

Bioinformatics

Human Genome Project

What are databases

- A database is an organised collection of structured information or data, typically stored in a computer system.
- Consists of basic units called records or entries.
- Protein entries → records
- Protein properties → fields

Three types of datatypes Management system :-

- Flat file indexing system
 - Can be used useful by searching and ordering
 - PPB used flat file indexing system
 - Used in System Retrieval System (SRS)
 - drawbacks : (i) Indexing → the index must be consistent
 - (ii) storage.
- Relational DBMS
 - Collection of tables
- Object-oriented DBMS

Relational DBMS

→ 2D format storing

unique columns in the short table called identifiers → these are called primary key

[candidate key → when more than 1 unique column]
 ↗ out of these 1 is primary key and others
 [alternate key]

primary key becomes foreign key for another table.

Characteristics : Scalability

③ Object Oriented DBMS

in the form of Objects as used in Object oriented programming

Database \rightarrow Databank \rightarrow Data warehouse

Biological database classified as

① Primary (archival) \rightarrow stores info without manipulation
 ↳ Raw Data
 ↳ Priority from research scientists

② Secondary database. \rightarrow only data of specific category
 (curated)
 e.g. mutants.

Composite Databases \rightarrow (Swiss Prot + TrEMBL) \rightarrow UniProtKB

Persons use such dataset \rightarrow ① Contain homologous Sequence or not.

② Repeat elements or regulatory seq. in non-coding DNA
 stretches

DNA vs Protein sequence \rightarrow There are very different DNA sequences that code for similar protein sequences.

\rightarrow DNA \rightarrow 4 letters \rightarrow A, T, C, G
 Protein \rightarrow 20 letters

\rightarrow ① degeneracy
 ② Randomness

Three data retrieval systems in molecular biology - ① SRS

② Entrez

③ DBGET

What can be discovered about a gene by database search

- Evolutionary information
- genomic information
- structural info.
- expression info.
- functional info.

Biological Databases

- (i) Nucleic Acid → GenBank
- (ii) Protein Sequence → UniProtKB
- (iii) Structure → PDB

Metabase

↳ a database of databases

BIOCYC

↳ NCBI is not a Metabase, it's an organization

- (i) GenBank (primary database)
Gene Sequence Database (GSDB) 1979, NCBI ↳ First DNA Sequence Database
ENA → European Nucleotide Archive
DDBJ → DNA Databank of Japan ↳ 19.8 trillion

* Submission of data to GenBank

- (i) Online web server → BankIt
- (ii) Stand alone sequin
- (iii) Batch Submitter → Tbl2asn
TBL2ASN

Types of sequences which can be deposited?

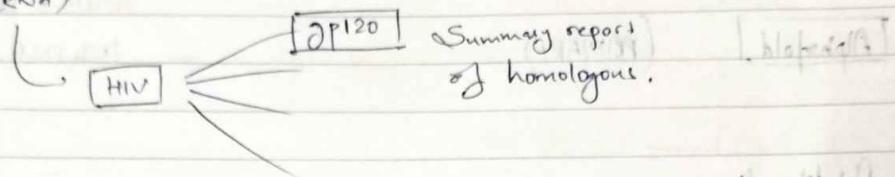
- (i) Whole cell genome
- (ii) EST → Expressed sequence tags cDNA
- (iii) GSS → Genome Sequence Survey
- (iv) Shotgun Sequence

What can not be deposited?

- (i) < 200 nt
- (ii) mixture of DNA + RNA Sequence
- (iii) Primers

Bac, VIR, MVV, PLT, ...

GenBank → Genomes
Retrovirus ← VIRAL Genomes
(RNA)



Primates. Eukaryotes
(*Homosapiens*)

↓
??
Past infection

[Genome sizes]

(ii) UniProt KB → Protein Sequences database (PRIMARY)

[DNA → translate into protein]

UniProt KB → Swiss Prot
→ trEMBL

(iii) Proteomics

(iv) 90% Archival 25%, 50%, 90%

(3) PDB : Structure Database (TextFile)
 → experimentally determined structures

~ 2,00,000

1978

models

AlphaFold

(PRIMARY)

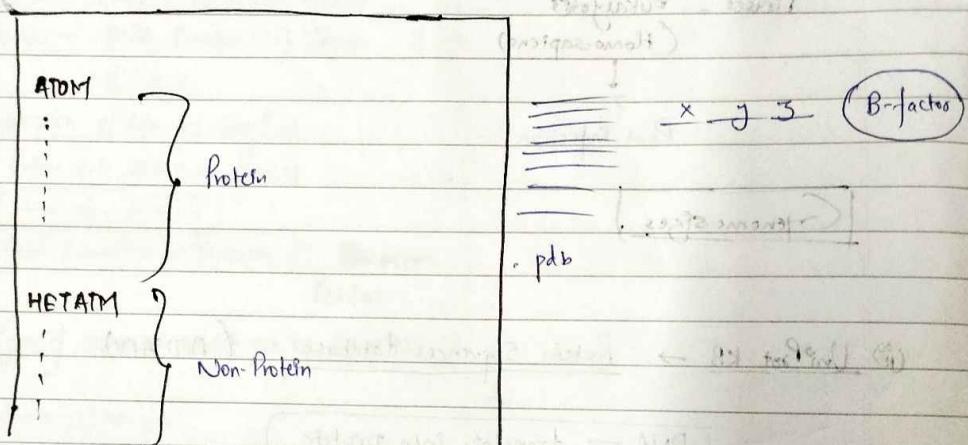
NDB, CSD

Protein

Protein + Complexes,

DNA, RNA, Sugars.

Flat file atoms



Homology Richard Owen (1843)

"Same Organ in different animals under every variety of form and function"



Homology

Orthology

Paralogy

Xenology

Homoplasy

↳ The counterpart of homology is usually considered to be homoplasy.

↳ A homoplasious trait is a similarity among organisms that was not inherited from the common ancestor of these organisms.

How to relate homology with 'gene'

- Homology →
 - ① Organism level
 - ② Trait-Specific
 - ③ Gene-Specific

Any gene present in two diff. species is called Orthology

Any gene present ~~within the same species~~ within the same species that has arisen through the gene duplication events. As a result of duplication, it can have different function while sharing a common ancestor → Paralogy

Relationship b/w two gene

Analogous

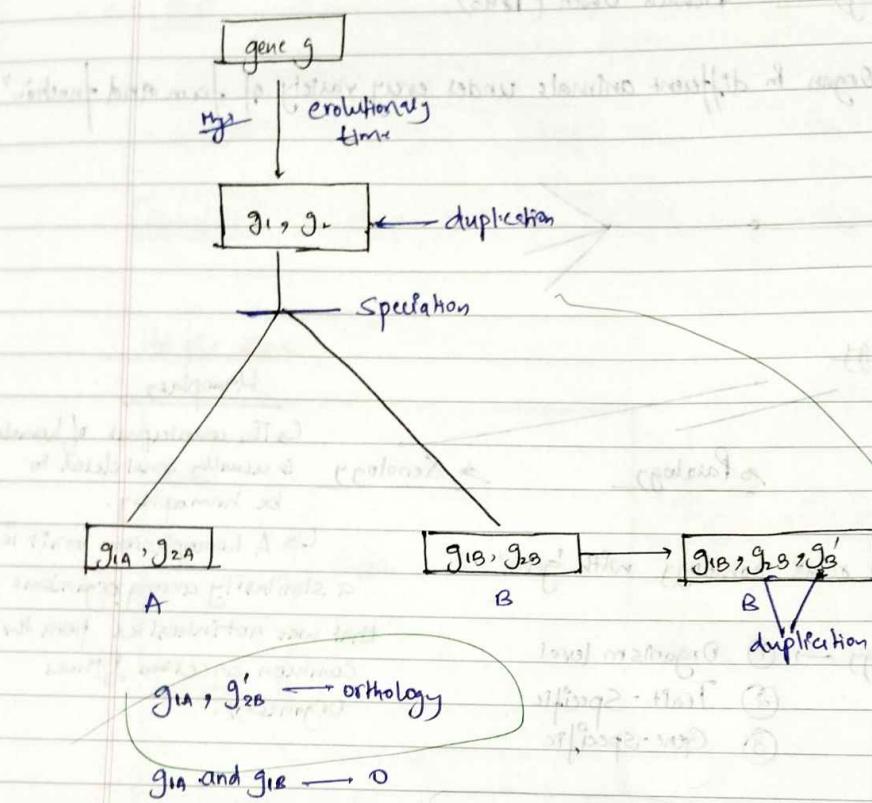
Xenology: Xenology is the formation of different sequences due to horizontal gene transfer ~~between~~ across species b/w diff. species. HGT is the transfer of genetic material from one org. to another through plasmid exchange or viral injection

Inparalogy →

Outparalogy → follows speciation after duplication in different species.

