

Lecture 2 Sequencing

Date

No. L

A sequence of DNA is denoted by $x_1 x_2 \dots x_k$ where $x_i \in \{A, T, G, C\}$. If $S = x_1 x_2 \dots x_k$, we denote $x_i = S(i)$.

The first problem we study about DNA sequences is the comparison of two sequences and determine how "similar" they are.

In general, the sequences we compare may not of the same length, i.e., two sequences may have different number of terms, but not too "different" from each other. So, we may use the so-called alignment to adjust their length (the number of terms).

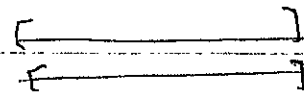
e.g.

AGACCTAG

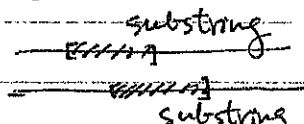
GCACCTGCAG

Types

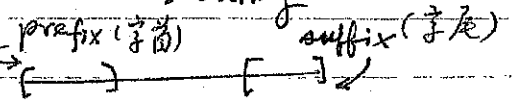
① Global alignment



② Local alignment



③ Semi-global alignment



(只存在字首或字尾作调整长度的动作。)

(*) Definition (Alignment)

The alignment is "the comparison of two or more nucleotide or protein sequences (strings) to determine the degree of similarity. $x_i \in \{\text{a set of twenty letters}\}$

(*) Commonly used to deduct functional or evolutionary relationships between genes and proteins.

(*) Global alignments attempt to align every base (or amino acid) in each aligned sequence (string).

(*) Local alignments will align only similar regions between sequences (strings), and leave regions with too many differences unaligned. (相差太多的部分没有调整的必要!)

(*) Multiple Sequence alignment refers to the process of aligning three or more nucleotide or protein sequences to identify similarities between the sequences. (among)

(*) Many algorithms are known so far!

More details

3

Definition (Alignment)

An alignment of two sequences S_1 and S_2 is a pair of sequences (S'_1, S'_2) obtained by insertion of spaces in S_1 and S_2 (respectively) such that

(1) $|S'_1| = |S'_2|$, and

(2) $\forall i, S'_1(i)$ is aligned with $S'_2(i)$ and either

$S'_1(i)$ or $S'_2(i)$ is not a space. (不可以同时为 space!)

(Note) S_1 can be viewed as a subsequence of S'_1 if (S_2) (S'_2) we consider the insertions as elements.

Example

$$S_1 = AGAC \leftarrow \langle A, G, A, C \rangle$$

$$S_2 = TACCC \leftarrow \langle T, A, G, C, C \rangle$$

$$\begin{array}{l|l} S'_1 = *AGAC & \text{or} \\ S'_2 = TACCC & \end{array} \quad \begin{array}{l} S'_1 = *AGAC* \\ \downarrow \quad \downarrow \downarrow \\ S'_2 = TACCC* \end{array}$$

Score of Alignment

$$\begin{array}{c} x_1 x_2 \dots x_k \\ y_1 y_2 \dots y_k \end{array}$$

(After alignment)

$$\text{For } i, j, k \left\{ \begin{array}{l} x_i = y_j \rightarrow p > 0 \\ x_i \neq y_j \rightarrow q < 0 \\ x_i \sim * \rightarrow r < 0 \\ * \sim y_j \rightarrow r < 0 \end{array} \right\} \text{ 负值}$$

two sequences

4

Similarity of S_1 and S_2

$$|S_1| = |S_2|$$

$$\text{Sim}(S_1, S_2) = \max_{(S'_1, S'_2)} \left\{ \sigma(S'_1, S'_2) = \sum_{i=1}^k \sigma(S'_1(i), S'_2(i)) \right\}$$

all alignments

$$\sigma(x_i, y_i) = \begin{cases} p > 0 & \text{if } x_i = y_i \\ q < 0 & \text{if } x_i \neq y_i \text{ and} \\ r < 0 & \text{if one of } x_i, y_i \\ & \text{is a space *} \end{cases}$$

Definition (Optimal Alignments)

The pair (S'_1, S'_2) with largest score is called an optimal alignment.

For example: (For AGC and AAAC) Let $p=1, q=-1$ and $r=-2$.

* AGC	A*GC	AG*C
AAAC	AAAC	AAAC
(-1)	(-1)	(-1)

are optimal alignments with $\text{Sim}(AGC, AAAC) = -1$.

(*) 在 "Alignment", 原来序列的前后顺序不作改变!

Algorithm

目標: Find $\text{Sim}(S_1, S_2)$ for input two sequences S_1 and S_2 .

1. Consider all possible alignments (S'_1, S'_2) .
2. Find $\sigma(S'_1, S'_2)$.

e.g. One insertion of AGC, \rightarrow *AGC, A*GC, AG*C, AGC*

$$S_2 = S'_2 \Rightarrow \sigma(S'_{1,1}, S'_2) = -1, \sigma(S'_{1,2}, S'_2) = -1$$

$$\sigma(S'_{1,3}, S'_2) = -1, \sigma(S'_{1,4}, S'_2) = -3$$

Theorem The time to accomplish "finding optimal alignment" is $O(|S_1| |S_2|)$ and the space we need is also $O(|S_1| |S_2|)$.

