru producerea de cancer, **p53**. Mutațiile au fost selectate aleator din baza de secvențe creată utilizând descrierea oferită de International Agency for Research on Cancer (IARC) Lyon, Franța³³. Baza de date IARC folosită este sub formă de fișier în format Excel și conține 18 585 intrări organizate în 42 coloane. Fiecare linie (coloană) reprezintă o singură mutație căreia îi este asignat un singur număr unic de identificare. Un număr unic de identificare este de asemenea atribuit mostrei de tumoră și pacientului. Din fiecare înregistrare se folosesc câmpurile care definesc acea mutație:

- numărul de identificare unic al mutației;
- locația mutației în intron sau exon³⁴ în gena p53 ;

³³ <http://www.iarc.fr/p53/Somatic.html>

³⁴ **Exonii** sunt regiuni ale AND-ului dintr-o genă care nu sunt înlăturate prin transcrierea ARN și sunt păstrate în molecula finală mesager de ARN(mRNA)[WIKI'05g].

- pentru mutații în exon, numărul codon-ului la care este locată mutația (1-393);
- poziția nucleotidelor mutației pe baza intrării în banca genetică, X54156;
- valorile de Yes (True) sau No (False) pentru a indica dacă poziția mutației cade într-o secvență CpG (Cytosine-phosphor-Guanine),
- valorile Yes (True) sau No (False) pentru a indica dacă poziția mutației cade într-o locație de conexiune ("splice site");
- natura mutației;
- pentru mutații în exoni, secvența de bază a codon-ului în care are loc mutația;
- secvența bazei mutate;
- tipul de amino acid (wild-type amino acid) codificat în codon-ul în care are loc mutația (abreviere din trei litere a amino acidului);
- amino acidul mutat codificat în codon-ul în care a avut loc mutația (abreviația de trei litere a amino acidului) și
- efectul mutației și al numărului de codon la care poate apare un codon de STOP.

Din baza de date creată au fost luate aleator 50 de mutații care produc *missense mutation* (când modificarea nucleotidelor produce modificare în codon, atunci amino acidul se modifică și în final și proteina), *nonsense mutations* (când modificarea în nucleotide determină modificarea codonului în TAA, TAG sau TGA – care sunt codon de STOP și prin urmare nu se va produce nici un amino acid, rezultând într-o proteină mai scurtă) și *silent mutations* (când schimbarea în nucleotide determină modificarea codon-ului dar nu se modifică nici amino acizii și prin urmare nici proteina).

Acest set de 50 de proteine, numit în acest context grupul p53, se așteaptă să formeze un grup compact de forme textuale în spațiul semnificațiilor biologice. În mod contrar, restul de 50 de proteine trebuie să apară ca forme textuale în același spațiu care nu numai să difere de un altul dar cu precădere de grupul p53. Setul complet de secvențe poate fi consultat în ANEXA 4.

5.3.2. Descrierea experimentelor. Rezultate.

Cu scopul de a ilustra cele două variante ale strategiei propuse, în primul rând au fost urmați niște pași clasici ai analizei exploratorii a datelor (*Exploratory Data Analysis*). Matricea conținând toate măsurile de *disimilaritate* posibile obținute $D(S_i, S_j)$, $i, j=1, 2, \dots, N$ pentru secvențele date este ilustrată în Figura 5.3.2-1, ca imagine la scara gri, pentru ambele variante algoritmice ale noii metode și trei modele n -gram diferite. În schema de vizualizare adoptată toate matricele prezentate (după o normalizare adecvată), partajează o scală comună în care 1 (alb) corespunde distanței maxime în fiecare matrice. Merită menționat faptul că reprezentarea spațială 'ideală' în acest caz este o matrice albă cu un singur segment negru în colțul din dreapta jos. Este prin urmare evident din Figura 5.3.2-1 faptul că modelarea 4-gram urmată de versiunea 'alternantă' a noului algoritm atinge o performanță aproape excelentă în căutarea în baza de secvențe date.

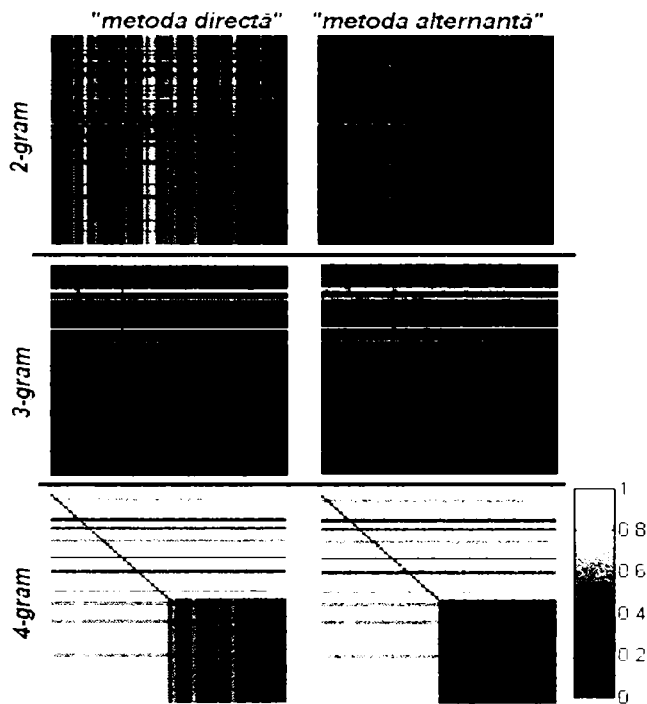


Figura 5.3.2-1. Vizualizarea matricelor conținând disimilitățile tuturor perechilor determinate de cele 100 de proteine din baza de date utilizată.

Datorită separării evidente a secvențelor în baza de date s-a trecut la determinarea apartenenței secvențelor vizualizate ca strâns grupate în colțul negru dreapta-jos al Figura 5.3.2-1. Astfel, cu scopul investigației lor, rezultatele obținute în urma aplicării metodei alternante pentru modelul 4-gram sunt proiectate într-un spațiu dimensional redus, reprezentat în Figura 5.3.2-2. Proteinele necunoscute sunt marcate cu albastru(formele definite de contur) iar cele provenind din mutațiile genei p53 în culoarea roșu(formele reprezentate cu continut hasurat). Este evident faptul că în acest caz s-a obținut o foarte bună soluție la problema identificării/separării a două clase.

Mențiune: Metoda de reprezentare vizuală a matricelor de disimilaritate îi aparține lui Dr. Nikos Laskaris, Universitatea Aristotel din Thessaloniki, Grecia, cu al cărui consimțământ și îndrumare a fost folosită.

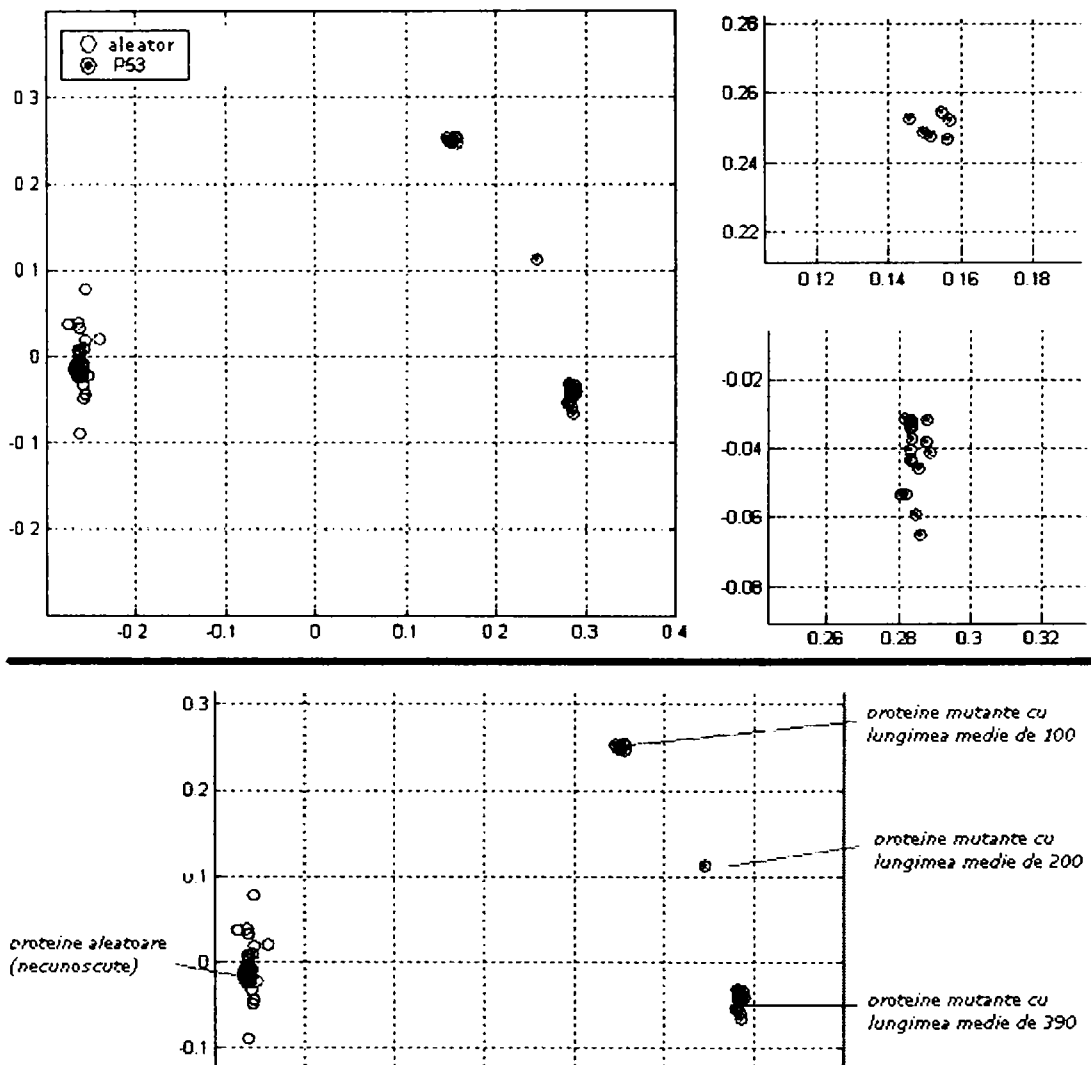


Figura 5.3.2-2. Reprezentarea în spațiu dimensional redus a secvențelor de proteine folosind măsura disimilarității pentru setul experimental de 100 secvențe

În al doilea rând, cu scopul furnizării unei măsuri cantitative a performanței celor două variante, s-a adoptat un index al acurateții care este derivat din metoda ROC (Receiver Operating Characteristic) și care a câștigat popularitate în validarea rezultatelor obținute din căutarea în baze de date (Liao and Noble [LIAO'03], Schäffer et al [SCHÄ'01]). Acest index, de obicei referit ca scor ROC trunchiat, este raportul dintre suprafața aflată sub curba ROC (în proiectarea valorilor adevărate și pozitive versus false și pozitive, pentru diferite praguri de disimilaritate). Mai explicit, după cum este menționat și în [SCHÄ'01]), pentru un număr T de cazuri adevărate și pozitive posibil de a fi găsite și un număr fix de cazuri false și pozitive n , acest index este proporția dată de dreptunghiul $[0, T] \times [0, n]$ care se află sub curba de sensibilitate. El ia valori în intervalul $[0-1]$, cu 1 corespunzând celei mai

ridicate performanțe. Acest scor ROC a fost introdus în Tabel 5.3.2-1 pentru diferite modele n -gram și ambele versiuni ale metodei iar în Figura 5.3.2-3 este o reprezentare grafică a traiectoriei urmate de fiecare metodă în circumstanțele impuse de modelul n -gram folosit și diverse valori prag pentru evaluare.

Tabel 5.3.2-1 Scorul ROC obținut pentru ambele metode pentru diferite valori prag.

model n -gram	Aria normalizata aflata sub curba ROC	
	<i>Metoda Directă</i>	<i>Metoda Alternantă</i>
2-gram	0.589	0.680
3-gram	0.723	0.817
4-gram	0.900	0.978

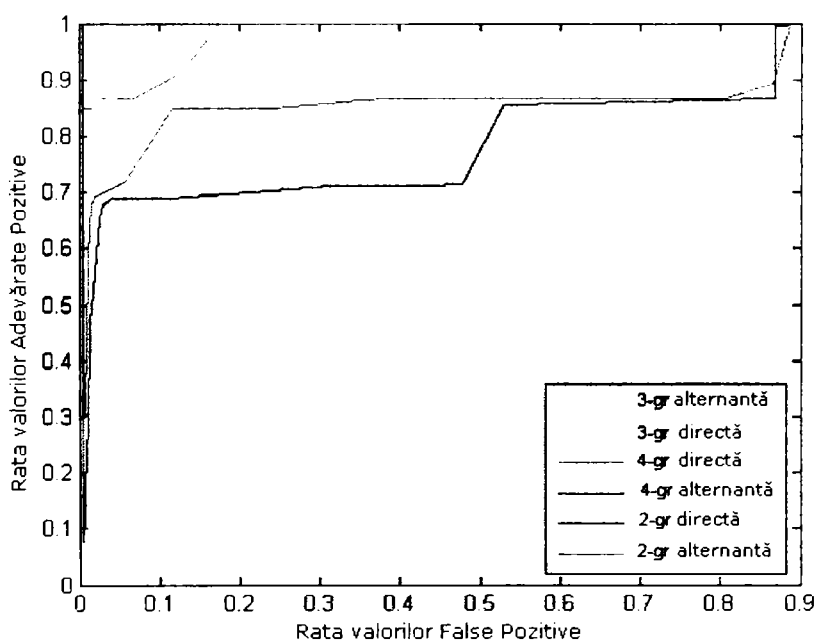


Figura 5.3.2-3. Reprezentarea evoluției valorilor raporturilor adevărate și pozitive versus false și pozitive

5.3.3. Analiza și interpretarea rezultatelor

Metoda prezentată constituie un prim pas în investigarea implicării modelării lingvistice în caracterizarea, manevrarea și înțelegerea datelor biologice sub formă de secvențe. În mod special, s-a studiat folosirea conceptului de cross-entropie aplicat modelelor n -gram cu scopul căutării eficiente în baze de date. Rezultatele experimentale au indicat soliditatea noii strategii algoritmice în exprimarea similarității dintre proteine [BOGA'05].

Fiind dată simplitatea conceptuală a abordării acestei metode noi, ea apare ca o alternativă de dorit tehnicilor bine consolidate de până acum.

Considerând problema generală de separare dintre măsurile de similaritate "globală" și "locală" dintre proteine, trebuie menționat faptul că noua abordare aparține primei categorii.

În timpul evaluării acestei noi metode s-a observat că dintre cele două variante introduse cea mai bună performanță este asociată cu a doua (metoda alternantă). Această remarcă semnifică nevoia de a veni cu îmbunătățiri care să depășească posibilele identificări eronate ale secvențelor similare datorate diferenței semnificative în lungime. În cazul excepțional când toate secvențele comparate ar avea aceeași lungime, *metoda directă* este echivalentă cu *metoda alternantă* și se comportă foarte bine.

În ce privește ordinul modelului n -gram folosit, după testarea ordinelor 2,3,4,5 s-a observat că performanța metodei crește odată cu ordinul modelului până la 4. După ordinul 5, datorită lipsei datelor, estimatorii probabilității maxime devin nerezonabil de uniformi și mici. Aceasta stabilește o limită superioară pentru modelul ales pentru această bază de date (probabil că un model de ordin ceva mai mare ar putea funcționa în alte baze de date).

Înainte de a se trece la munca de îmbunătățire a acestei metode trebuie remarcat faptul că este o tehnică statistică la bază. Ea ar putea deveni mai performantă prin incorporarea de cunoștințe biologice cum ar fi lucrul cu grupe funcționale de amino acizi

În final, un alt aspect care merită a fi luat în considerare este de a testa comportamentul acestei metode vis-a-vis de dimensiunea bazelor de date de secvențe.

5.4. Determinarea similarității secvențelor dintr-o bază structurată de proteine (SCOP). Experimente.

5.4.1. Baza de secvențe

Aceeași strategie propusă pentru măsurarea similarității proteinelor a fost demonstrată și validată de data aceasta folosind un set de 1460 proteine extrase din baza de resurse de secvențe Astral SCOP, versiunea 1.67. În conformitate cu organizarea structurală a secvențelor biologice în baza de secvențe SCOP [SCOP], și strategia sa de căutare, toate proteinele au fost selectate din clasa a (care include toate alfa proteinele caracterizate prin aceeași structură α -helix) astfel încât să prezinte grade diferite de similaritate. Din corpusul original disponibil au fost selectate și incluse în baza de date experimentală numai acele familii de secvențe care conțin cel puțin 10 secvențe (această restricție va fi apreciată mai târziu deoarece a fost dictată de măsura de **precizie** adoptată pentru evaluarea metodei).

În acest fel, au fost selectate 31 de familii diferite, populate inegal (a se vedea în Tabel 5.4.1-1). Reamintesc aici că adnotarea bazei de date experimentale păstrează adnotarea originală care se bazează pe semnificația biologică a conceptului de similaritate și ca urmare poate fi considerată ca « valoare de

adevăr » pentru clasificarea proteinelor și măsurile de similaritate testate. În consecință, se așteaptă ca toate proteinele aparținând aceleiași familii să apară ca un grup compact de forme textuale și fiind dată o măsură adecvată de similaritate am putea diferenția diverse familii.

Baza de date (de 1460 proteine) a fost organizată în trei seturi diferite, deoarece rezultatele experimentale obținute prin noua metodă se doresc a fi comparate cu cele furnizate de metoda adoptată de un alt program care poate avea limitări privind numărul de secvențe procesate (poate accepta la intrare doar până la 500 secvențe). În Tabel 5.4.1-2 sunt incluse toate informațiile necesare înțelegerii organizării adoptate. Fiecare set conține un număr diferit de familii (depinde de numărul de secvențe conținut de fiecare familie) astfel încât numărul total de secvențe să nu depășească 500 per set. Baza de date completă (organizată în 3 subseturi) este disponibilă la cerere și va putea fi accesată public de pe pagina web a laboratorului AIIA (Artificial Intelligence and Information Analysis, Universitatea Aristotel, Salonic, Grecia).

Tabel 5.4.1-1. Numărul de superfamilii și familii folosite în experimente.

	Nr. Superfamii	Nr. Familii
Set1	5	(1,5), (2,3), (3,3), (4,2), (5,2) Total =15
Set2	4	(1,2), (2,4), (3,2), (4,2) Total=10
Set3	5	(1,1), (2,2), (3,1), (4,1), (5,1) Total=6

Pentru fiecare set s-a menționat numărul superfamiliiilor (în prima coloană). În a doua coloană este prezentat în mod explicit sub formă de perechi de forma (x,y), numărul de familii (y) care aparțin superfamiliei corespunzătoare (x).

Tabel 5.4.1-2. Familiile SCOP incluse în experimente^a

Set1		Set2		Set3	
ID-ul familiei de proteine	Nr. Secv.	ID-ul familiei de proteine	Nr. secv.	ID-ul familiei de proteine	Nr. Secv.
a.39.1.1	10	a.1.1.2	28	a.123.1.1	109
a.39.1.2	30	a.1.1.3	213	a.22.1.1	103
a.39.1.4	14	a.4.1.1	21	a.22.1.3	15
a.39.1.5	119	a.4.1.2	19	a.45.1.1	85
a.39.1.8	16	a.4.1.12	67	a.93.1.1	78
a.26.1.1	32	a.4.1.3	11	a.133.1.2	76
a.26.1.2	45	a.25.1.1	52		
a.26.1.3	30	a.25.1.2	37		466
a.35.1.1	10	a.138.1.1	22		
a.35.1.2	26	a.138.1.3	27		
a.35.1.5	13				
a.3.1.1	85		497		
a.3.1.2	21				
a.118.1.1	34				
a.118.1.14	12				
	497				

^aPentru fiecare familie numele complet al fiecărei proteine este disponibil la: <http://astral.berkeley.edu/scopseq-1.67.html>.

Simbolul fiecărei proteine menționate conține patru domenii reprezentând complet clasificarea structurală biologică a proteinelor. Spre exemplu a.26.1.2 înseamnă clasa 'a', plierea '26', superfamilia '1', și familia '2'. 'Nr.seq' indică numărul de secvențe de proteine conținute de fiecare familie.

5.4.2. Descrierea experimentelor.

Cele două variante ale strategiei propuse în analiza similarității sunt testate folosind noua bază de secvențe extrasă din SCOP[BOGA'06d]. Pentru aceasta au fost urmați aceiași pași clasici ai analizei exploratorii a datelor care au fost menționați și în experimentele anterioare. Matricea conținând toate măsurile de disimilaritate posibile obținute $D(S_i, S_j)$, $i, j=1, 2, \dots, N$ pentru seturile 1-3 este ilustrată în Figura 5.4.2- respectiv, ca imagine la scara gri, pentru ambele variante algoritmice ale noii metode și trei modele n -gram diferite. În schema de vizualizare adoptată toate matricele obținute (după o normalizare adecvată), partajează o scală comună în care 1(alb) corespunde distanței maxime în fiecare matrice. Se reamintește faptul că reprezentarea spațială 'ideală' este o matrice albă cu niște segmente negre în jurul liniei diagonale. Din aceste trei figuri este evident faptul că modelarea bazată pe modelul 4-gram urmtă de ambele versiuni ale noului algoritm are o foarte bună performanță în compararea secvențelor din baza de proteine dată. Cu scopul de a furniza măsuri de performanță cantitative pentru cele două variante, s-a adoptat și un index al acurateții căutării care este derivat din măsura **precizie** [BAEZ'99]. Acest index este raportul calculat prin împărțirea numărului de proteine corect clasificate (identificate de algoritm ca primele 10 cele mai similare) cu 10 (numărul minim de secvențe din fiecare familie). Și anume, fiecare proteină a fost la un moment dat considerată query și s-a măsurat acuratețea primelor 10 secvențe identificate în cadrul setului ca fiind cele mai similare cu secvența de query. Cu alte cuvinte, luând în considerare eticheta de clasă/familie a fiecărei proteine s-au identificat proteinele care au aceeași etichetă cu secvența de query (i.e. a number from 1 to 10). Procedura a fost repetată pentru toate proteinele în seturi individuale iar la final s-a făcut media părților estimate cu scopul de a furniza un scor de precizie totală pentru fiecare set în parte. Pentru completitudine, s-au repetat măsurările preciziei pentru cazul în care algoritmul nou s-ar fi aplicat întregului număr de 1460 de secvențe din baza de proteine. Valorile obținute astfel nu au fost semnificativ diferite față de valorile corespunzătoare celor trei seturi, dovedind o evidență în plus a robusteții noii metode și indică faptul că performanța sa este bună raportat la mărimea bazei de date.

În bioinformatică, algoritmi de clasificare încearcă să grupeze secvențele care sunt cumva înrudite. La modul general, algoritmi de clasificare sunt bazați pe legături singulare construind o închidere tranzitivă a secvențelor cu o similaritate care depășește un anumit prag. Scorul de similaritate este adesea bazat pe aliniamentul de secvențe. De obicei clasificarea secvențelor este folosită pentru a produce un set non-redundant de secvențe reprezentative și clasele de secvențe sunt adesea sinonime (dar nu identice) cu familiile de proteine. Determinarea unei structuri reprezentative pentru fiecare grupare este scopul multor inițiative genomice. Obiectivul general de a grupa proteinele în familii conduce la detectarea mai sensibilă de noi membri și o discriminare îmbunătățită împotriva potrivirilor false, pe baza caracteristicilor esențiale conservate într-o familie.

Cea mai evidentă măsură a similarității (sau disimilarității) dintre două secvențe este distanța dintre ele. O cale de a începe investigarea modului de grupare este de a defini o metrică adecvată și de a calcula matricea distanțelor dintre toate perechile de secvențe. Dacă distanța este o măsură bună a disimilarității, atunci se așteaptă ca distanțele dintre secvențele aflate în aceeași grupare să fie semnificativ mai mici decât distanțele față de secvențele aflate în alte grupări [BOGA'06b].

Ca urmare a observației din experimentele realizate asupra setului de 100 de secvențe, unde s-a observat potențialul noii metode de a clasifica secvențele în două categorii evidente (secvențe mutante ale proteinei generate de gena p53 și restul) s-a propus o analiză mai aprofundată a acestui caz considerând rezultatele obținute din baza structurată de proteine. Astfel, folosind statisticile Hubert [LASK'02] în determinarea factorilor de corelație dintre grupările de secvențe evidențiate în reprezentarea vizuală și valoarea lor de adevăr furnizată de codificarea din baza de date SCOP, s-au obținut valorile tabulate în Tabel 5.4.2-2. Statisticile lui Hubert reprezintă o tehnică de corelare și măsoară gradul de corespondență liniară dintre matricea distanțelor unor modele și matricea etichetelor categoriilor corespunzătoare. Descrieri detaliate ale acestei metode pot fi găsite în [JAIN'88].

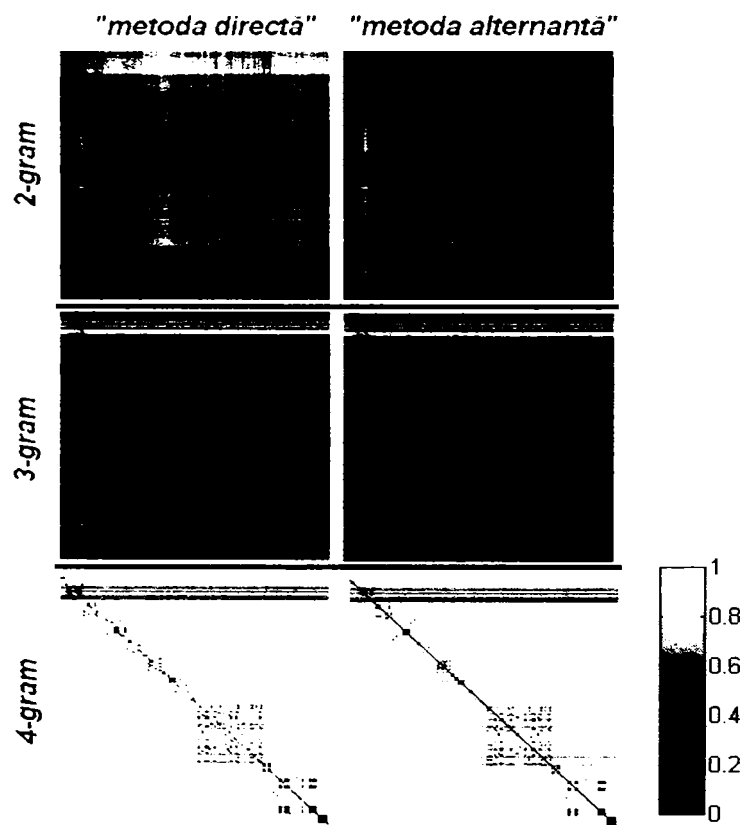


Figura 5.4.2-1.

Figura 5.4.2-1. Vizualizarea matricilor care conțin toate disimilaritățile posibile pentru perechile de proteine ale Set1 format din 497 secvențe pentru modele 2,3,4-gram.

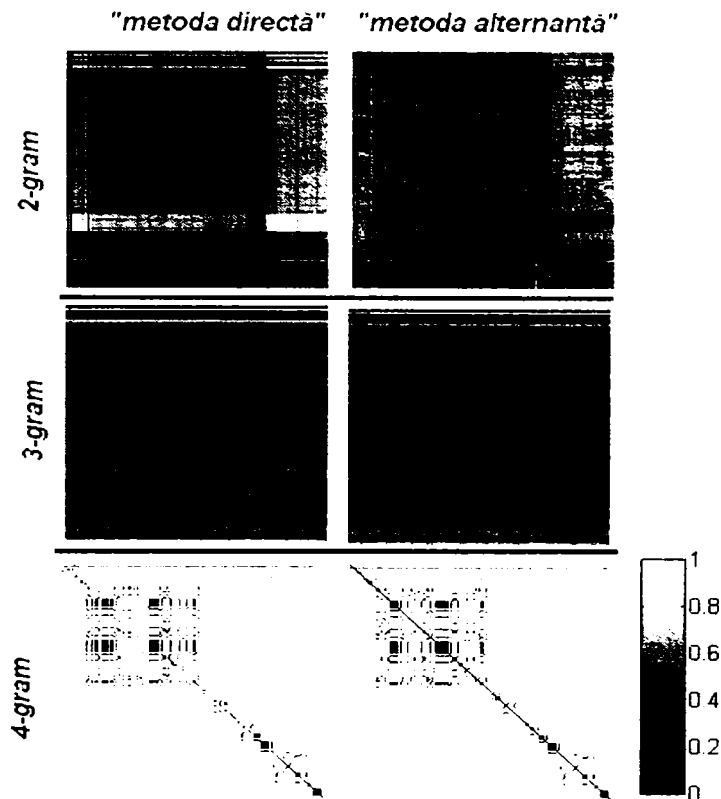


Figura 5.4.2-2.