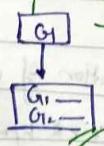
classmate

Date _____



Paralogy

- ① In (gen) paralog, → after speciation
- ② Out (Allo) paralog, → Before Speciation
- ③ Ortho paralog →

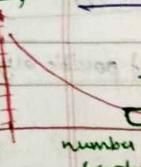


Gene that is result of a duplication of a whole genome.

Analogy

- due to
- convergent evolution
- parallel evolution
- reversal mutation

Identify



Com

① Al

Si

C

A

Date _____
Page _____

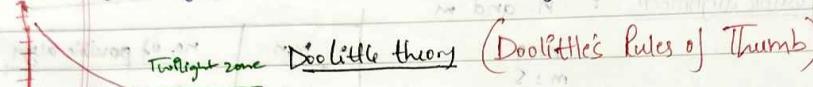
proteins that have similar active sites but different backbone separate

Analogy due to convergent evolution → Convergent Evolution : It refers to the phenomenon in which the physical traits of the two diff. species are similar but they do not have any common ancestor. The similarity b/w the species is due to the adaptation to similar environment or functional requirements.

→ Pseudo-Homology → Shows similarity b/w diff. ancestors.

→ Homology

Identity → Similarity (%) of genes



Doolittle theory (Doolittle's Rules of Thumb)

number of residues $g_1 \& g_2$

- $\geq 30 \rightarrow$ homologous
- $\leq 15-30 \rightarrow$ may/may not be homologous
- $\leq 15 \rightarrow$ not homologous

If two genes g_1 & g_2 are similar $\not\rightarrow$ homologous. } Similarity vs Homology
 If two genes g_1 & g_2 are homologous $\not\rightarrow$ similar.

Comparison (pattern 1 and pattern 2)

DNA : A / C / G / T

① Alignment

$\begin{array}{c} \checkmark \checkmark \checkmark \checkmark \\ \text{A} \text{G} \text{G} \text{T} \\ \text{A} \text{A} \text{G} \text{T} \end{array}$

$$\begin{aligned} \text{Similarity} &= \frac{3}{4} \\ \text{similarity} &= \frac{3}{4} \rightarrow \end{aligned}$$

S₁ : G G G A S₂ : AT G G G T

G G G A - - -
A T G G G A

→ global alignment

$\begin{array}{c} \text{C} \text{G} \text{C} \text{Y} \\ \text{G} \text{C} \text{G} \text{C} \end{array}$

→ local alignment → max of no. better matches without a gap.

- Match
- Mismatch
- Gap
- Global aligned
- Local aligned

pair wise sequence alignment.
Fission yeast (*Saccharomyces cerevisiae*) & humans
common ancestor 1 billion years ago.

classmate

Red

Page

S1 : G T A T T A G C A

Cyclin-dependent CDK-CDK2

S2 : GTACC

(a clock gene which controls the cell cycle)

T A G C
TA C C

Highly conserved region of two sequences

local alignment

G T A T T A G C A
G T A - - - C C -

covered the whole length of the sequences. → Global Alignment

No. of possible alignment : n and m

n = 9

m = 5

without gap = $n^m + m^n$

with gap = $n^m \cdot m^n$

n, m	No. of possible align.
1, 1	3
2, 2	13
3, 3	63
4, 4	321
5, 5	1683
	8989

A C G G C A

A A G C -

+2 +1 +2 +2 -2 = 5

A - G G C A

A A G - C -

+2 -2 +2 -2 +2 -2 = 0

Scoring scheme

Assume that M : +2

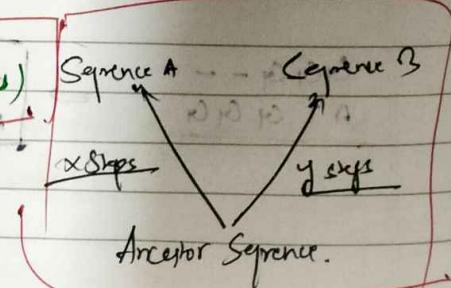
Mis : +1

Gap : -2

Higher score → Better Alignment

Optimal Alignment : $\max_{\theta} S^{\theta}(x, y)$

$d(S_1, S_2) = \max_{\text{alignment of } S_1, S_2} \text{Score (alignment)}$



PIK-CDK2

Significance of sequence alignment

- Sequence alignment is useful for discovering functional, structural and evolutionary information about the biological sequences.
- Can locate subsequences in DNA which can be useful to identify regulatory elements.
- Can locate DNA sequences that might overlap. DNA sequencing

Methods : (i) Dot Plot Method

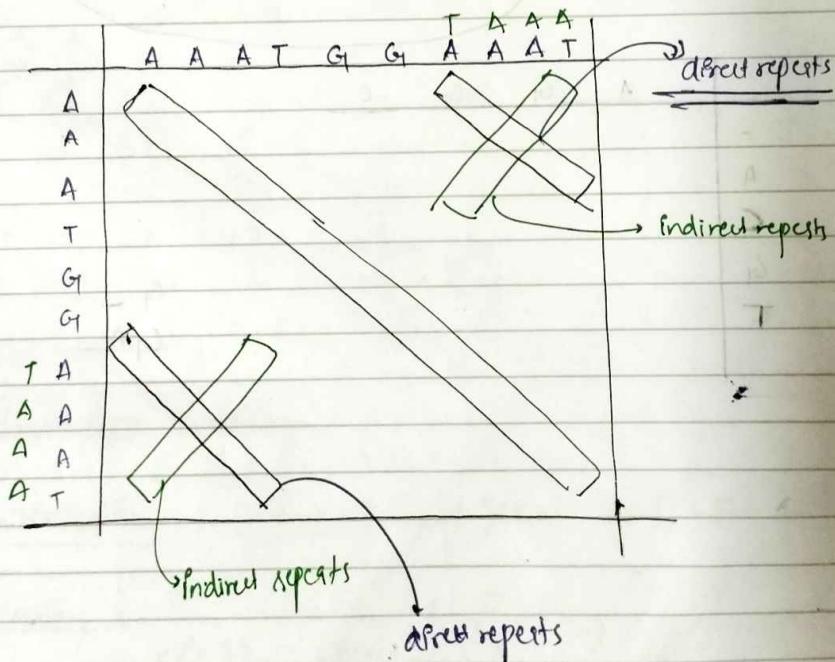
two sequences

Alignment

possible align.



