

Sequence Alignment & Similarity Search.

Motivations:

1. A new gene \rightarrow what is its function?

Compare with the database of known genes
& based on similarity, formulate hypothesis.

e.g. 1984 \rightarrow newly discovered cancer-causing
viral oncogene (resides in a virus) with all known genes & found
a close match with a gene involved in
growth & development of cell.

2. A segment of chromosome identified for
certain function \rightarrow what
gene resides there? \Rightarrow Compare sequenced
segment with the gene database

e.g. cystic fibrosis gene in chromosome 7 (human)
in 1989, matches with a gene
responsible for coding ATP binding proteins.

① DNA "sequence similarity" is more than a "string matching" problem.

- Evolutionary processes: Mutation, replication error causes insertion & deletion & substitution.
- Exon-intron chains.

e.g. ATATAT → Similar from the angle of genome analysis.
TATA TA →

⊛ Edit distance: (Vladimir Levenshtein, 1966)

"Edit distance" between two strings is the "minimum no." of "editing operations" needed to transform from one string into another, where the edit operations are insertion of a symbol, deletion of a symbol & substitution of one symbol for another.

e.g.

TGCATAT
 ↓ delete T
 TGCATATA
 ↓ delete last A
 TGCATA
 ↓ insert A at the front
 ATGCATA
 ↓ substitute G for C in the 3rd. pos.
 ATCCAT
 ↓ insert a G before the last A
 ATCCGAT

∴ edit dist. (TGCATAT, ATCCGAT) = 5

TGCATAT
 ↓ insert A at the front
 ATGCATAT
 ↓ delete T in the 6th pos.
 ATGC AAT
 ↓ subst. G for A in the 5th pos.
 ATGCCAAT
 ↓ substitute C for G in the 3rd pos.
 ATCCGAT

no. edit dist.
 (TGCATAT, ATCCGAT) = 4

∴ edit dist. ⇒ 4

→ 5 Edit operator:

TGCATAT → ~~TGC~~ATAT
 A?C?GAT--

I ⇒ Insert
 D ⇒ Delete
 S ⇒ Subst.

→ 4- Edit oper.

TGCATAT
 ATCCGAT

⇒ string matching problem with insertion of gaps (indels) in any of the strings (with the constraint that simultaneous insertions at any position are not allowed) and also substitution of symbols. (a cost factors are generally associated with all these operations).

Given m -string & n -string, S_m & S_n edit operators. \Rightarrow $\frac{m+n}{2}$ (only insertion & deletion)

Minimum edit oper. = $(n-m)$, $n > m$. (only deletion & subst.)

$\therefore (n-m) < \text{edit-distance}(S_m, S_n) < \frac{m+n}{2}$
 fighter based \hookrightarrow max. (m, n) .
 However if subst. is allowed, then $\frac{m+n}{2}$ is the answer.

⊕ Modeling edit operations:

- Two pointers movement (for each string in one pointer)
 - ~~start~~ Initialize both the pointers in front of starting alphabet
 - move either one of them from its current position to the next alphabet (as gap or indel is introduced in front of the static pointer.) A gap in the first string is insertion else (target string) it is deletion
 - move both the pointers (substitution / matching operation)
 - Stop when both the pointers arrived at the last alphabet.

Increased pointer.
 e.g. Ref: ~~P~~ T G C A T A T

Target: P_t A T C C G A T

→ Step 1. move P_t:
 P_r T G C A T A T
 A T C C G A T
 ↑ P_t
 Edit operation: I

Step 2. move P_r & P_t:
 P_r ↓
 T G C A T A T
 A T C C G A T
 ↑ P_t
 Edit operation: M

Step 3: move P_r & P_t:
 T G C A T A T
 A T C C G A T
 ↓ ↓ ↓
 Edit operation: S, M, S

Step 4 move P_r:
 T G C A T A T
 A T C C G A T
 ↓
 Edit operation: D

Step 7 & 8. move P_r & P_t:
 T G C A T A T
 A T C C G A T
 ↑ ↓ ↓
 ↑ ↑
 Edit operation: M, M.