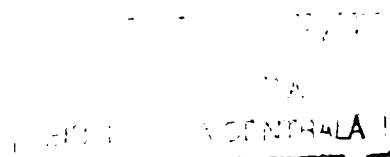
Figura 5.4.2-2. Vizualizarea matricilor care conțin toate disimilaritățile posibile pentru perechile de proteine ale Set2 format din 497 secvențe pentru modele 2,3,4-gram.

În Tabel 5.4.2-1 sunt incluse valorile de precizie furnizate de cele două abordări ale noului algoritm pentru diferite modele n -gram.

Tabel 5.4.2-1. Valorile preciziei

Set	Metoda Directă			Metoda Alternantă		
	2-gram	3-gram	4-gram	2-gram	3-gram	4-gram
1	0.439	0.662	0.830	0.471	0.646	0.823
2	0.446	0.650	0.874	0.439	0.605	0.860
3	0.534	0.865	0.931	0.574	0.828	0.919

Sunt prezentate valorile **preciziei** obținute din rezultatele oferite de noua măsură de similaritate pentru ambele abordări testând modele 2,3,4-gram pentru setul de date structurate obținut din baza de proteine SCOP.



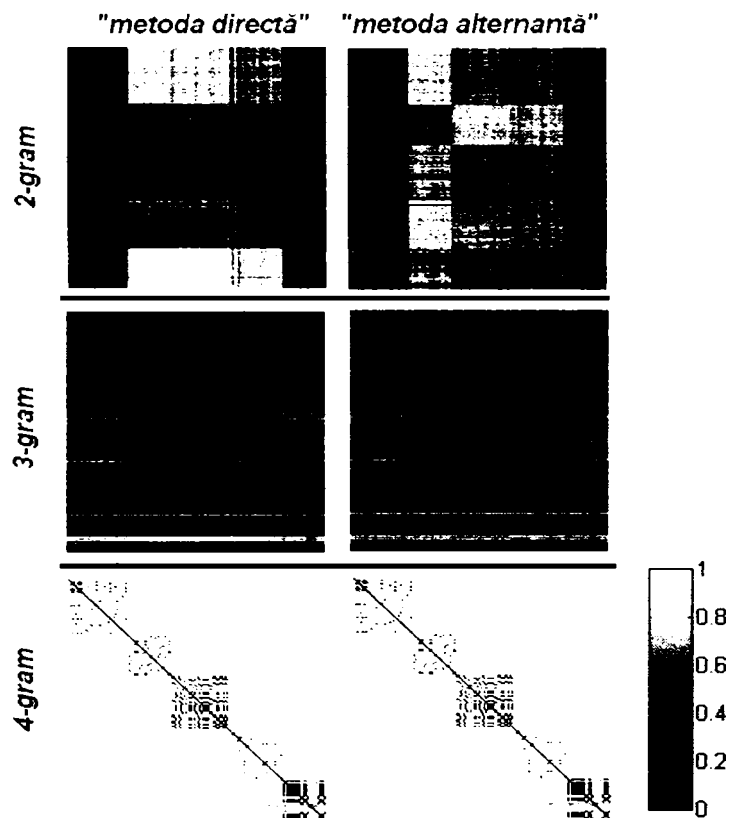


Figura 5.4.2-3. Vizualizarea matricilor care conțin toate disimilaritățile posibile pentru perechile de proteine ale Set3 format din 466 secvențe pentru modele 2,3,4-gram.

Tabel 5.4.2-2. Valorile coeficienților de corelație pentru Set1, 2 și 3 folosind statisticile Hubert.

Set	Metoda Directă	Metoda Alternantă
1	0.1230	0.1110
2	0.0347	0.0354
3	0.3339	0.3622

În interpretarea valorilor coeficienților de corelație sunt considerate bune valorile mari și nu este cazul atins cu presupunerea cu care s-a pornit (posibila identificare a clasificării biologice structurale a secvențelor în familii/superfamilii). În aceste circumstanțe, intervenția raționamentului biologului ajută la elucidarea semnificației grupurilor obținute. Explicația vine din faptul că de multe ori, secvențele de lungimi diferite și cu un conținut parțial identic pot aparține unei aceleiași familii biologice. Astfel, grupurile identificate în acest caz sunt secvențele similare în reprezentare textuală. Această concluzie este deja justificată de

rezultatele obținute în identificarea grupului proteinelor mutante din baza de date anterioară.

5.4.3. Analiza și interpretarea rezultatelor

În timpul evaluării noii metode folosind baza de date structurată, s-a observat performanța apropiată a celor două variante algoritmice indicând lungimea apropiată a secvențelor care se așteaptă să aparțină aceleiași familii. Această remarcă induce observația conform căreia acele secvențe (de lungime apropiată și care par a fi asemanătoare) au un grad ridicat de similaritate. Dacă a doua metodă (*alternantă*) ar fi dat rezultate mai bune, ar fi trebuit să ne așteptăm să avem diferențe semnificative de lungime între secvențele clasificate drept similare și prin urmare aparținând aceleiași familii. Aceasta presupune că este important să se vină cu îmbunătățiri care să depășească posibilele identificări eronate ale secvențelor similare datorate diferenței semnificative în lungime. Așa după cum s-a observat și în cazul primelor experimente pentru ordinul modelului n -gram folosit, după testarea ordinilor 2,3,4 s-a observat în Tabel 5.4.2-1. Valorile preciziei ce conține valorile preciziei și Figurile 5.4.2-1, 2, 3 că performanța metodei crește odată cu ordinul modelului până la 4. După ordinul 5, datorită lipsei datelor, estimatorii probabilității maxime devin nerezonabil de uniformi și mici.

Dacă acordăm mai multă atenție reprezentării vizuale din rezultatele obținute (formele rezultate de-a lungul diagonalei principale în prezentarea 2D corespund unor grupuri bine definite de proteine, în special în cazul modelării cu 4-gram), putem considera că structura relevată prin folosirea noii metode de similaritate poartă o semnificație biologică. Mai ales, presupunerea că fiecare grup definit este indicativ pentru existența unei familii/superfamilii de proteine. Cea mai clară evidență este oferită de Fig. 3, unde al treilea set de date conține un număr de 6 familii aparținând la 5 superfamilii diferite (conform definiției biologice) iar această informație coincide cu numărul de 5 grupuri identificate, care constituie o informație indusă de date (*data-driven information*).

Această prezumție este demontată de calcularea coeficienților de corelație dintre etichetele corespunzătoare fiecărei secvențe din setul original de date și gradul de similaritate indicat de grupările rezultate. Chiar dacă precizia indică un grad ridicat de performanță a metodei nu se poate spune același lucru despre clasificarea structurală. Totuși, este de apreciat faptul că noua metodă de similaritate poate oferi o măsură de grupare pentru secvențe similare din punct de vedere al compoziției lor textuale.

5.5. Concluzii

Metoda prezentată în această lucrare constituie un prim pas în investigarea implicării modelării lingvistice în caracterizarea, manevrarea și înțelegerea datelor biologice sub formă de secvențe. În mod special s-a studiat folosirea conceptului de cross-entropie aplicat modelelor n -gram cu scopul căutării eficiente în baze de date. Rezultatele experimentale au indicat soliditatea noii strategii algoritmice în exprimarea similarității dintre proteine. Fiind dată simplitatea conceptuală a acestei noi abordări, ea apare ca o alternativă de dorit tehnicilor bine consolidate de până acum. Așa după cum au fost remarcate în secțiunile de analiză a rezultatelor

experimentale, principalele concluzii ale acestui capitol sunt rezumate în cele ce urmează.

- Considerând problema generală de separare dintre măsurile de similaritate "globală" și "locală" dintre proteine, trebuie menționat faptul că noua abordare aparține primei categorii;
- Conceptul de n -gram chiar dacă a fost folosit în lucrări anterioare în domeniul bioinformaticii, cum ar fi [GANA'04b], [ERHA'80], [KARL'91], [KARL'96], pentru abordarea actuală el are meritul de a fi prima încercare de a adopta acest pas dual în compararea secvențelor biologice. Prin urmare, au fost necesare unele experimente pentru a descoperi cel mai bun mod în care aceste idei pot fi valorificate în domeniul specific de aplicabilitate;
- Rezultatele obținute în urma aplicării celor două variante ale noi metode în ambele experimente releva un grad ridicat de sensibilitate la lungimea secvențelor. Dacă în primul caz metoda alternantă furnizează cele mai bune rezultate, în cazul bazei de date structurate performanța lor este oarecum apropiată. Această remarcă subliniază nevoia de a veni cu îmbunătățiri care să depășească posibilele identificări eronate ale secvențelor similare datorate diferenței semnificative în lungime.
- În cazul excepțional când toate secvențele comparate ar avea aceeași lungime, *metoda directă* este echivalentă cu *metoda alternantă* și se comportă foarte bine.
- Pentru ordinul modelului n -gram folosit, după testarea ordinelor 2,3,4,5 în ambele experimente, s-a observat că performanța metodei crește odată cu ordinul modelului până la 4. După ordinul 5, datorită lipsei unei cantități mai mari de date, estimatorii probabilității maxime devin nerezonabil de uniformi și mici. Aceasta stabilește o limită superioară pentru modelul ales pentru această bază de date. Este foarte probabilă existența unor grupări de 4 amino acizi cu importanță biologică ce participă în mod decisiv la caracterizarea secvențelor de proteine.
- Folosirea rezultatelor furnizate de detectarea similarității secvențelor în formarea de clusteri este evidentă doar că trebuie făcută precizarea că această grupare se bazează pe similaritatea compoziției textuale. Desigur că intervenind cu raționamente biologice se poate obține mult mai mult decât o clasificare textuală.
- Înainte de a se trece la munca de îmbunătățire a acestei metode trebuie remarcat faptul că aceasta este o tehnică pur statistică. Ea ar putea fi îmbunătățită prin incorporarea de cunoștințe biologice cum ar fi lucrul cu grupe funcționale de amino acizi. De asemenea, se are în vedere căutarea unei metode de ajustare a valorilor pentru menținerea evenimentelor care nu participă la evaluarea secvențelor și care ar putea contribui la o mai mare consistență a evaluării.

Ca o constatare finală, această metodă oferă un mod eficient de a captura caracteristicile comune ale secvențelor comparate și în același timp de a evita misiunea mai puțin plăcută de a alege parametrii, funcții suplimentare sau metode de evaluare. Performanța ridicată și caracterul facil de implementare, luate împreună cu eficiența computațională fac din această metodă o alternativă promițătoare la măsurile sofisticate, bine cunoscute, utilizate în compararea secvențelor de proteine.

MENTIUNI:

Aceste experimente au fost realizate în cadrul proiectului EU Biopattern: Computational Intelligence for biopattern analysis in Support of eHealthcare, Network of Excellence Project No. 508803. Adresez pe această cale mulțumiri d-lui dr.ing. Nikos Laskaris pentru sprijinul acordat în realizarea și interpretarea rezultatelor experimentale obținute.

6. STUDIU COMPARATIV AL ALGORITMULUI DE SIMILARITATE PROPUȘ

Pentru a evalua performanța noii metode comparativ cu a altor algoritmi existenți în domeniul bioinformaticii este necesară o abordare comună atât a datelor de testare cât și a funcțiilor executate.

În cazul de față s-a dorit obținerea valorilor de similaritate pe baza comparațiilor multiple de perechi de secvențe. Aceasta este realizată folosind unul dintre cele mai populare programe de aliniament multiplu, Clustal W accesibil la <http://www.ebi.ac.uk/clustalw/>, site-ul oficial al Institutului European de Bioinformatică. Mai există și alte programe care atestă o performanță ridicată sau apropiată de cea furnizată de Clustal W, în special din punct de vedere computațional dar în cazul de față ele nu furnizează valori concrete de similaritate decât realizează aliniamentul multiplu al secvențelor identificând regiunile comune. Acesta este însă un aspect mai puțin folositor în analiza rezultatelor obținute cu ajutorul noii metode propuse, deoarece în prezenta fază a dezvoltării aplicației se urmărește mai puțin alinierea secvențelor cât identificarea gradului de similaritate și eventuala grupare corectă a secvențelor similare.

Un alt aspect care se constituie într-un impediment în testarea noii metode în comparație cu alte programe este faptul că majoritatea acestora determină similaritatea unei secvențe doar raportat la una din bazele de date standard prevăzute de aplicație. Aceasta limitează o analiză mai detaliată a rezultatelor obținute și a posibilelor observații suplimentare. Prin urmare, Clustal W permite o analiză comparativă pentru o bază de date comună, propusă de utilizator.

6.1. Algoritmul CLUSTAL W

Pentru alinierea a două secvențe practica standard este de a folosi programarea dinamică [Need/Wunsch]. Aceasta ar garanta un aliniament optimal din punct de vedere matematic, fiind dată o tabelă de valori pentru potriviri și nepotriviri între amino acizi sau nucleotide (ex. Matricele PAM250 sau BLOSUM62) și penalități pentru inserții sau ștergeri de diferite lungimi [THOM'94]. Încercările de generalizare a programării dinamice la aliniamente multiple sunt limitate la un număr redus de secvențe scurte [THOM'94]. Pentru mai mult de 8 proteine de lungime medie problema este dificil de calculat, depinzând mult de puterea computațională disponibilă. Prin urmare, toate metodele capabile să lucreze cu date mai mari, în timp computațional rezonabil, folosesc euristici.

Clustal W este considerat unul dintre cele mai populare programe folosite pentru aliniamentul multiplu de secvențe biologice și a beneficiat de multe îmbunătățiri față de versiunea inițială Clustal V. El realizează un aliniament multiplu progresiv iar metoda de aliniament constă din trei mari etape: (i) toate perechile de

secvențe sunt aliniate separat cu scopul de a calcula o matrice a distanțelor ce furnizează divergența fiecărei perechi de secvențe; (ii) un arbore ghid este calculat pe baza matricei de distanțe; (iii) secvențele sunt aliniate progresiv în conformitate cu gradul de ramificare din arborele ghid. Un exemplu ce folosește 7 secvențe de globină cu structura terțiară cunoscută este ilustrat în Figura 3.4.5-2 din capitoul 3 alături de alte observații generale referitoare la metoda de aliniament pe care o folosește.

Aliniamentul se aplică pentru secvențe cu diverse grade de divergență și are o semnificație biologică, calculând cele mai bune potriviri pentru secvențele selectate. Ele sunt aliniate astfel încât identitățile, similaritățile și diferențele pot fi văzute. Relațiile evoluționare pot fi văzute via cladograme³⁵ sau filograme³⁶ ce pot fi opțional generate.

Programul oferă posibilitatea alegerii modului de aliniament dintre o metodă rapidă de aproximare [BASH'87] ce permite să fie aliniate un număr foarte mare de secvențe și o alta mai înceată dar cu un grad mai mare de precizie. În cazul primei metode, scorurile sunt calculate folosind numărul de k-tupluri de potriviri (înseamnă reziduuri identice, de obicei de lungime 1 sau 2 pentru proteine sau de lungime 2-4 pentru secvențe de nucleotide) în cel mai bun aliniament dintre două secvențe minus o penalitate fixă pentru fiecare gap. A doua metodă, după cum o afirmă autorii, oferă o mai mare acuratețe a scorurilor, fiind complet bazată pe programarea dinamică și folosește două penalități pentru gap-uri (pentru deschidere sau extindere) și o matrice de ponderi a amino acizilor. Scorurile în acest caz sunt calculate ca număr de identități în cel mai bun aliniament împărțit la numărul de reziduuri comparate (pozițiile de gap sunt excluse).

Scorurile în ambele metode sunt calculate inițial ca procentaj de scor identitate și sunt convertite în distanțe prin împărțirea cu 100 și apoi scăzând din 1.0. În aceste distanțe inițiale nu se aplică nici o corecție pentru substituții multiple.

În ce privește complexitatea acestui algoritm interesul este orientat spre complexitatea algoritmului de estimare a similarității dintre secvențe. Astfel, conform unui studiu realizat în BMC Bioinformatics [EDGA'04] măsura de similaritate este identitatea fracțională calculată dintr-un aliniament global. Aliniamentul global al unei perechi de secvențe sau profile este calculat folosind algoritmul spațiului linear al lui Myers-Miller [37 din Edgar, Muscle] care este spațiul $O(L)$ și timpul $O(L^2)$ în secvența tipică de lungime L . Fiind date N secvențe și astfel $N(N-1)/2 = O(N^2)$ perechi, este nevoie, prin urmare, de timpul $O(N^2L^2)$ și spațiul $O(N^2+L)$ pentru construirea matricei distanță. Mai departe, operațiile realizate cresc gradul de complexitate în timp și spațiu.

Interfața preluată din programul Clustal W accesibil online pe site-ul oficial al Institutului European de Bioinformatică (<http://www.ebi.ac.uk/clustalw/>), cu toate opțiunile de prelucrare posibile, este prezentată în Figura 6.1-1.

³⁵ Cladogramele sunt un arbore de diagrame folosiți în clasificarea biologică a membrilor din punct de vedere al evoluției pentru a ilustra relațiile filogenetice [WIKI'06a]

³⁶ Filogramele sunt arbore filogenetici care indică relațiile dintre diverse specii de plante sau animale grupate în funcție de relațiile dintre ele și de asemenea exprimă un sens al timpului sau ratei de evoluție. Aspectul temporal al filogramei lipsește dintr-o cladogramă [WIKI'06a]

YOUR EMAIL	ALIGNMENT TITLE	RESULTS	ALIGNMENT	CPU MODE
<input type="text"/>	Sequence	interactive ▾	full ▾	single ▾
KTUP (WORD SIZE)	WINDOW LENGTH	SCORE TYPE	TOPDIAG	PAIRGAP
def ▾	def ▾	percent ▾	def ▾	def ▾
MATRIX	GAP OPEN	END GAPS	GAP EXTENSION	GAP DISTANCES
def ▾	def ▾	def ▾	def ▾	def ▾

OUTPUT		PHYLOGENETIC TREE		
OUTPUT FORMAT	OUTPUT ORDER	TREE TYPE	CORRECT DIST.	IGNORE GAPS
aln w/numbers ▾	aligned ▾	none ▾	off ▾	off ▾

Enter or Paste a set of Sequences in any supported format Help

```

LENLEMLDKDGHKIKITDFGLCKREGIKDGRATEKTFCGIPEYLAPEVLE
DNDYGRAVDWUWGLGVVYEMMCGRLPFYNQDHEKLFELILMEEIRF
PRTLGPPEAKSLLSGLLKKDPKQRLGGGSEDAKEIMQHRFFAGIVWQ
HVYEKKLSPPFKPQVTSETDTRYFDEEFTAQMITITPPDQDSMEC
VDSERRPHFPQFSYSASGTA

>AAH84541.txt
MAV FVLLALVAGVLGNEFSILKSPGSVVF RGNWNP IGERIPDVA
ALSMGFSVKEDLSUPGLAVGNL FHRPRATVMVMVKGVNKLALPPGS
VISYPLENAVPPSLDSVANSIHS LFSEETPVVLQLAPSEERVYMG
KANSVEEDI SVTLDQLDMLLEQEMSVLSSLDLNSLGDRIEVDLLEL

```

Upload a file:

Figura 5.4.3-1 Interfața oferită de programul Clustal W

6.2. Experimente

Cu scopul realizării comparației rezultatelor furnizate de Clustal W și cele ale noii strategii propuse s-a folosit baza de proteine prezentată în subcapitolul 5.4.1. Aceasta răspunde necesității programului Clustal W de a considera un număr mai mic sau egal cu 500 de secvențe la o execuție. Acest tool presupune un set de parametri de intrare și în acest scop s-a decis utilizarea valorilor parametrilor implicați cum ar fi : Protein Gap Open Penalty = 10.0, Protein Gap Extension Penalty = 0.2, Protein matrix = Gonnet, oferiți de European Molecular Biology Laboratory și European Bioinformatics Institute (EMBL-EBI) la adresa oficială de web <http://www.ebi.ac.uk/>

Pentru validare s-a folosit aceeași măsură dată de *precizie*[32] după cum a fost deja prezentată în 5.4.2. Astfel, pentru cele trei seturi de secvențe și metoda de similaritate aplicată de Clustal W au fost obținute valorile prezentate în Tabel 6.2-1, alături de cele obținute prin noua metodă propusă.

Merită menționat faptul că noua strategie algoritmică atinge (în cazul modelului 4-gram al setului 3) performanța ridicată a metodei CLUSTAL W.

Tabel 6.2-1 Valorile **preciziei** obținute pe baza rezultatelor de similaritate furnizate de programul CLUSTAL W sunt în coloana 'CLUST.W' urmate de coloanele corespunzătoare celor două metode noi de similaritate folosind modele 2,3,4-gram pentru celei trei seturi de date.

Set	Clust.W	Metoda Directă			Metoda Alternantă		
		2-gram	3-gram	4-gram	2-gram	3-gram	4-gram
1	0.872	0.439	0.662	0.830	0.471	0.646	0.823
2	0.921	0.446	0.650	0.874	0.439	0.605	0.860
3	0.932	0.534	0.865	0.931	0.574	0.828	0.919

De asemenea, s-a aplicat aceeași metodă de verificare a gradului de corelare dintre grupurile obținute utilizând similaritatea furnizată cu metoda din Clustal W și gruparea pe baza căreia este structurată baza de secvențe SCOP [BOGA'06e]. Astfel, folosind statisticile lui Hubert, rezultatele finale sunt prezentate în Tabel 6.2-2 alături de cele obținute cu noua metodă.

Tabel 6.2-2. Valorile de corelație pentru set1, 2 și 3 ale bazei de secvențe obținute din Astral SCOP folosind statisticile Huberts

Set	Clustal W	Metoda Directă	Metoda Alternantă
1	0.202	0.123	0.111
2	0.127	0.034	0.035
3	0.297	0.334	0.362

6.3. Interpretarea rezultatelor

Compararea rezultatelor obținute cu Clustal W și noua metodă arată că în termeni de performanță a determinării similarității, noua metodă atinge performanța obținută de Clustal W pentru modelul 4-gram. În ce privește rezultatele aplicării tehnicii de corelare, după cum se observă, cele mai bune rezultate sunt obținute pentru setul 3 de secvențe. Astfel, recurgând la o investigație vizuală a secvențelor implicate s-a observat că în setul 1 și 2 multe secvențe sunt de lungime relativ mică și deoarece noua metoda funcționează bine cu precădere pe secvențe lungi, aceasta poate fi o explicație a diferențelor de valori obținute pentru aceste seturi de secvențe. S-a observat că valoarea de adevăr a grupurilor biologice oferite de baza de secvențe SCOP se potrivește strict pentru funcțiile biologice ale proteinelor

acordând mai puțină importanță aspectelor reziduale și lungimii secvențelor de proteine. Cu alte cuvinte, uneori interesul biologic nu este strans legat de toate aspectele care caracterizează secvențele biologice.

6.5. Concluzii

Cele mai evidente concluzii la care s-a ajuns în urma investigațiilor din acest capitol sunt cele ce urmează.

- În compararea dintre noua metodă de similaritate și cea realizată de programul Clustal W, un prim aspect îl reprezintă complexitatea algoritmică. Astfel, dacă în cazul noii metode vorbim de o complexitate maximă de timp situată în jur de $O(N \times M)$ pentru ambele abordări, cu N și M lungimile a două secvențe comparate), în ce privește complexitatea spațiului stocat pentru *metoda directă* se estimează o complexitate $O(N)$, cu N lungimea secvenței procesate cu rol de query, iar pentru *metoda alternantă* complexitatea spațiului poate fi aproximată de $O(\min\{N, M\})$. Așa după cum a fost menționat în cazul algoritmului Clustal W pentru o pereche de secvențe de lungime L s-a estimat o complexitate a timpului de $O(L^2)$ și a spațiului de $O(L)$. Astfel dacă e vorba de S secvențe, complexitatea este calculată după formula $S(S-1)/2 = O(S^2)$ perechi, prin urmare se estimează timpul $O(S^2 L^2)$ și spațiul $O(S^2 + L)$ pentru construirea matricei distanță. În aceleași circumstanțe noua metodă ar presupune o complexitate a timpului $O(S^2 N M)$ iar a spațiului de $O(S^2 + N)$ pentru *metoda directă* și respectiv $O(S^2 + \min\{N, M\})$ pentru *metoda alternantă*. Din acest punct de vedere, este evidentă o ușoară superioritate a noii metode propuse.
- Un alt atu al noii metode este pus în valoare de faptul că în termeni de performanță rezultatele obținute pentru modelul 4-gram sunt foarte apropiate de cele ale Clustal W.
- Considerând simplitatea algoritmică și eficiența computațională a noii metode este justificată sugerarea ei ca primă alegere atunci când se realizează o căutare în baze mari de secvențe.
- Există o motivare întemeiată de a perfecționa această nouă metodă dar în acest moment trebuie remarcat faptul că în ce privește determinarea similarității noua tehnică folosită atinge performanțele metodei CLUSTAL W.

Mențiuni

Această parte de cercetare a fost susținută de proiectul European Biopattern: Computational Intelligence for biopattern analysis in Support of eHealthcare, Network of Excellence Project No. 508803.

7. CONCLUZII FINALE

Progresul rapid în cercetare legat de secvențe genomice/proteomice motivează nevoia de conceptualizare și analiză a acestora care îmbină cunoștințele biologiei moleculare cu cele matematice și ale sistemelor automate de calcul pentru eficientizarea determinării de răspunsuri la diverse întrebări biologice. Avantajele cooperării între aceste domenii ale științei valorificate în această teză sunt:

- + Păstrarea informației biologice sub forma organizată a bazelor de secvențe;
- + Posibilitatea de explorare a unei cantități imense de informații în timp scurt;
- + Extragerea automată de caracteristici pe baza comparației secvențelor, operațiune de altfel dificilă și mare consumatoare de timp prin analiza cu metode biologice sau manuale.

Pentru analiza și compararea secvențelor biologice este necesară folosirea de metode și tehnici algoritmice adecvate. Dintre metodele de analiză abordate cu preponderență în această teză se conturează două tendințe:

- + Metode de aliniament bazate pe programarea dinamică,
- + Metode orientate spre folosirea modelelor Markov ascunse.

Datorită complexității problemei de comparare a secvențelor cu scopul determinării similarității sau al identificării de segmente cu rol biologic important au fost puse în valoare următoarele aspecte:

- + Pentru determinarea similarității sunt necesare funcții pentru evaluarea aliniamentelor în vederea stabilirii distanței dintre secvențe, care implică folosirea matricelor de substituție sau tehnici algoritmice optimizate independente de aliniamentul automat al secvențelor;
- + Aliniamentele multiple folosesc de obicei tehnicile pentru compararea de perechi de secvențe pe care le încorporează în algoritmi de parcurgere și evaluare a tuturor secvențelor implicate în procesul de comparare;
- + Diversitatea de metode implică identificarea unei mulțimi distincte de determinare a similarității, stabilirea celui mai performant bazându-se pe minimizarea timpului și spațiului computațional;
- + Metodele de aliniament și determinare a similarității propuse urmăresc îmbunătățirea continuă a performanțelor sistemului de analiză urmărind performanțe ridicate la costuri cât mai mici;
- + Teza subliniază implicarea estimărilor statistice în evaluarea rapoartelor de asociere pentru ponderea elementelor în matricele de substituție și în

- * funcțiile de evaluare ale aliniamentelor, cu precădere pentru modelele Markov ascunse.

Metodele identificate se constituie în nuclee care stau la baza programelor comerciale destinate comparării, analizei și extragerii automate de informații relevante din baze mari de secvențe.

Tehnicile de recunoaștere au o largă aplicabilitate în domeniul bioinformaticii, iar obiectivul principal al tezei este de a veni cu o nouă soluție la problemele de analiză și căutare în baze de secvențe biologice mari. S-a realizat o identificare a acestor tehnici în funcție de scopul la care servesc și s-a constatat, în concordanță cu alte aprecieri ale specialiștilor, că până în momentul de față nu se cunoaște vreo evaluare a performanțelor comparative a mai multor metode care să decidă asupra superiorității totale a uneia sau alteia. Explicația vine din aspecte ca:

- * Tipul de format al datelor prelucrate. Nu este același lucru a se obține secvențe similare prin realizarea de aliniamente de secvențe ca și prin interogarea lingvistică a unor descrieri textuale despre aceste secvențe;
- + Scopul urmărit de utilizator. Scopul urmărit prin analiza și/sau căutarea de secvențe poate să fie de la unul simplu de a identifica obiectul de interes până la cel de a obține detalii suplimentare despre respectivul „obiect”.