- Draw {
- As the length of sequences increases, matrix size increases.
 - database search



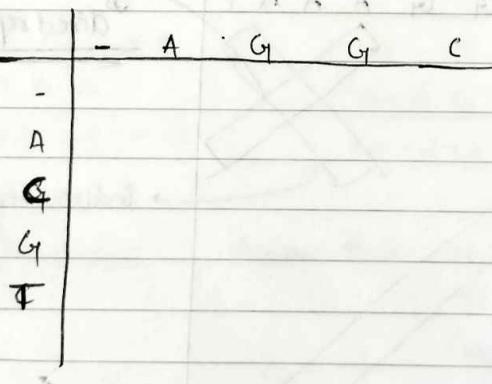
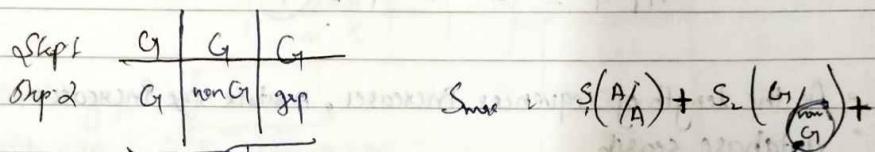
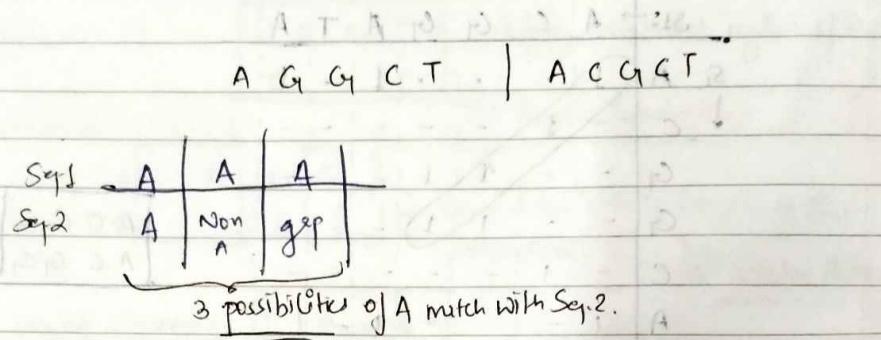
(GN) Global Alignment : Needleman - Wunsch Algorithm
 (LS) Local Alignment : Smith - Waterman Algorithm

classmate

Date _____

Page _____

2. Dynamic Programming (DP method) \rightarrow best method to alignment.



G -
G + A

Global Alignment

(Needle Man Wunsch Algorithm)

$$S(0,0) = 0$$

$$S(i,j) = \max \left\{ \begin{array}{l} S(i-1, j-1) + w(x_i, y_j) \\ S(i-1, j) + \text{Gap-penalty} \\ S(i, j-1) + \text{Gap-penalty} \end{array} \right\}$$

Time complexitySpace : $O(mn)$ Time : $O(mn)$

- Filling the matrix $O(mn)$
- Backtrace $O(mn)$

$$M = +1, MM = -1, \text{ Gap} = -3$$

-	G	A	A	G	A	
-	0	-3	-6	-9	-12	-15
G	3	1	-2	-5	-8	-11
A	-6	2	1	-1	-4	-7
A	-9	-5	-1	0	-3	
A	-12	-8	-4	0	-1	
T	-15	-11	-7	-3		
T	-18	-14	-10	-6	-4	

$$\begin{aligned} 0+1 &= 1 \\ &= -6 \\ &= -6 \end{aligned}$$

Match (m)
Mismatch (mm)

$$S(i,j) = \max \left\{ \begin{array}{l} S(i-1, j-1) + \delta_{x,y} \\ S(i-1, j) + \delta_{gap} \\ S(i, j-1) + \delta_{gap} \end{array} \right\}$$

Local Alignment(i) Initialization

$$S(0,0) = 0, S(i,0) = S(0,i) = 0 \quad \checkmark i, j$$

(ii) Iteration

$$S(i,j) = \max \left\{ \begin{array}{l} S(i-1, j-1) + \delta_{x,y} \\ S(i-1, j) + \delta_{gap} \\ S(i, j-1) + \delta_{gap} \end{array} \right\}$$

	-	<u>G A A</u>	G A	
-	-	0 0 0 0 0 0		Max Value (1)
<u>G</u>	0	1 0 0 1 0	<u>G A A</u>	<u>G A A</u>
A	0 0	2 1 0 2	<u>L1 + L1 = 3</u>	(Smith Waterman)
A	0 0 1	3 0 1		
4	0 0 1 2 2 1			
T	0 0 0 0 1 1			
T	0 0 0 0 0 0			

Overlapping1. Initialization

$$S(0,0) = 0, S(i,0) = S(0,i) = 0$$

Similar to Global Alignment except
gap penalties observe boundary

2) Iteration

$$S(i,j) : \max \left\{ \begin{array}{l} S(i-1, j-1) + S_{match} \\ S(i, j-1) + S_{mismatch} \\ S(i-1, j) + S_{mismatch} \end{array} \right\}$$

maximum value either
in last row or last column

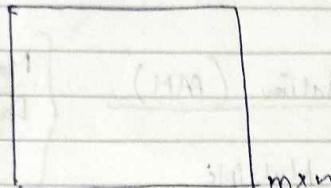
3) Trace back

Maximum value either in
last row or last column

17/08/23

10^{10} operations/c

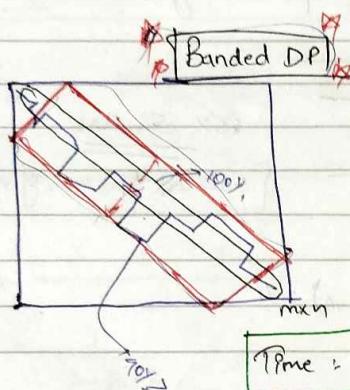
Disadvantages : time and space complexity.



Time taken : $O(mn)$

Tracing back : $O(m+n) \approx O(10^3)$

$$\begin{cases} m \approx 10^3 \\ n \approx 10^3 \end{cases}$$



$$S_{1m}(S_1, S_2) \geq 90\%$$

$$S_{1m}(S_1, S_2) = 100\%$$

$$\text{Time} = O(kn) = O(n)$$

Only a diagonal band of the dynamic programming matrix is computed and filled. Reduces time & memory requirements.

If the diagonal band consists of K diagonals (width K), then dynamic programming takes $O(Kn)$, much faster than $O(n^2)$ of standard DP.

$$\begin{array}{l} AGGT \\ AA CGT \\ -4 +2 -1 +2 +2 \\ -1 +3 -2 +3 +3 \end{array} \quad \begin{array}{c} \text{Match} : +2 \\ M-M : -1 \\ \text{Gap} : -7 \end{array} \quad \begin{array}{r} +3 \\ -2 \\ -1 \end{array} \quad \begin{array}{r} 15 \\ +4 \\ -6 \end{array}$$

* Hamming Distance

A M A N
A M I N
① → Hamming Distance

(3/4) * 100

* Levenshtein Distance

A M I A I N (Relationship)
A M A I I N

Substitution / Addition / Deletion
to make similar.

Substitution (or Scaling) Matrices

PAM / BLOSUM

③ PAM (or Percent) Accepted mutation (PAM)

Margaret Dayhoff Jr
1928

84 families of proteins → Multiple Sequence Alignment.



85% 8 miles

{ count the no. of mutations }

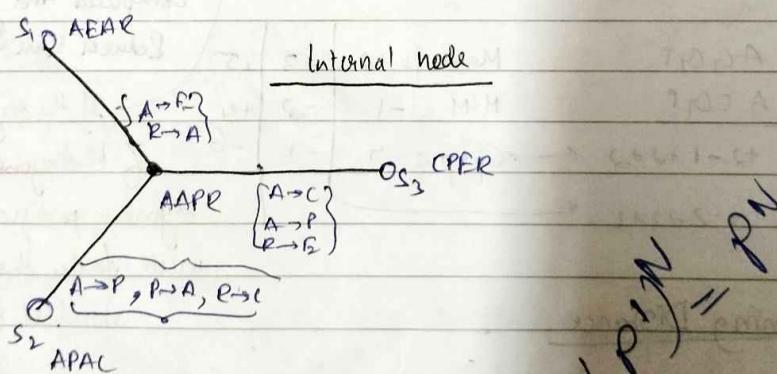
A E A R

A P A C

C P E R

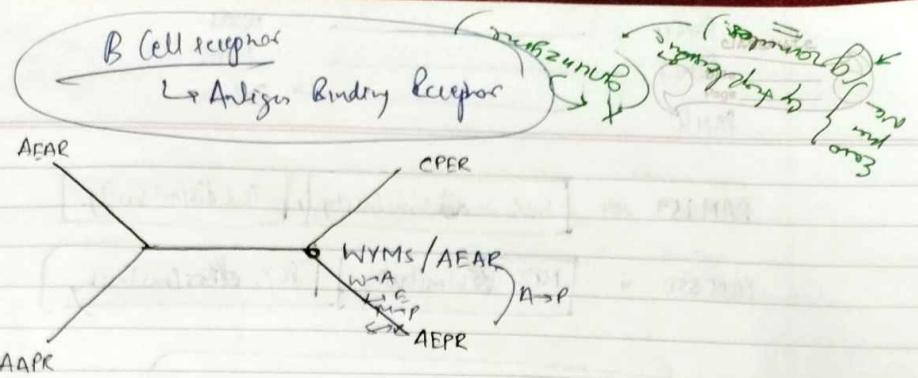
$$A \rightarrow A \quad (M_{A \rightarrow A}) = 4$$

$$M_{\mathrm{Ac}} > 2$$



Rooted tree \rightarrow

Unrooted tree \rightarrow



- Minimum no. of mutation
- Pseudomonas halice.