Proof. Since constant work is required per entry in the matrix, the proof follows. ■

Note No algorithm is known that uses asymptotically less time and has the same generality. Although, there are algorithms for more specific problems.

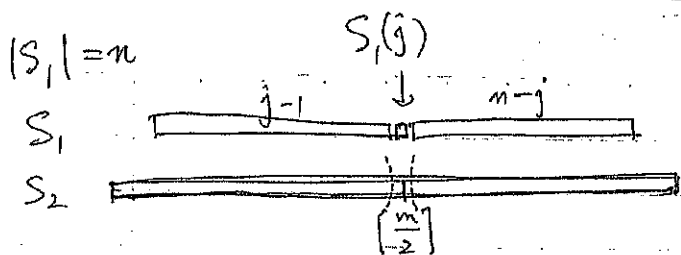
Note 儲存兩序列的 sequences 並不難; 但是, 如果每個 sequence 都有 10,000 字母, 則寫出矩陣就要填入 $(10,000)^2 = 100,000,000$ 壹佰萬個數, 是龐大的空間。

(*) We are aiming to find the "optimal alignments" instead of "presenting the matrix". (不斷地調整長度 & 比較當然很花時間!)

Bonnie Berger's Idea

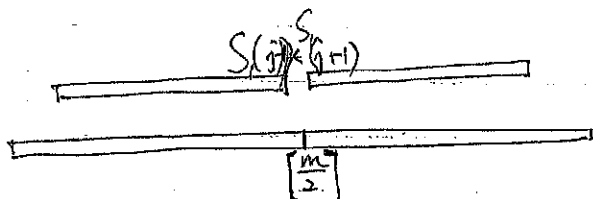
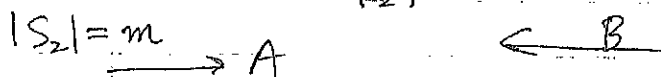
(*) $A(i, j)$: Score of optimal alignment of $S_1[1, i]$ and $S_2[1, j]$.

6



$$\text{Sim}(S_1, S_2)$$

$$= A(|S_1|, |S_2|)$$



固定 $\lfloor \frac{m}{2} \rfloor$, 变动 "j" (in S_1).

$\rightarrow A$ $\leftarrow B$ (把两 sequences 都倒过来看)

Step 1 选择中间的字母 (从 S_2 中).

(这个字母对应的位置不外乎是 $S_1(j-1)$ 或 $S_1(j)$ 或 $S_1(j+1)$)

中间的空格) (For some j)

Step 2 计算 ($\forall j$)

$$\max \left\{ \begin{aligned} & (A(\lfloor \frac{m}{2} \rfloor - 1, j-1) \pm 1) + B(\lfloor \frac{m}{2} \rfloor, n-j); \text{ and} \\ & (A(\lfloor \frac{m}{2} \rfloor - 1, j) - 2) + B(\lfloor \frac{m}{2} \rfloor, n-j). \end{aligned} \right.$$

Step 3

The maximum tell you the best score for aligning the $\frac{m}{2}$ (th) character of S_2 with some character of S_1 or gap. ($2n$ values!)

Time we need

$T(m, n) = 2cmn$

$T(m, 1) = m$

$T(1, n) = n$

$T(m, n) \leq cmn + \max_{j=1}^n (2c(\lceil \frac{m}{2} \rceil - 1)j + 2c\lfloor \frac{m}{2} \rfloor(n-j))$

By induction

$\leq cmn + \max_{j=1}^n (cmj + cm(n-j))$

$\leq cmn + cmn$

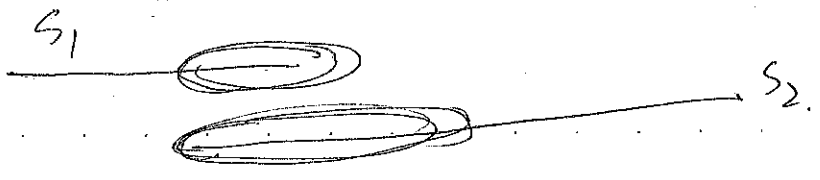
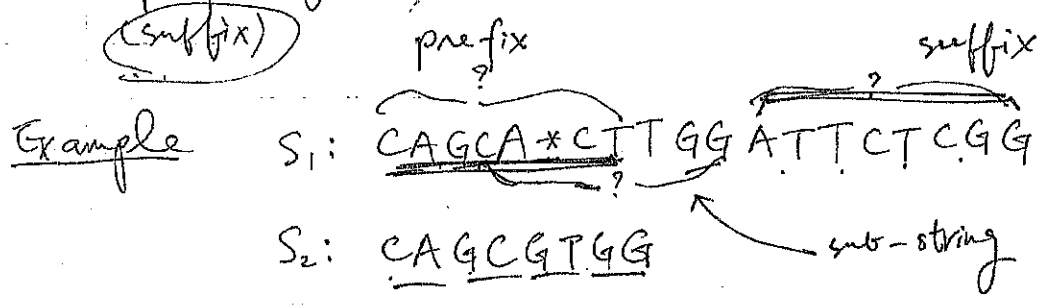
$\leq 2cmn$