Pornind strict de la aceleași premise și urmărind aceleași obiective se poate concluziona că vor fi alese metodele cu timpi computaționali cât mai reduși, consum minim de resurse și gradul ridicat de acuratețe al rezultatelor.

Un aspect caracteristic tuturor metodelor de comparație și analiză investigate în această teză este gradul ridicat de dependență de estimări de diverși parametri, sau estimări biologice colectate din diverse surse.

- + La metodele de aliniament este evidentă necesitatea implicării matricei de substituție și a evaluărilor de penalități de gap (pentru marcarea absențelor de elemente) sau inserții și ștergeri de amino acizi;
- + Folosind modelele Markov ascunse trebuie avută în vedere calcularea de diverse ponderi de tranziție de la o stare la alta sau extragerea de „profile” din cadrul familiilor de secvențe pentru a putea evalua gradul de apartenență al unei secvențe la o familie, folosind algoritmi bazati pe principiile implementării automatelor finite deterministe (AFD).

Prin analizarea metodelor existente în compararea secvențelor biologice și a experimentelor folosind principiile modelării lingvistice aplicate pentru sistemele de recunoaștere a vorbirii s-a dorit realizarea unui cadru informațional pentru a motiva și a atrage atenția asupra contribuției personale a autoarei la dezvoltarea unei noi metode eficiente de comparație, analiză și căutare în baze de secvențe biologice.

Abordări ale conceptelor folosite în sistemele de recunoaștere a vorbirii sunt întâlnite în literatura biologică de specialitate orientate cu precădere către analiza secvențelor, construirea de profile pentru gruparea lor sau generarea de aliniamente de secvențe.

Din experimentele realizate de autoare, se poate concluziona că modelele lingvistice statistice se pot identifica drept un domeniu de importanță majoră în

analiza și exploatarea unei cantități mari de informație fiind totodată ușor de manipulat și cu rezultate eficiente în aplicarea lor.

Metoda nou propusă în această teză constituie un prim pas în investigarea implicării modelării lingvistice în caracterizarea, manevrarea și înțelegerea datelor biologice sub formă de secvențe.

În mod special s-a studiat folosirea conceptului de entropie încrucișată (cross entropie) aplicat modelelor n -gram cu scopul căutării eficiente în baze de date. Rezultatele experimentale au indicat soliditatea noii strategii algoritmice în exprimarea similarității dintre proteine.

Fiind dată simplitatea conceptuală a acestei noi abordări, ea apare ca o alternativă de dorit tehnicilor studiate și utilizate până în prezent. Pe baza analizei rezultatelor experimentale referitoare la eficiența metodei propuse se evidențiază aspectele principale:

- + Considerând problema generală de separare dintre măsurile de similaritate "globală" și "locală" dintre proteine, noua abordare aparține primei categorii;
- + Conceptul de n -gram, folosit în lucrări anterioare în domeniul bioinformaticii, pentru abordarea actuală are meritul de a fi prima încercare de a adopta acest pas dual în compararea secvențelor biologice.
- + Rezultatele obținute în urma aplicării celor două variante ale noi metode (metoda directă și metoda alternantă) în ambele experimente relevă un grad ridicat de sensibilitate la lungimea secvențelor.
 - Dacă în cazul bazei de secvențe cu mutații ale genei P53 metoda alternantă furnizează cele mai bune rezultate, în cazul bazei de date structurate derivate din SCOP performanța lor este oarecum apropiată.
 - În cazul excepțional când toate secvențele comparate ar avea aceeași lungime, *metoda directă* este echivalentă cu *metoda alternantă* și se comportă foarte bine.
- + Pentru ordinul modelului n -gram folosit, după testarea ordinilor 2,3,4,5 în ambele experimente, s-a observat că performanța metodei crește odată cu ordinul modelului până la 4. După ordinul 5, datorită lipsei datelor, estimatorii probabilității maxime devin nerezonabil de uniformi și mici.
 - Aceasta stabilește o limită superioară pentru modelul ales în această bază de secvențe (este sugerată observarea existenței unor grupări de 4 amino acizi cu importanță biologică ce pot participa în mod decisiv la caracterizarea secvențelor de proteine).
- + Rezultatele furnizate de reprezentările grafice ale experimentelor pentru detectarea similarității secvențelor în formarea de clusteri este evidentă dar

cu precizarea că această grupare se bazează pe similaritatea compoziției textuale.

- Intervenind cu raționamente biologice se poate obține mult mai mult decât o clasificare textuală.
- + Metoda propusă folosește tehnici statistice de evaluare dar ea poate fi dezvoltată și îmbunătățită prin incorporarea de cunoștințe biologice cum ar fi lucrul cu grupe funcționale de amino acizi și impunerea unor ajustări în cazul evenimentelor singulare, purtătoare de valoarea zero în momentul evaluării.

Ca o remarcă finală, această metodă oferă un mod eficient de a captura caracteristicile comune ale secvențelor comparate și în același timp de a evita misiunea de a alege parametri, funcții suplimentare sau metode de evaluare. Performanța ridicată și caracterul facil de implementare, împreună cu eficiența algoritmică fac din această metodă o alternativă promițătoare la măsurile sofisticate, bine cunoscute, utilizate în compararea secvențelor de proteine.

Compararea performanțelor noii metode de determinare a similarității propuse și Clustal W, unul dintre cele mai populare programe de aliniament multiplu de secvențe, relevă:

- + Din punct de vedere al complexității algoritmice vorbim de o complexitate maximă de timp situată în jur de $O(N \times M)$ pentru ambele abordări, cu N și M lungimile a două secvențe comparate.
- + În ce privește complexitatea spațiului stocat:
 - pentru *metoda directă* se estimează o complexitate $O(N)$, cu N lungimea secvenței procesate cu rol de query (referință în căutare),
 - iar pentru *metoda alternantă* complexitatea spațiului poate fi aproximată de $O(\min\{N, M\})$.
- + Așa după cum a fost menționat în cazul algoritmului Clustal W, pentru o singură secvență de lungime L s-a estimat o complexitate a timpului de $O(L^2)$ și a spațiului de $O(L)$. Dacă e vorba de N secvențe, se observă o superioritate a noii metode propuse.

Un alt avantaj al noii metode este pus în valoare de faptul că în termeni de performanță rezultatele obținute pentru modelul 4-gram sunt foarte apropiate de cele ale Clustal W.

Datorită simplității algoritmice și eficienței sistemului automat de determinare a similarității este justificată sugerarea ei ca primă alegere atunci când se realizează o căutare în baze mari de secvențe.

Există o motivare întemeiată de a perfecționa această nouă metodă propusă, ea fiind în acest format o primă variantă de identificare a similarității folosind într-o mare măsură tehnici de evaluare statistică.

Contribuții personale

Partea de contribuții personale ale autoarei este structurată pe capitole pentru a facilita identificarea lor rapidă în conținutul tezei.

În capitolul 3 s-a realizat o investigație a tehnicilor de determinare automată a similarității secvențelor biologice propuse relativ recent. Ele vin ca alternativă la metodele de aliniament fundamentale studiate și folosite în prezent și oferă o imagine de ansamblu asupra noilor alternative propuse/experimentate pentru determinarea similarității secvențelor (reprezentate ca structuri primare în cazul proteinelor).

În capitolul 4 aportul de contribuții personale se evidențiază în cadrul experimentelor care pun în valoare potențialul modelelor lingvistice statistice în analiza textuală. Se pornește de la implementarea programelor de analiză folosite, continuă cu perspectiva din care este abordată analiza modelelor lingvistice statistice folosite în sistemele automate de recunoaștere a vorbirii și se finalizează în interpretările și concluziile ce vor deveni bază de pornire în propunerea noii măsuri de evaluare a contextului secvențelor biologice.

O prezentare sintetică a acestor aspecte este următoarea:

- ÷ Dezvoltarea aplicației în limbajul Perl 5.0 care permite o analiză flexibilă a modelelor lingvistice statistice;
- + Experimente cu modele lingvistice de dimensiuni 2, 3, ..., 9, 12, 15, 18, 20 care depășesc evaluările uzuale limitate de obicei la gradul de 4 și foarte rar 5.
- ÷ Se vine cu propunerea unui indicator al evaluării gradului de acuratețe al estimării și anume, valoarea cross entropiei obținute când setul de testare este acoperit în mod strict.
- + Este realizat un studiu asupra influenței evenimentelor n -gram care nu participă la evaluarea entropică. Astfel, atenția a fost direcționată spre frecvența apariției acestor evenimente constatând că evoluția indicatorilor de frecvență poate oferi o referință asupra conținutului setului de testare.

În capitolul 5, atât noua metodologie propusă pentru determinarea similarității dintre secvențele biologice, discutată pe larg în capitolele corespunzătoare, cât și implementarea ei întregesc partea de contribuții personale aduse la elaborarea acestei teze de doctorat. În cele ce urmează sunt remarcate aspectele principale ale acestor contribuții:

- ÷ Un fapt evident este transferul de cunoștințe din domeniul sistemelor automate de recunoaștere în cel al biologiei moleculare.
- + Tehnica aplicată în modelarea lingvistică pentru sistemele automate de recunoaștere a vorbirii a suferit modificări în sensul restrângerii numărului de evenimente considerat. Mai precis de la cel posibil, oferit de toate combinațiile de lungimea n -gram ce ar putea fi întâlnite, la cel al evenimentelor existente în secvențele curente analizate.

- + Realizarea profilului de referință al secvenței de căutare, adică a *modelului perfect* pe baza căruia se estimează distanța dintre secvențele comparate este un concept nou care fixează un nivel al așteptărilor distribuției datelor (conceptualizat ca scor perfect) pentru secvența cu care se face comparația.
- + Determinarea unei formule de calcul al similarității dintre două secvențe pe baza evaluării statistice a conținutului secvențelor. Distanța dintre scorul perfect al secvenței de query și cel obținut prin proiecția acesteia asupra celei cu care se compară furnizează gradul de disimilaritate dintre cele două secvențe.
- + Abordarea *metodei alternante* care vine să evidențieze nevoia de simetrie în aplicarea metodei, fapt exploatat însă în evidențierea diferenței de lungime a secvențelor. Din această perspectivă se poate considera că noua metodă de similaritate realizează o "proiecție" a secvenței de query asupra secvenței cu care este comparată.
- + Rezultatele obținute în cazul comparației unei secvențe de căutare cu mai multe secvențe existente într-o bază de date pot să furnizeze și un criteriu de ordonare al acestora în funcție de valoarea scorului de similaritate/disimilaritate. Prin urmare, aceasta poate fi considerată și o metodă de data mining.

O recunoaștere a caracterului inovativ propus prin aceasta nouă metodă este susținută de publicațiile naționale și internaționale în care este prezentată și analizată sub diverse aspecte (i.e. metodă de similaritate pentru secvențe biologice, metodă de identificarea rapidă a secvențelor mutante, folosirea metodei în procesul de grupare de secvențe, pentru data mining, în compararea de secvențe text în general).

În capitolul 6 se realizează o analiză și evaluare comparativă a performanțelor metodei propuse cu cele ale programului Clustal W. Aceasta este realizată din punct de vedere al preciziei valorilor experimentale și a complexității algoritmului care stă la baza generării valorilor de similaritate.

Destinată domeniului medical și în special biologiei meculare, teza are un caracter esențialmente interdisciplinar combinând cunoștințe de: biologie, automatizări, matematică, informatică, calculatoare. Încadrarea ei în domeniul automatizării este susținută de scopul pe care îl deservește independent de asistența omului („*automatos*” în limba greacă înseamnă „functionare de la sine”) precum și de suportul teoretic al acestei specializări. Deoarece automatica este preocupată de automatizarea celor mai diferite (prin natura lor) procese (tehnice, energetice, militare etc.), tot astfel și preocupările cuprinse în această teză – analiza automată a similarității secvențelor biologice – urmărește realizarea automată a scopului propus, ceea ce o poate încadra în domeniul preocupărilor specializării de „Automatică”.

Din punct de vedere al perspectivelor în dezvoltarea noii metode propuse și analizate în această teză autoarea menționează:

-
- ÷ Se poate considera că această nouă abordare a problematicii comparației secvențelor biologice este în esență valabilă pentru orice alt tip de secvențe textuale. Desigur că există posibilitatea de a-i crește gradul de performanță intervenind în abordarea statistică folosită.
 - + Ușor de remarcat este sensibilitatea la lungimea secvențelor comparate, aspect care poate fi folosit ca indicator al gradului de divergență în reprezentarea secvențelor.
 - ÷ Gruparea fiind o acțiune realizată tot pe baza determinării de relații de similaritate/disimilaritate dintre obiecte poate beneficia de o nouă strategie de determinare a similarității pentru cazul obiectelor reprezentate textual.
 - + Elaborarea unui produs IT performant folosind Oracle RDBMS (Relational DataBase Management System), bazat pe crearea unei structuri relaționale poate concretiza utilitatea noii metode în practică. Acesta va pune în valoare atât relațiile de similaritate dintre secvențele biologice cât și alte caracteristici ale secvențelor implicate, venind semnificativ în sprijinul cercetării bio-medicale.

Astfel, se conturează un potențial semnificativ de folosire și dezvoltare a acestei noi metode propuse în cazuri greu accesibile strategiilor existente.

ANEXA 1

Codurile acceptate de acizi nucleici sunt:

Nucleic Acid Code	Meaning
A	A denosine
C	C ytidine
G	G uanine
T	T hymidine
U	U racil
R	G A (pu R ine)
Y	T C (p Y rimidine)
K	G T (K etone)
M	A C (a M ino group)
S	G C (S trong interaction)
W	A T (W eak interaction)
B	G T C (not A) (B comes after A)
D	G A T (not C) (D comes after C)
H	A C T (not G) (H comes after G)
V	G C A (not T, not U) (V comes after U)
N	A G C T (a N y)
-	gap of indeterminate length

Codurile acceptate de amino acizi sunt:

Amino Acid Code	Meaning
A	Alanine
B	Aspartic acid or Asparagine
C	Cysteine
D	Aspartate
E	Glutamate
F	Phenylalanine
G	Glycine
H	Histidine
I	Isoleucine
K	Lysine
L	Leucine

M	Methionine
N	Asparagine
P	Proline
Q	Glutamine
R	Arginine
S	Serine
T	Threonine
U	Selenocysteine
V	Valine
W	Tryptophan
Y	Tyrosine
Z	Glutamate or Glutamine
X	Any
*	translation stop
-	gap de of indeterminate length

- **ALN/ClustalW format**

CLUSTAL W (1.82) multiple sequence alignment

```
FOSB_MOUSE
MFQAFPGDYDSGSRCS SSPSAESQYLSSVDSFGSPPTAAASQECAGLGEMPGSFVPTVTA 60
FOSB_HUMAN
MFQAFPGDYDSGSRCS SSPSAESQYLSSVDSFGSPPTAAASQECAGLGEMPGSFVPTVTA 60
```

```
FOSB_MOUSE
ITTSQDLQWLVQPTLISSMAQSQGQPLASQPPAVDPYDMPGTSYSTPGLSAYSTGGASGS
120
FOSB_HUMAN
ITTSQDLQWLVQPTLISSMAQSQGQPLASQPPVDPYDMPGTSYSTPGMSGYSSGGASGS
120
```

*****.*****.*.***.*****

- **AMPS Block file format**

```
>Seq_1
>A0231
>HAHU
>Four_Alpha
>Globin
>GLobin_C
*
```

ARNDLQ
AAAAAA
PPPPPP
PP PPP
WW WWW
LLLLLL
IIVVLL
*

- **Codata**

```
ENTRY      IXI_234
SEQUENCE
      5      10      15      20      25      30
  1 TSPASIRPPAGPSSRPAMVSSRRTRPSPPG
  31 PRRPTGRPCCSAAPRRPQATGGWKTCSGTC
  61 TTSTSTRHRGRSGWSARTTTAACLRASKS
  91 MRAACSR SAGSRPNRFAP TLMSSCITSTTG
 121 PPAWAGDRSHE
```

///

```
ENTRY      IXI_235
SEQUENCE
      5      10      15      20      25      30
  1 TSPASIRPPAGPSSR-----RPSPPG
  31 PRRPTGRPCCSAAPRRPQATGGWKTCSGTC
  61 TTSTSTRHRGRSGW-----RASRKS
  91 MRAACSR SAGSRPNRFAP TLMSSCITSTTG
 121 PPAWAGDRSHE
```

///

- **EMBL**

```
ID  MMFOSB   standard; RNA; MUS; 4145 BP.
XX
AC  X14897;
XX
SV  X14897.1
XX
DT  23-NOV-1989 (Rel. 21, Created)
DT  12-SEP-1993 (Rel. 36, Last updated, Version 2)
XX
DE  Mouse fosB mRNA
XX
KW  fos cellular oncogene; fosB oncogene; oncogene.
XX
OS  Mus musculus (house mouse)
OC  Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
Mammalia;
OC  Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Mus.
XX
```

RN [1]
 RP 1-4145
 RX MEDLINE; 89251612.
 RA Zerial M., Toschi L., Ryseck R.P., Schuermann M., Mueller R., Bravo R.;
 RT "The product of a novel growth factor activated gene, fos B, interacts with
 RT JUN proteins enhancing their DNA binding activity";
 RL EMBO J. 8:805-813(1989).
 XX
 DR MGD; MGI:95575; Fosb.
 DR SWISS-PROT; P13346; FOSB_MOUSE.
 DR TRANSFAC; T00291; T00291.
 XX
 CC clone=AC113-1; cell line=NIH3T3;
 XX
 FH Key Location/Qualifiers
 FH
 FT source 1..4145
 FT /db_xref="taxon:10090"
 FT /organism="Mus musculus"
 FT CDS 1202..2218
 FT /db_xref="SWISS-PROT:P13346"
 FT /note="fosB protein (AA 1-338)"
 FT /protein_id="CAA33026.1"
 FT
 /translation="MFQAFPGDYDSGSRCSSSPSAESQYLSSVDSFGSPPTAAASQECA
 FT
 GLGEMPGSFVPTVTAITTSQDLQWLQPTLISSMAQSQGQPLASQPPAVDPYDMPGTSY
 FT
 STPGLSAYSTGGASGSGGPSTSTTTSGPVSARPARARPRRPREETLTPEEEEEKRRVRRE
 FT
 RNKLAALKCRNRRRELTDRLQAETDQLEEEKAELESEIAELQKEKERLEFVLVAHKPGC
 FT
 KIPYEEGPGGPLAEVRDLPGSTSAKEDGFGWLLPPPPPPPLPFQSSRDAPPNLTASLF
 FT
 THSEVQVLGDPFPVSPSYTSSFVLTCEVSAFAGAQRRTSGSEQPSDPLNSPLLAL"
 XX
 SQ Sequence 4145 BP; 960 A; 1186 C; 1007 G; 991 T; 1 other;

```

ataaattctt atttgacac tcacaaaat agtcacctgg aaaaccgct tttgtgaca 60
aagtacagaa ggcttggtca catttaaact actgagaact agagagaaat actatcgcaa 120
actgtaatag acattacatc cataaaagt tccccagtcc ttattgtaat attgcacagt 180
gcaattgcta catggcaaac tagttagca tagaagtcaa agcaaaaaca aaccaaagaa 240
aggagccaca agagtaaac tgttcaacag ttaatagttc aaactaagcc attgaatcta 300
tcattgggat cgtaaaatg aatcttcta cacctgcag tgtatgattt aactttaca 360
gaacacaagc caagttaaa atcagcagta gagatattaa aatgaaaagg tttgctaata 420
gagtaacatt aaatacctg aaggaaaaaa aacctaaata tcaaaataac tgattaaat 480
tcacttgcaa attagcacac gaatatgcaa cttggaatc atgcagtgtt ttattaaga 540

```

```

ctctcaatga ctctgatctc cggntngtct gtaattctg gattgtcgg ggacatgcaa 4080
tttacttct gtaagtaagt gtgactgggt ggtagatctt ttacaatcta tatcggtgag 4140
aattc 4145
//

```

- **GCG/MSF**

```

DNA_MULTIPLE_ALIGNMENT 1.0
Four anthropoidea
MSF: 50 Type: N Check: 2666 ..
//

```

```

Name: Homo_sapiens Len: 50 Check: 8318 Weight: 1.00
Name: Pan_paniscus Len: 50 Check: 7854 Weight: 1.00
Name: Gorilla_gorilla Len: 50 Check: 7778 Weight: 1.00
Name: Pongo_pigmaeus Len: 50 Check: 8716 Weight: 1.00
//

```

```

Homo_sapiens AGUCGAGUC...GCAGAAAC
Pan_paniscus AGUCGCGUCG..GCAGAAAC
Gorilla_gorilla AGUCGCGUCG..GCAGAUAC
Pongo_pigmaeus AGUCGCGUCGAAGCAGA..C
Homo_sapiens GCAUGAC.GACCACAUUUU.
Pan_paniscus GCAUGACGGACCACAUCAU.
Gorilla_gorilla GCAUCACGGAC.ACAUCAUC
Pongo_pigmaeus GCAUGACGGACCACAUCAUC
Homo_sapiens CCUUGCAAAG
Pan_paniscus CCUUGCAAAG
Gorilla_gorilla CCUCGCAGAG
Pongo_pigmaeus CCUUGCAGAG

```

- **GDE**

```

{
name "Short name for sequence"
longname "Long (more descriptive) name for sequence"
sequence-ID "Unique ID number"
creation-date "mm/dd/yy hh:mm:ss"
direction [-1|1]
strandedness [1|2]
type [DNA|RNA|PROTEIN|TEXT|MASK]
offset (-999999,999999)
group-ID (0,999)
creator "Author's name"
descrip "Verbose description"
comments "Lines of comments that can be fairly arbitrary text about a
sequence. Return characters are allowed, but no internal double quotes
or brace characters. Remember to close with a double quote"
sequence "gctagctagctagctagctcttagctgtagctgtagctgatgctagct
gatgctagctagctagctagctgatcgatgctagctgatcgtagctgacg"

```



```
gactgatgctagctagctagctagctgtctagtgtcgtagtgcttattgc"
}
```

• **Genebank**

```
LOCUS      MMFOSB          4145 bp  mRNA  linear  ROD 12-SEP-1993
DEFINITION Mouse fosB mRNA.
ACCESSION  X14897
VERSION    X14897.1  GI:50991
KEYWORDS   fos cellular oncogene; fosB oncogene; oncogene.
SOURCE     Mus musculus.
  ORGANISM Mus musculus
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Mus.
REFERENCE  1 (bases 1 to 4145)
AUTHORS    Zerial,M., Toschi,L., Ryseck,R.P., Schuermann,M., Muller,R. and
            Bravo,R.
TITLE      The product of a novel growth factor activated gene, fos B,
            interacts with JUN proteins enhancing their DNA binding activity
JOURNAL    EMBO J. 8 (3), 805-813 (1989)
MEDLINE    89251612
PUBMED     2498083
COMMENT    clone=AC113-1; cell line=NIH3T3.
FEATURES   Location/Qualifiers
  source    1..4145
            /organism="Mus musculus"
            /db_xref="taxon:10090"
  CDS       1202..2218
            /note="fosB protein (AA 1-338)"
            /codon_start=1
            /protein_id="CAA33026.1"
            /db_xref="GI:50992"
            /db_xref="MGD:95575"
            /db_xref="SWISS-PROT:P13346"

/translation="MFQAFPGDYDSGSRCS SSPSAESQYLSSVDSFGSPPTAAASQEC
AGLGEMPGSFVPTVTAITTSQDLQWLQVQPTLISSMAQSQQQPLASQPPAVDPYDMPGT
SYSTPGLSAYSTGGASGSGGPSTSTTTSGPVSARPARARPRRPREETLTPEEEEEKRRV
RRERNKLAALKCRNRRRELTDRQLQAETDQLEEEKAELESEIAELQKEKERLEFVLVAH
KPGCKIPYEEGPGGPLAEVRDLPGSTSAKEDGFGWLLPPPPPPPLPFQSSRDAPPNL
TASLFTHSEVQVLGDPFPVSPSYTSSFLTCPEVSAFAGAQR TSGSEQPSDPLNSPS
LLAL"
BASE COUNT  960 a  1186 c  1007 g  991 t   1 others
ORIGIN
1 ataaattctt attttgacac tcacccaaat agtcacctgg aaaaccgct tttgtgaca
```

```

61 aagtacagaa ggcttggtca cattaaatc actgagaact agagagaaat actatcгаа
121 actgtaatag acattacatc cataaaagtt tccccagtc ttattgtaat attgcacagt
181 gcaattgcta catggcaaac tagttagca tagaagtaa agcaaaaaca aaccaaagaa
241 aggagccaca agagtaaaac tgttcaacag ttaatagttc aaactaagcc attgaatcta
.
.
.
1021 gcagagggaa cttgcatcga aacttgggca gttctccгаа cgggagacta agttccccg
1081 agcagcgcac tttggagacg tgtccggtct actccggact cgcatctcat tccactcggc
1141 catagccttg gcttccccgc gacctagcg tggtcacagg ggccccctg tgcccagggg
1201 aatgtttcaa gcttttcccg gagactacga ctccgggtcc cgggtgtagct catcacctc
1261 cgccgagtct cagtacctgt

```

- **NBRF/PIR**

```

>P1;FOSB_MOUSEFOSB_MOUSE 338 bases
MFQAFPGDYD SGSRCSSSPS AESQYLSSVD SFGSPPTAAA SQECAGLGEM
PGSFVPTVTA ITTSQDLQWL VQPTLISSMA QSQGQPLASQ PPAVDPYDMP
GTSYSTPGLS AYSTGGASGS GGPSTSTTTS GPVSARPARA RPRRPREETL
TPEEEKRRV RRERNKLAAL KCRNRRRELT DRLQAETDQL EEEKAELESE
IAELQKEKER LEFVLVAHKP GCKIPYEEGP GPGPLAEVRD LPGSTSAKED
GFGWLLPPPP PPPLPFQSSR DAPPNLASL FTHSEVQVLG DPFVPSY
TSSFVLTCPV VSAFAGAQRV SGSEQPSDPL NSPSSLAL*

```

- **PDB format**

```

COMPND  MOL_ID: 1;
COMPND  2 MOLECULE: S-ADENOSYLMETHIONINE SYNTHETASE;
COMPND  3 CHAIN: A, B;
COMPND  4 SYNONYM: MAT, ATP\;L-METHIONINE S-ADENOSYLTRANSFERASE;
COMPND  5 EC: 2.5.1.6;
COMPND  6 ENGINEERED: YES;
COMPND  7 BIOLOGICAL_UNIT: TETRAMER;
COMPND  8 OTHER_DETAILS: TETRAGONAL MODIFICATION

```

- **Pfam/Stockholm format**

```

# STOCKHOLM 1.0
#=GF ID CBS
#=GF AC PF00571
#=GF DE CBS domain
#=GF AU Bateman A
#=GF CC CBS domains are small intracellular modules mostly found
#=GF CC in 2 or four copies within a protein.
#=GF SQ 67
#=GS O31698/18-71 AC O31698

```

```

#=GS O83071/192-246 AC O83071
#=GS O83071/259-312 AC O83071
#=GS O31698/88-139 AC O31698
#=GS O31698/88-139 OS Bacillus subtilis
O83071/192-246      MTCRAQLIAVPRASSLAE..AIACAQKM....RVSrvpvyers
#=GR O83071/192-246 SA 999887756453524252..55152525....36463774777
O83071/259-312      MQHVSAPVVFVFECTRLAY..VQHKLRAH....SRAVAIVLDEY
#=GR O83071/259-312 SS CCCCCHHHHHHHHHHHHHHH..EEEEEEEE....EEEEEEEEEE
O31698/18-71        MIEADKVAHVQVGNLEH..ALLVLTKT....GYTAIPVLDPS
#=GR O31698/18-71 SS CCCHHHHHHHHHHHHHHHHH..EEEEEEEE....EEEEEEEEHHH
O31698/88-139       EVMLTDIPRLHINDPIMK..GFGMVINN.....GFVcVende
#=GR O31698/88-139 SS CCCCCCHHHHHHHHHHHHH..HEEEEEEE....EEEEEEEEEEH
#=GC SS_cons        CCCCCHHHHHHHHHHHHHHH..EEEEEEEE....EEEEEEEEEEH
O31699/88-139       EVMLTDIPRLHINDPIMK..GFGMVINN.....GFVcVende
#=GR O31699/88-139 AS _____*_____
#=GR_O31699/88-139_IN _____1_____2_____0_____
//