A_{xy}

$$P_{xy} = \frac{m_x}{\sum_{h \neq x} A_{yh}}$$

for non-diagonal values of matrix

Diagonal value of matrix

$$P_{xx} = 1 - m_x$$

m_x : mutability of x

$$m_x = \frac{1}{L - p_x \cdot z} \sum_{h \neq x} A_{xh}$$

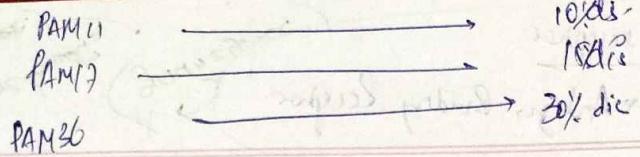
L: length of alignment

p_x : background prob. or a priori prob.
 z : scaling factor.

21/08/13

PAM₂ : PAM₁ × PAM₂ \Rightarrow 2% accepted mutation

PAM_n : PAM₁ × PAM_(n) : (PAM₂)ⁿ



PAM 159 \rightarrow [30% dissimilarity] & [70% dissimilarity]

PAM 350 \rightarrow [14% dissimilarity] & [86% dissimilarity]

1990, Henikoff and Henikoff, BLOSUM

n Blocks (500 BLOCKS)

Blocks Substitution Matrix.

$q_{ij} = \text{no. of times } i^{\text{th}} \text{ changed to } j^{\text{th}}$ for $i, j \in \text{AA}$

$p_{ij}^{(o)} = \frac{q_{ij}^{(o)}}{\sum_{j=1}^{20} \sum_{i=1}^{20} q_{ij}^{(o)}}$

p_{ij}^(o) is the probability that a mutation occurs from amino acid i to amino acid j

Marginal Prob $p_i^{(o)} = \frac{\sum_{j=1}^{20} q_{ij}^{(o)}}{\sum_{i=1}^{20} \sum_{j=1}^{20} q_{ij}^{(o)}}$

Log Odd ratio : $S_{ij}^{(o)}$

$$S_{ij}^{(o)} = 2 \log \frac{p_{ij}^{(o)}}{p_i^{(o)} \times p_j^{(o)}}$$

$$P[X,Y] = P[X] \times P[Y] \rightarrow X, Y \rightarrow \text{independent}$$

$$P[X,Y] \neq P[Y/X] \cdot P[X] \rightarrow X, Y \rightarrow \text{dependent}$$

$$\begin{array}{c|ccccc} & P & C & V & A & P \\ \hline P & & C & A & A & C \\ C & P & & V & A & P \\ P & A & C & A & V & \end{array}$$

	A	C	P	V	
A	6x2 + 3	1	1	2	18
C	3	1x2 + 7	3		15
P	1	7	4x1 + 2		18
V	2	3	2	1x1	9
					60

$$P_A = \frac{18}{60}, P_C = \frac{15}{60}, P_P = \frac{18}{60}, P_V = \frac{9}{60}$$

$$P_{ij}^o = \frac{9}{60}$$

$$S_{ij}^o = 2 \log_2 \frac{P_{ij}^o}{P_i^o \times P_j^o}$$

$$S_{AA} = 2 \log_2 \frac{P_{AA}}{P_A \times P_A}$$

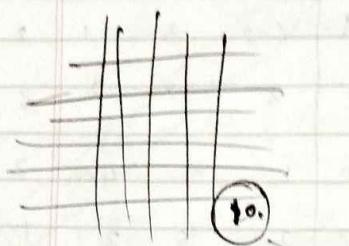
$$S_{AA} = 2 \log_2 \frac{\frac{12}{60}}{\frac{18}{60} \times \frac{18}{60}}$$

BLOSUM45 : 45% simi \rightarrow PAM100

BLOSUM60 : 60% simi

BLOSUM80 : 80% simi

↓ divergence
Higher index are appropriate for more conserved proteins



2^{\log_2}

$2^{(10/2)} = 2^5 = 32$ means the matrix is 32% more efficient for alignments.

standard in protein database "scoring". Higher numbers better alignment

BLOSUM 62, 2008 Nature Biotechnology

RBLUM 62 (incorrect)

RBLUM 62 (correct)

still BLOSUM62 performs better than RBLUM62

With PAM matrices (Percentage difference)

The score indicates the percentage of substitutions per position \Rightarrow higher index are appropriate for more distant proteins

With BLOSUM matrices (percentage similarity)

The score indicates the percentage of conservation \Rightarrow higher index are appropriate for more conserved proteins

Scoring of MSA / Multiple sequence alignment

Gaps

$$\text{Total gap cost} = \text{Gap open cost} + (n-1) \text{gap ext. cost}$$

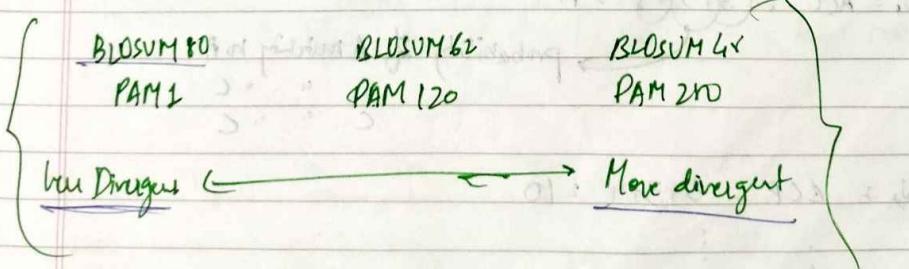
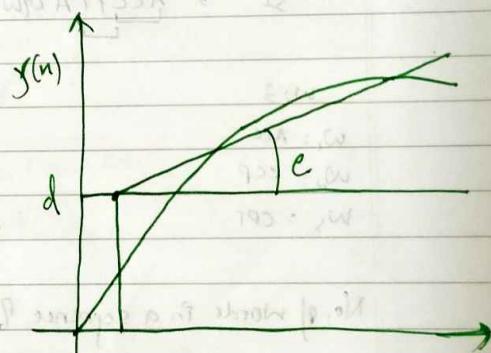
n gaps

A P A A C P
A P Q - - P
open gap Internal gaps

A P A A C C P C
P - - C C P C
End gaps -10 > -25

$$y(n) = d + (n-1)e$$

↓ ↓
 gap open gap extend



20/08/23

classmate

Date _____

Page _____

Pairwise Sequence Alignment

- Methods of Sequence Alignment

- ↳ Dot Plot Method.
- ↳ D.P.
- ↳ Substitution Matrices.
- ↳ Heuristic Methods.

Pattern Search

BLAST - Basic local alignment Search Tool, Altschul, 1990

- (i) Let us divide query sequence into small words

SI \rightarrow A C C P T A G W Y

$w = 3$

$w_1 = ACC$

$w_2 = CCP$

$w_3 = CPR$

No. of words in a sequence of length $L \Rightarrow L-w+1$

- (ii) Assign the score for each words.

$$w_1 = ACC \cdot (5+3+3) = 11$$

probability of a matching to A

C " " C
C " " C

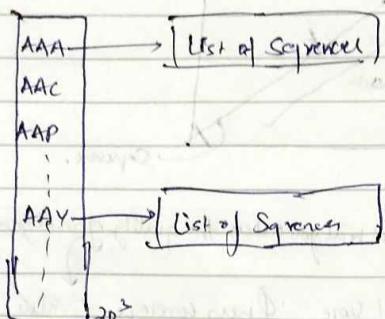
$$w_2 = CCP = 3+3+4 = 10$$

: : :

Let us decide a cut-off score $\rightarrow 10$

(iii) Scan the database sequence using these high scoring words.

Table of words vs Sequences in DB.