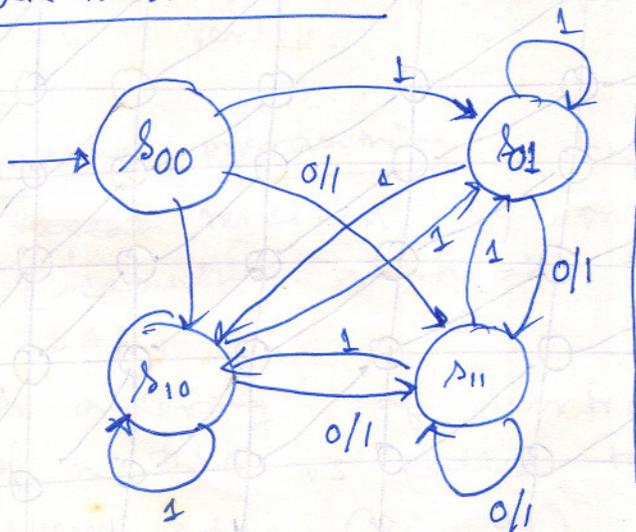
∴ Total no. of edit operations I + 4M + 2S + D
 ⇒ 4, M=0.

N.B. (1) Movement is allowed from one alphabet to the next (jumping index)

(ii) No. of $\downarrow P_r$ movements = n (length of ref. string)
 No. of $\uparrow P_t$ " " = m (" of target ")

(iii) 00 \Rightarrow No movement of $\downarrow P_r$ & $\uparrow P_t$ (only at initial state)
 01 \Rightarrow Movement of $\downarrow P_r$
 10 \Rightarrow " " $\uparrow P_t$
 11 \Rightarrow " " $\downarrow P_r$ & $\uparrow P_t$

A State transition model:



0/1 \rightarrow Cost for matching 0 else 1.

stop after when

$$\sum 1 = n + m$$

Let n_{ij} = no. of times S_{ij} occurred.

$$\therefore n_{01} + n_{10} + 2 \cdot n_{11} = m + n; \quad n_{01} + n_{11} = n; \quad n_{10} + n_{11} = m;$$

(no. of integer solutions of the above eqn.)

Give the no. of ~~cost~~ valid sequences.

\Rightarrow constraints:

$$n_{01} + n_{11} = n;$$

$$n_{10} + n_{11} = m;$$

$$n_{01}, n_{10}, n_{11} \geq 0;$$

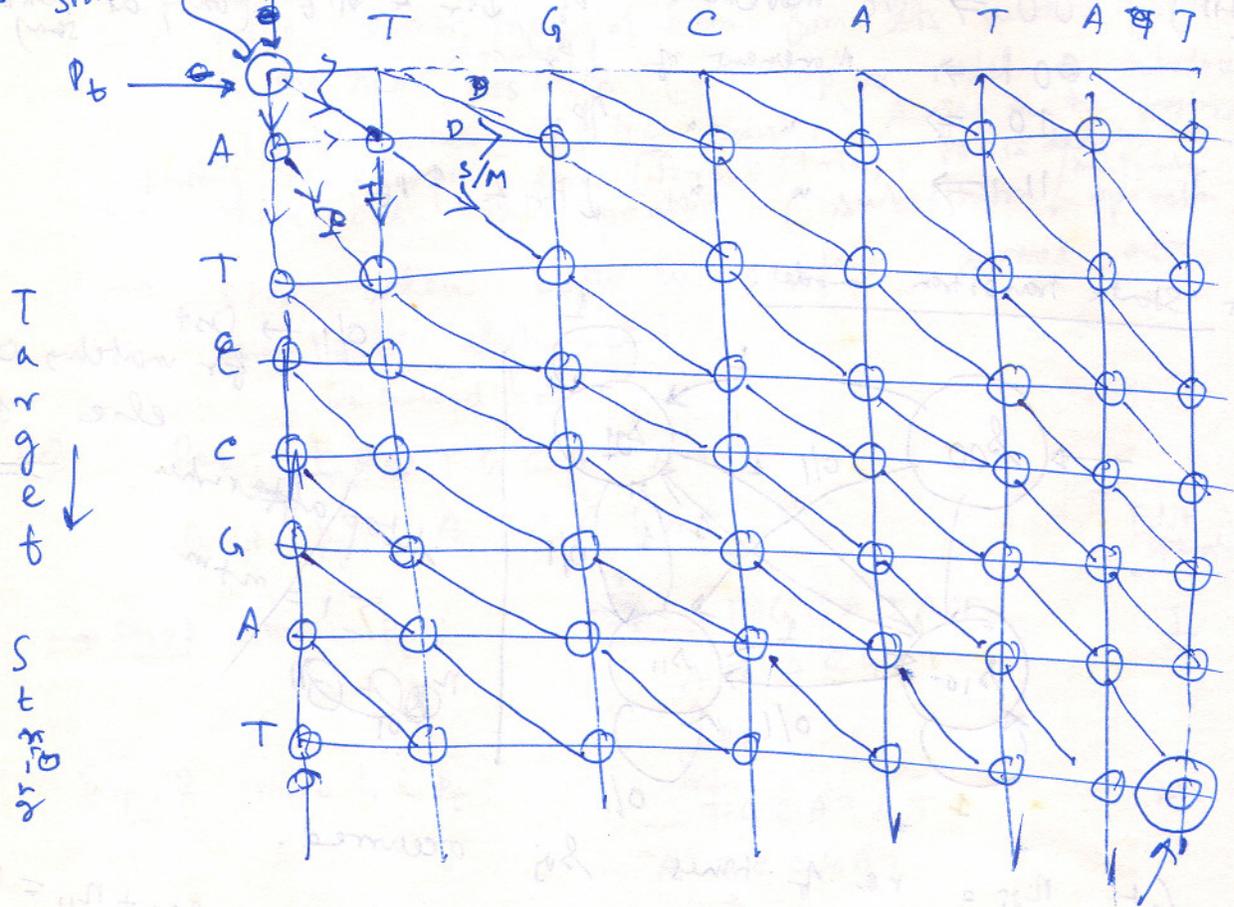
Content free transitions are not able to capture the ~~to~~ a concept of valid string.