```

Phylip

```

1 338 I
FOSB_MOUSE MFQAFPGDYD SGSRCSSSPS AESQYLSSVD SFGSPPTAAA
SQECAGLGEM
PGSFVPTVTA ITTSQDLQWL VQPTLISSMA QSQGQPLASQ PPAVDPYDMP
GTSYSTPGLS AYSTGGASGS GGPSTSTTTS GPVSARPARA RPRRPREETL
TPEEEEKRRV RRERNKLAAA KCRNRRRELT DRLQAETDQL EEEKAELESE
IAELQKEKER LEFVLVAHKP GCKIPYEEGP GPGPLAEVRD LPGSTSAKED
GFGWLLPPP PPPLPFQSSR DAPPNLASL FTHSEVQVLG DPFPVVSPSY
TSSFVLTCPV VSAFAGAQR TSGSEQSDPL NSPSSLAL

```

Raw

```

ataaattcttatttgacactcacaaaatagtcacctggaaaacccgctttttgtgaca
aagtacagaaggcttggtcacatttaaatcactgagaactagagagaatactatcgcaa
actgtaatagacattacatccataaaagtttcccagtccttattgtaatattgcacagt
gcaattgctacatggcaaaactagtgtagcatagaagtcaaagcaaaaacaaacaaagaa
aggagccacaagagtaaaactgttcaacagttaatagttcaaactaagccattgaatcta
tcattgggatcgtaaataaatgaatcttctcacctgcagtgtagtattactttaca

```

RSF(Rich Sequence Format)

```

!!RICH_SEQUENCE 1.0
..
{
name chkhba
type DNA
longname chkhba
checksum 980
creation-date 4/15/98 16:42:47

```

```
strand 1
sequence
  ACACAGAGGTGCAACCATGGTGCTGTCCGCTGCTGACAAGAACAACGTCAAGGGCATCTT
  CACCAAATCGCCGGCCATGCTGAGGAGTATGGCGCCGAGACCTTGGAAAGGATGTTTCCAC
  CACCTACCCCCCAACCAAGACCTACTTCCCCCACTTCGATCTGTACACGGCTCCGCTCA
  ...
}
{
name davagl
type DNA
longname davagl
checksum 7399
creation-date 4/15/98 16:42:47
strand 1
sequence
  GTGCTCTCGGATGCTGACAAGACTCACGTGAAAGCCATCTGGGGTAAGGTGGGAGGCCAC
  GCCGGTGCCTACGCAGCTGAAGCTCTTGCCAGAACCTTCTCTCCTTCCCCACTACCAAA
  ...
}
```

- **UniProtKB/Swiss-Prot**

```
ID FOSB_MOUSE STANDARD; PRT; 338 AA.
AC P13346;
DT 01-JAN-1990 (Rel. 13, Created)
DT 01-JAN-1990 (Rel. 13, Last sequence update)
DT 15-JUN-2002 (Rel. 41, Last annotation update)
DE Protein fosB.
GN FOSB.
OS Mus musculus (Mouse).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Mus.
OX NCBI_TaxID=10090;
RN [1]
RP SEQUENCE FROM N.A.
RX MEDLINE=89251612; PubMed=2498083;
RA Zerial M., Toschi L., Ryseck R.-P., Schuermann M., Mueller R.,
RA Bravo R.;
RT "The product of a novel growth factor activated gene, fos B, interacts
RT with JUN proteins enhancing their DNA binding activity.";
RL EMBO J. 8:805-813(1989).
RN [2]
RP SEQUENCE FROM N.A.
RX MEDLINE=92158623; PubMed=1741260;
RA Lazo P.S., Dorfman K., Noguchi T., Mattei M.-G., Bravo R.;
RT "Structure and mapping of the fosB gene. FosB downregulates the
```

RT activity of the fosB promoter.";
RL Nucleic Acids Res. 20:343-350(1992).
CC -!- FUNCTION: FOSB INTERACTS WITH JUN PROTEINS ENHANCING THEIR DNA
CC BINDING ACTIVITY.
CC -!- SUBUNIT: HETERODIMER (BY SIMILARITY).
CC -!- SUBCELLULAR LOCATION: NUCLEAR.
CC -!- INDUCTION: BY GROWTH FACTORS.
CC -!- SIMILARITY: BELONGS TO THE BZIP FAMILY. FOS SUBFAMILY.
CC -----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See <http://www.isb-sib.ch/announce/>
CC or send an email to license@isb-sib.ch).
CC -----
DR EMBL; X14897; CAA33026.1; -.
DR EMBL; AF093624; AAD13196.1; -.
DR PIR; S04108; TVMSFB.
DR PIR; S35477; S35477.
DR HSSP; P01100; 1FOS.
DR TRANSFAC; T00291; -.
DR MGD; MGI:95575; Fosb.
DR InterPro; IPR000837; Leuzip_Fos.
DR InterPro; IPR004827; TF_bZIP.
DR Pfam; PF00170; bZIP; 1.
DR PRINTS; PR00042; LEUZIPPRFOS.
DR SMART; SM00338; BRLZ; 1.
DR PROSITE; PS00036; BZIP_BASIC; 1.
KW Nuclear protein; DNA-binding.
FT DNA_BIND 161 179 BASIC MOTIF.
FT DOMAIN 183 211 LEUCINE-ZIPPER.
SQ SEQUENCE 338 AA; 35976 MW; E9D031A4BEAE48EC CRC64;
MFQAFPGDYD SGSRCSSSPS AESQYLSSVD SFGSPPTAAA SQECAGLGEM PGSFVPTVTA
ITTSQDLQWL VQPTLISSMA QSQQGPLASQ PPAVDPYDMP GTSYSTPGLS AYSTGGASGS
GGPSTSTTTS GPVSARPARA RPRRPREETL TPEEEKRRV RRERNKLA AA KCRNRRRELT
DRLQAETDQL EEEKAELESE IAELQKEKER LEFVLVAHKP GCKIPYEEGP GPGPLAEVRD
LPGSTSAKED GFGWLLPPPP PPPLPFQSSR DAPPNLASL FTHSEVQVLG DPFVPSY
TSSFVLTCPE VSAFAGAQRT SGSEQPSDPL NSPSLLAL
//

ANEXA 2

Aliniament de secvențe³⁷

Algoritmul de aliniament al secvențelor folosind o abordare a programării dinamice încearcă întotdeauna să urmeze cel mai bun rezultat până la un moment dat.

Se încearcă alinierea a două secvențe prin inserarea unor gap-uri la diferite locații necunoscute (i.e. obținerea a două noi secvențe cu goluri astfel încât să fie maximizat scorul acestui aliniament). Măsurarea scorului poate fi autodefinit. El este determinat prin "**premiul de potrivire**", "**penalitatea de nepotrivire**" și "**penalitatea de gap**". Cu cât scorul este mai mare cu atât este mai bun aliniamentul. Dacă ambele penalități sunt setate la 0 scopul este de a găsi întotdeauna un aliniament cu **potriviri maxime** până la un moment. **Potrivire maximă** = cel mai mare număr de potriviri pentru o secvență care se poate obține prin permiterea tuturor ștergerilor posibile unei alte secvențe.

Se obișnuiește compararea a două secvențe de ADN sau Proteine pentru predicția de similaritate a funcționalității lor.

Există trei algoritmi fundamentali larg răspândiți: **Needleman-Wunsch (1970)**, **Sellers (1974)**, **Smith-Waterman (1981)**

Aliniamentul global: primii doi algoritmi sunt folosiți pentru aliniamentul global.

Se presupune că se aliniază și se potrivesc cele două secvențe de la stânga la dreapta. Fiecare celulă din matricea de scoruri denotă un sub-aliniament între cele două subsecvențe care formează o sub-matrice. În acest moment există trei moduri de a înainta, "*inserează un gap secvenței X*"-- mergând în jos în matrice, "*inserarea unui gap în secvența Y*"-- mergând la dreapta în matrice, "*încearcă o potrivire*" -- ambele secvențe X și Y merg mai departe în mod diagonal. Pentru același motiv, pentru fiecare celulă în matrice, există trei sub-matrice, (i.e. trei perechi de subsecvențe conduc la acest sub-aliniament). Alegând cel mai bun scor dintre aceste trei căi, se obține scorul optimal pentru sub-matricea curentă.

³⁷ Conform (<http://www2.cs.uh.edu/~zhzhzhao/Review/alignment.htm>) University of Houston, Texas

Needleman Wunsch algorithm with gap

	A	B	C	N
A	2	1	0	-1
J	1	1	0	-1
C	0	0	3	2

scored and just for (N&C) Note

Target

ABC&AJC

	A	B	C	N
A	2	1	0	-1
J	1	1	0	-1
C	0	0	3	2

ABCN&AJ

	A	B	C	N
A				
J				
C	0	0	3	2

ABCN&AJC

Target Cell_Score

current match (one alignment)

ABC
A__

ABC
AJC

after match

ABC_N
A__C

ABCN
AJC_

So: $M(i,j) = M(i-1,j-1) + S(i,j)$
 $M(i,j)$: for ABCN & AJ_C
 $M(i-1,j-1)$: for ABC & AJ_
 $M(i,j) = M(i-1,j-1) + Compare(N&C)$
 while $Compare(N&C) = S(i,j)$

So: $M(i,j) = M(i,j) - Gap$
 instead of $M(i,j) = M(i,j) - Gap + S(i,j)$
 Because: $M(i,j)$ is the score for ABCN_ & AJ__C
 $M(i,j) = M(i,j) - Compare(N&C)$
 while $Compare(_&C) = -Gap$, or $W(r)$ --a gap function

So: from 'ABC & A' to 'ABCN & A__C' essentially, it follows 'ABC & A__' and 'ABC_ & A__J' and 'ABC_N & A__JC' i.e. always go down to $M(i-1,j-1)$, then step forward together. So as to the other jumps. Gap penalty will be applied in the above way. i.e. r in $(1,j-1)$, s in $(1,i-1)$. $W(r)$ and $V(s)$ are penalty functions

make more sense: A & AJC have score 2-1-1=0, instead of -1

Penalitatea de gap:

O funcție generală de penalitate gap are forma $W(k) = C_{open} + C_{length} * k$, unde k este lungimea gap-ului, C_{open} este constanta de penalitate de deschidere a gap-ului (gap-open penalty constant), și C_{length} constanta de penalitate a lungimii de gap (gap-length penalty constant), de asemenea, uneori C_{end} va fi aplicată.

Pot să existe diferite formule pentru calculul scorului.

Problema:

Pentru alinierea a două secvențe: ABCNJRQCLCRPM, AJCJNRCKCRBP

Algoritmul Needleman Wunsch : $M[i,j]=\text{Max}\{S[i,j]+M[i-1,j-1],M[i-1,j-k],M[i-r,j-1]\}$, nu poate face față cazului cu gap.

Condiție:

Fiind dat: Premiul de potrivire: $P_a=1$ puncte, nici o penalitate pentru nepotrivire și gap.

Proces:

Pas 1: pentru completarea matricei S cu 1 dacă e potrivire, altfel 0.

	A	B	C	N	J	R	Q	C	L	C	R	P	M
A	1	0	0	0	0	0	0	0	0	0	0	0	0
J	0	0	0	0	1	0	0	0	0	0	0	0	0
C	0	0	1	0	0	0	0	1	0	1	0	0	0
J	0	0	0	0	1	0	0	0	0	0	0	0	0
N	0	0	0	1	0	0	0	0	0	0	0	0	0
R	0	0	0	0	0	1	0	0	0	0	1	0	0
C	0	0	1	0	0	0	0	1	0	1	0	0	0
K	0	0	0	0	0	0	0	0	0	0	0	0	0
C	0	0	1	0	0	0	0	1	0	1	0	0	0
R	0	0	0	0	0	1	0	0	0	0	1	0	0
B	0	1	0	0	0	0	0	0	0	0	0	0	0
P	0	0	0	0	0	0	0	0	0	0	0	1	0

Pas 2: calculează scorul pe baza formulei $M[i,j]=\text{Max}\{S[i,j]+M[i-1,j-1],M[i-1,j-k],M[i-r,j-1]\}$. Următorul caz pornește invers, din dreapta jos, astfel formula anterioară ar trebui modificată relativ. De asemenea, există o presupunere: linia (-1 prima linie) și coloana (-1 prima coloană) care nu sunt afișate în matrice, dar vor fi folosite pentru calcularea primei linii (linia 0) și coloane (coloana 0), sunt toate 0, ceea ce înseamnă un șir gol potrivit unei secvențe în cazul absenței penalității de gap.

	A	B	C	N	J	R	Q	C	L	C	R	P	M
A	1	0	0	0	0	0	0	0	0	0	0	0	0
J	0	0	0	0	1	0	0	0	0	0	0	0	0
C	0	0	1	0	0	0	0	1	0	1	0	0	0
J	0	0	0	0	1	0	0	0	0	0	0	0	0
N	0	0	0	1	0	0	0	0	0	0	0	0	0
R	0	0	0	0	0	1	4	3	3	2	2	0	0
C	3	3	4	3	3	3	3	4	3	3	1	0	0
K	3	3	3	3	3	3	3	3	3	2	1	0	0
C	2	2	3	2	2	2	2	3	2	3	1	0	0
R	2	1	1	1	1	2	1	1	1	1	2	0	0
B	1	2	1	1	1	1	1	1	1	1	1	0	0
P	0	0	0	0	0	0	0	0	0	0	0	1	0

Pas 3: se parcurge înapoi din colțul stânga sus, și întotdeauna se selectează valoarea maximă din cea mai îndepărtată coloană și linie, după care se sare la următorul maximum în următoarea linie sau coloană.

	A	B	C	N	J	R	Q	C	L	C	R	P	M
A	8	7	6	6	5	4	4	3	3	2	1	0	0
J	7	7	6	6	6	4	4	3	3	2	1	0	0
C	6	6	7	6	5	4	4	4	3	3	1	0	0
J	6	6	6	5	6	4	4	3	3	2	1	0	0
N	5	5	5	6	5	4	4	3	3	2	1	0	0
R	4	4	4	4	4	5	4	3	3	2	2	0	0
C	3	3	4	3	3	3	3	4	3	3	1	0	0
K	3	3	3	3	3	3	3	3	3	2	1	0	0
C	2	2	3	2	2	2	2	3	2	3	1	0	0
R	2	1	1	1	1	2	1	1	1	1	2	0	0
B	1	2	1	1	1	1	1	1	1	1	1	0	0
P	0	0	0	0	0	0	0	0	0	0	0	1	0

Pas 4: prezintă aliniamentul ca algoritmul al lui Sellers arătat mai jos.

Algoritmul lui Sellers:

$M[i,j] = \text{Max}\{S[i,j] + M[i-1,j-1], M[i,j-1] + W(1), M[i-1,j] + W(1)\}$. //1 este numărul unu și nu litera 'l'.

Condiție:

Fiind dat: Premiul de Potrivire: $P_a = 2$ puncte, Penalitatea de Nepotrivire: $P_m = -1$ puncte, Penalitatea de Gap: $P_g = -1$ point.

Proces:

Primul pas: se crează matricea inițială, $(N+1) \times (M+1)$, M este pentru axa X, N pentru Y, și această matrice este numită S.

Aceasta este, se ia partea albă a următoarei matrice completând P_a și P_m relativ. Dacă nu există potrivire, nici nepotrivire și nici gap, valoarea celulei este implicit 0.

Prima linie (linia -1) și prima coloană (coloana -1) care sunt marcate cu gri înseamnă potrivirea unei secvențe cu una goală/un gap. Strict vorbind, ele aparțin celui de-al doilea pas. În timp ce în acest pas ele toate ar trebui să fie **Pg** (ca și -1 sau 0), exceptând intersecția dintre linia -1 și coloana -1 care este 0.

Oricum, este de sens comun de a le avea disponibile în mod direct.

	A	B	C	N	J	R	Q	C	L	C	R	P	M
A	2	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
J	-1	-1	-1	-1	2	-1	-1	-1	-1	-1	-1	-1	-1
C	-1	-1	2	-1	-1	-1	-1	2	-1	2	-1	-1	-1
J	-1	-1	-1	-1	2	-1	-1	-1	-1	-1	-1	-1	-1
N	-1	-1	-1	2	-1	-1	-1	-1	-1	-1	-1	-1	-1
R	-1	-1	-1	-1	-1	2	-1	-1	-1	-1	2	-1	-1
C	-1	-1	2	-1	-1	-1	-1	2	-1	2	-1	-1	-1
K	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
C	-1	-1	2	-1	-1	-1	-1	2	-1	2	-1	-1	-1
R	-1	-1	-1	-1	-1	2	-1	-1	-1	-1	2	-1	-1
B	-1	2	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
P	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	2	-1

Pasul secund: se crează matricea de scoruri numită **M**, din stânga-sus spre direcția dreapta-jos.

D: înseamnă din diagonală apropiată, L: vecinătatea stângă, U: vecinul superior pe baza formulei:

$$M[i,j] = \text{Max}\{S[i,j] + M[i-1,j-1], M[i-1,j] - Pg, M[i,j-1] - Pg\},$$

$1 \leq i \leq M+1, 1 \leq j \leq N+1$, -Pg ar putea fi mai general ca: și amintește **săgeata** dinspre celula/celulele care este/sunt $\text{Max}\{S[i,j] + M[i-1,j-1], M[i-1,j] - Pg, M[i,j-1] - Pg\}$, la $M[i,j]$.

		A	B	C	N	J	R	Q	C	L	C	R	P	M
	0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-10	-11	-12	-13
A	-1	D2	L1	L0	L-1	L-2	L-3	L-4	L-5	L-6	L-7	L-8	L-9	L-10
J	-2	U1	D1	DL0	DL-1	D1	L0	L-1	L-2	L-3	L-4	L-5	L-6	L-7
C	-3	U0	DU0	D3	L2	L1	DL0	DL-1	D1	L0	DL-1	L-2	L-3	L-4
J	-4	U-1	DU-1	U2	D2	D4	L3	L2	L1	DL0	DL-1	DL-2	DL-3	DL-4
N	-5	U-2	DU-2	U1	D4	LU3	D3	DL2	DL1	DL0	DL-1	DL-2	DL-3	DL-4
R	-6	U-3	DU-3	U0	U3	D3	D5	L4	L3	L2	L1	D1	L0	L-1
C	-7	U-4	DU-4	DU-1	U2	DU2	U4	D4	D6	L5	DL4	L3	L2	L1
K	-8	U-5	DU-5	U-2	U1	DU1	U3	DU3	U5	D5	DL4	DL3	DL2	DL1
C	-9	U-6	DU-6	DU-3	U0	DU0	U2	DU2	D5	DLU4	D7	L6	L5	L4
R	-10	U-7	DU-7	U-4	U-1	DU-1	D2	DLU1	U4	D4	U6	D9	L8	L7
B	-11	U-8	D-5	U-5	U-2	DU-2	U1	D1	U3	DU3	U5	U8	D8	DL7
P	-12	U-9	U-6	DU-6	U-3	DU-3	U0	DU0	U2	DU2	U4	U7	D10	L9

Al treilea pas: parcurgerea înapoi din colțul dreapta-jos în modul următor:

Parcurgerea simplă: fără reamintirea săgeților, adecvate potrivirilor maxime, altfel se pot introduce erori.

//întotdeauna se pornește de la cea mai mare linie și ce mai îndepărtată și cea mai îndepărtată coloană.

```

for(i=M+1,j=N+1;i>0&&j>0;)
  {
    highlight the cell P(x,y) with biggest number in column P.x in 0~j and
row P.y in 0~i;
    if(ties happen) trace every P(x,y);
    i=P.x; j=P.y;
    i--; j--;
  }

```

	i	A	B	C	H	J	R	Q	C	L	C	R	P	H
J														
A			-1	-2	-3	-4	-5	-6	-7	-8	-9	-10	-11	-12
J		1		0	-1	1	0	-1	-2	-3	-4	-5	-6	-7
C		0	0		2	1	0	-1	1	0	-1	-2	-3	-4
J		-1	-1	2	2		3	2	1	0	-1	-2	-3	-4
H		-2	-2	1		3	3	2	1	0	-1	-2	-3	-4
R		-3	-3	0	3	3		4	3	2	1	1	0	-1
C		-4	-4	-1	2	2	4	4		5	4	3	2	1
K		-5	-5	-2	1	1	3	3	5		4	3	2	1
C		-6	-6	-3	0	0	2	2	5	4		6	5	4
R		-7	-7	-4	-1	-1	2	1	4	4	6		8	7
B		-8	-8	-5	-2	-2	1	0	3	3	5	8		7
P		-9	-9	-6	-3	-3	0	-1	2	2	4	7		9 ← start

Parcurea standard:

se pornește din **start**, se urmează săgețile înapoi la **end**.

	i	A	B	C	H	J	R	Q	C	L	C	R	P	H
J														
A		2	1	0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-10
J		1	1	0	-1	1	0	-1	-2	-3	-4	-5	-6	-7
C		0	0	2	2	1	0	-1	1	0	-1	-2	-3	-4
J		-1	-1	2	2	4	3	2	1	0	-1	-2	-3	-4
H		-2	-2	1	4	3	3	2	1	0	-1	-2	-3	-4
R		-3	-3	0	3	3	5	4	3	2	1	1	0	-1
C		-4	-4	-1	2	2	4	4	6	5	4	3	2	1
K		-5	-5	-2	1	1	3	3	5	5	4	3	2	1
C		-6	-6	-3	0	0	2	2	5	4	7	6	5	4
R		-7	-7	-4	-1	-1	2	1	4	4	6	9	8	7
B		-8	-8	-5	-2	-2	1	0	3	3	5	8	8	7
P		-9	-9	-6	-3	-3	0	-1	2	2	4	7	10 ← 9 ← start	

Al patrulea pas: se inserează gap-uri și se calculează scorul.

Pornind din colțul dreapta-sus (end), se citesc secvențele X și respectiv Y. Dacă există o săgeată în jos, se inserează un gap pentru secvența X, o săgeată la dreapta, inserează un gap pentru secvența Y.

Rezultate:

$$2-1+2-1+2-1+2-1+2-1+2+2-1+2-1=9$$

A	B	C	N	J		R	Q	C	L	C	R		P	M
A	J	C		J	N	R		C	K	C	R	B	P	
A	B	C		N	J	R	Q	C	L	C	R		P	M
A	J	C	J	N		R		C	K	C	R	B	P	

$$2-1+2-1+2-1+2-1+2-1+2+2-1+2-1=9$$

Aliniament Local: Smith-Waterman,

$$M[i,j]=\text{Max}(M[i-1,j-1]+S[i,j], \text{Max}(M[i,r]+W(r)), \text{Max}(M[s,j]+W(s)), 0)$$

Idea vine din aliniamentul global și este similară în mod special cu algoritmul lui Sellers (1974). A se vedea următoarele:

	H	E	A	G	A	W	G	H	E	E
P	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
A	0.00	0.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00
W	0.00	0.00	0.00	0.00	0.00	2.00	0.67	0.33	0.00	0.00
H	0.00	0.00	0.00	0.00	0.00	0.67	1.67	?		
E										
A										
E										

Diferența dintre aliniamentul local și global constă în următoarea observație: aliniamentul local setează o celulă la 0 când ea devine negativă, ceea ce înseamnă că **o potrivire va fi repornită din această celulă** (deoarece valoarea unei celule vine din trei direcții, $M[i,j]$ atingând 0 denotă faptul că scorurile din toate cele trei direcții sunt balansate a fi 'fără valoare'), în timp ce aliniamentul global va considera rezultatul de la început care este cumulativ la potrivirea curentă. De asemenea, când se parcurge înapoi, aliniamentul local va porni din valoarea maximă, deși aliniamentul global va parcurge întreaga lungime a celor două secvențe. Un

aliniament local ar trebui să fie metoda aleasă atunci când omologia a două secvențe este distantă sau secvențele sunt prea lungi

Proces:

Primul pas: se inițializează matricea **S** pentru a fi completată cu 0.

	H	E	A	G	A	W	G	H	E	E
P	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
A	0.00									
W	0.00									
H	0.00									
E	0.00									
A	0.00									
E	0.00									

Pasul doi: se continuă inițializarea matricei **S** prin adăugarea de puncte de Potrivire/Nepotrivire la matricea obținută anterior. A se nota, linia (-1 prima) și coloana (-1 prima) nu sunt afișate, deoarece ele nu vor avea efect pe linia 0 și coloana 0 după aplicarea la formulă.

Al treilea pas: se calculează matricea de scoruri și se parcurge înapoi de la cea mai ridicată valoare dintr-o celulă din matrice și urmând săgețile până la celula care conține valoarea 0 (care nu se include).

	H	E	A	G	A	W	G	H	E	E
P	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
A	0.00	0.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00
W	0.00	0.00	0.00	0.67	0.00	2.00	0.67	0.33	0.00	0.00
H	0.00	0.00	0.00	0.00	0.33	0.67	1.67	1.67	0.33	0.00
E	0.00	1.00	0.00	0.00	0.00	0.33	0.33	1.33	2.67	1.33
A	0.00	0.00	2.00	0.67	1.00	0.00	0.00	0.00	1.33	2.33
E	0.00	1.00	0.67	1.67	0.33	0.67	0.00	0.00	1.00	2.33

Al patrulea pas: se prezintă rezultatul ca în algoritmul Sellers de mai sus.

ANEXA 3

123D, (<http://123d.ncifcrf.gov/123D+.html>). 123D-este un program care combină profile de secvențe, predicția structurii secundare și potențialele capacități de contact pentru a descrie o secvență de proteine prin setul de structuri. Pentru o secvență oarecare de proteină selectată din baza de date a NCBI rezultatele furnizate de acest program sunt de următoarea formă :

Database of **4381** domains
 Matrix: **gonnet**
 Secondary structure prediction
 Contact Capacity Potentials
 Contact Capacity Potentials
 Penalty for Gap opening: **-10.000000**
 Penalty for Gap extension: **-1.000000**

Q_Seq: **musculus** length: **437**a.a.
 20 best hits:

PDB ID	raw score	Z-score	a.a. aligned	% identities	SCOP family
<u>1vhh</u> 162 a.a.	6499	33.62	162	100	<u>d.65.1.2</u>
<u>1at0</u> 145 a.a.	4911	25.00	145	35	<u>b.86.1.1</u>
<u>1cyda</u> 244 a.a.	1036	3.98	236	14	<u>c.2.1.2</u>
<u>1dq3a1</u> 168 a.a.	995	3.76	135	15	<u>b.86.1.2</u>
<u>1eqja2</u> 274 a.a.	952	3.53	232	13	<u>c.76.1.3</u>
<u>2ae2a</u> 260 a.a.	948	3.51	250	14	<u>c.2.1.2</u>
<u>1qg6a</u> 257 a.a.	919	3.35	246	13	<u>c.2.1.2</u>
<u>1e6wa</u> 260 a.a.	916	3.33	251	11	<u>c.2.1.2</u>

....

Alignments:

```
seq 1vhh_ 6499 437 0.00 -100000.00 162 0
s_m = 6191 s_ss = 124 s_cc = 184 go = 0 ge = 0
HHHHHHHHHHHEH HE H HH HHHHH
MLLLLARCFLVILASSLLVCPGLACGPRGFGKRRHPKLTPLAYKQFIPNVAEKTGAS
*****
.....RRHPKLTPLAYKQFIPNVAEKTGAS
UUUUU TTTEETT TTTTTTT

EEE EEE H EEEE HHHHHHHHHHHHEEHEEH
GRYEGKITRNSERFKELTPNYNPDIIFKDEENTGADRLMTQRCKDKLNALISVMNQWPG
*****
```

GRYEGKITRNSERFKELTPNYNPDIIFKDEENTGADRLMTQRCKDKLNALAISVMNQWPG
T TTTTTGGG EE TTEEE TTTT GGEE HHHHHHHHHHHHHHHHTTT

EEEE E HHH EEEEE HHHHHHHHHH EEEE
VKLRVTEGWDEDGHHSEESLHYEGRAVDITTSRDRSKYGMLARLAVEAGFDWVYYESKA

VKLRVTEGWDEDGHHSEESLHYEGRAVDITTSRDRSKYGMLARLAVEAGFDWVYYESKA
T EEEEE TTTT TTTTGGGG EEEEEETT GGGHHHHHHHHHHH TTEEEEEETT

...

FFAS (Fold & Function Assignment System) un program utilizat pentru analiza secvențelor biologice din punct de vedere structural.