$\varphi = \text{ACCEPTAGWY}$

ACCPY

AGWY

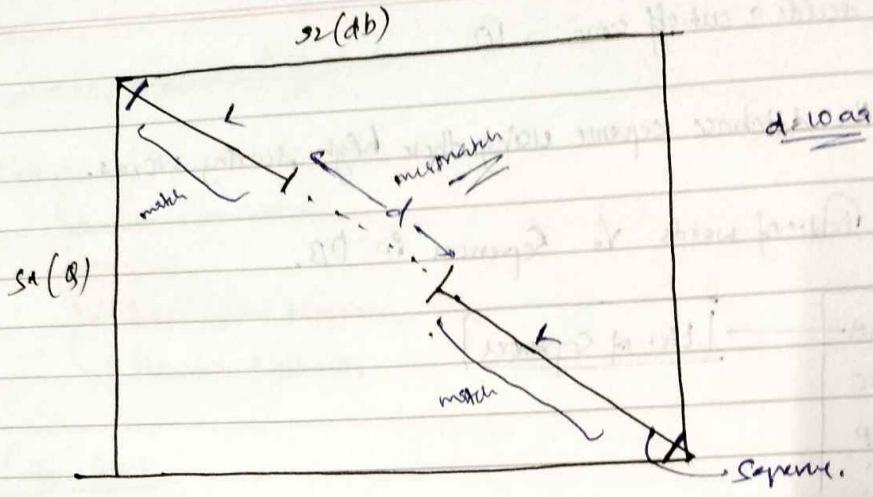
Sequence 1

Sequence 2

100aa

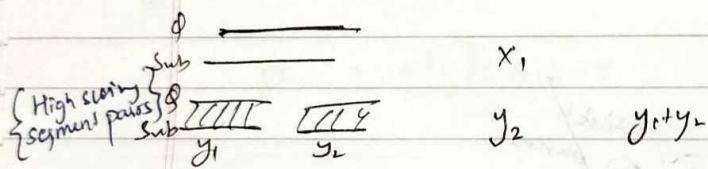
100aa

total 1000aa



BLAST

Raw Score . Total Score . Query coverage e-value SI



S is the aligned size of a sequence database in which the current motif would be found.

$S' = \frac{AS - \ln(Kmn)}{\ln 2}$

S' : Bit score.

$$S' = \frac{AS - \ln(Kmn)}{\ln 2}$$

$\lambda, K \rightarrow$ Constants

m, n - length of sequences being aligned

$$m = n = ?$$

$$P(\text{Score} \geq S) = \frac{1}{2^S}$$

Prob that we get a score similar to 'S' by just random alignment.

$$E\text{-value} = N \times P\text{-Value}$$

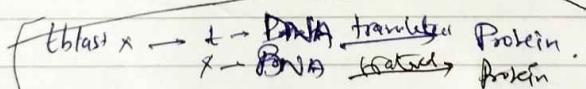
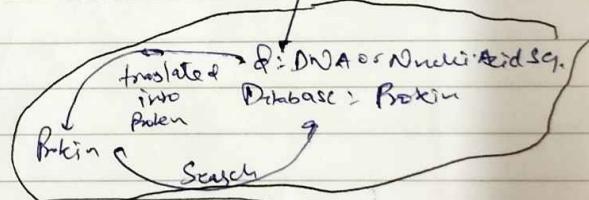
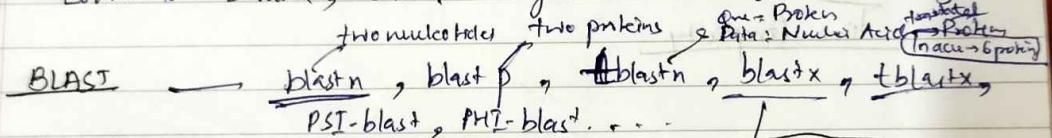
$$E_{\text{value}} = \frac{N}{2^S} \quad (N = n \times m)$$

query length

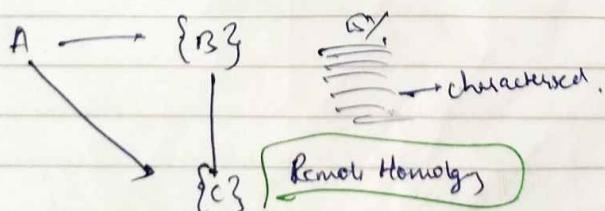
m: database length

Prob. of getting a better alignment by just random alignment.

Lower the E-Value, better the alignment m = database length



PSI-BLAST — Position specific Iterative BLAST



PSI-BLAST
Pattern Hit
Initiated

for detection of Remote homology

28/08/23

classmate

Date _____

Page _____

Comparison of Sequences

↪ Alignment

- ① Whole genome sequence comparison : Duplication / Translocations / Insertion / deletion.
- ② Differences obtained from NCIs : Some regions of the genome called pseudo genes.
- ③ Non Coding Regions
- ④

Time Complexity

27/08/21

Dynamic Programming

$$S_i^p - S_j^p \quad S_j^p \quad \left. \begin{array}{l} S_i^p \\ - S_j^p \\ S_j^p \end{array} \right\} \rightarrow \text{pair Segm}$$

 S_1, S_2, S_3 S_i^p
 S_j^p
 S_k^p Multiple Segm
alignment

$$\begin{array}{c|c|c} S_i^p & S_j^p & - \\ - & S_j^p & - \\ - & - & S_k^p \end{array} \quad \begin{array}{c|c|c} S_i^p & - & S_j^p \\ S_j^p & S_j^p & - \\ - & S_k^p & S_k^p \end{array}$$

? alignment

$$S(S_i, S_j, S_k) = (\max) \quad \left. \begin{array}{c} S_{i+1}, S_{j+1}, S_{k+1} + E_{i,j,k} \\ \vdots \\ \vdots \end{array} \right\}$$

$$\left. \begin{array}{l} \text{no. of sequences} = k \\ \text{size of sequence} = h \end{array} \right\}$$

Time Complexity $O((2^{k-1})^h \cdot (n^k))$

Heuristic Methods(i) Star Alignment

$S_1 : ACT$ → Start Sequence

$S_2 = TCT$

Δ

$S_3 : CT$

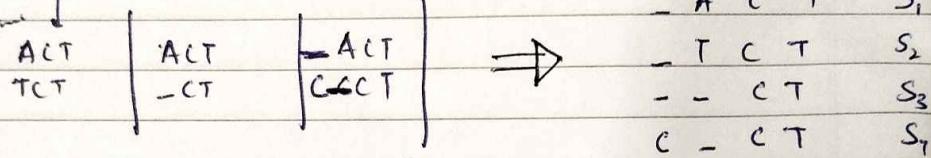
Aligning S_2, S_3, S_4 with S_1

$S_4 = CCT$

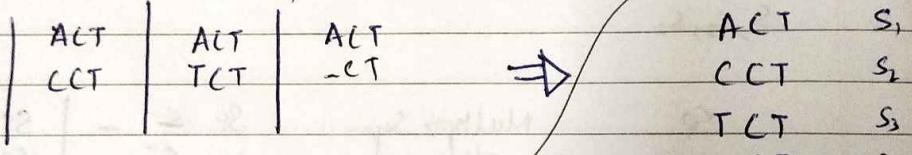


(ii) One of the sequences from multiple representations

$(S_1, S_2), (S_1, S_3), (S_1, S_4)$



$(S_1, S_4), (S_1, S_2), (S_1, S_3)$

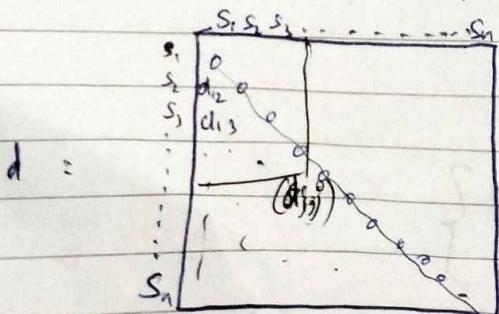


different

How to fix this order problem

Progressive Method

(i) Calculate a distance matrix b/w Sequences



$$d_{i,j} = \frac{s_{i,j}}{l_{i,j}}$$

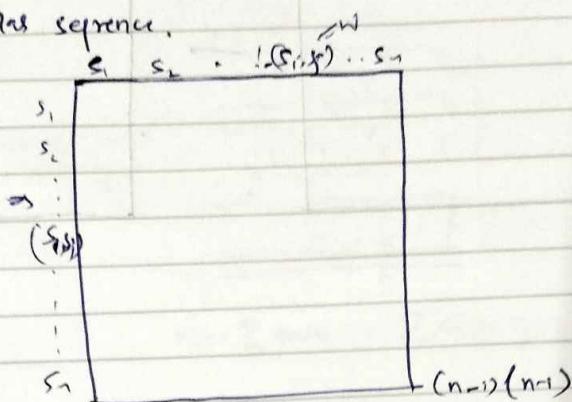
$s_{i,j}$ = no. of substitution

$l_{i,j}$ = length of the alignment

gaps are not counted.