It is quite difficult to ~~consider~~ ^{enumerate} all possible valid sequences & get the minimum out of them. States should be defined from the context.

Context Sensitive PSM: (Position of $\downarrow P_r$ & $\uparrow P_t$ define a State).

Ref. sm

e.g. Starts state, P_r



Transitions are possible only to right ^{stop state} neighbor, bottom / diagonal neighbor for D/I/S/M operations.

⇒ Computation of edit distance is the shortest path from the starting state to the stop state following the neighborhood definition. ^{equivalent to}

⇒ Directed Acyclic Graph (DAG).

- Alg.:
- ① Start from the starting vertex and initialise distances to its neighbours vertices.
 - ① For every vertex u ,
 - if its predecessors are visited, mark it ~~visited~~ visited & compute its dist. from the source as $\min \{ \text{dist. (predecessor)} + \text{cost between } (u, \text{pred.}) \}$ & also store the predecessor info. ^{iteratively}
 - ② Perform step 1 until you reach destination vertex.

Dynamic programming: Ensure no vertex is ~~revisited~~ revisited by proper traversal order (scanning order) \Rightarrow topological ^{sorted} order for a DAG

An ordering of vertices v_1, v_2, \dots, v_n of a DAG is called topological if every edge (v_i, v_j) of the DAG connects a vertex with a smaller index to a vertex with a larger index, i.e., $i < j$.

\rightarrow In this case row-wise / column-wise scanning / zigzag scanning should be ok. $\Rightarrow O(m \cdot n)$

\rightarrow For every vertex store also the predecessor from which it is the shortest path from the source. This helps in retracing the path back.

\rightarrow If weights are inversely ^{related} proportional to dist., then ~~maximum~~ ^{longest} path will provide the solution (maximisation problem). \Rightarrow Used for sequence alignment.

* Longest Common Subsequence problem:

A subsequence of a string v is simply an (ordered) sequence of alphabets (not necessarily consecutive) from v .

e.g. $v = \text{ATTGCTA}$

$A, G, C, A, A, T, T, A \in$ ^{set of} $\text{subsequences of } v$
 but $T, G, T, T, T, C, G \notin$