Aliniamentele profil-la-profil au avantajul informației prezente în secvențele omologe de proteine pentru a amplifica modelul secvenței care definește familia. Ca un rezultat, ele sunt în stare să detecteze omologiile îndepărtate dincolo de utilizarea altor metode de comparare a secvențelor.

O variantă îmbunătățită a acestui tool este serverul FFAS03 care acceptă secvențe de proteine furnizate de utilizator și generează automat un profil care este comparat după aceea cu câteva biblioteci de profile ale secvențelor incluzând PDB proteins, COG, PFAM, SCOP și o bibliotecă de obiective genomice structurale. Serverul permite, de asemenea, accesul la analiza aliniamentului, aliniament multiplu și tooluri de modelare comparativă.

Descrierea metodei se poate face pe scurt menționând pașii principali. În primul pas este pregătit aliniamentul multiplu al secvențelor folosind PSI-BLAST. Sunt rulate cinci execuții pe baza de secvențe de proteine NR85S. În al doilea pas, toate secvențele găsite de PSI-BLAST sunt utilizate pentru a construi profilul. Ponderile sunt asignate secvențelor pe baza similarității lor cu toate celelalte secvențe incluse în același rezultat al PSI-BLAST. În plus, FFAS aliază două profile folosind un algoritm local de programare dinamică. Valoarea scorului de comparație dintre pozițiile n și m din cele două profile este calculată ca și produs (dot product) al coloanei a n -a din primul profil și a m -a coloană din al doilea profil. După asocierea de valori tuturor pozițiilor matricea obținută este normalizată. Aliniamentul optimal este calculat de un algoritm de programare dinamică. În ultimul pas, scorul elementar astfel obținut, este translatat în scorul final FFAS comparându-l cu distribuția scorurilor primare obținute prin compararea secvențelor neînrudite.

Datele de intrare pentru serverul FFAS sunt sub forma secvențelor de amino acizi în format FASTA. Serverul acceptă secvențe de lungime cuprinsă între 25 și 2000 reziduuri dar algoritmul a fost optimizat pentru secvențe de lungime cuprinsă între 50 și 500 reziduuri.

Exemplu :

Prezintă cel mai bun model din fiecare bază de date și toate mostrele cu scoruri semnificative acoperind diverse fragmente căutate. O parte demonstrativă din rezultatul execuției pentru o secvență de proteine idicată (mus

musculus selectată aleator din baza de secvențe locată pe serverul NCBI) este ilustrată de tabelul următor.

# Query	Length	Result vs. Range	Score	%id	Covered by template(s)
1 sonic hedgehog [Mus musculus] (11/16/05)	430	<u>PDB0205</u> 34-191	-85.700	94	<u>1vhh</u> mol:protein length:162 Sonic Hedgehog
		193-344	-26.900	34	<u>1at0</u> mol:protein-het length:145 17- Hedgehog
2 gi 18203507 sp Q9ULZ3 ASC_HUMAN Apoptosis-associated speck-like protein containing a CARD (hASC) (PYD and CARD domain containing protein) (Target of methylation-induced silencing 1) (Caspase recruitment domain protein 5) (11/16/05)	195	<u>PDB0205</u> 1-91	-49.300	100	<u>1ucp_A</u> mol:protein length:91 Apoptosis-Associated Speck-Like Protein Conta
		103-195	-28.600	19	<u>1cww_A</u> mol:protein length:102 Apoptotic Protease Activating Factor 1

3 (11/16/05)	195 <u>PfamA160</u>	5-91	-38.900	21 <u>PF02758.5</u> ; AIM2_HUMAN/3-87; PAAD/DAPIN/Pyrin domain
		112-195	-28.400	11 <u>PF00619.9</u> ; CED3_CAEVU/7-91; Caspase recruitment domain
4 gi 18203507 sp Q9ULZ3 ASC_HUMAN Apoptosis-associated speck-like protein containing a CARD (hASC) (PYD and CARD domain containing protein) (Target of methylation-induced silencing 1) (Caspase recruitment domain protein 5) (11/16/05)	195 <u>PfamA160</u>	5-91	-38.900	21 <u>PF02758.5</u> ; AIM2_HUMAN/3-87; PAAD/DAPIN/Pyrin domain
		112-195	-28.400	11 <u>PF00619.9</u> ; CED3_CAEVU/7-91; Caspase recruitment domain

3D PSSM Fold Recognition Server , este o metodă web rapidă pentru recunoașterea plicurilor proteinelor folosind profilele secvențelor reprezentate 1D sau 3D cuplate cu structura secundară și Informația Potențială de Solvatare (Solvation Potential Information).

Serverul 3D-PSSM este proiectat să preia secvența de proteină supusă investigației și să încerce să prezică structura tridimensională și funcția sa probabilă. Se utilizează o bibliotecă de structuri de proteine deja cunoscute față de care este descompusă și evaluată pentru compatibilitate proteina investigată. Sunt utilizate o varietate de componente de evaluare : 1D-PSSMs (profile de secvențe construite din omologii apropiate) ; 3D-PSSMs (profile mai generale conținând mai multe omologii îndepărtate), potrivirea elementelor structurii secundare și tendințele reziduurilor în secvențe de query pentru a ocupa diverse nivele ale accesibilității solventului.

Exemplu de rezultate oferite în urma analizei pentru o secvență de proteine oarecare, TRAD_HUMAN_domain sunt ilustrate în tabelul următor.

Description: TRAD_HUMAN_domain
 Remote Host: larch.lif.icnet.uk
 Remote Address: 143.65.53.159
 E-mail Address: demo@icrf.icnet.uk
 Input Format: single
 Query Length: 100

E-value Key (% Certainty)
 95% 90% 80% 70% 50%

[View Multiple Sequence Alignment](#)

View Alignment	SCOP	Template Lengths	Model	PSSM E-value	PSSM Type	E-value
<input checked="" type="checkbox"/> View Alignment	1ntr-0 26% id	85		0.062	R1D	0.653 <input checked="" type="checkbox"/> bind.ptm REPEAT.PHOSPHORYLATION contains
<input checked="" type="checkbox"/> View Alignment	1d3f-0 17% id	127		1.1	R1D	0.872 <input checked="" type="checkbox"/> REPEAT binding contains membrane.subcellular
<input checked="" type="checkbox"/> View Alignment	1sknF0 20% id	74		2.24	R1D	0.946 <input checked="" type="checkbox"/> contains subcellular location
<input checked="" type="checkbox"/> View Alignment	1etp_A2 22% id	92		2.9	1D	0.746 <input checked="" type="checkbox"/> group.ptm subcellular location

[View Model](#)

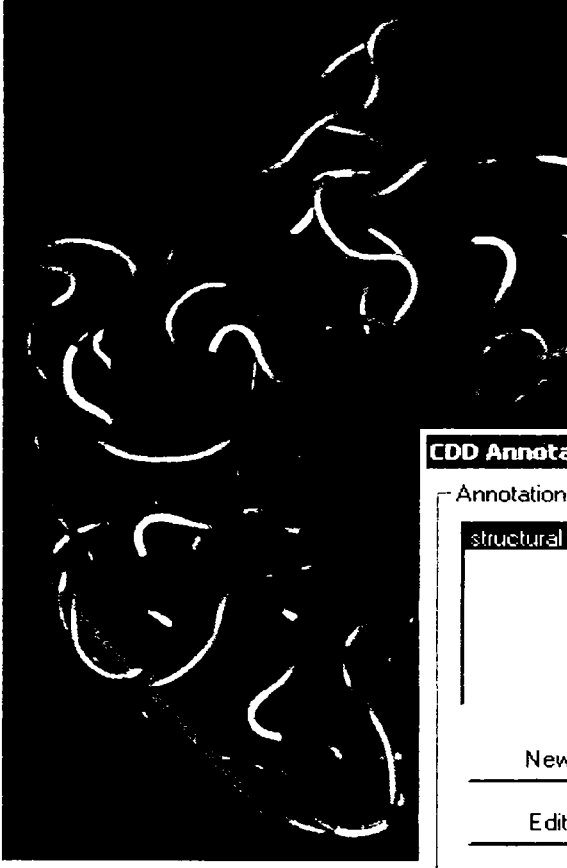
Query	
2ndary Str. Pred.	
Reliability	

Cn3D, este o aplicație ajutătoare pentru căutarea web care permite vizualizarea structurilor 3D din baze de secvențe gestionate de NCBI Entrez.

În cele ce urmează este făcută o exemplificare simplă a caracteristicilor afișate pentru domeniul WD40. Cn3D prezintă o structură de proteine reprezentativă, aliniamentul familiei și tablouri de adnotare cu informații despre caracteristicile adnotate ale acestei familii de proteine. Reziduurile conservate într-o formă caracteristică acestui domeniu sunt marcate luminos.

WD40 - Cn3D 4.1

File View Show/Hide Style Window CDD Help



CDD Descriptive Item

Name: WD40

WD40 domain, found in a number of e a wide variety of functions including ac in signal transduction, pre-mRNA proc assembly; typically contains a GH dipe its N-terminus and the WD dipeptide a 40 residues long, hence the name WC a conserved core; serves as a stable p

Show Annotations Panel | Show

CDD Annotations

Annotations

structural tetrad

New | Highlight | Move Up

Edit | Delete | Move Down

WD40 - Sequence/Alignment Viewer

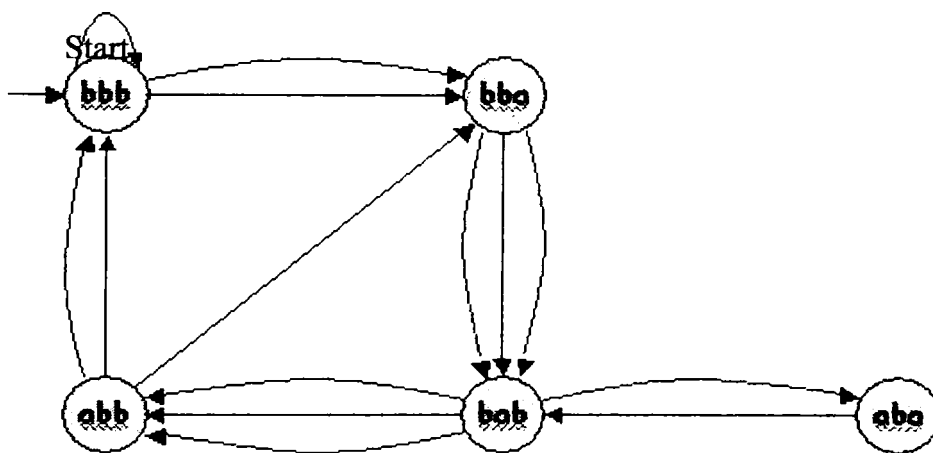
View Edit Mouse Mode Unaligned Justification Imports

ATDQ_A	~ ~ Q I V T S S ~ ~ ~ ~ ~ D T T C A L W D ~ ~ ~ ~ ~
consensus	~ ~ g R I L S S S s ~ ~ ~ ~ ~ r D K T I K V W D v e ~ ~ ~ ~ ~
IJV2_A	~ ~ R V L L G G ~ ~ ~ ~ ~ Q G Q L I S D Q ~ ~ ~ ~ ~
gi 1346108	~ ~ a H I L T A S g ~ ~ ~ ~ ~ D M T C A L W D i p ~ ~ ~ ~ ~
gi 2506476	~ ~ l P W I I S A s ~ ~ ~ ~ ~ d D Q T I R I W N w q ~ ~ ~ ~ ~
gi 461848	~ ~ d F I L V G T q ~ ~ ~ ~ ~ H P T L R L Y D i n ~ ~ ~ ~ ~
gi 3122853	s p g L Y F L T A g ~ ~ ~ ~ ~ d Q G T L R V W E a a s ~ ~ ~ ~ ~
gi 3023307	~ ~ h H L L L A T s ~ ~ ~ ~ ~ s D C F L K L W D i n ~ ~ ~ ~ ~
gi 1176045	~ ~

ANEXA 4

Simularea unui lanț Markov și generarea pseudo-aleatoare de text

Reprezentat grafic, lanțul Markov este compus dintr-o mulțime de stări cu o stare inițială numită stare de start și un set de tranziții de la o stare la alta. Figura de mai jos ilustrează un lanț Markov cu 5 stări și 14 tranziții.



Simularea unui lanț Markov. Pentru a simula traiectoria prin lanțul Markov se începe la starea de Start. La fiecare pas se selectează unul din arcele de plecare la alegere și se trece la starea vecină. Spre exemplu, dacă lanțul Markov este în starea "bab", atunci va traza în starea "abb" cu probabilitatea $\frac{3}{4}$ și în starea "aba" cu probabilitatea $\frac{1}{4}$. Mai jos sunt descriși primii 10 pași ai unei traiectorii posibile pornind din "bbb":

bbb -> bbb -> bba -> bab -> abb -> bbb -> bbb -> bba -> bab -> abb

Generarea de text. Pentru generarea unui text pseudo-aleator în conformitate cu un model Markov de ordinul k construind un lanț Markov dintr-o secvență de text se pot parcurge următoarele etape: Se include o stare pentru fiecare k -gram distinct care apare în text; de asemenea, se includ N tranziții, una pentru fiecare k -gram în text, care trece la următoarea k -gram (suprapusă); se tratează textul ca un șir ciclic pentru a atinge marginile. Spre exemplu, dacă $k = 3$ și secvența de text este

"bbbabbabbbba", atunci primele două tranziții sunt de la "bbb" la "bba" și de la "bba" la "bab" iar ultima (ciclic) este de la "abb" la "bbb". Lanțul Markov complet este ilustrat în figura anterioară.

vocabulary components total=74	occurrence frequency total frequency=621087	vocabulary components total=74	occurrence frequency total frequency=621087
!	18	a	41264
#	8	b	6490
\$	180	c	13827
%	27	d	19810
&	29	e	60113
'	1779	f	9749
(6370	g	8830
)	6	h	22864
0	1556	i	34080
1	34	j	664
2	1569	k	3728
3	1275	l	19678
4	1254	m	11381
5	622	n	34709
6	652	o	36601
7	532	p	9306
8	821	q	552
9	940	r	32847
A	1368	s	33433
B	500	t	41192
C	302	u	13006
D	696	v	4783
E	1785	w	8393
F	777	x	890
G	475	y	7999
H	831	z	459
I	26	ξ	5167
J	765	Π	5167
K	1950	1	1220
L	2081	2	640
M	336	3	358
N	177	4	387
O	909	5	534
P	9	6	332
Q	170	7	357
R	33	8	400
S	98407	9	608

Vocabularul activ utilizat în experimente

Detalii despre modul în care a fost creată baza de secvențe mutante

Un set de date format din 50 de secvențe de proteine aparținând genei p53 [X54156 din NCBI GenBank] și 50 secvențe aleatoare au fost folosite pentru a evalua algoritmul propus pentru similaritatea proteinelor.

Mutațiile au fost selectate aleator din baza de secvențe creată utilizând descrierea oferită de International Agency for Research on Cancer (IARC) Lyon, Franța³⁸. Baza de date IARC folosită este sub formă de fișier în format Excel și conține 18 585 întrări organizate în 42 coloane. Fiecare linie (coloană) reprezintă o singură mutație căreia îi este asignat un singur număr unic de identificare. Un număr unic de identificare este de asemenea atribuit mostrei de tumoră și pacientului. Din fiecare înregistrare se folosesc câmpurile care definesc acea mutație:

- numărul de identificare unic al mutației;
- locația mutației în intron sau exon³⁹ în gena p53 ;
- pentru mutații în exon, numărul codon-ului la care este locată mutația (1-393);
- poziția nucleotidelor mutației pe baza intrării în banca genetică, X54156;
- valorile de Yes (True) sau No (False) pentru a indica dacă poziția mutației cade într-o secvență CpG (Cytosine-phosphor-Guanine),
- valorile Yes (True) sau No (False) pentru a indica dacă poziția mutației cade într-o locație de conexiune ("splice site");
- natura mutației;
- pentru mutații în exoni, secvența de bază a codon-ului în care are loc mutația;
- secvența bazei mutate;
- tipul de amino acid (wild-type amino acid) codificat în codon-ul în care are loc mutația (abreviere din trei litere a amino acidului);
- amino acidul mutat codificat în codon-ul în care a avut loc mutația (abreviația de trei litere a amino acidului) și
- efectul mutației și al numărului de codon la care poate apare un codon de STOP.

Din baza de date creată au fost luate aleator 50 de mutații care produc *missense mutation* (când modificarea nucleotidelor produce modificare în codon, atunci amino acidul se modifică și în final și proteina), *nonsense mutations* (când modificarea în nucleotide determină modificarea codonului în TAA, TAG sau TGA – care sunt codon de STOP și prin urmare nu se va produce nici un amino acid, rezultând într-o proteină mai scurtă) și *silent mutations* (când schimbarea în nucleotide determină modificarea codon-ului dar nu se modifică nici amino acizii și prin urmare nici proteina). Setul de secvențe este următorul :

>AAH84538.txt

```
MSDVAIVKEGWLHKRGEYIKTWRPRYFLLKNDGTFIGYKERPQDQDQREAPLNNFSVAQCQL
MKTERPRPNTFIIRCLQWTTVIERTFHVETPEEREWTTAIQTVADGLKKQEEEEEMDFRSGSPS
```

³⁸ <http://www.iarc.fr/p53/Somatic.html>

³⁹ **Exonii** sunt regiuni ale AND-ului dintr-o genă care nu sunt înlăturate prin transcrierea ARN și sunt păstrate în molecula finală mesager de ARN(mRNA)[WIKI'05g].

DNSGAEEMEVSLAKPKHRVTMNEFEYLKLLGKGTFGKVILVKEKATGRYYAMKILKKEVIVAKD
EVAHTLTENRVLQNSRHPFLTALKYSFQTHDRLCFVMEYANGGELFFHLSRERVFSEDRARFYG
AEIVSALDYLHSEKNVVYRDLKLENMLDKDGHKIDTDFGLCKEIKDGATMKTFCGTPEYLAP
EVLEDNDYGRAVDWWGLGVVMEYEMMCGRLPFYNQDHEKLFELILMEEIRFPRTLGPPEAKSLLS
GLLKKDPKQRLGGGSEDAKEIMQHRFFAGIVWQHVVYEKLSPPFKPQVTSETDTRYFDEEFTA
QMITITPPDQDDSMECVDSERRPHFPQFSYSASGTA

>AAH84541.txt

MAVFVLLALVAGVLGNEFSILKSPGSVFRNGNWPIPERIPDVAALSMGFSVKEDLSWPGL
AVGNLFHRPRATVMVMVKGVNKLALPPGSVISYPLENAVPFSLDSVANSIHSLSFSEETPVVLQL
APSEERVYMGKANSVFEDLSVTLRQLRNRLFQENSVLSSLPLNSLSRNNEVDLLFSELQVLH
DISSLLSRHKHLAKDHSPDLYSLELAGLDEIGKRYGEDSEQFRDASKILVDALQKFADDMYSLY
GGNAVVELVTVKSFDTSLIRKTRTILEAKQAKNPASPYNLAYKYNFEYSVVFNMVLWIMIALALA
VIITSYNIWNMDPGYDSIIYRMTNQKIRMD

>AAH84544.txt

GPNARRVGSPPGRGVPGGKSAVGAPRPRCPGAMAMVTGGWGGPGGDTNGVDKAGGYPRAA
EDDSASPPGAASDAEPGDEERPLQVDCVCGDKSSGKHYGVFTCEGCKSFFKRSIRRNLISY
TCRSNRDCQIDQHHRNQCYCRLKKCFRVGMRKEAVQRGRIPHSLPGAVAASSGSPPGSALA
AVASGGDLFPGQPVSELIAQLLRAEPYPAAAGRFGAGGGAAGAVLGIDNVCELAARLLFSTVE
WARHAPFFPELPVADQVALLRLSWSELFVLNAAQAALPLHTAPLLAAAGLHAAPMAAERAVAFM
DQVRAFQEVDKLGRLQVDSAEGCLKAIALFTPDACGLSDPAHVESLQEKALQVALTEYVRAQ
YPSQPQRFGRLLLRLPALRAVPASLISQLFFMRLVGKTPIETLIRDMLLSGSTFNWVYGGGQ

>AAH84554.txt

MASGSAGKPTGEAASPAPASAIGGASSQPRKRLVSVCDHCKGKMLVADLLLLSSEARPVLF
GPASSGAGAESFEQCRDTIARTKGLSILTHDVQSQLNMGRFGEAGDSLVELGDLVVSLETC
AHAAYLAAVATPGAQPAQPLVDYRVTRCRHEVEQCAVLRATPLADMTPLQLLEVSQGLSR
NLKFLTDACALASDKSRDRFSREQFKLVKCMSTSASALLACVREVKVAPSELARSRCALFSGP
LVQAVSALVGFATEPQFLGRAAAVSAEGKAVQTAILGGAMSVWSACVLLTQCLRDLAQHPDGG
AKMSDHRERLRNSACAVSEGCTLLSQALRERSSPRTLPPVNSNSVN

>AAH84612.txt

LLLLVAILLGSAVQTKTFDRIIGGEECVPHSQPWQVALYYFSDYLCGGILIDEWVLTAAHCNQ
SNLQVLLGAHNRTKPTDHKQYTYAVKICPRCDFPVTYNNDIMLLKLASKANMNCHVKTIQLA
SDLVEDNTECLASGWGTTTSPEENYPDKLQCVNLSTVSNSECQACYPEDDITDNLMLCAGNMA
GGKDTCKGDSGGPLVCNGELHGITSWGHYICGLPNKPGVFTKVFNIDWISDIMQENENPCCY
ETT

>AAH84617.txt

MAVPPKSTKERLESLLDDLEVLRSREVIETLALSRSQKLSQPGEENQILELLIQKDGEFQELMKVA
FSQGKTHQEMQVLEKEVEKRDSDIQQLQKQLKEAEHILATAVYQAKEKLSIDKARKGVISSE
ELIKYAHRISASNAVCAPLTWVPGDPRRPPYPTDLEMRSGLLGQMSNLPTNGVNGHLPGDALAA
GRLPDVLAQPYPWQSSDMSMNMLPPNHSNEFLMESLGPKNKENEEDVEVMSTDSSSSSSDSD

>AAH84660.txt

MVHTARLSRRLCSAGGSMELGQLGHSLGLNGQQLERLVEVEEGRSLVLRVQEGDQKVAVYRS
DLRLCQTPKCAGDCGQLHLRCFYILGSCNRSPCKFNHCIRNAHNSVLQKFHLD SMPIDELRQ
LLLQNDPNLLPDICSHYNRGDGPHGSCYKNCNKLHVCQHFLQGDCKFGEKCKRCHDLSEE
ETLKKLTKWQLSDSLLPGLLETYTNAHTLKSSCDRPPRNVVKSAEKTKPGPSKSGTSQEIEEI
CLYFIRNSCSFKEKCVREHFHPYRWQVYTNGTWKDLDNMEMIEKSYCDPNSRAVLVNLDFA

MTYQSNKVKRLSTPSSASKPPHFVFTTDWKWYWMDEYNKWVEYGTESDRHANSTICSSDLE
 NFYQSNQTADVFKAGKHKYLSLFSKMNQVQRNLQYETKRRVCRRPVFVSGKDVKKRSSTSEP
 SKEDKNTPEHWDKGQTPDLGYKLVLLSPLSEEYSKIEAMFCRTLRTIRIHSIERIQNLALWEVYQ
 W

>AAH84676.txt

MTFEPDPADLALSSIPGHETFDPRRHRFSEELKQPIMKKARKVQVPEEQKDEKYWSRRYKN
 NEAAKRSRDARRLKENQISVRAAFLEKENALLRQEVVAVRQELSHYRAVLRSRYQAQHGTL

>AAH84693.txt

MEGTQEALSGKMRLLFPAARTSLLMLRLNEAALRALQECQQQVVPVIAFQGGQRYLRLPGP
 GWSCLFSFIVSQCGQEGGLDLVYQRLGRSGPNCLHCLGSLRERLTIWAAMDTIPAPLLAQEH
 LTEGTRESESWQDSEDEPEGHPQMALQEVSDPLASNHEQSLPGSSSEPMQWEVRNHTYLS
 NREPDQPLSSASQKRLDKKRSAPITTEEPEEKRPALPLASSPLQGLSNQDSPEEQDWGQDA
 DGDSRLEQSLSVQSASESPSEVPDYLLQYSTIHSAEQQAYEQDFETDYAEYRILHARVGA
 ASQRFTELGAIEKRLQRGTPEHKVLEDKIVQYKFRKRYPSYSEEKRRCEYLHEKLSHIKGLIL
 EFEENRGS

>AAH84726.txt

MTTSQKHRDFVAEPMGEKPVGSLAGIGDALGKRLEERGFDAKAYVVLGQFLVLKKDEDLFREWL
 KDTCGANAKQSRDCFGLREWCDLFL

>BC008979.txt

MSDSWVPNSASGQDPGGRRRAWAELLAGRVKREKYNPERAQLKESAVRLLRSHQDLNALLL
 EVEGPLCKKLSLSKVIDCDSSEAYANHSSFFIGSALQDQASRLGVPVIGLSAGMVASSVGGIC
 TAPAETSHPVLLTVEQRKLSLLEFAQYLLAHSMFSLFCQELWKIQSSLLLEAVWHLHVQG
 IVSLQELLESHPMHAVGSLFRNLCCLEQMEASCQHADVARAMLSDVFQMFVLRGFGQKNS
 DLRRTVEPEKMPQVAVDVLRMLIFALDALAAGVQEEESTHKIVRC

>BC036531.txt

MFSFVDLRLLLLLAATALLTHGQEEGQVEGQDEDIPPITCVQNGRLRYHDRDVWVKPEPCRICVC
 DNGKVLCDVICDETKNCPGAEVPEGECCVCPDGSSEPTDQETTGVGPKGDTGPRGPRGP
 AGPPGRDGIPGQPLGPPGPPGPPGPPGLGGNFAPQLSYGYDEKSTGGISVPGPMGSPGPRG
 LPGPPGAPGPQGFQGGPEGEPGASGPMGPRGPPGPPGKNGDDGEAGKPRPGERGPPGP
 QGARGPLGTAGLPGMKGRGFSGLDGAKGADGAPGPKGEPGSPGENGAPGQMGPRGLPGE
 RGRPGAPGPAGARGNDGATGAAGPPGPTGPAGPPGFPGAVGAKGEAGPQGPGRGSEGPQGV
 GEPGPPGAGAAGPAGNPGADGQPGAKGANGAPGIAGAPGFPGARGSPGQGGPPGPKG
 NSGEPGAPGSKGDTGAKGEPGVPVQGGPPGAGEEGKRGARGEPGPTGLPGGPGERGGPGS
 RGFPGADGVAGPKGPAGERGSPGAPGKSPGEGAGRPGEAGLPGAKGLTGSPPGSPGPDGKT
 GPPGPAGQDGRPGPPGPPGARGQAGVMGFPGPKGAAGEPGKAGERGVPPGAVGPAGKDG
 EAGAQQPPGAPGAGERGEQGPAGSPGFQGLPGPAGPPGEAGKPEEQGVPGDLGAPGPGA
 RGERGFGERGVQGGPPGAPRGANGAPNDGAKGDAGAPGAPGSQGAPGLQGMPPGERGA
 AGLPMPKGDGRDAGPKGADGSPGKDGVRGLTGPVGGPAGAPGDKGESGSPGAPGPTGAR
 GAPGDRGEPGPPGAPGAGPPGADGQPGAKGEPGDAGAKGDAGPPGAPGAPGPPGPIGNV
 APGAKGARGSGAPGATGFPGAAGRVPVGGPSGNAGPPGPPGAPGKEGGKGRGETGPAGR
 PGEVGGPPGPPGAGEKGSPPGADGAPAGPTGPPQGIAGQRGVGLPGQRGERGFPLGPPSG
 EPKQGPSGASGERGPPGPMGPPGLAGPPGESGREGAPGAEGSPGRDGSAGKDRGETGP
 AGPPGAPGAPGAPGVPVGPAGKSGDRGETGPAGAPGVPVGPARGPAGPQGPGRDKGETGEQ
 GDRGIKGRGFSGLQGGPPGPPGPPGPPGPPSAGFDFFLQPPQEKAHHDGGYRADDANVVRD
 PGPRGRTGDAGPVGGPPGPPGPPGPPGPPSAGFDFFLQPPQEKAHHDGGYRADDANVVRD
 RDLEVDTTLKLSLQIENIRSPEGSRKNPARTCRDLKMHSDWKSGEYWIDPNQGCNLDAIK