(ii) To find the most similar Sequence.

(iii) Merge the most similar sequence.



(iv) PGMA (Pair Group Method - Arithmetic Mean)

$$d(w, x) = \frac{d(u, x) + d(v, x)}{2}$$

$$u = s^p$$

$$v = s^p$$

(v) WPGMA

$$d(w, x) = w_1 d(u, x) + w_2 d(v, x)$$

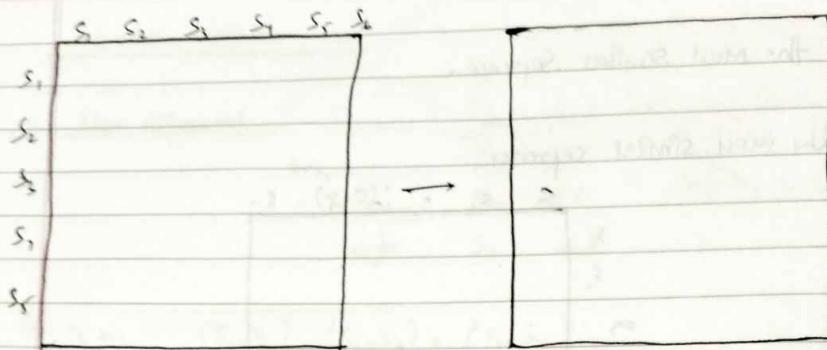
(vi) OverWeighted PGMA (UPGMA)

$$m(u) = \text{no. of nodes under } u$$

$$d(w, x) = a(u)d(u, x) + b(v)d(v, x)$$

$$a(u) = \frac{m(u)}{m(u) + m(v)}$$

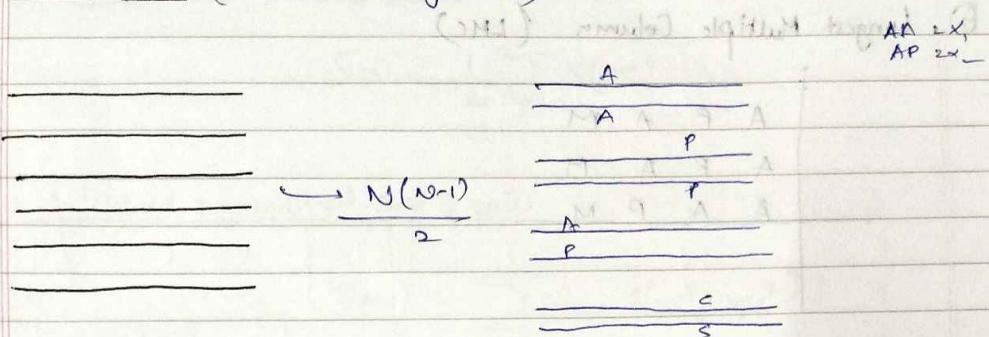
$$b(v) = \frac{m(v)}{m(u) + m(v)}$$



Multiple Sequence Alignment

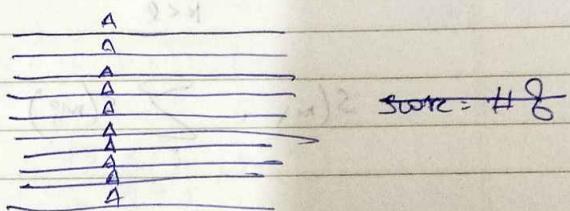
- ↳ How to align multiple sequences.
- (a) Dynamic Programming and its limitation $O(n^2)$
- (b) Heuristic method \rightarrow Progressive method of MSA (clustering)

T-COFFEE (Evaluate the alignment)



$$\text{No. 8 pairs} = \sum x_i$$

Score: # 8 constant pairs
total no. 8 pairs



For

$S_1 = APAM$

$S_2 = PAPA$

$S_3 = APPA$

$(S_1, S_2) = APAM$

$PAPA$

$(S_1, S_3) =$

$\begin{matrix} A & 2 & 3 \\ P & 3 & \\ M & 1 & \end{matrix}$

$(S_2, S_3) =$

$(A, A)^2 + (A, P)^2 + (P, P)^2 + (P, A)^2$

$= 1^2 + 1^2 + 1^2 + 1^2 = 4$

How to Score the MSA?

① Longest Multiple Column (LMC)

A	P	A	M
A	P	A	M
A	A	P	M

(2) Sum of Pairs

$$S(m_i) = \sum_{k < l} s(m_i^k, m_i^l)$$

$$s(m_i) = \sum_i s(m_i)$$

$s(m_i)$: Score at the i^{th} position

$s(m_i^k, m_i^l)$: Score b/w pairs of two sequence b/w k & l

$$s(m_1) = s(A, A) + s(A, A) + s(A, A) \quad (12) \\ 3s(A, A)$$

$$s(m_2) = s(P, P) + s(P, A) + s(P, A)$$

$$s(m) = s(m_1) + s(m_2) + s(m_3) + s(m_4) \quad \text{in}$$

Normalised sum of pairs :

$$S(m) = \frac{1}{n} \sum_i s(m_i)$$

Weighted Normalized sum of pairs

$$s(m_i) = \sum_{k=1}^K w_{k,i} s(m_i^k, m_k)$$

$$\text{MSA}_1 = 100$$

$$\text{MSA}_2 = 80$$

Star Consensus Sum of pairs

$$s(m_i) = \sum_{k=1}^K \delta(m_i^k, m_k)$$

Consensus sequence.

MSA
 A P A M M
 P A P A M
 A P A M M
A P A M M
 Consensus sequence

$$s(m_i) = s(A, A) + s(A, P) + s(A, M)$$

How to find Consensus sequence

A P A F
 A A P F
 S P S A
 Consensus sequence → { A P A / P S F }

Entropy Measurement (Calulation)

$$\text{Shannon's Entropy} = - \sum_{x \in \text{lecters}} \sum_{i=1}^L p_x \log_2 p_x$$

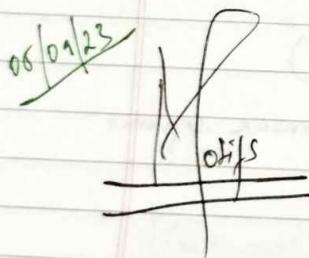
A C T G
 A C G G
 A T G C
 A C T A

$$E_1 = - \left[p_A \log_2 p_A + 0 \cdot 0 + 0 \right] = -4 \times (1 \times \log_2 1)$$

$$E_2 = - \left[3 \times (p_A \log_2 p_A) + (p_T \log_2 p_T) \right]$$

$$E_3 = - \left[2 \times (p_A \log_2 p_A) + 2 \times (p_T \log_2 p_T) \right]$$

$$+ \left[3 \times \left(\frac{3}{4} \log_2 \frac{3}{4} \right) + \frac{1}{4} \log_2 \frac{1}{4} \right]$$



concerned ; short sequence segment
having similar functions.

Signature / Patterns

① DNA : TATAAT → RE GAATTCT → TF CACGTC

② RNA : RRS (AGYGAGYGV) → Splicing of pre-mRNA

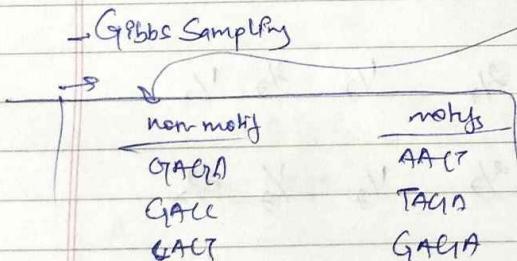
③ Protein : PKKKRKV (NLS), X-C-X-C-X-H-X-H
 ↴ Zinc finger Motifs

- ① Genes, patterns \rightarrow Pattern Matching,
- ② Genes, do not know the pattern \Rightarrow Discover the pattern.
- ③ Patterns don't know the set of genes \Rightarrow Classification

How to find motifs?