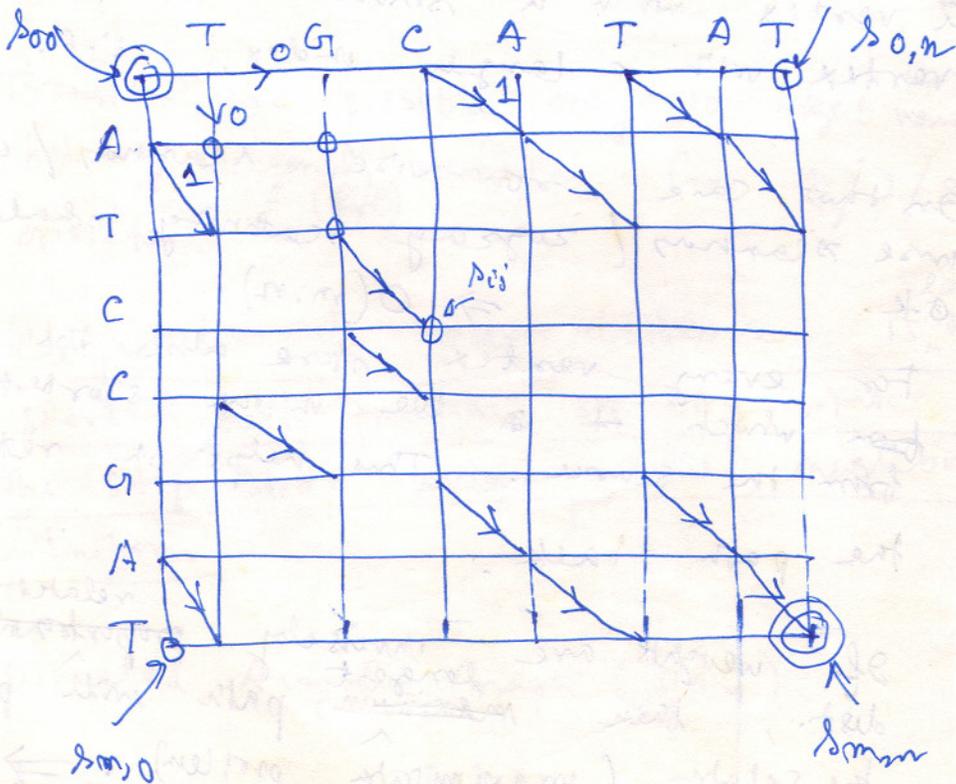
Problem Statement:

Let $v = v_1 v_2 \dots v_n$ & $w = w_1 w_2 \dots w_m$ be two strings. Find the longest common subsequence of v & w .

s.t. $v_{i_t} = w_{j_t}$

& $0 < i_1 < i_2 < \dots < i_n, 1 < j_1 < j_2 < \dots < j_m$

\Rightarrow No mutation or substitution allowed.



$$\begin{cases} \delta_{i,0} = \delta_{0,j} = 0, & 1 \leq i \leq m, \quad 1 \leq j \leq n \\ \delta_{i,j} = \max \begin{cases} \delta_{i-1,j} + 0 \\ \delta_{i,j-1} + 0 \\ \delta_{i-1,j-1} + 1, & \text{if } v_i = w_j \end{cases} \end{cases}$$

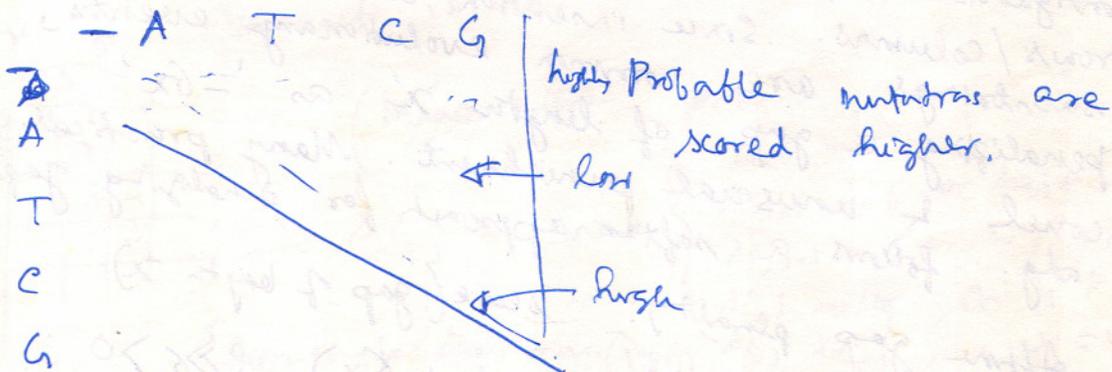
* B. (c) whose backtracking pointers $\leftarrow, \uparrow, \swarrow$ at every vertex.