VFCNMETGETCVYPTQPSVAQKNWYISKNPDKDRHVWFGESMTDGFQFEYGGQGS DPADVA
IQLTFLRLMSTEASQNTYHCNSVAYMDQQTGNLKKALLLQGSNEIEIRAEGNSRFTYSVTVD
GCTSHTGAWGKTVIEYKTKTSRLRIIDVAPLDVVGAPDQEFQFDFVGHVCFL

>M17398.txt

MEPFVVLVLCLSFMLLFSLWRQSCRRRKLPPGPPTPLPIIGNMLQIDVKDICKSFTNFSKVYGPV
TVYFGMNPVVFHGYEAVKEALIDNGEEFSGRGNPISQRITKGLGISSNGKRWKEIRRFSLT
NLRNFGMGKRSIEDRVQEEAHCLVEELRKTASPCDPTFILGCAPCNVICSVVFQKRFDYKQD
NFLTLMKRFNENFRILNSPWIQVCNFPLLIDCFPGTHNKVLKNVALTRSYIREKVKEHQASLD
VNNPRDFMDCFLIKMEQEKDNQKSEFNENLVGTVADLFVAGTETTSTTLRYGLLLLLKHPEVT
AKVQEEIDHVIGRHRSPCMQDRSHMPYTDVAVVHEIQRYSDLVPTGVPHAVTTDTKFRNYLIPK
SFDNKIMLA

>NM0003972.txt

MAVSRLDRFLILLDTGTTPVTRKAAAQQLGEVVKLHPHELNNLLSKVLIYLRSANWDTRIAAGQ
AVEAIVKNVPEWNPVPRTRQEPTSESSMEDSPTTERLNDFRFDICRLLQH GASLLGSAGAEFEV
QDEKSGEVDPKERARQRKLLQKGLNMGEAIGMSTEELFNDEDLDYTPTSASFVNKQPTLQ
AAELIDSEFRAGMSNRQKNKAKRMAKLFKQSRDAVETNEKSNDSTDGEPEEKRRKIANVVI
NQSANDSKVLIDNIPDSSSLIETNEWPLESFCEELCNDLFNPSWEVRHGAGTGLREILKAHG
KSGGKMGDSTLEEMIQQHQEWLEDLVIRLLCVFALDRFGDFVSDEVVAPVRETCAQTLGVVLK
HMNETGVHKTVDVLLKLTQEWEVRHGGLGIKALAVRQDVINTLLPKVLRTRIEGLQDLDD
DVRVAAAASLVPVVESLVYLQTKVPFIINTLWDALLEDDLTA STNSIMTLLSSLLTYPQVQQC
SIQQSLTVLVRVWPFLHHTISSVRAALETFTLLSTQDQNSSSWLIPILPDMLRHIFQFCVLE
SSQEILDLIHKVWMELLSKASVQYVVAACPWMGAWLCLMMQPSHLPIDLNMLLEVKARAKE
KTGGKVRQGGQSNKEVLQEYIAGADTIMEDPATRDFVVMRARMMAAKLLGALCCICDPGVN
VVTQEIKPAESLGQLLFLHNSKSALQRISVALVICWAALQKECKAVTLAVQPRLLDILSEHLY
YDEIAPFTRMQNECKQLISSLADVHIEVGNRVNNVLTIDQASDLVTTVFNEATSSFDLNPQV
LQQLDSKRQQVQMTVTETNQEWQVLQLRVHTFAACAVVSLQQLPEKLNPIIKPLMETIKKEEN
TLVQNYAAQCIKLLQCTTRTPCPNSKIIKNLCSL CVDPYLTPCVTPVPTQSGQENSKGST
SEKDGMMHHTVTKHRGIITLYRHQKAFAITSSRRGPTPKAVKAQIADLPAGSSGNILVELDEAQK
PYLVQRRGAEFALTTIVKHFGGEMAVKPLHLWDAMVGPLRNTIDINNFDGKSLLDKGDSPAQE
LVNSLQVFETAASMDSELHPLL VQHLP HLYMCLQYPSTAVRHMAARCVGMVMSKIATMETMNI
FLEKVL PWLGAIDDSVKQEGAEALACVMEQLDVGIVPYIVLLVVPVLGRMSDQTD SVRFMAT
QCFATLIRLMPLEAGIPDPNMSAELIQLKAKERHFLEQLLDGKKLENYKIPVPINAELRKYQQD
GVNWLAFNLKYKLHGILCDDMGLGKTLQ SICILAGDHCHRAQEYARSKLAECMPLPSLVVCP
TLTGHVWDEVGKFC SREYLNPLHYTGPPTERIRLQH QVKRHN LIVASYDVVRNDIDFFRNIKFN
YCILDEGHV IKNGKTKLSKAVKQLTANYRIILSGTPIQNNVLELW SFLDFLMPGFLGTERQFAAR
YGKPIASRDARSSSREQEAGVLAMDALHRQVLPFLLRRMKEDVLQDLPPKIIQDYCTLSPLQ
VQLYEDFAKSRACDQVDET VSSATLSEETEKPKLKATGHVFQALQYLRKLCNHPALVLT PQHPE
FKTTAEKLA VQNSSLHDIQHAPKLSALKQLLLDCGLGNGSTSESGTESVVAQHRI LIFCQLKSM
LDIVEHDLLKPHLPSVTYLRDLGSI PPQRHSIVSRFNNDPSIDVLLLTHVGGGLGNLTGADTV
VFVEHDWNP MRDLQAMDRAHRIGQKR VNVYRLITRGTLEEKIMGLQKFKMNIANTVISQEN
SSLQSMGTDQLLLDFTLDKDGKAEKADTSTSGKASMKSILENLSDLWDQEYDSEYSLENFM
HSLK

>NM000729.txt

MNSGVCLCVLMAVLAAGALTQVPPADPAGSGLQRAEEAPRRQLRVSQRTDGESRAHLGALL
ARYIQQARKAPSGRMSIVKNLQNLDP SHRISDRDYMGMDFGRRSAEEYEYPS

>NM001895.txt

MSGPVPSRARVYTDVNTHRPREYWDYESHVVEWGNQDDYQLVRKLGKGYSEVFEAINITNN
 EKVVVKILKPVKIKREIKILENLRGGPNITLADIVKDPVSRTPALVFEHVNTDFKQLYQTL
 TDYDIRFMYEILKALDYCHSMGIMHRDVKPHNVIMIDHEHRKLRLLIDWGLAEFYHPGQYENVR
 VASRYFKGPELLVDYQMYDYSLDMWSLGCMLASMIFRKEPFFHGHNDYDQLVRIAKVLGTEDL
 YDYIDKYNIELDPRFNDILGRHSRKRWERFVHSENQHLVSPEALDFDKLLRYDHSRSLTAREA
 MEHPYFYTVMKQARMGSSSPGGSTPVSSANMMSGISSVPTPSPLGPLAGSPVIAAANPLG
 MPVAAAAGAQQ

>NM004028.txt

MVAFKGVWTQAFWKAVTAEFLAMLIFVLLSLGSTINWGGTEKPLPVDMLISLFCGLSIATMVQ
 CFGHISGGHINPAVTAMVCTRKISIAKSVFYIAAQCLGAIIGAGILYLVTPPSVVGGLGVTMVH
 GNLTAGHGLLVELIITFQLVFTIFASCDSKRTDVTGSIALAIGFSVAIGHLFAINYTGASMPARS
 FGPAVIMGNWENHWIYVWGPIIGAVLAGGLYEYVFCPDVEFKRRFKEAFSKAAQQTGKSYMEV
 EDNRSQVETDDLILKPGVVHVIDVDRGEEKKGDQSGEVLSSV

>NM005143.txt

MSALGAVIALLLWGQLFAVDSGNDVTDIADDGCPKPEIAHGYVEHSVRYQCKNYYKLRTEGD
 GVYTLNDKKQWINKAVGDKLPECEADDGCPKPEIAHGYVEHSVRYQCKNYYKLRTEGDGVY
 TLNNEKQWINKAVGDKLPECEAVCGKPKNPANPVQRILGGHLDAGSFPWQAKMVSHHNLTT
 GATLINEQWLLTTAKNFLNHNENATAKDIAPTLTYVGGKQLVEIEKVVLPNYSQVDIGLIK
 KQKVSVNERVMPICLPSKDYAEVGRVGVSWGWRNANFKFTDHLKYVMLPVADQDQCIRHYE
 GSTVPEKTPKSPVGVQPILNEHTFCAGMSKYQEDTCYGDAGSAFAVHLEEDTWYATGILSF
 DKSCAVAEGVYVKVTSIQDWVQKTIEN

>NM018684.txt

MADEQEIMCKLESIKEIRNKTLMKIKARLKAFFEALSEERHLKEYKQEMDLLLQEKMAHVE
 ELRLIHADINVMENTIKQSENDLNKLESTRRLHDEYKPLKEHVDALRMTLGLQLPDLCEEEEEK
 LSLDYFEKQKAEWQTEPQEPPIPESLAAAAAAQQLQVARKQDTRQTATFRQQPPPMKACLSC
 HQQIHRNAPICPLCKAKSRSRNPKPKRKQDE

>NM024596.txt

MNEHLSSLIKKRKCMQPKDFNFKTPENDKRFQKKEKMAKELQRQKTNLDDVDPILLFESNG
 SLIYPTTIEINSSHHSAMEKRLQEMKEKRENLSPTSSQMIQQSHDNPSNSLCEAPLNISRDTLC
 SDEYFAGGLHSSFDLDCGNSGCGNQRKLEGSINDIKSDVCISSLVLKANNIHSPPSFTHLDK
 SSPQKFLSNLSKEEINLQRNIAGKVVTPHQKQAAGMSQETFEKYRLSPTLSSTKGHLLIHSRP
 RSSSVKRKRKRVSHGSHSPPKEKCKRKRSTRRSIMPRLQLCRSEGRQLQHVAGPALEALSCGESSY
 DDYFSPDNLKERYSENLPPESQLPSSPAQLSCRSLSKKERTSIFEMSDFSCVGGKTRTVDITNF
 TAKTISSPRKTGNTEGRATSSCVTSAPPEALRCCRQAGKEDACPEGNGFSYTIEDPALPKGHD
 DDLTPLEGSLEEMKEAVGLKSTQNKGTTSKISNSSEGEAQSEHEPCFIVDCNMETSTEEKENLP
 GGYSGSVKNRPTRHDVLDSDCDGFKDLIKPHEELKKSGRGKPTRTLVMTSMPSEKQNVVIQ
 VVDKLGFSIAPDVCETTTHVLSGKPLRTLNVLLGIARGCWVLSYDWVWLSLELGHWISEEPF
 ELSHHFAAPLCRSECHLSAGPYRGTLFADQPMFVSPASSPPVAKLCELVHLCGGRVSVQVPR
 QASIVIGPYSGKKKATVKYLSEKWVLDSTQHKVCAPENYLLSQ

>O76273 .txt

MKGISKILSASIMVMKLGNVYSAVPLCSNTYDPSQQQPSYVLIPSTPEAITNCAYSPPKNAYVPS
 SPTTSSSTPGTNNNDNETSPTTEDVGTCKISVVKHCDTPGASSTPCEPEQTIPAQPVTMATVTPA
 IIASVQTPSVVSVIPVTQKVIQPATMIVPPSSIIPGYYPNGTPAAPGQQGQILSGSVLAPGASSC
 QLVPGNTPGQMLPGMTPGVSPCLPTQGGGDSNQITIPGIVYPCQPGQGGSGSNQITIPGIVISPC
 QPGQGGSGSNQITIPGIVYPCQPGQGGSGSNQITIPGIVISPCQPGQGGSGSNQITIPGIVYPCQ
 QNGDGSNQITIPGIVISPCQPGQGGNGGTTGQPGQCVSVPQTPNPIAMPPISGISGNGYPTS

TTYTQSLGQLGPCIDVQKPTSSCESQTNEKSTMQYAMEACAAPTPTVVIGNSEYLVGPGMYSS
LTSPCNCCCQC

>P62242.txt

MGISRDNWHKRRKTGGKRKPYHKKRKYELGRPAANTKIGPRRIHTVRVRGGNKKYRALRLDV
GNFSWGSECCTRKTRIIDVVYNASNELVRTKTLVKN CIVLIDSTPYRQWYESHYALPLGRKKG
AKLTPEEEEEILNKKRSKKIQKKYDERKKNAKISSLLEEQQGKLLACIASRPGQCGRADGYVL
EGKELEFYLRKIKARKGK

>P63210 .txt

MPVINIEDLTEKDKLKMEVDQLKKEVTLERMLVSKCCEEVRDYVEERSGEDPLVKGIPEDKNPF
KELKGGCVIS

>P63234.txt

MSLLTEVETYVLSIVPSGPLKAEIAQRLEDVFAGKNTDLEALMEWLKTRPILSPLTKGILGFVFTL
TVPSERGLQRRRFVQNALNGNDPNMDRAVKLYRKLKREITFHGAKEIALSYSAGALASCMG
LIYNRMGAVTTEVAFLVCATCEQIADSQHRSHRQMVATTNPLIRHENRMVLASTTAKAMEQM
AGSSEQAAEAMEVASQARQMVQAMRAIGTHPSSSAGLKDDLLENLQAYQKRMGVQMQRFK

>P63274.txt

MGRVRTKTVKKAARVIEKYTRLGNDFHNTKRVCIEIAIIPSKKLRNKIAGYVTHLMKRIQRGP
VRGISIKLQEEERERRDNVPEVSALDQEIIEVDPDTKEMLKLLDFGSLSNLQVTQPTVGMNFK
TPRGAV

>P63289 .txt

HSDAVFTDNYTRLRKQMAVKKYLNSILN

>P63302.txt

MALAVRVVYCGACGYKSKYLQKKKLEDEFPGRLDICGEGTPQATGFFEVMVAGKLIHSKKKG
DGYVDTESKFLKLVAAIKAALAQG

>P67778.txt

MAAKVFESIGKFLALAVAGGVVNSALYNVDAGHRAVIFDRFRGVQDIVVGEGTHFLIPWVQK
PIIFDCRSRPRNVPVITGSKDLQNVNITLRILFRPVASQLPRIYTSIGEDYDERVLPSITTEILKSV
VARFDAGELITQRELVSRQVDDLTERAATFGLILDDVSLTHLTFGKEFTEAVEAKQVAQQEAE
RARFVVEKAEQQKAAIISAEGDSKAAELIANSLATAGDGLIELRKLEAAEDIAYQLSRSRNITY
LPAGQSVLLQLPQ

>P67780.txt

MAYPFQLGLQDATSPIMEELLHFHDHTLMIVFLISSLVLYIISLMLTTKLHTSTMDAQEVETVW
TILPAIILILIALPSLRILYMMDEINNPSLTVKTMGHQWYWSYEYTDYEDLNFDSYMIPQELKPG
ELRLLEVDNRVVLPMEMTIRMLISSEDLVLSWAVPSLGLKTD AIPGRLNQTTLMAMRPGLYYGQ
CSEICGSNHSFMPIVLEMVPLSYFETWSALMV

>P67809.txt

MSSEAETQQPPAAPPAAPALSAADTKPGTTGSGAGSGGPGGLTSAAPAGGDKKVIATKVLGTV
KWFNVRNGYGFINRNDTKEDVVFHQTAIKKNNPRKYLRVGDGETVEFDVVEGEKGAEANV
TGPGGVVPVQGSKYAADRNHYRRYPRRRGPPRNYQQNYQNSSESGEKNEGSESAPEGQAQQRR
PYRRRRFPYYMRRPYGRRPQYSNPPVQGEVMEGADNQGAGEQGRPVRQNMYRGRPRFRR
GPPRQRQPREDGNEEDKENQGETQQQPPQRRYRRNFYRRRRPENPKPQDGKETKAADP
PAENSSAPEAEQGGAE

>P67811.txt

MLSLDFLDDVRRMNRQLYYQVLNFGMIVSSALMIWKGLMVITGSESPIVVVLSGSMEPAFHR
GDLLFLTNRVEDPIRVGEIVVFRIEGREIPIVHRVLKIHEKQNGHIKFLTkgdnnavddrGLYKQ
GQHWLEKKDVVGRARGFVPHYIGIVTILMNDYPKFKYAVLFLGLFVLVHRE

>P67997.txt

MANLGCWMLVLFVATWSDLGCKKRPKPGGWNTGGSRYPGQGSPGGNRYPPQGGGGWGQ
PHGGGWGQPHGGGWGQPHGGGWGQPHGGGWGQGGGTHNQWHKPSKPKTSMKHMAGA
AAAGAVVGGGLGGYMLGSAMSRPLIHFGNDYEDRYRENMYRYPNQVYRPPVDQYSNQNNFV
HDCVNITIKQHTVTTTTKGENFTETDVKMMERVVEQMCITQYEKESQAYYQRGSSMVLFSPP
VILLISFLIFLIVG

>P68037.txt

MAASRRLMKELEEIRKCGMKNFRNIQVDEANLLTWQGLIVPDNPPYDKGAFRIEINFPAEYPFK
PPKITFKTKIYHPNIDEKGQVCLPVISAENWKPATKTDQVIQSLIALVNDPQPEHPLRADLAEY
SKDRKKFCKNAEEFTKKYGEKRPVD

>Q6A548.txt

MPRAPRCRAVRALLRGRYREVLPLATFLRRLGPPGRLLVRRGDPAAFRALVAQCLVCVPWGARP
PPAAPCFRQVSCLKELVARVVQRLCERGARNVLAFGFALLDGARGGPPVAFTTSVRSYLPNTVT
ETLRGSGAWGLLLRRVGDVLTLLARCALYLLVAPSCAYQVCGPPLYDLCAPASLPLPAPGLP
GLPGLPGLGAGAGASADLRPTRQAQNSGARRRRGSPGSGVPLAKRPRRSVASEPERGAHRSF
PRAQQPPVSEAPAVTPAVAASPAASWEGGPPGTRPTPAWHPYPGPQGVPHDPAHPETKRFLY
CSGGRERLRPSFLLSALPPTLSGARKLVETIFLGSAPQKPGAARRMRRLPARYWRMRPLFQELL
GNHARCPYRALLRTHCPLRAMAAKEGSGNQAHRGVVICPLERPVAAPQEQTDSLRLVQLLRQ
HSSPWQVYAFLRACLWLVPTGLWGSRHNQRRFLRNVKFFISLKGKAKLSLQELTWKMKVRD
CTWLHGNGACCVPAAEHRREEILARFLVLVDGHIYVVKLLRSFFVYVTTETTFQKNRLLFFYRKS
WSQLQSIGIRQLFNSVHLRELSEAEVRRHREARPALTSRLRFLPKPSGLRPVINDYIMGART
FHRDKKVQHLTSQLKTLFVNLNERARRPSLLGASMLGMDDIHRAWRTFVLRIRAQNPAPQLY
FVKVDVTGAYDALPQDLRVEIVANVIRPQESTYCVRHAYVQRTARGHVRKAFKRHVSTFADL
QPYMRQFVERLQETSLLRDAVVIEQSSSLNEAGSSLFHLFLRLVHNHVVIRIGGKSYIQCGVQP
GSILSTLLCSLCYGMERRLFPGIEQDGVLLRLVDDFLLVTPHLTQAQAFRLTLVKGVPYGCRA
NLQKTAVNFPVEDGALGSAAPLQLPAHCLFPWCGLLLDTRTLEVSCDYSSYAHTSIRASLTFSQ
GAKPGRNMRRKLLAVLRLKCCALFDLQVNGIHTVYMNVYKIFLLQAYRFHACVLQLPFNQPV
KNPSFFLRVIADTASCCYLLKARNAGLSLGAKGASGLFPSEARWLCLHAFLLKLAHHSHTYR
CLLALQAQAAHLRSRQLPRGTAALEAAADPSLTADFKTILD

>Q8EUC5 .txt

MAIKSVLYDKYKDVTKSLMKEFNKSTMEIPRLEKIVINSGLGDATSDSKIVEIGLKLHLITG
QKPIATKSKKSIATFKLREGQAIGAKVTLRRENMWNFLTKLISIAIPRIRDFRGLSTKSFNGN
YTFGIKEQIIFPEIVYDDVKKLRGFDITIVTSAKTDEAMFLLKELGMPFVKTEAK

>Q8EUK2 .txt

MAVPQRKVTHSRKAKRGSHLHLSIPTLVACKRCGKKITPHRVCNSCGYYKNKKVPQIEA

>Q8EUM5.txt

MENKNQKHNEFHEKNQSQKDNVVKKENLHEDQSDLNDANFDDGGKKNKLNKNDIKLN
HLVDSLKNENNNKQAKIESLERQINLLNENFKSEVIKKASEAQTKLDEKIKFQAKYETELKHA
KKYALKSSAIELIDIVSNFELAVNSKVTNPEIANYLKGFQMFANMFKNYFQQNGITEIPVNLNDD
FNAEVMQAFETQKAPNTQPNKVIKIKKGYKLHDIVLVPATVIVSE

>Q8EW28.txt

MEFNFEPSNEKDILKKYSRNLNEEVSANKLNKIIGREQEIRRVEILSRKEKNNPVLIGEPGVGK
TAIVEGFVQKIIISKEVPENLVHCVVYEVNLSLSLIAGTFLQGEFEKRLNALIKEAKQNNGAVILFID
EIHQLMGMGKAGNNSGMDAANILKPIARGEIKIIGATTSNEYRQYIEKDGALERRFQKILVEE
PTPEEALTIMRGLKEKWEIYHKVRIQDNALVASVKLSERYISDKYLPDKAIDLIDEAAAIAKITEA
HTSPSEPINKKIFYLETEKIALSKEEGTNQKERINEIEIELEKQERDLVEKEWKEQKEQVAL
NKIKKEIEKNNWDVERYQNQGEYTEASKILYSVLPKLLKLEAIEKSISENKKVLIKDYISEIDVA
EISRITKIPLNKIFEKEQDKLLNLFNLLKRVKQGQDEALKLVSDTVLKNRVGINNPNRPIGSFLF
VGPTGVGKTEVAKSLAENLFNTEKAIVRINMSEYMEKHSISRLIGAPPGYIGYEQAGELSEQIRR
KPYSVLLDEIEKAHPDILNVLLQVLDGTLKDNQGRNINFKNTIIMTSNVGALYLMENKEDMF
ERELKTSFKPEFLNRIDEIIFNSITMEFAKEIAQKMLDDLEKRLKDNNYQITFDKSVVEYVAKN
GYSKEYGARPINRFIQKTIENTEFITESILKNELIKDRSTIINFANNKLAAILKSN

>Q8EWR9.txt

MKIIKITPRGFCCKGVVDAYATCKKIAKLYPNHEKYLIGWLVHNKEIIELEELGIQTKDDKNHSR
SEIDSIEIKDKNNPPIVIFSAHGTDQKTIDKAREKGLVFDTTTCIYVTKTHDLIKEIEQGYQIFY
IGVNNHPETISTLSIDKSIILIETVNDIENIRTESEKPIFVTNQTTSIYEFEEITELSKYKNIEFK
NDICNAAKDRQDAVINMPSEVDLLLVDGDIKSNNSKLLVEIGIKQIESHLIWNKNIKDEWFI
NKKCLAITSGCSTPTWLANYVIIIFLEKKLGTAND

>Q8EX21 .txt

MATIAQLIRHDRKDKFKKSKSPALMYTNSLKKKRTYNPSPYKRGVCTRVGTMTPKKPNSALR
KYAKVKLTNGYEVLAYIPGEGHNLQEHSVVLIEGGRVKDLPGVRYHIVRGTLDTSGVEKRRQQ
RSGYGAKRPKEKKE

>Q8EX25.txt

MLLVAFKDPRIKAVNLIGGQAKPGPQLASIGINMGEFTKQFNDQTKDKNGEVIPCITAFKD
KSFTFEIKTTPVTMLLKAANIKVGAKNSKTETVATISREKALEIAKTKLVDTNANDEEAVLRMV
AGSAKQMGIKIEGVDPVVHKDGKKK

>Q8MKI5.txt

MENTENSVDAKSFKNAETKILHGSKSMDSGMSFDNSYKMDYPEMGLCIIINNKNFHKSTGMA
PRSGTDVDAANLRETFTNLKYEVRNKNDLTCEEIILEMNSVSKEDHRSKSSFCVLLSHGDEGI
IFGTNGPVDLRKVTGFFRGDYCRSLTGKPKLFIQACRGTELDGDIETDSGIEDDMACQKIPVE
ADFLYAYSTAPGYYSWRNSKDGSWFIQSLCAMLKLYAHKLEFMHILTRVNRKVATEFESFSLD
SAFHGKKQIPCIVSMLTKELYLYH

>Q8NB59.txt

MAIEGGERTCGVHELICIRKVSPEAVGFLSAVGVFIIIMLLFLYINKKFCFENVGGFPDLGSEYS
TRKNSQDKIYNSYMDKDEHGSSSESEDEALGKYHEALSRTHNSRLPLADSRQRNYAWETRQK
YSPLSAEYDGYSSSEASIDEGNCIQMRRTPLDELQPPPYQDDSGSPHLSCTPSEIGDSKCEFS
HCSNSPRCSYNKCPSEGSTGHEIESFHNGGYEEDVPSDSTAVLSPEDMSAQGSSSQLPKPFD
EPEAKYGLDVTDFDYDSQEQLLVTVTAVTDIPTYNRTGGNSWQVHLVLLPIKKQRAKTSIQRG
PCPVFTETFKFNHVESEMIGNYAVRFRLYGVHRMKKEKIVGEKIFYLTKLNLQGKMSLPVILEPS
YNHSGCDSQMSVSEMSCSESTSSCQSLEHGSVPEILIGLLYNATTGRLSAEVIKGSFKNLAA
NRPPNTYVKLTLNLSMGQEMSKCKTSIRRQPNPVYKETFVVFQVALFQLSDVTLILSVYNKRSM
KREEMIGWISLGLNSSGEEELNHWTEMKESKGGQVCRWHALLES

>Q99J21.txt

MATPAGRRASETERLLTPNPGYGTQVGTSPAPTTPTTEEDLRRRLKYFFMSPCDKFRAKGRKPC
KLMLQVVKILVVTVQLILFGLSNQLVVTFREENTIAFRHLFLLGYSDGSDDTFAAYTQEQLYQAI

FYAVDQYLILPEISLGRYAYVRGGGGPWANGSALALCQRYYHRGHVDPANDTFDIDPRVVTDC
 IQVDPDRPPDIPSEDLDFLDGASASYKNLTLKFHKLINVTIHFQLKTINLQSLINNEIPDCYTF SIL
 ITFDNKAHSGRIPIRLETKTHIQECKHPSVSRHGDNSFRLLFDVVVILTCSLSFLLCARSLLRGFL
 LQNEFVVMWRRRRGREISLWERLEFVNGWYILLVTSVLTISGTMKIGIEAKNLA SYDVCSIL
 LGTSTLLVWGVIRYLTFHFKYNI LIATLRVALPSVMRFCCCVAVIYLYGFCGWIWLGYPYHVKF
 RLSMVSECLFSLINGDDMFVTF AAMQAQQGHSSLVWLFSQLYLYSFISLFIYMVLSLFIALITG
 AYDTIKHPGGTGT EKSELQAYIEQCQDSPTS GKFRRGSGSACSLFCCCGRDSPEDHSLLVN//

>Q9NUX5.txt

MSLVPATNYIYTPLNQLKGGTIVNVYGVVKKFPYLSKGTDYCSVVTIVDQTNVKTCLLFSGN
 YEALPIIYKNGDIVRFHRLKIQVYKQETQGITSSGFASLT FEGTLGAPIIPRTSSKYFNFTTEDHK
 MVEALRVWASTHMSPSWTLLKCDVQPMQYFDLTCQLLGKAEVDGASFLKLVWDGTRTPFPS
 WRVLIQDLVLEGDLSHIHRLQNLTIIDILVYDNHVVHVARSLKVGSLRISLHTKLQSMNSENQT
 MSLFEHLHGGSYGRGIRVLPESNSDQVLDKLDLESANLTANQHS DVICQSEPDDSFSSGS
 VSLYEVEQCQLSATILTDHQLERTPLCAILKQKAPQQYRIRAKLRSYKPRRLFQSVKLHCPKC
 HLLQVEPHEGDLDIIFQDGATKTPDVKLQNTSLYDSKIWTTKNQKGRKVAVHFVNNGILPLS
 NECLLIEGGTLSEICKLSNKFNSVIPVRS GHEDLELLDLSAPFLIQGTIHHYGCKQCSSLRSIQN
 LNSLVDKT SWIPSSVAEALGIVPLQYVFMFTLDDGTGVLEAYLMDSDKFFQIPASEVLMDD
 LQKSVDMIMDMFCPPGIKIDAYPWLECFIKSYNVTNGTDNQICYQIFDTTVAEDVI