- ① Enumeration (word-based) methods:
- ② Probabilistic methods
- ③ Nature-Inspired Algorithms

Probabilistic Methods



~~S₁~~ ACGA GATA
~~S₂~~ GAGA ACTA
~~S₃~~ TAGA GACC
~~S₄~~ CACT GAGA

$$p_{c,0} = \frac{n_{c,0} + k_c}{(N-1)(L-W) + k_b}$$

$n_{c,0}$:= # of times 'c' has appeared in the non-motif region

k_c, k_b - pseudocounts

$$k_c = 1, k_b = 4 \text{ (DNA)} \\ = 20 \text{ (AA)}$$

L = length of sequence

N = No. of Sequence

W = motif length

$$P_{C,k} = \frac{N_C k + K_C}{(L-w)(N-1) + K_b}$$

$$P_{A,k} = \frac{N_A k + K_A}{(L-w)(N-1) + K_b}$$

(1)

Non-Motif

G A G A
 G A C C
 C A C T

k=0

A: $\frac{5}{16}$

A A C T
 T A G A
 G A G A

k=1 k=2 k=3 k=4

$\frac{2}{16}$ $\frac{4}{16}$ $\frac{1}{16}$ $\frac{3}{16}$

C: $\frac{5}{16}$

$\frac{1}{16}$ $\frac{1}{16}$ $\frac{2}{16}$ $\frac{1}{16}$

G: $\frac{4}{16}$

$\frac{2}{16}$ $\frac{1}{16}$ $\frac{3}{16}$ $\frac{1}{16}$

T: $\frac{2}{16}$

$\frac{2}{16}$ $\frac{1}{16}$ $\frac{1}{16}$ $\frac{2}{16}$

S₁ = ACGAGATA

	<u>motif</u>	<u>non-motif</u>

S1 = ACGA GATA

	Motif	Non-Motif	(Motif)
ACGA	$\frac{2}{4} \times \frac{1}{3} \times \frac{2}{3} \times \frac{3}{4}$	$\frac{5}{8} \times \frac{1}{16} \times \frac{9}{16} \times \frac{5}{16}$	0.18
CGAG	$\frac{1}{4} \times \frac{1}{3} \times \frac{1}{3} \times \frac{1}{3}$		= 0.068
GAGA	$\frac{2}{4} \times \frac{4}{3} \times \frac{3}{3} \times \frac{3}{4}$	$\frac{1}{4} \times \frac{7}{16} \times \frac{5}{16} \times \frac{8}{16}$	$\boxed{4.91}$
ACTAT	$\frac{2}{4} \times \frac{1}{3} \times \frac{1}{3} \times \frac{2}{3}$	$\frac{5}{16} \times \frac{4}{16} \times \frac{5}{16} \times \frac{2}{16}$	= 0.64
GATA	$\frac{2}{4} \times \frac{4}{3} \times \frac{1}{3} \times \frac{3}{4}$	$\frac{4}{16} \times \frac{5}{16} \times \frac{2}{16} \times \frac{5}{16}$	= 3.24

Non-motif (removing S1) = CACT GAGA

CACT	GA/G(A)
CA CC	TAG(A)
GAGA	AACT
<u>Non-motif</u>	motif

A

C

G

T

Artificial Bee Colony Algorithm

Position =

4	3	1	5
---	---	---	---

$$\begin{array}{r}
 A \quad G \quad A \quad T \\
 G \quad A \quad A \quad C \\
 T \quad A \quad G \quad A \\
 G \quad A \quad G \quad A \\
 \hline
 P(GAAGA) = \dots
 \end{array}$$

$$\begin{array}{r}
 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8 \\
 S_1 \quad A \quad C \quad G \quad A \quad G \quad A \quad T \quad A \\
 S_2 \quad G \quad A \quad G \quad A \quad A \quad C \quad T \quad A \\
 S_3 \quad T \quad A \quad G \quad A \quad G \quad A \quad C \quad C \\
 S_4 \quad C \quad A \quad C \quad T \quad G \quad A \quad G \quad A \\
 \hline
 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8
 \end{array}$$

$$\begin{array}{r}
 1 \quad 2 \quad 3 \quad 4 \\
 \hline
 \end{array}$$

$$A \quad C \quad G \quad A$$

$$A \quad G \quad A \quad A$$

$$G \quad A \quad G \quad A$$

$$T \quad G \quad A \quad G$$

$$P(A.GA.GA) =$$

no word length

$$\text{Score} \times \sum \frac{\max(b_i)}{w \times N}$$

$$\frac{2/4 + 3/4 + 2/4 + 2/4}{4 \times 4}$$

How to find Motifs ??

How to represent motif?

① Consensus pattern

P A T A A T

A T T A A T

C A T A A G

C A T A A T

C A T A A T

NOPAC

A and G = R

C and T = Y

A and T = W

G and C = S

A and C = M

G and T = K

A, C and G = not T = V

Count ~~W~~ T
next ~~Y~~ U
↓

② Regular Expression

= [ACG,T] - [AT] - T - A - A - [GT]

C, G & T = not A = B

• [ACGAT] - {Cg} = T - A - A - {AC}

A, C, G, T = N → Any

A, G & T = not C = D

∴ [ACGAT] - {Cg} - T - A - {AC}

↳ This can also have a range

$A^{(n,m)}$ → maximum no. of times
min no.
of times

$A^{(n, \cdot)}$, $A^{(\cdot, m)}$

e.g. $A^{(2, \infty)}$ $A^{(\cdot, 2)}$

2 to ∞ ∞ to 2

(3) Position Weight Matrix

		Count Matrix					
		1	2	3	4	5	6
A		1	3	0	4	4	0
C		1	0	0	0	0	0
G		1	0	0	0	0	1
T		1	1	4	6	6	2

frequency Matrix

$$\begin{bmatrix} Y_A & 3/4 & 0 & 1 & 2 \\ Y_A & 0 & 0 & 0 & 0 \\ Y_A & 0 & 0 & 0 & 0 \\ Y_A & Y_A & 1 & 0 & 0 \end{bmatrix} \quad f_{ij} = \frac{n_{pij}}{\sum n_{pj}}$$

Position Weight Matrix

$$W_{ij}^o = \ln \left(\frac{f_{ij}^o}{p_i^o} \right) \quad p_i^o = \text{prior or background prob. of letter } (A, C, G, T)$$

$$W_{ij}' = \ln \left(\frac{f_{ij}'}{p_i'} \right)$$

$$f_{ij}' = \frac{n_{pij} + k}{\sum n_{pj} + k} \quad k: \text{pseudo Constant}$$

Most Probable Motif

(Pst-Blast)

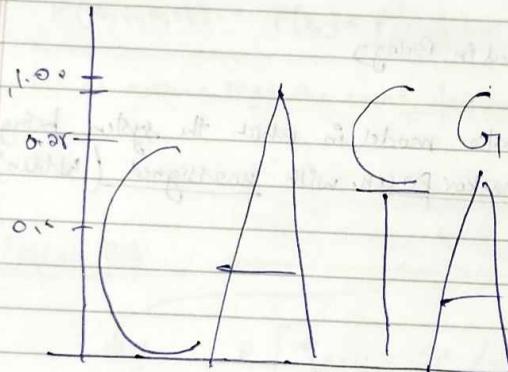
C A C A . .

$$0.93 + 1.25 + 1.25 + 0.93$$

WM - (ATAA)

	1	2	3	4
A	0.43	1.25	0.93	1.3
C	0.73	-0.5	1.25	0.55
G	-0.58	-1.3	-0.95	0.25
T	1.25	0.6	0.60	-0.6

(4) Web logo / Sequence logo



PSSM → how to measure the quality of PSSM ??