(b) Computation of edit distance:

$$d_{i,0} = i; \quad d_{0,j} = j; \quad 1 \leq i \leq m, \quad 1 \leq j \leq n;$$

$$d_{i,j} = \min \begin{cases} d_{i-1,j} + 1 \\ d_{i,j-1} + 1 \\ d_{i-1,j-1} + 0 & \text{if } v_i = w_j \end{cases}$$

* Global Sequence Alignment: Needleman & Wunsch (1970)
 (The use of context sensitive FSM with all ^{edges} transitions).
 scoring matrix ($\delta(\)$)



For protein comparison 21x21 scoring matrix.:

PAM (Point accepted Mutation) \rightarrow

BLOSUM (Block substitution) \rightarrow

$$D_{i,j} = \max \begin{cases} D_{i-1,j} + \delta(v_i, -) \\ D_{i,j-1} + \delta(-, w_j) \\ D_{i-1,j-1} + \delta(v_i, w_j) \end{cases}$$

def: $\delta(v_i, -) = \delta(-, w_j) = -\sigma$ [insertion & deletion]

$$\delta(v_i, w_j) = -\mu \text{ (Mismatch), for } v_i \neq w_j \\ = 1 \text{ (match)}$$

$$D_{i,j} = \max \begin{cases} D_{i-1,j} - \sigma \\ D_{i,j-1} - \sigma \\ D_{i-1,j-1} - \mu, \text{ if } v_i \neq w_j \\ D_{i-1,j-1} + 1, \text{ if } v_i = w_j \end{cases}$$

∴ Resulting score = no. of matches $- \mu \times$ no. of mismatches $- \sigma \times$ no. of insertion & deletion.

Alignment with 'Gap' penalties:

A gap in an alignment is defined as a contiguous sequence of spaces (indels) in one of the rows/columns. Since insertion and deletions of substrings are common evolutionary events, penalizing a gap of length 'x' as '- σx ' is cruel & unusual punishment. Many practical alg. follows a softer approach for penalizing gaps.

Affine gap penalty: score(gap of length x)

$$= -(p + \sigma x), \text{ } p, \sigma > 0$$

$$\downarrow$$

$$D_{i,j} = \max \begin{cases} D_{i,j-1} - \sigma \\ D_{i-1,j} - (p + \sigma) \end{cases}$$

$$\uparrow$$

$$D_{i,j} = \max \begin{cases} D_{i,j-1} - \sigma \\ D_{i-1,j} - (p + \sigma) \end{cases}$$

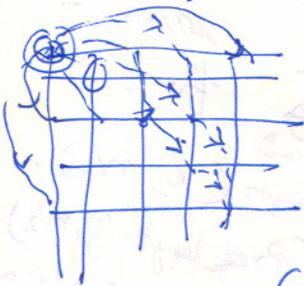
$$\delta_{i,j} = \max \begin{cases} \delta_{i-1,j-1} + \delta(v_i, w_j) \\ \delta_{i,j} \\ \delta_{i,j} \end{cases}$$

N.B. [For every (i,j) the vertex maintain $\delta_{i,j}$, $\delta_{i,j}$ & $\delta_{i,j}$ separately].

* Local Sequence Alignment problem: Smith-Waterman (1981)

Problem statement: Given strings v & w and a scoring matrix δ , compute substring of v & w whose global alignment, as defined by δ , is maximal among all global alignments of all substrings of v & w .

Model: Additional edges from the starting state (vertex) to each state (vertex) with "zero" cost & also edges from the each vertex to the destination vertex.



(Maximisation problem)

$$\therefore \delta_{i,j} = \max \begin{cases} 0 \\ \delta_{i-1,j} - \delta \text{ or } \delta_{i,j} + \delta(v_i, -) \\ \delta_{i,j-1} - \delta \text{ or } \delta_{i,j-1} + \delta(-, w_j) \\ \delta_{i,j} - \delta_{ij}, \text{ if } v_i \neq w_j \\ \delta_{i,j} + 1, \text{ if } v_i = w_j \end{cases}$$

$\delta(v_i, w_j)$

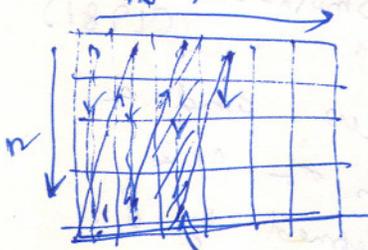
(Smith-Waterman (1981))

* Space Efficient Alignment Algo.:

Divide & Conquer Strategy:

→ Usual Dynamic Algo.: Time Complexity: $O(nm)$ (no. of edges)
 Storage: $O(nm)$ (no. of vertices)
 [For storing back-back pointers].