>Q9NWW6.txt

MKTFIIGISGVTNSGKTLAKNLQKHLPNCSVISQDDFFKPESEIETDKNGFLQYDVLEALNME
 KMMSAISCWME SARHSVST DQESAEIPIIIEGFLFNKPLDTIWNRSYFLTIPYEECKRRR
 STRVYQPPDSPGYFDGHVWPMYLYRQEMQDITWEVVYLDGTKSEEDLFLQVYEDLIQELAKQ
 KCLQVTA

>Q9NXK6.txt

MLSLKLPRLFSIDQIPQVFHEQGILFGYRHPQSSATA CILSLFQMTNETLNIWTHLLPFWFFAWR
 FVTALYMTDIKNDSSWPM LVYMCTSCVYPLVSSCAHTFSSMSKNARHICYFLDYGAVNLFSL
 GSAIAYSAYTFDALMCTTFHDYVALAVLNTILSTGLSCYSRFLEVQKPRLCKVIRVLAFAYPYT
 WDSLPIFYRLFLFPGESAQNEATSYHQHMIMTLLASFLYSAHLPERLAPGRFDYIGHSHQLFH
 VCVILATHMQMEAILLDKTLRKEWLLATSKPFSFSQIAGAILLCIIFSLSNIIYFSAALYRIPKPEL
 HKKET

>Q9UKH8.txt

MGQTKSKIKSKYASYLSFIKILLKRGVVKVSTKNLIKLFQIIEQFCPWFPEQGTLDLKDWKRIGK
 ELKQAGRKGNIIP LTVWNDWAIKAALPEFQTEEDSVSVSDAPGSCLIDCNEKTRKKSQKEME
 GLHCEYVAEPVMAQSTQNV DYNQLQEVYIPETLKLEGKGPVLPLESKP

>Q9Y3S2.txt

MPKKKTGARKKAENRREREKQLRASRSTIDLAKHPCNASMECDKQCRRQKNRAFCYFCNSVQ
 KLPICACQCGKTKCMMKSSDCVIKHAGVYSTGLAMVGAICDFCEAWVCHGRKCLSTHACACPL
 TDAECVE CERGVWDHGGRIFSCSFCHNFLCEDDQFEHQASCQVLEAETFKCVSCNRLGQHSC
 LRCKACFDDHTRSKVFKQEKGKQPPCKCGHETQETKDLMSSTRSLKFRQTGGEEGDGAS
 GYDAYWKNLSSDKYGDTSYHDEEEDEYEAEDEEEDEEGRKDSDESSDLFTNLNLGRTYAS
 GYAHYEEQEN

>Q9Y6I0.txt

WASQVSENRPVCKAIIQGKQFEGLVDTGADVSI IALNQWPKNWPKQKAVTGLVGIGTASEVY
 QSTEILHCLGPDNQESTVQPMITSIPLNLWGRDLLQQWGAEITMPAPSYSPSTSQKIMTKMGYIP
 GKGLGKNEDGIKIPVEAKINQEREGINPC

>protmutt1.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRGLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNRRPIL
TIITLEDSSGNLLGRNSFEVRVCACPGDRDRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

>protmutt10.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKT

>protmutt11.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTX

>protmutt1176.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRGLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNRRPIL
TIITLEGSSGNLLGRNSFEVRVCACPGDRDRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

>protmutt1179.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRGLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNRRPIL
TIITLEESSGNLLGRNSFEVRVCACPGDRDRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

>protmutt1181.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRGLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNRRPIL
TIITLEDTSGNLLGRNSFEVRVCACPGDRDRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

>protmutt1184.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRGLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNRRPIL
TIITLEDYSGNLLGRNSFEVRVCACPGDRDRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP

QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS

>protmutt1187.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS

>protmutt1192.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS

>protmutt12.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS

>protmutt13.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKT

>protmutt14.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTY

>protmutt15.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYLGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS

>protmutt1501.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPPHHERCSDSDGLA

PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNM CNSSCMGGMNRRPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEP GGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS D

> protmutt1505.txt

MEEPQSDPSVEPPLSQETFS DLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKT CPVQLWVDSTPPPGTRVRAMAIYKQS QHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNM CNSSCMGGMNRRPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPSGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEP GGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS D

> protmutt1510.txt

MEEPQSDPSVEPPLSQETFS DLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKT CPVQLWVDSTPPPGTRVRAMAIYKQS QHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNM CNSSCMGGMNRRPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPRSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEP GGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS D

> protmutt1540.txt

MEEPQSDPSVEPPLSQETFS DLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKT CPVQLWVDSTPPPGTRVRAMAIYKQS QHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNM CNSSCMGGMNRRPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEP GGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS D

> protmutt16.txt

MEEPQSDPSVEPPLSQETFS DLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYHGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKT CPVQLWVDSTPPPGTRVRAMAIYKQS QHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNM CNSSCMGGMNRRPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEP GGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS D

> protmutt17.txt

MEEPQSDPSVEPPLSQETFS DLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYHGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKT CPVQLWVDSTPPPGTRVRAMAIYKQS QHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNM CNSSCMGGMNRRPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEP GGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS D

> protmutt18.txt

MEEPQSDPSVEPPLSQETFSDLWKLLENVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQRSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHILIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNRRPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPXXXPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS

>protmutt19.txt

MEEPQSDPSVEPPLSQETFSDLWKLLENVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQCSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHILIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNRRPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPXXXPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS

>protmutt2.txt

MEEPQSDPSVEPPLSQETFSDLWKLLENVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPS

>protmutt20.txt

MEEPQSDPSVEPPLSQETFSDLWKLLENVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQVSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHILIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNRRPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPXXXPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS

>protmutt21.txt

MEEPQSDPSVEPPLSQETFSDLWKLLENVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHILIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNRRPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPXXXPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDS

>protmutt22.txt

MEEPQSDPSVEPPLSQETFSDLWKLLENVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGGYGFRLGFLHSGTAKSVTCTY
SPALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGL
APPQHILIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNRRPI
LTIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSS
PQPXXXPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTS
RHKKLMFKTEGPDS

>protmutt23.txt

MEEPQSDPSVEPPLSQETFSDLWKLLENVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGRYGFRLGFLHSGTAKSVTCTYS

PALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGSDCCTTIHYNMCMNSSCMGGMNRRLPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRRDRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

> protmutt24.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGS

> protmutt25.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGS

> protmutt26.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGSDCCTTIHYNMCMNSSCMGGMNRRLPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRRDRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

> protmutt27.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGSDCCTTIHYNMCMNSSCMGGMNRRLPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRRDRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

> protmutt28.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYDFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGSDCCTTIHYNMCMNSSCMGGMNRRLPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRRDRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

> protmutt29.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGSDCCTTIHYNMCMNSSCMGGMNRRLPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRRDRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

> protmutt3.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
 PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVSPSRKTYQGSYGFRLGFLHSGTAKSVTCTYS
 PALNKMFCQLAKTQCPVQLWVDSTPPPQTRVRAMAIYKQSQHMTEVVRRCPPHHERCSDSDGLA
 PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
 TIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
 QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTS
 HKKLMFKTEGPDS

>protmutt30.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
 PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVSPKTYQGSYGLRGLHSGTAKSVTCTYS
 PALNKMFCQLAKTQCPVQLWVDSTPPPQTRVRAMAIYKQSQHMTEVVRRCPPHHERCSDSDGLA
 PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
 TIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
 QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTS
 HKKLMFKTEGPDS

>protmutt31.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
 PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVSPKTYQGSYGSRLGFLHSGTAKSVTCTY
 SPALNKMFCQLAKTQCPVQLWVDSTPPPQTRVRAMAIYKQSQHMTEVVRRCPPHHERCSDSDGL
 APPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR PI
 LTIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSS
 PPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTS
 RHKLMFKTEGPDS

>protmutt32.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
 PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVSPKTYQGSYGCRLGFLHSGTAKSVTCTY
 SPALNKMFCQLAKTQCPVQLWVDSTPPPQTRVRAMAIYKQSQHMTEVVRRCPPHHERCSDSDGL
 APPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR PI
 LTIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSS
 PPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTS
 RHKLMFKTEGPDS

>protmutt33.txt

MEEPQSDPSVKPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
 PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVSPKTYQGSYGFRLGFLHSGTAKSVTCTYS
 PALNKMFCQLAKTQCPVQLWVDSTPPPQTRVRAMAIYKQSQHMTEVVRRCPPHHERCSDSDGLA
 PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
 TIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
 QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTS
 HKKLMFKTEGPDS

>protmutt4.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
 PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVSPKTYQGSYGFRLGFLHSGTAKSVTCTYS
 PALNKMFCQLAKTQCPVQLWVDSTPPPQTRVRAMAIYKQSQHMTEVVRRCPPHHERCSDSDGLA
 PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
 TIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP

QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

> protmutt5.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQ

> protmutt6.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQRTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

> protmutt7.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQNTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

> protmutt740.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDVRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
TIITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

> protmutt742.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDERNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPILT
IITLEDSSGNLLGRNSFEVRVCACPGRDRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSPQ
PKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSRH
KKLMFKTEGPDSD

> protmutt745.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVVRRCPHHERCSDSDGLA
PPQHLIRVEGNLRVEYLDD

> protmutt748.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
 PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
 PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVRRCPHHERCSDSDGLA
 PPQHLIRVEGNLRVEYLDDINTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPILT
 IITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSPQ
 PKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSRH
 KKLMFKTEGPDS

>protmutt750.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
 PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
 PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVRRCPHHERCSDSDGLA
 PPQHLIRVEGNLRVEYLDDRHTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
 TIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
 QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
 HKKLMFKTEGPDS

>protmutt753.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
 PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
 PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVRRCPHHERCSDSDGLA
 PPQHLIRVEGNLRVEYLDDRITFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPILT
 IITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSPQ
 PKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSRH
 KKLMFKTEGPDS

>protmutt755.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
 PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
 PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVRRCPHHERCSDSDGLA
 PPQHLIRVEGNLRVEYLDDRNAFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
 TIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
 QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
 HKKLMFKTEGPDS

>protmutt8.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
 PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKIYQGSYGFRLGFLHSGTAKSVTCTYS
 PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVRRCPHHERCSDSDGLA
 PPQHLIRVEGNLRVEYLDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
 TIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP
 QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
 HKKLMFKTEGPDS

>protmutt9.txt

MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEA
 PRMPEAAPRVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAKSVTCTYS
 PALNKMFCQLAKTQCPVQLWVDSTPPPGTRVRAMAIYKQSQHMTEVRRCPHHERCSDSDGLA
 PPQHLIRVEGNLRVEYLDRNTFRHSVVVPYEPPEVGS DCTTIHYNMCMNSSCMGGMNR RPIL
 TIITLEDSSGNLLGRNSFEVRVCACPGRRRTEENLRKKGEPHHELPPGSTKRALPNNTSSSP

QPKKKPLDGEYFTLQIRGRERFEMFRELNEALELKDAQAGKEPGGSRAHSSHLKSKKGQSTSR
HKKLMFKTEGPDSD

BIBLIOGRAFIE

- [AHO'74] **Aho, V.A., Hopcroft, J.E., Ullman, J.D.:** „ *The design and Analysis of Computer Algorithms*”, Addison-Wesley, Menlo Park, California, 1974;
- [ALTS'90] **Altschul, S., Gish, W., Miller, W., Myers, E., Lipman, D.J.:** "*Basic local alignment search tool*", J. Mol. Biol., vol. 215, no.3, pp:403-410, 1990.
- [ALTS'97] **Altschul, S., Madden, T., Schaffer, A., Zhang, J., Zhang, Y., Miller, W., Lipman, D.:** "*Gapped Blast and PSI-Blast: a new generation of protein database search programs*", Nucleic Acids Research, vol.25, no.17, pp: 3389-3402 1997.
- [ANAS'03] **Anastasiadis, A.D., Magoulas, G.D., Liu, X.:**"*Classification of Protein Localisation Patterns via Supervised Neural Network Learning*", Proceedings of the Fifth Symposium on Intelligent Data Analysis (IDA-03), Berlin, Germany, 2003.
- [ARRA'84a] **Arratia, R.A.,and Watterman, M.S.:** "An Erdos-Renyi Low with Shifts", Adv. Appl. Math.(in press).
- [ARRA'84b] **Arratia, R.A., Gordon, L., and Watterman, M.S.:** "An Extreme Value Distribution for Sequence Matching", manuscript, 1984.
- [ARSL'01] **Arslan, N.A., Egecioglu Ö., Pevzner, P.A.:** "*A new approach to sequence comparison: normalized sequence alignment*", Bioinformatics, vol.17 no.4 , pp: 327-337, 2001.
- [BACO'86] **Bacon, D.J., and Anderson, W.F.:** "*Multiple sequence alignment*", J. Mol.Biol. vol. 191, pg: 153-161, 1986.
- [BAEZ'99] Baeza-Yates, R., Ribeiro-Neto, B., in "Retrieval Evaluation", *Modern Information Retrieval*, Ed. Addison Wesley, 1999, pp:75-81.
- [BAIL'94] **Bailey, T.L., Elkan C.:** "*Fitting a mixture model by expectation maximization to discover motifs in biopolymers*", Proceedings of International Conference on Intelligent Systems for Molecular Biology, PubMed, 2:28-36, 1994.
- [BALD'94] **Baldi, P., Chauvin, Y., Hunkapiller, T., McClure, M.A.:** "*Hidden Markov models of biological primary sequence information*", Proc. Natl.Acad.Sci.USA 91(3):1059-1036, 1994.

- [BART'96] **Barton, G.J.:** *"Protein Sequence Alignment and Database Scanning"* in M.J. E. Sternberg, *Protein Structure Prediction - a practical approach*, IRL Press at Oxford University Press, 1996.
- [BASH'87] **Bashford, D., Chothia, C. and Lesk, A.M.:** *"Determinants of a protein fold: Unique features of the globin amino acid sequences"*, *J.Mol.Biol.* 196:199-216, 1987.
- [BATE'04] **Bateman, A., Coin, L., Durbin, R., Finn, R.D., Hollich, V., Griffiths-Jones, S., Khanna, A., Marshall, M., Moxon, S., Sonnhammer, E.L., Studholme, D.J., Yeats, C., Eddy, S.R.:** *"The Pfam protein families database"*, *Nucleic Acids Research*, 32, 2004.
- [BELL'00] **Bellagarda, J.:** *"Exploiting latent semantic information in statistical language modeling,"* *Proceedings IEEE*, vol. 88, no.8, pp:1279-1296, 2000.
- [BENS'99] **Benson, G.:** *"Tandem repeats finder: a program to analyze DNA sequences"*, *Nucleic Acids Res.*, 27(2):573-80, 1999.
- [BERG'85] **Berger, J.:** *"Statistical Decision Theory and Bayesian Analysis"*, Springer-Verlag, New York, 1985.
- [BIG'05] **BioInformatics Glossary**,
<http://big.mcw.edu/display.php/2241.html>
- [BLAN'00] **Blanchette, M., Schwikowski, B., Tompa, M.:** *"An exact algorithm to identify motifs in orthologous sequences from multiple species"*, *Proceedings of the 8th International Conference on Intelligent Systems for Molecular Biology (ISMB)*,: 37-45, 2000.
- [BOGA'03a] **Bogan-Marta, A., Wambacq, P.:** *"From Latent Semantic Indexing to Text Categorization?"*, published in the Journal "Analele Universitatii Bucuresti", No. LII, seria Matematica- Informatica", vol. 1, 2003, ISSN 1224-7170, pg. 87-96.
- [BOGA'03b] **Bogan-Marta, A., Robu, N.:** *"Investigations Over ``add one'' Smoothing Method for Speech Recognition"*, *Proceedings of EMES'03*, May 24-26, Oradea, Romania, pp.93-97.
- [BOGA'03c] **Bogan-Marta, A., Robu, N.:** *"Indicators for Text Relevance Using Techniques from Information Theory Field"*, *Proceedings of SINTES 11th Edition*, Section Computer Science & Engineering, Craiova, Romania, 23-24 October 2003, ISBN 973-8043-415-5, ISBN 973-8043-415-5 pg. 286-290.

- [BOGA'04a] **Bogan-Marta, A., Robu, N.:** "*The Use of Statistical Models For Large Text Data Analysis*", Analele Universitatii din Craiova, Seria Inginerie Electrica, no. 27, vol.1, ISSN 1223-530X, pp.87-92, 2004.
- [BOGA'04b] **Bogan-Marta,A., Győrödi, R., Győrödi,C.:** "*Implications of syntactic and semantic categories of grammars on a speech recognition tool*", in Proceedings of The Sixth International Scientific Conference, ECI 2004, Kosice-Herľany, Slovakia, pp. 86-90, 2004.
- [BOGA'05a] **Bogan-Marta,A., Gavrielides,M.A., Pitas, I., Lyroudia,K.:** "*A New Statistical Measure of Protein Similarity based on Language Modeling*", GENSIPS 05, IEEE International Workshop on Genomic Signal Processing and Statistics , Newport, Rhode Island, ISBN 9-812-38846-X, 2005.
- [BOGA'05b] **Bogan-Marta,A., Laskaris,N., Gavrielides,M.A., Pitas, I., Lyroudia,K.:** "*A novel efficient protein similarity measure based on n-gram modeling*", CIMED 2005, IEEE, IEE, Second International Conference on Intelligence in Medicine and Healthcare, Costa da Caparica, Lisbon, Portugal, ISBN:0-86341-520-2, pp. 122-127,2005.
- [BOGA'06c] **Bogan-Marta, A., Robu, N., Pater, M.:** "*String comparison in terms of statistical evaluation applied on biological sequences*", IEEE, Proceedings of ICCCC'06, Baile Felix Spa-Oradea, pp.86-91, 2006.
- [BOGA'06d] **Bogan-Marta,A., Pitas, I., Lyroudia,K.:** "*Statistical Method of Context Evaluation for Biological Sequence Similarity*", IFIP 2006, Santiago de Chile, published in 'Artificial Intelligence in Theory and Practice', Springer , 2006.
- [BOGA'06e] **Bogan-Marta, Robu,N.:** "*A study of Sequence Clustering on Protein's Primary Structure using a Statistical Method*", Acta Polytechnica Hungarica, Journal of Applied Sciences at Budapest Tech, Hungary, vol3, issue3, ISSN 1785-8860, pp.17-27, 2006.
- [BOGA'06a] **Bogan-Marta,A., Robu,N.:** "*Markow chains for a New Datamining Techique applied on Biological Sequences*", CONTI'2006, The 7th International Conference on Technical Informatics, 8-9 June 2006, Timisoara, Romania, pp.227-230.
- [BOGA'06b] **Bogan-Marta,A., Robu,N.:** "*Clustering Sequences with a Statistical Content Evaluation Method*", SACI2006, 3rd Romanian-Hungarian Joint Symposium on Applied Computational Intelligence, Timișoara, Romania, May 25-26, 2006, pp.59-69.
- [BORD'91] **Bordo, D., and Argos, P.:** "*Suggestions for "safe" residue substitutions in site-directed mutagenesis.*", J Mol Biol., Feb 20; 217(4), pp:721-9, 1991.
- [BOSW'84] **Boswell, D.R., and MacLachlan, A.D.:** "*Sequence Comparison by Exponentially-damped Allignment*", Nucleic Acid Res., 12, 457-467, 1984.

- [BREJ'00] **Brejová, B., DiMarco, C., Vinař, T., Hidalgo, S.R., Holguin, G., Patten, C.:** "*Finding Patterns in Biological Sequences*", Project report for CS798g, Fall 2000.
- [BROW'92] **Brown, P.F., Della Pietra, A. S., Della Pietra, V.J., Mercer Robert, L.R., and Jennifer, C.L.:** "An estimation of an upper bound for the entropy of English", in *Association for Computational Linguistics*, Yorktown Heights, NY 10598, P.O. Box 704, 1992.
- [BYER'84] **Byers, T.H., and Waterman, M.S.:** "*Determining All Optimal and Near-optimal Solutions When Solving Shortest Path Problems by Dynamic Programming*", *Operat. Res.* (in press), 1984.
- [CAMO'03] **Camoglu, O., Kahveci, T., Singh, A.K.:** "*PSI: indexing protein structures for fast similarity search*". *PubMed, Bioinformatics*, vol. 19 Suppl. 1, pp: i81-i83, 2003.
- [CHIA'04] **Chiang, J-H., Yu, H-C., Hsu, H-J.:** "*Gis: a biomedical text-mining system for gene information discovery*", *Bioinformatics*, 20(1): 120-121, 2004.
- [CLAV'93] **Claverie, J-M.:** "*Detecting frame shifts by amino acid sequence comparison*", *J. Molecular Biology* 234: 1140-1157, 1993.
- [COHE'75] **Cohen, D.N., Reichert, T.A., and Wong, A.K.C.:** "*Matching Code Sequences Utilizing Context Free Quality Measures*", *Math. Biosci.*, 24, pp.25-30, 1975.
- [COLL'84] **Collins, J.F., and Coulson, A.F.W.:** "*Applications of Paralell Processing Algorithms for DNA Sequence Analysis*", *Nucleic Acids Research.*, 12, 181-192,
- [COWA'98] **Coward E., Drablos, F.:** "*Detecting periodic patterns in biological sequences*", *Bioinformatics*, 14: 498-507, 1998.
- [DAVI'92] **Jones, D.T, Taylor, W.R., and Thornton, J.M.:** "*The rapid generation of mutation data matrices from protein sequences*", *CABIOS* 8: 275-282, 1992.
- [DAYH'78] **Dayhoff, M.O.:** "*Atlas of Protein Sequence and Structure*", *Natl. Biomed. Res. Found.*, Washington, Vol. 5, Suppl. 3, pp. 345-352, 1978.
- [DELC'75] **Delcoigne, A., and Hansen, P.:** "*Sequence Comparison by Dynamic Programming*", *Biometrika*, 62, 661-664, 1975.
- [DEMA'95] **De Marcken, K. :** "*On the unsupervised induction of phrase-structure grammars*", *SIGDAT*, 1995,
<http://www.demarcken.org/carl/papers/sigdat.pdf>
- [DESP'02] **M. Deshpande, M., Karypis, G.:** "*Evaluation of techniques for Classifying Biological Sequences*", *Proceedings of the 6th Pacific-Asia Conference on Advances in Knowledge Discovery and Data Mining*, : 417 - 431, 2002.
- [DINK'67] **Dinkelbach, W.:** "*On nonlinear fractional programming*", *Management Science*, vol. 18, no.7, pp: 492-498, 1967.

- [DUDA'01] **Duda, R.O., Hart, P. E., Stork., D. G.,:** " *Pattern classification*", Information Theory, John Wiley&Sons, Inc., 2nd edition, 2001.
- [DUMA'82] **Dumas, J.P., and Nino, J.,:** " *Efficient Algorithms for Folding and Comparing Nucleic Acid Sequences*", Nucleic Acids Res., 80, 197-206, 1082.
- [DUME'56] **Dumey, A.I.,:** " *Indexing for Rapid Random-access Memory*", Comput. Automat., 5, 6-8, 1956.
- [DURB'98] **Durbin, R., Eddy, S., Krogh, A., Mitchison, G.,:** " *Biological Sequence Analysis. Probabilistic models of proteins and nucleic acid,*" Cambridge University Press, 1998.
- [EDGA'04] **Edgar, R.C.,:** " *MUSCLE: multiple sequence alignment with high accuracy and high throughput*", Nucleic Acids Research 32(5), 1792-97, 2004.
- [ERHA'80] **Erhan,S., Marzolf,T., Cohen,L., :** " *Amino-acid neighborhood relationships in proteins: breakdown of amino-acid sequences into overlapping doublets, triplets and quadruplets.*" Int. J. Biomed Comput, vol. 11(1), pp:67-75, 1980.
- [ERIC'83] **Erickson, B.W., and Sellers, P.H.,:** " *In Time Warps, String Edits, and Macromolecules: the Theory and Practice of Sequence Comparison*", Eds. D. Sankoff and J.B. Kruskal, Addison-Wesley, London, pp. 55-90, 1983.
- [ESKI'01] **Eskin, E., Grundy, W.N., and Singer, Y.,:** " *Using mixtures of common ancestors for estimating the probabilities of discrete events in biological sequences,*" Bioinformatics, vol.17 Suppl 1, pp: 65-73, 2001.
- [ESTE'04] **Ester, M., Zhang, X.,:** " *A Top-Down Method for Mining Most-Specific Frequent Patterns in Biological Sequences*", Proceedings of the Fourth SIAM International Conference on Data Mining, Lake Buena Vista, Florida, USA, April 22-24, 2004.
- [FICK'84] **Fickett, J.W.,:** " *Fast Optimal Alignment*", Nucleic Acids Res., 12: 175-180, 1984.
- [FITC'67] **Fitch, W.M., and Margoliash, E.,:** " *Construction of Polygenetic Trees*", Science, 155, 279-284, 1967.
- [FITC'69] **Fitch, W.M.,:** " *Locating Gaps in Amino Acid Sequences to Optimize the Homology Between Two Proteins*", Biochem. Genet., 3, 99, 1969.
- [FITC'71] **Fitch, W.M.,:** " *Towards Defining the Course of Evolution: Minimum Change for a Specific Tree Topology*", Syst.Zool.20, 406-416, 1971.
- [FITC'83] **Fitch, W.M., and Smith, T.F.,:** " *Optimal Sequence Alignments*" , Proc., natn., Acad. Sci. U.S.A., 80, 1382-1386, 1983.
- [FLOR'96] **Floroiu, C.,:** " *Evaluarea produselor software*", arhiva revista Byte Romania, disponibil online la: <http://www.byte.ro/byte96-05/eva.html>, (accesibil la 1 Septembrie 2006).

- [GANA'04] **Ganapathiraju, M.K., Klein-Seetharaman, J., Balakrishnan, N., and Reddy, R.:** "Characterization of protein secondary structure-application of latent semantic analysis using different vocabulary," IEEE Signal Processing Magazine, vol. 21, no.3, pp: 78-87, 2004.
- [GANA'04b] **Ganapathiraju, M., Manoharan, V., and Klein-Seetharaman, J.:** "Statistical sequence analysis using n-grams", J. Appl. Bioinformatics, vol.3 (2), pp:193-200, 2004
- [GOAD'82] **Goad, W.B., Kanehisa, M.I.:** "Pattern Recognition in Nucleic Acid Sequences I. A General Method for Finding Local Homologies and Symmetries", Nucleic Acid Res., 10: 247-263, 1982.
- [GONN'92] **Gonnet, G.H., Cohen, M.A., Benner, S.A.:** "Exhaustive Matching of the Entire Protein Sequence Database", Science 256: 1443-1445, 1992.
- [GORD'73] **Gordon, A.D.:** "A Sequence Comparison Statistic and Algorithm", Biometrika, 60, 197-200, 1973.
- [GORO'97] **Gorodkin, J., Heyer, L. J., Stormo, G. D.:** "Finding the most significant common sequence and structure motifs in a set of RNA sequences", Nucleic Acids Research, 25(18): 3724-3732, 1997.
- [GOTO'82] **Gotoh, O.:** "An improved algorithm for matching biological sequences." J. Mol. Biol. vol.162, pp:705-708, 1982.
- [GRIBS'87] **Gribkov, M., McLachlan, A.D., and Eisenberg, D.:** "Profile analysis: detection of distantly related proteins." Proceedings of the National Academy of Sciences of the USA, 84:4355-4358, 1987.
- [HENI'92] **Henikoff, S., Henikoff, J.G.:** "Amino acid substitution matrices from protein blocks", in Proceedings of Natural Academic Sciences, 89:10915--10919, 1992.
- [HERT'90] **Hertz, G.Z., Hartzell, G.W., III, Stormo, G.D.:** "Identification of consensus patterns in unaligned DNA sequences known to be functionally related", Comput. Appl. Biosci. 6, 81-92.
- [HIGG'94] **Higgins, D., Thompson J., Gibson, T., Thompson, J.D., Higgins, D.G., Gibson, T.J.:** "CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice", Nucleic Acids Res. 22:4673-4680, 1994.
- [HUGH'96] **Hughey, R., and Krogh, A.:** "Hidden Markov models for sequence analysis: extension and analysis of the basic method." Computer Applications in the Biosciences, 12: 95-107.
- [JAIN'88] **Jain, A.K., Dubes, R.C.:** "Algorithms for Clustering Data", Prentice Hall, Englewood Cliffs, New Jersey, 1988.
- [JAKK'00] **Jaakkola, T., Diekhans, M., and Haussler, D.:** "A discriminative framework for detecting remote protein homologies, " J. Computational Biology, vol. 7, pp: 95-114, 2000.