Position-Specific Scoring Matrix

Information content of PSSM

$$I_{pj} = f_{pj} \log_2 \left(\frac{f_{pj}}{p_i} \right)$$

$$\text{Intrinsic} = \sum_j \sum_i I_{pj}$$

better Quality

Hidden Markov Model

$$P(AAAA) = \left(\frac{1}{4}\right)^4$$

$$P(AATA) = \left(\frac{1}{4}\right)^4$$

$$P(AATA) = \left(\frac{1}{4}\right)^4$$

$$\vdots (4^4)$$

$$P(TTTT) = \left(\frac{1}{4}\right)^4$$

DNA

$$P(AAAA) = \left(\frac{1}{2^4}\right)^4$$

⋮

$$(2^4)^4$$

$$P(YYYY) = \left(\frac{1}{2^4}\right)^4$$

protein

- Implemented in Electronics, used in Biology

→ Defined as statistical markov model in which the system being modeled is assumed to follow markov process, with unassigned (hidden) states.

Markov process

↳ A process which follows statistical Markov properties.

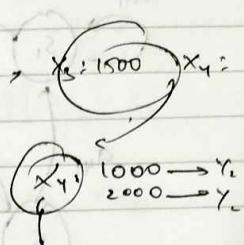
Markov Property

↳ If anyone is able to predict the future given only the present

Dollars vs Pakistan 1000 bet

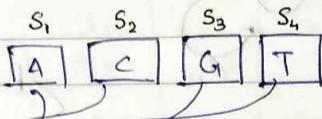
INR 1000

$$X_0 = 1000 \text{ }, X_1 = 1500 \text{ }, X_2 = 2000 \text{ }, X_3 = 1500 \text{ }, X_4 = \dots$$



X_4 is predicted given X_3 .

This is called Markov property.



$P(CTAG)$:

$$P(S_0, S_1, S_2, S_3) = P[S_0] \times P[S_1/S_0] \times P[S_2/S_1, S_0] \times P[S_3/S_2, S_1, S_0]$$

Let us assume that the event follows the Markov property.

$$= P[S_0] \times P[S_1/S_0] \times P[S_2/S_1] \times P[S_3/S_2]$$

$$\boxed{P[S_0] \times q_{24} \times q_{41} \times q_{13}}$$

Let's Define

$$\star \quad a_{ij} = P[q_{t+1} = S_j^o / q_t = S_i^o]$$

q_t : Represents a state at time t .

+ i, j

* Transition Prob

$$a_{ij} \geq 0 \quad \forall i, j$$

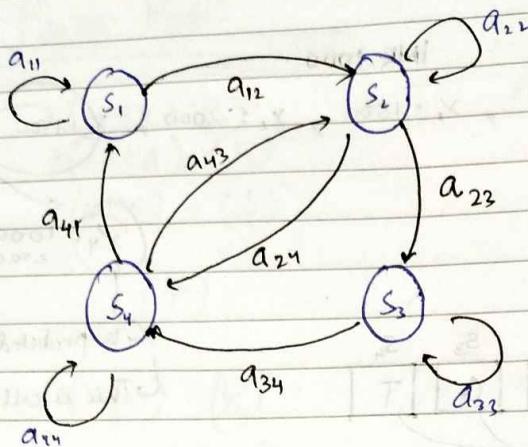
$$\sum a_{ij} = 1$$

Initial State Property

$$\star \quad \pi_i^o = P[q_1 = S_i^o] \quad \forall i$$

$$\sum \pi_i^o = 1 \quad 1 \leq i \leq N$$

$$= \boxed{\pi_2 \times q_{24} \times q_{41} \times q_{13}}$$



every node is connected.

Non-Ergotic Model vs Ergotic Model

If any one of the node is not directly connected then it is called non-ergotic model.

① No. of state $S = \{S_1, S_2, \dots, S_n\}$

② Transition prob. $= a_{ij}$

③ Initial State Prob. $= \pi_i^0$

Example Weather Prediction

S_1 : Sunny

S_2 : Rainy

S_3 : Cloudy

0 : $S_2 S_3 S_1 S_1 S_2 S_3 S_2$

$$P[\theta] = P[S_1 S_2 S_3 S_4 S_5 S_6]$$

$$= P[S_1] \times P[S_2/S_1] \times P[S_3/S_1 S_2] \times P[S_4/S_1 S_2 S_3] \times P[S_5/S_1 S_2 S_3 S_4] \times P[S_6/S_1 S_2 S_3 S_4 S_5]$$

$$\pi_1 \times a_{23} \times a_3 \times a_{11} \times a_{12} \times a_{23} \times a_{32}$$

(Given)

	S_1	S_2	S_3
S_1	0.4	0.3	0.3
S_2	0.2	0.6	0.4
S_3	0.1	0.7	0.2

$$P[\theta] = 1 \times 0.2 \times 0.1 \times 0.4 \times 0.3 \times 0.2 \times 0.7$$

Ans What is the prob. same weather in Cont. for d days

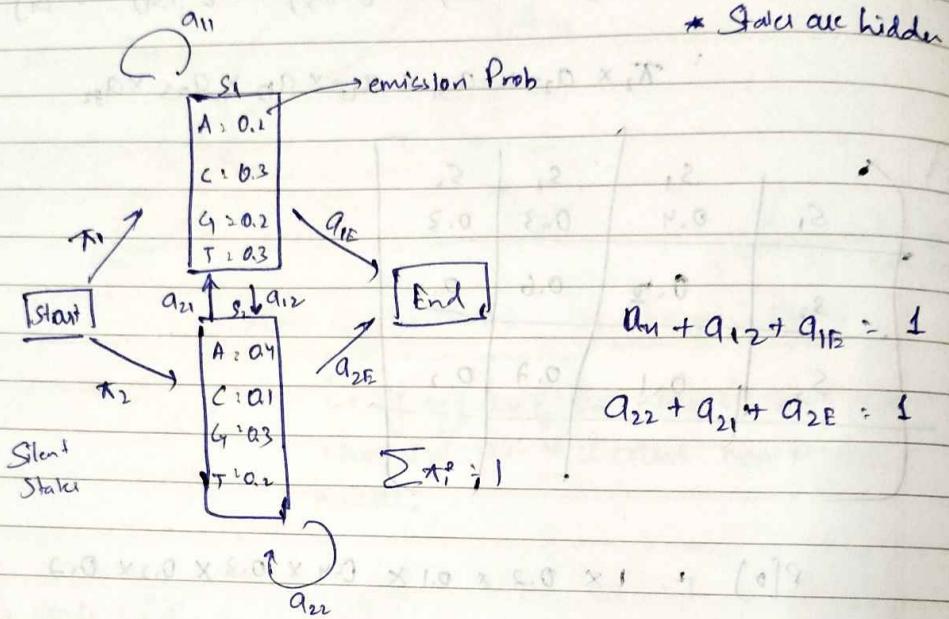
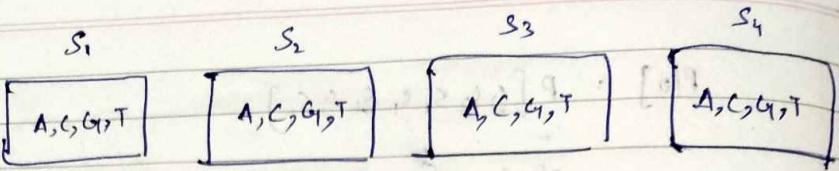
$$= \pi^e \times (a_{ii})^{d-1} (1-a_{ii}) \cdot p^e(d)$$

$$\bar{d}_i = \sum_{a=1}^6 a \cdot p^e(a) = \sum_{a=1}^6 d_i (a_{ii})^{d-1} (1-a_{ii}) \cdot \frac{1}{1-a_{ii}}$$

$$\bar{d}_1 = \frac{1}{1-0.4}$$

$$\bar{d}_2 = \frac{1}{1-0.6}$$

$$\bar{d}_3 = \frac{1}{1-0.2}$$



Observation $L=1$,

$$\{0, 1, E\}, \{0, 2, E\}, \{$$

Observation $L=2$

$$\{0, 1, 1, 3\}, \{0, 1, 2, 3\}, \{0, 2, 1, 3\}, \{0, 2, 2, 3\}$$

Possible Paths $\rightarrow 2^L \rightarrow N^L$