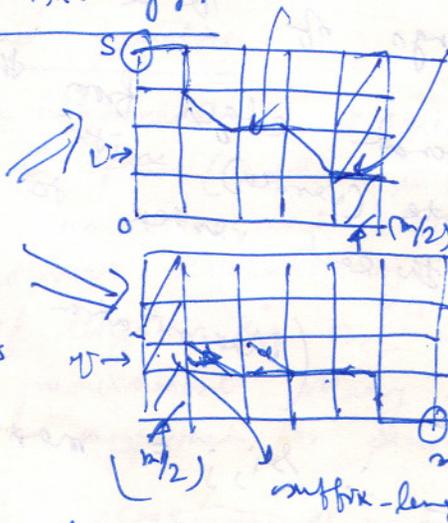
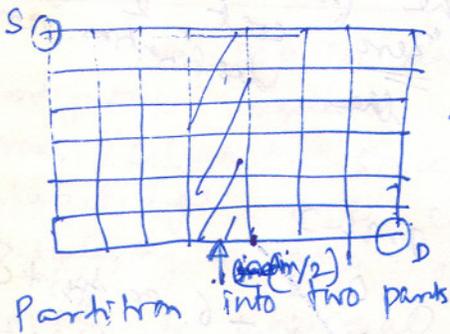
→ However, if only 'score' is required (not the aligned sequence), $O(n \text{ or } m)$ is sufficient. Store only previous row or column during scanning.



remember only these scores at next column.

Time complexity: $O(nm)$
 space: $O(n)$ } When only maximal score is required (or edit distance required).

→ Divide & Conquer Strategy:



Prefix-length $(i, m/2)$
 Compute the maximal score at every vertex of the shaded area (from S)

- same -
 for D.

∴ maximal length $(S, D) = \max(\text{Pref-length}(i, m/2) + \text{suffix-length}(i, m/2))$
 & the vertex $(i, m/2)$ will be on the path for which it is max.

i.e. $i = \text{argmax}(\dots)$

Computational requirement: $O(n \cdot m/2) + O(n \cdot m/2) \approx O(nm)$
 Storage: $O(n) + O(n)$.

6. Four Russian Speed Up: (Proposed by Vladimir Arlazarov, Effim Dinic, Michael Kronrod, & Igor Pasadzev in 1970, applied to sequence comparison by William Masek & Michael Paterson in 1980).

Store all optimal paths & score for all strings: (Storage space: $4^t \times 4^t$) \Rightarrow for nucleotides

$$\text{Let, } t = \frac{\log n}{4}$$

$$\therefore n = 2^{4t} = 4^{2t}$$

$$\therefore 4^t \times 4^t = n^{\frac{1}{2}} \times n^{\frac{1}{2}} = n$$

\therefore Storage space requirement: $O(n)$ for all optimal paths.

This makes computation of $\beta_{i,j}$ for access to string $\Rightarrow O(\frac{n}{t}) \Rightarrow O(\log n)$

$$\therefore \text{Computation of all } \beta_{i,j} \Rightarrow O\left(\frac{n}{t} \times \frac{n}{t}\right) \times \log n$$

$$\Rightarrow O\left(\frac{n^2}{t^2}\right) \times \log n$$

$$\Rightarrow O\left(\frac{n^2 (\log n)}{(\log n)^2}\right)$$

$$\Rightarrow O\left(\frac{n^2}{\log n}\right)$$

\Rightarrow Discuss lexon chaining problem

(i) Suffix based Alignment.

* Heuristic Similarity Search Techniques:

→ Dynamic Algorithm: $O(n^2)$ not suitable for searching through large genome database (10^{10} nucleotides) & query string ($\sim 10^3$)

→ Starting in the early 1990's biologists had no choice but to use fast heuristics as an alternative to quadratic sequence algorithm.

→ D. Lipman & W.R. Pearson (1985) → FASTA

→ S. Altschul, W. Gish, W. Miller, E. Myers and J. Lipman (1990) → Basic Local Alignment Search Tool (BLAST)

→ Modified in 1997.

• l -mer filtration technique for approximate pattern matching (with k -mismatches)

If an n -letter substring of a query approximately matches an n -letter substring of the text, then the two substrings share at least one l -mer for a sufficiently large value of l .