- [JOHN'86] **Johnson, M.S., and Doolittle, R.F.:** "A method for the simultaneous alignment of three or more amino acid sequences," *J. Mol. Evol.*, vol. 23, pp: 267-278, 1986.
- [JURA'00] **Jurafsky, D., Martin, J. H.:** "Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition." Prentice Hall, 1st edition, January 26, 2000.
- [KALT'04] **Kaltchenko, A.:** "Algorithms for Estimating Information Distance with Application to Bioinformatics and Linguistics", The Computing Research Repository disponibilă online la : <http://arxiv.org/abs/cs.CC/0404039> (la data de 8 Septembrie 2006), 2004.
- [KANE'82] **Kanehisa, M.I., and Goad, W.B.:** „Pattern Recognition in Nucleic Acid Sequences II. An Efficient Method for Finding Locally Stable Secondary Structures”, *Nucleic Acids Research*, 10, 265-277, 1982.
- [KARC'98] **Karchin R., Hughey, R.:** "Weighting hidden Markov models for maximum discrimination", *Bioinformatics*, 14(9): 772-782, 1998.
- [KARL'83] **Karlin, S., Ghandour, G., Ost, F., Tavaré, S., and Korn, L.J.:** „New Approaches for Computer Analysis of Nucleic Acid Sequences”, *Proc. Natn. Acad. Sci. U.S.A.*, 80, 5660-5664, 1983.
- [KARL'84] **Karlin, S., Ghandour, G., and Foulser, D.E.:** „Comparative Analysis of Human and Bovine Papillomaviruses”, *Mol.Biol.Evol.*, 1, 357-370, 1984.
- [KARL'91] **Karlin, S., Bucher, P., Brendel, V., and Altschul, S.F.:** "Statistical methods and insights for protein and DNA sequences", *Annu Rev Biophys Chem.*, vol 20, 175-203, 1991.
- [KARL'96] **Karlin, S., Bruge, C.:** "Trinucleotide repeats and long homopeptides in genes and proteins associated with nervous system disease and development", *Proc Natl Acad Sci USA*, vol. 93(4), pp:1560-1565, 1996.
- [KASI'99] **Kasif, S.:** "Datascop: Mining Biological Sequences", *IEEE Intelligent Systems*, : 38-43, Nov/Dec. 1999.
- [KATT'00] **Katti, M.V., Sami-Subbu, R., Ranjekar, P.K., and Gupta, V.S.:** "Amino acid repeat patterns in protein sequences: Their diversity and structural-functional implications". *Protein Sci.*, vol. 9, pp: 1203-1209, 2000.
- [KORN'77] **Korn, L.J., Queen, C.L., and Wegman M.N.:** „Computer Analysis of Nucleic Acid Regulatory Sequences”, *Proc. Natn. Acad. sci. U.S.A.*, 74, 4401-4405, 1977.
- [KRAS'04] **Krasnogor, N., Pelta, D.A.:** "Measuring the similarity of protein structures by means of the universal similarity metric," *Bioinformatics Advance Access*, vol. 20, pp: 1015-1021, 2004.

- [KROE'04] **Kroese, D. P., Rubinstein, R.Y., Taimre, T.:** "Application of the Cross-Entropy Method to Clustering and Vector Quantization", *Journal of Machine Learning Research*, 2004.
- [KROG'94] **Krogh, A., Brown, M., Mian, I., Sjolander, K., Haussler, D.:** "Hidden Markov models in computational biology: Application to protein modeling", *J. Molecular Biology*, 235:1501-1531, 1994.
- [KRUS'83] **Kruskal, J.B., and Sankoff, D.:** "In Time Warps, String Edits, and Macromolecules: the Theory and Practice of Sequence Comparison", Eds. D. Sankoff and J.B. Kruskal, Addison-Wesley, London, 1983.
- [LAND'98] **Landauer, T., Foltz, P., and Laham, D.:** "Introduction to latent semantic analysis," *Discourse Processes*, vol. 25, pp:259-284, 1998.
- [LASK'02] **Laskaris, N.A., Ioannides, A.A.:** "Semantic geodesic maps: a unifying geometrical approach for studying the structure and dynamics of Single trials evoked responses. *Electroenceph.*" *Clinical Neurophysiology*, 113:1209-1226, 2002.
- [LAQU'81] **Laquer, H.T.:** "Asymptotic Limits for a Two-dimensional Recursion", *Stud.appl.Math.* 46, 271-277, 1981.
- [LAWR'93] **Lawrence, C.E., Altschul, S.F., Boguski, M.S., Liu, J.S., Neuwald, A.F., Wootton, J.C.:** "Detecting subtle sequence signals: a Gibbs sampling strategy for multiple alignment", *Science*, 262(5131):208-214, 1993.
- [LEE'94] **Kowlakowski, L.F., and Rice, K.A.:** "Accepted-Mutation Parsimony Functionally Classifies G-Protein-Coupled Receptors", *Science (in press)*, 1994.
- [LESK'86] **Lesk, A., Levitt, M., Chothia, C.:**"Alignment of the amino acid sequences of distantly", *Protein Engineering*, 1(1):77-78, 1986.
- [LEVE'66] **Levenstein, V.I. :** "Binary Codes Capable of Correcting Deletions, Insertions, and Reversals", *Cybernet, Control Theory*,10, pp. 707-701.
- [LI'01] **Li, M., Badger, J.H., Chen, X., Kwon, S., Kearney, P., and Zangh, H.:** "An information based sequence distance and its application to whole mitochondrial genome phylogeny," *Bioinformatics*, vol.17, pg: 149-154, 2001.
- [LI'03] **Li, M., Chen, X., Li, X., Ma, B., and Vitanyi, P. :** "The similarity metric," *Proceedings of the 14th ACM-SIAM Symposium Discrete Algorithms*, 2003.
- [LI'97] **Li, M., and Vitany, P.:** "An introduction to Kolmogorov Complexity and Its Applications," Springer Verlag, 1997.
- [LIAO'03] **Liao, L., Noble, W.S.:** "Combining pairwise sequence similarity and support vector machines for remote protein homology detection," *J. Comp. Biol.*, vol.10, no.6, pp.857-868, 2003.

- [LIU'04] **Liu, X., Croft, W.B.,** : "*Statistical Language Modeling For Information Retrieval*", Center for Intelligent Information Retrieval, Univ. of Massachusetts, <http://ciir.cs.umass.edu/pubfiles/ir-318.pdf>
- [LORD'03] **Lord, P.W., Stevens, R.D., Brass, A., Goble, C.A.,**: "*Semantic similarity measures as tools for exploring the gene ontology*," PubMed, Pac. Symp. Biocomput., pp: 601-612, 2003.
- [MADE'95] **Madej, T., Gibrat, J.F., and Bryant, S.H.,**: "*Threading a database of protein cores*," Proteins: Structure, Function and Genetics, vol. 23, pp: 356-369, 1995.
- [MANN'00] **Manning, C.D., Schütze, H.,**: "*Foundations of statistical natural language processing*", Massachusetts Institute of Technology Press, Cambridge, Massachusetts London, England, pag:60-60;554 - 556;557 - 588, 2000.
- [MARK'03] **Markowitz, F., Edler, L., Vingron, M.,**: "*Support Vector Machines for Protein Fold Class Prediction*", Biometrical Journal, 45(3): 377 - 389, 2003.
- [MARK'71] **Markov, A.A.**: "*Extension of the limit theorems of probability theory to a sum of variables connected in a chain.*", reprinted in Appendix B of: R. Howard. *Dynamic Probabilistic Systems, volume 1: Markov Chains*. John Wiley and Sons, 1971.
- [MART'80] **Martinez, H.M.,**: " A New Algorithm for Calculating RNA Secondary Structure", Manuscript, .
- [MART'83] **Martinez, H.M.,**: "An Efficient Method for Finding Repeats in Molecular Sequences", Nucleic Acids Res., 11, 4629-4634, 1983.
- [MITC'97] **Mitchell, T.,**: "*Machine learning, WCB McGraw-Hill*", Boston, 1997.
- [MURZ'95] **Murzin, A.G., Brenner, S.E., Hubbard, T., Chithia, C.,**: "*SCOP: A structural classification of protein database for the investigation of sequences and structures*", J.Mol.Biol., 247, 1995.
- [NCBI'05] National Center for Biotechnology Information, accesibil online la adresa: (<http://www.ncbi.nlm.nih.gov/>), (accesat la 5 Iunie 2005).
- [NEED'70] **Needelman, S.B., Wunsch, C.D.,**: "*A General Method Applicable to the Search for Similarities in Amino Acids Sequences of Two Proteins*", J. Mol. Biol, 48, pp.444-453, 1970.
- [NOTR'00] **Notredame, C., Higgins, and D., Heringa J.,**: "*T-Coffee: A novel method for multiple sequence alignments*" Journal of Molecular Biology, 302: 205-217, 2000.
- [PARK'98] **Park, J., Karplus, K., Barrett, C., Hughey, R., Haussler, D., Hubbard, T., and Chotia, C.,**: "*Sequence comparisons using multiple sequences detect three times as many remote homologues as pairwise methods*", J. Mol. Biol., vol. 284 no.4, pp: 1201-1210, 1998.

- [PEAR'88] **Pearson, W.R., and Lipman, D.G.:** *"Improved tools for biological sequence analysis"*, Proc. Natl Acad. Sci. USA, vol. 85, pp: 2444-2448, 1988.
- [PEAR'96] **Pearson, W.R.:** *"Effective protein sequence comparison"*, Methods Enzymol., vol. 266, pp: 227-259, 1996.
- [PEAR'97] **Pearson, W.R.:** *"Identifying distantly related protein sequences"*, Comput. Applic. Biosci., vol.13, pp: 325-332, 1997.
- [PEAR'98] **Pearson, W.R.:** *"Empirical statistical estimates for sequence similarity searches,"* J. Mol. Biol., vol. 276, pp: 71-84, 1998.
- [PEVZ'00] **Pevzner, P.A., Sze, S.H.:** *"Combinatorial approaches to finding subtle signals in DNA sequences"*, Proceedings of the Eighth International Conference on Intelligent Systems for Molecular Biology (ISMB), pp: 269 - 278, 2000.
- [PYLE'99] **Pyle, D.:** *"Data Preparation for Data Mining"*, Morgan Kaufmann, Bk&CD Rom edition, March 15, 1999.
- [QUEE'82] **Queen, C.M., Wegman, N., and Korn, L.J.:** *"Improvements to a Program for DNA Analysis: a procedure to find Homologies Among Many Sequences"*, Nucleic Acids Res., 10, 449-456, 1982.
- [REBH'98] **Rebhan, M., Chalifa-Caspi, V., Priluski, V., Lancet, D.:** *"GeneCards: a novel functional genomics compendium with automated data mining and query reformulation support"*, Bioinformatics, 14, no.8, pp: 656-664, 1998.
- [REIC'73] **Reichert, T.A., Cohen, D.N., and Wong, A.K.C.:** *"An Application of Information Theory to Genetic Mutations and Matching of Polypeptide Sequences"*, J. Theor. Biol., 42, 245-261, 1973.
- [RIGO'98] **Rigoutsos, I., Floratos, A.:** *"Motif discovery without alignment or enumeration"*, Proceedings of the second annual international conference on Computational Molecular Biology (RECOMB), New York, pp:221-227, 1998.
- [ROSE'00] **Rosenfeld, R.:** *"Two decades of statistical language modeling: where do we go from here?"* Proceedings of the IEEE, 88(8), 2000.
- [ROSE'94] **Rosenfeld, R.:** *"Adaptive Statistical Language Modeling: A Maximum Entropy Approach"*. Ph.D. Thesis. Carneige Mellon University, April 1994.
- [RYCH'00] **Rychlewski, L., Jaroszewski, L., Li, W., Godzik, A.** *Protein Science* 9:232-241, 2000.
- [SAIG'04] **Saigo, H., Vert, J-P., Ueda, N., Akutsu, T.:** *"Protein homology detection using string alignment kernels,"* J. Bioinformatics, vol.20 no.11, pp:1682-1689, 2004.
- [SANK'72] **Sankoff, D.:** *"Matching Sequences Under Deletion-Insertion Constraints"*, Proc. Natn. Acad. Sci. U.S.A., 68, 4-6, 1972 .

- [SANK'73] **Sankoff, D., and Sellers, P.H.:** "Shortcuts, Diversions and Maximal Chains in Partially Ordered Sets", *Discrete Math.*, 4, 287-293, 1973.
- [SANK'76] **Sankoff, D., Cedergren, R.J., and Lapalme, G.:** "Frequency of Insertion-Deletion, Transversion, and Transition in the Evolution of 5S Ribosomal RNA," *J. Mol. Evol.* vol.7, pp: 133-149, 1976.
- [SANK'83] **Sankoff, D., and Kruskal, J.B.:** "Time Warps, String Edits, and Macromolecules: the Theory and Practice of Sequence Comparison", Addison-Wesley, London, 1983.
- [SANK'85] **Sankoff, D.:** "Simultaneous Solution of the RNA Folding; Alignment and Protosequence Problem", *SIAM, J. Applied Mathematics*, 45, pp: 810-825, 1985.
- [SANT'89] **T.G. Santner and D.E. Duffy:** "The Statistical Analysis of Discrete Data," Springer-Verlag, New York, 1989.
- [SARK'02] **Sarkar, I.N., Phil, M., and Rindfleisch, T.C.:** "Discovering Protein Similarity using Natural Language Processing," Cognitive Science Branch of Lister Hill National Center for Biomedical Communications, National Library of medicine, 2002.
- [SCHÄ'01] **Schäffer, A., Aravind, L., Madden, L., Shavirin, S., Spouge, J., Wolf, Y., Koonin, E., Altschul, S.,** *Nucleic Acids Research*, 29, (14), pp: 2994-3005, 2001.
- [SCOP] **Structural Classification Of Proteins**, site-ul oficial disponibil online la: <http://scop.mrc-lmb.cam.ac.uk/scop/>, (accesat la 2 Mai 2005).
- [SEAR'95] **Searls, D.B., and Murphy, K.P.:** "Automata-theoretic models of mutation and alignment", *Proceedings of the third International Conference on Intelligent Systems for Molecular Biology*, pp: 341-349, 1995.
- [SELL'74a] **Sellers, P.:** "An Algorithm for the Distance Between Two Finite Sequences", *Comb. Theoriz*, 16, 253-258, 1974.
- [SELL'74b] **Sellers, P.:** "On the Theory and Computation of Evolutionary Distances", *SIAM J. Appl. Math.*, 26, pp.787-793, 1974.
- [SELL'79] **Sellers, P.:** "Pattern Recognition in Genetic Sequences", *Proc. Natn. Acad. Sci. U.S.A.*, 76, pp.3041, 1979.
- [SELL'80] **Sellers, P.:** "The theory and Computation of Evolutionary Distances: Pattern Recognition", *J. Algorithms*, 1, pp.359-373, 1980.
- [SHAN'51] **Shanon, C.E.:** "Prediction and Entropy of Printed English", *Bell Systems Technical Journal*, 30: 50-64, 1951.
- [SHEP'80] **Shepard, R.N.:** "Multidimensional Scaling, Tree-Fitting, and Clustering", *Science*, 210, pp.390-398, 1980.
- [SJOL'96] **Sjolander, K., Karplus, K., Brown, M., Hughey, R., Krogh, A., Mian, I.S. and Haussler, D.:** "Dirichlet mixtures: a method for improved detection of weak but significant protein sequence homology," *J. Bioinformatics*, Vol 12, pp: 327-345, 1996.

- [SMIT'81] **Smith, T.F., Waterman, M.S., Fitch, W.M.:** *"Comparative Biosequence Metrics"*, J. Mol. Evol, 18, 38-46, 1981.
- [SMIT'81a] **Smith, T.F., and Waterman, M.S.:** *"Identification of Common Molecular Subsequences"*, J. Mol. Biol., 147, 195-197, 1981.
- [SMIT'81b] **Smith, T.F., Waterman, M.S., Fitch, W.M.:** *"Comparison of Biosequences"*, Adv. Appl. Math., 2, 482-489, 1981.
- [SMIT'84] **Smith, T.F., Waterman, and Burks, C.:** *"The Statistical Distribution of Nucleic Acids Similarities"*, In prep.
- [SMIT'90] **Smith, H.O., Annau, T.M., Chandrasegaran, S.:** *"Finding sequence motifs in groups of functional related proteins"*, Proceedings of the National Academy of Sciences of the United States of America, 87(2): 826-830, 1990.
- [SSEA'05] **SSEARCH**, descrierea aplicației disponibilă online la data de 9 iunie 2005 la adresa: <http://ori.nibb.ac.jp/SIT/SSEARCH.html>.
- [STAN'70] **Stanton, R.G., and Cowan, D.D.:** *"Note on a 'Square Functional' Equation"*, SIAM Rev.12,277-297, 1970.
- [STUD'78] **Studnicka, G., rahn, G., Cummings, I., and salser, W.:** *"Computer Method for Predicting the Secondary Structure of Single Stranded RNA"*, Nucleic Acids Res., 5, 3365-3387, 1978.
- [TAYL'84] **Taylor, P.:** *"A fast Homology Program for Aligning Biological Sequences"*, Nucleic Acid Res., 12, 447-455, 1984.
- [THOM'94] **Thompson, J.D., Higgins, D.G., and Gibson, T.J.:** *"CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, positions-specific gap penalties and weight matrix choice."* Nucleic Acids Research, 22:4673-4680, 1994.
- [TOMP'99] **Tompa, M.:** *"An exact method for finding short motifs in sequences, with application to the ribosome binding site problem"*, in Proceedings of the 7th International Conference on Intelligent Systems for Molecular Biology (ISMB), pp: 262-271, 1999.
- [TUER'90] **Tuerk, C.,Gold, L.:** *"Systematic evolution of ligands by exponential enrichment: RNA ligands to bacteriophage T4 DNA polymerase"*, Science, 249: 505-510, 1990.
- [UKKO'83] **Ukkonen, E.:** *"On Approximate String Matching"*, Proc. Int. Conf. Found. Comp.Theor.Lectures Notes in Comp. Sci., 158, 487-496, 1984.
- [UKKO'84] **Ukkonen, E.:** *"Algorithms for Approximate String Matching"*, Informat. Control (in press), 1984.
- [ULAM'72] **Ulam, S.M.:** *"In Applications of Number Theory to Numerical Analysis"*, Ed. S.K. Zaremba, Academic Press New York, pp.1-3, 1972.
- [VANC'03] **Van Compennolle, D.:** *"Spoken Language Science and Technology"*,2003, http://www.esat.kuleuven.ac.be/compil/pub/spoken_language/TOC.htm

- [VANH'98] **van Helden, J., Andre, B., Collado-Vides, J.,**: "Extracting regulatory sites from the upstream region of yeast genes by computational analysis of oligonucleotides frequencies", *J. of Molecular Biology*, 281(5), pp:827-832, 1998.
- [WAGN'83] **Wagner, R.H.,**: "*In Time Warps, Sring Edits, and Macromolecules: the Theory and Practice of Sequence Comparison*", Eds. D. Sankoff and J.B. Kruskal, Addison-Wesley, London, pp. 215-235, 1983.
- [WANG'00a] **Wang, J.T.L., Ma,Q., Shasha, D., Wu, C.H.,**: "*Application of Neural Networks to Biological Data Mining: A Case Study in Protein Sequence Classification*", Proceedings of the sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Boston, Massachusetts, United States, 2000.
- [WATE'76] **Waterman et al., 1976 Waterman, M. S., Smith, T. F. , Beyer, W. A.,**: "*Some biological sequence metrics.*", *Advances in Mathematics*, 20, 367-387, 1976.
- [WATE'83] **Waterman, M.S.:** "*Sequence Alignment in the Neighborhood of the Optimum with General Applications to Dynamic Programming*", *Proc. Natn. Acad. Sci. U.S.A.*, 80, 3123-3124, 1983.
- [WATE'84] **Waterman, M.S.:** "*General methods of sequence comparison*", *Bulletin of Methematical Biology*, Vol. 46, No. 4, pp. 473-500, printed in Great Britain, 1984.
http://www.cmb.usc.edu/papers/msw_papers/msw-054.pdf
- [WATE'86] **Waterman, M.S.:** "*Multiple sequence alignment by consensus*", *Nucleic Acids Res.*,vol. 14, no.22, pp: 9095-9102, 1986.
- [WHEE'96] **Wheeler, D.:** "*Weight Matrices for Sequence Similarity Scoring*", sintezã, <http://www.icp.ucl.ac.be/~opperd/private/matrices.html>;
- [WIKI'05a] **Wikipedia,** The free encyclopedia.Nucleotide.
<http://en.wikipedia.org/wiki/Nucleotide>
- [WIKI'05b] **Wikipedia,** Sequencing,
http://en.wikipedia.org/wiki/Biological_sequence;
- [WIKI'05c] **Wikipedia,** The free enciclopedia. Protein folding.
http://en.wikipedia.org/wiki/Protein_folding;
- [WIKI'05d] **Wikipedia,** The free enciclopedia. Peptide.
<http://en.wikipedia.org/wiki/Peptide>;
- [WIKI'05e] **Wikipedia,** Biological database
http://en.wikipedia.org/wiki/Biological_database;
- [WIKI'05f] **Wikipedia,** The free enciclopedia. Multiple alignment.
http://en.wikipedia.org/wiki/Sequence_alignment
- [WIKI'05g] **Wikipedia,**The free enciclopedia. Exon.
<http://en.wikipedia.org/wiki/Exon>

- [WIKI'06a] **Wikipedia**, The free encyclopedia. Cladistics.
<http://en.wikipedia.org/wiki/Cladistics>
- [WILB'83] **Wilbur, W.J., and Lipman, D.J.**: "*Rapid Similarity Searches of Nucleic Acid and Protein Data Bank*", Proc. Natn.Acad.Sci.U.S.A., 80, 726-730, 1983.
- [WILB'84] **Wilbur, W.J., and Lipman, D.J.**: "*The Context Dependent Comparison of Biological Sequences*", SIAM, J. appl. Math., in press, 1984.
- [WONG'74] **Wong, A.K.C., Reichert, T.A., Cohen, D.N., and Ayyun, B.O.**: "*A Generalized Method for Matching Informational Macromolecular Code Sequences*", Comput., Biol., Med., 4, pp.43-57, 1974.
- [WU'00] **Wu, C.H., McLarty, J.**: "*Neural Networks and Genome Informatics*", Elsevier Science, 2000.
- [WU'03] **Wu, K-P., Lin, H-N., Sung T-Y., and Su, W-L.**: "*A new Similarity Measure among Protein Sequences*," IEEE Computer Society Bioinformatics Conference (CSB'03), 2003.
- [YAO'04] **Yao, Q., Huang, X., An, A.**: "*Applying Language Modeling to Session Identification from Database Trace Logs*" work in progress, <http://www.cs.yorku.ca/~qingsong/research/ngram.pdf>
- [YOUN'97] **Young, S., Bloothoof, G.**: "*Corpus-Based Methods in Language and Speech Processing*", Kluwer Academic Publishers, 1997.
- [ZHAN'99] **Zhang, Z., Berman, P., Wiehe, T., and Miller, W.**: "*Post-processing long pairwise alignments*," Bioinformatics, vol.15, no.12, pp: 1012-1019, 1999.

INDEX

A

acid nucleic, 24
 ADN, 9, 11, 13, 14, 21, 24, 25, 26,
 31, 43, 44, 49, 53, 73, 82, 84, 85,
 89, 94, 160
 ajustarea (smoothing), 97
 aliniament local, 26, 27, 38, 41, 43,
 65, 66, 168
 aliniament optimal, 43, 47, 53, 74,
 138
 aliniamentele structurale, 27
 aliniamentul global, 26, 27, 139, 160
 aliniamentul multiplu, 66, 79, 88, 90,
 138, 170
 aliniamentul multiplu, 26
 aliniamentul secvențelor, 12, 26, 27,
 32, 64, 68, 72, 90
 amino acizi, 11, 15, 16, 21, 23, 25,
 26, 31, 32, 33, 34, 37, 39, 41, 86,
 113, 128, 136, 138, 144, 145, 146,
 150, 170

C

caracteristicile ROC, 41
 căutarea iterativă, 79, 87
 clasificatorului k-NN, 70
 codon, 15, 124, 155, 177
 codul genetic, 14, 15
 complexitate de timp, 122
 complexitatea Kolmogorov, 40
 complexitatea spațiului, 122, 142,
 146
 cross-entropiei, 99, 115, 116

D

densități de potrivire, 54
 distanță de editare, 29
 distanță evoluționară, 35

E

entropia, 95, 97, 98, 99, 100, 104,
 112, 114, 115, 116

entropia condiționată, 98
 entropia încrucișată, 98
 estimarea entropiei, 102, 104, 108,
 111
 evaluarea entropiei, 101, 103, 107,
 111
 evaluarea modelelor lingvistice, 95,
 98, 100, 112, 114
 explorarea ontologică a genelor, 39

F

formate de secvențe, 22
 formula lui Stirling, 43

G

gene, 14, 25, 26, 41, 82, 83, 84, 153,
 155, 158, 198, 203
 grad de asociere, 32

H

helix, 16, 17, 41, 128

I

indels, 30, 46, 47, 48, 49, 50, 53, 57

L

lanț Markov, 70, 95, 96, 175
 lanțurile Markov, 70, 71, 95
 LSA, 9, 41

M

măsură a similarității, 31
 măsurile de disimilaritate, 124, 130
 matrice de ponderi, 31, 89, 90, 139
 matricea cu ponderi, 81
 matricea de scoruri, 32, 160, 164,
 168
 matricea de substituție, 32, 86
 matricelor de scoruri, 30
 metoda LSA, 41
 metoda ROC, 38, 126

metoda VSM (Vector Space Model, 41
 metodele de aliniament, 44
 metrica Lee, 42
 model lingvistic statistic, 96
 model Markov, 36, 78, 96, 103, 175
 model matematic Markov, 96
 model probabilistic, 31, 60, 62, 65,
 78, 81
 modelele Markov ascunse, 81
 modelele Markov selective, 71
 modelelor n-gram, 97, 99, 121, 127,
 135, 145
 Motiv al unei secvențe, 25
 motive, 82, 84, 87
 mutații, 14, 25, 30, 32, 35, 37, 44,
 45, 123, 124, 145, 177

N

normalizării textului, 101
 nucleotide, 13, 15, 21, 24, 25, 32,
 33, 54, 56, 71, 73, 75, 84, 85, 86,
 89, 124, 138, 139, 177

O

omologia, 25, 84, 92, 168

P

penalitățile, 90
 perechi de HMM, 58
 perechi HMM, 62, 63, 64, 65
 perplexitatea, 95, 97, 100, 111
 plieri, 15, 21, 92
 probabilitate de tranziție, 78
 probabilităților de tranziție, 34, 60
 profil HMM, 37
 profile de HMM, 58

programare dinamică, 43, 44, 45, 46,
 49, 57, 58, 59, 60, 63, 64, 65, 66,
 73, 74, 76, 85, 86, 92, 170
 proteine, 11, 12, 14, 15, 16, 17, 18,
 19, 20, 21, 24, 25, 26, 27, 31,
 32, 33, 34, 35, 37, 38, 39, 40, 41,
 43, 44, 58, 66, 72, 79, 84, 85, 86,
 87, 89, 90, 91, 92, 94, 113, 114,
 115, 116, 117, 118, 120, 121, 123,
 124, 125, 126, 127, 128, 129, 130,
 131, 133, 134, 135, 136, 138, 139,
 140, 142, 145, 146, 169, 170, 172,
 173, 177

R

regiuni, 16, 25, 27, 38, 39, 51, 52,
 53, 54, 57, 66, 74, 75, 76, 84, 86,
 90, 93, 123, 177
 repetițiilor tandem, 83
 reziduu, 24, 31, 35, 36, 63, 90

S

scorului de similaritate, 44, 148
 secvențelor biologice, 11, 12, 13, 21,
 25, 27, 29, 41, 43, 53, 57, 67, 69,
 72, 73, 83, 87, 92, 93, 94, 114,
 115, 123, 128, 136, 143, 144, 145,
 147, 148, 149, 170
 secvențiere, 13, 14, 24, 37
 similaritatea compoziției textuale,
 136, 146
 similaritatea funcțională, 20
 statisticile lui Hubert, 131
 structura primară, 13, 15, 18

V

vectorizarea algoritmilor, 75