Th.: If the strings x_1, \dots, x_n & y_1, \dots, y_n match with at most k mismatches, then they share an l -mer for $l = \lfloor \frac{n}{k+1} \rfloor$, i.e. $x_{i+1} \dots x_{i+l} = y_{i+1} \dots y_{i+l}$ for some $1 \leq i \leq n-l+1$

Proof: Partition into $(k+1)$ groups, $\{1 \dots l\}, \{l+1, \dots, 2l\}, \dots$

— where $l = \lfloor \frac{n}{k+1} \rfloor$;

As there are k mismatches, one of this set will not have any mismatch. \square

This theorem motivates the following l -mer filtration algorithm for query matching with k -mismatches:

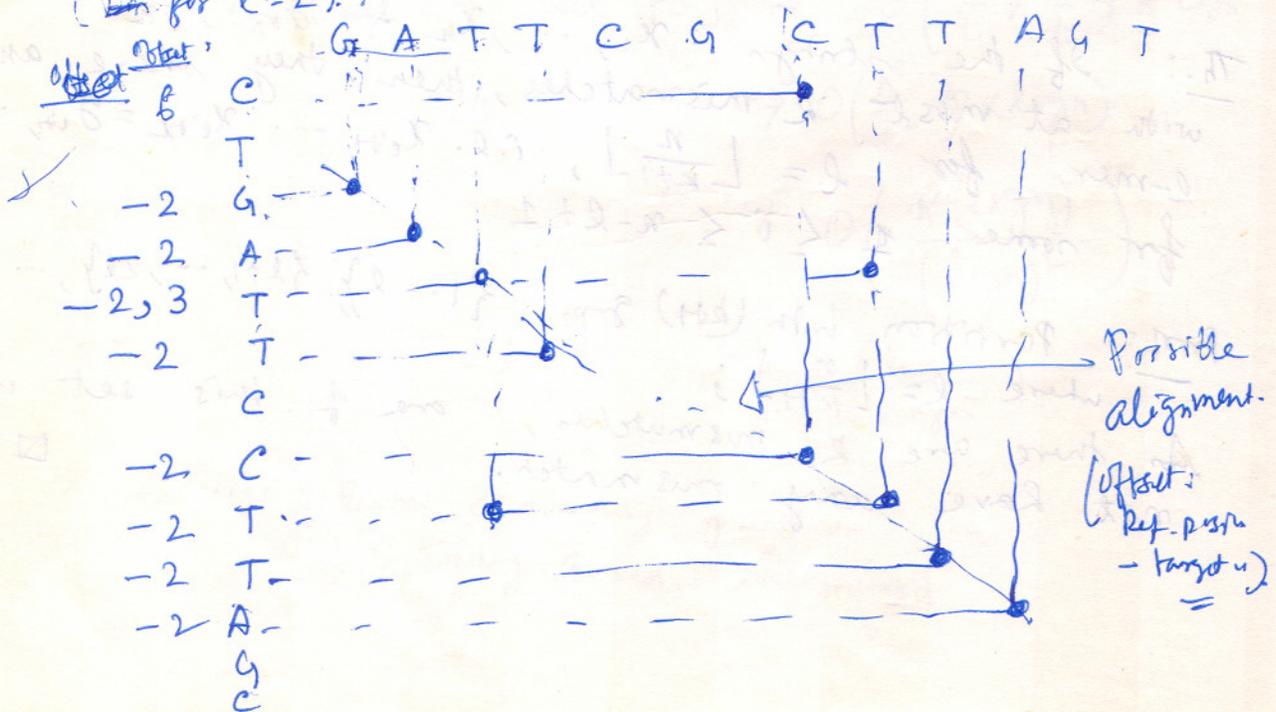
- Potential match detection: Find all matches of l -mers in both the query & the text for $l = \lfloor \frac{n}{k+1} \rfloor$. [Use of hashing or suffix tree]
- Potential match verification: Verify each potential match by extending it to the left & to the right until the first $k+1$ mismatches are found or the beginning/end of the query/text is found.

• FASTA:

→ FASTA search begins by breaking the search sequence into words (for nucleotides words of length 4 to 6, for proteins 1 to 2).

→ Find the positions of l -mer matches & represent in the form of dot matrices.

(In for $l=2$):



→ From this dot matrix, it selects some diagonal with a high concentration of 1's and groups these runs of ^{diagonal} 1's into longer runs. [~~Dyn~~, Local sequence alignment on these segments by dynamic programming).

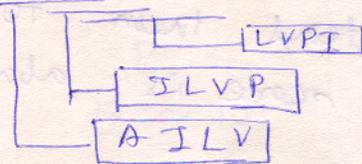
→ 'dots' are noted as offsets from the ref. positions. - Adjacent Alphabets with similar offsets are grouped together.

• BLAST:

Input sequence: AILVPTV

1. Break the query sequence into words.

A I L V P T V I G C T



} select those words which are biologically more relevant e.g. sequences with frequently observed aminoacids are discarded.

2. Search for exact word matches (also called high-scoring pairs or HSPs) in the database sequence.

A I L V P T V
M V Q G W A L Y D F L K C R A I L V G T V I A M L ...

(in both directions)

3. Extend the match - until the local alignment score falls below a fixed threshold (the most recent version of BLAST allows gaps in the extended match).

→ For any two l -mers x_1, \dots, x_l & y_1, \dots, y_l ,
 BLAST defines the segment score as

$$\sum_{i=1}^l \delta(x_i, y_i)$$
 where $\delta(x, y)$ is the similarity score
 between x & y .

→ A segment pair is locally maximal if
 its score cannot be improved by shifting
 it to any direction. BLAST attempts to
 find all locally maximal segment pairs in
 the query & database sequences with scores
 above some set threshold. If the threshold
 is high enough, then the set of all l -mers
 that have score above a threshold is not
too large. This becomes a multiple-pattern
 matching problem (use to keyword tree / suffix tree).

→ After the potential matches are located,
 BLAST attempts to extend them to see
 whether the resulting score is above the
 threshold.

→ The choice of the threshold in BLAST
 is guided by the Altschul-Dembo-Karlin
 statistics.

$$E(Q) = \text{Number of matches with scores above } Q$$

$$= K \cdot m \cdot n \cdot l^{-\lambda}$$

$m, n \Rightarrow$ length of two strings,
 $K \Rightarrow$ constant
 $\lambda \Rightarrow$ parameter

λ is a free root of the eqn.:

$$\sum_{x, y \in A} p_x p_y e^{\lambda \delta(x, y)} = 1$$

$p_x, p_y \Rightarrow$ frequencies of x, y .

The probability that there is a match of score greater than 0 betw. two "random" sequences of length n & m is computed as

$$1 - e^{-E(\theta)} \quad (\text{as it is } 1 - e^{-E(\theta)})$$

* Alignment Scores and Statistical Significance of Database Searches:

E -score(S): expected no. of sequences of score $\geq S$, would have been found randomly

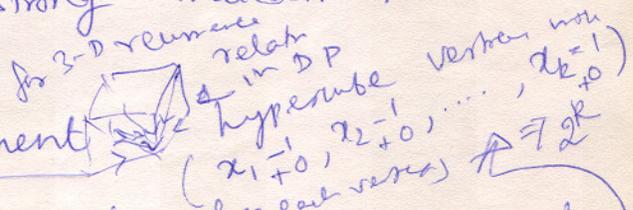
P -score(S): prob. that one or more sequences of score $\geq S$ would have been found randomly.

Low values of E & P indicate the search result was unlikely to have been obtained by random chance & thus is likely to bear evolutionary relationship to the query sequence.

e.g. $E < 10^{-3} \Rightarrow$ indicative of stat. signif.

$E \sim 10^{-50} \Rightarrow$ strong indication.

* Multiple Sequence Alignment



\rightarrow Extend pairwise sequence alignment for R multiple sequences with dynamic programming in a k -dimensional space $\Rightarrow O((2^n)^k)$

\rightarrow Heuristic method: (i) $\binom{R}{2}$ alignment & combine them
 or (ii) get the best alignment pair & merge them into a new string (once a gap is always a gap).
 Progressively, merge the best possible alignment with the combined string.
 e.g. CLUSTAL (Higgins, Thompson & Gibson (1996))

• CLUSTAL (Higgins & Sharp, 1988) ;

→ Constructs a phylogenetic tree to determine the degrees of similarity among the sequences being aligned.

→ Using this tree as a guide, closely related sequences are aligned two at a time using DP of the pairwise alignments.

→ Sensitive to initial choice of pairs.

→ CLUSTALW (1996) uses sequence weighting, position specific gap penalties & weight matrix choice depend on how the sequences are divergent.