Global Alignment

Semi-global Comparison :

To find the optimal alignment between suffix of S_1

with prefix of S_2



Multiple Sequences Alignment (MSA)

Definition (Sum-of-pairs, SP-score)

Let α be a collection of alignments of S_1, S_2, \dots, S_n :
 S'_1, S'_2, \dots, S'_n . Then, the SP-score of α , denoted by
 $SP(\alpha) = \sum_{1 \leq i < j \leq n} \alpha(S'_i, S'_j)$, $\alpha(S'_i, S'_j)$ is called the score
of α_{ij} , denoted by $s(\alpha_{ij})$

(*) Since MSAs are applied in finding the similarity of protein sequences, we use characters of proteins for examples.

S_1 : RCTLEE

S'_1 : RCTLEE

S_2 : RCLEE

$\alpha \Rightarrow$

S'_2 : RC*LEE

S_3 : CTLEE

S'_3 : *CTLEE

S_4 : CTEE

S'_4 : *CT*EE

(With $p = 1$, $q = -1$ and $r = -2$)

* \leftrightarrow *
 A LOSE 0

$$SP(\alpha) = A(\alpha_{1,2}) + A(\alpha_{1,3}) + A(\alpha_{1,4}) + A(\alpha_{2,3}) + A(\alpha_{2,4}) + A(\alpha_{3,4})$$

$$= 3 + 3 + 0 + 0 + (-3) + 2 = 5.$$

蛋白质的代表符号 (20种)

		缩写
1.	A : Alanine	Ala
2.	C : Cysteine	Cys
3.	D : Aspartic Acid	Asp
4.	E : Glutamic Acid	Glu
5.	F : Phenylalanine	phe
6.	G : Glycine	Gly
7.	H : Histidine	His
8.	I : Isoleucine	Ile
9.	K : Lysine	Lys
10.	L : Leucine	Leu
11.	M : Methionine	Met
12.	N : Asparagine	Asn
13.	P : Proline	Pro
14.	Q : Glutamine	Gln
15.	R : Arginine	Arg
16.	S : Serine	Ser
17.	T : Threonine	Thr
18.	V : Valine	Val
19.	W : Tryptophan	Trp
20.	Y : Tyrosine	Tyr

MSA Problem

Given k sequences, find the optimal alignment of these k sequences.

Dynamic Programming

Time-complexity = $O(k^2 n^k)$

(Exponential in k !)

觀察: 每一次調整都要和剩下的 sequences 對一次!

研究進展

① MSA problem is NP-complete!

Wang and Jiang (1994) + Bonizzoni and Vedova (2001)

2. Gusfield (1993): $(2 - \frac{2}{k})$ -approximation algorithm

(If $k=2$, then the algorithm is working.)

3. Perzner (1992): $(2 - \frac{3}{k})$ -approx. algorithm. (Works for $k=3$!)

(*) 4. Bafna, Lawler and Perzner (1997): Approximation algorithms for MSA, Theoretical Computer Science, 182, 233-244.

$(2 - \frac{l}{k})$ -approx. algorithm for $l < k$.

[Sequence Alignment Terms Explained](#)

[Sequence Alignment Software](#)

[Sequence Alignment Web Resources](#)

[Glossary](#)

[About](#)

[Contact](#)

Sequence Alignment Terms

Alignment

The comparison of two or more nucleotide or protein sequences to determine the degree of similarity. Commonly used to deduct functional or evolutionary relationships between genes and proteins.

Assembly

The process of combining short DNA sequence fragments into larger units by looking for overlaps between different fragments. Often required because the length of the genes studied exceeds the length of the sequence fragments produced by DNA sequencing machines. Also used to combine several fragments that cover the same region, for example in forward and reverse direction, with the goal to reduce errors in the consensus sequence.

Consensus Sequence

A single sequence generated from an alignment or assembly of sequence fragments that is the "best fit" for the given sequences. Historically, majority ("vote based") and inclusive methods were most commonly used to determine consensus sequence. For sequence assemblies, these methods have often been replaced by quality-based consensus methods. Quality-based consensus sequences are typically more accurate than majority-based sequences, and can reduce the need for manual editing of sequence assemblies drastically.

Contig

The result of a sequence assembly or alignment that shows the arrangement of the fragments to form a contiguous large sequence.

Dynamic Programming

A computer-science based method to find the optimal alignment between sequences. For two sequences, this algorithm creates a two-dimensional matrix based on identity or similarity of bases (or amino acids) in both sequences, and then finds the highest-scoring path to obtain the alignment. A commonly used dynamic programming method is the Needleman-Wunsch algorithm. A nice graphical display of the dynamic programming methods for sequence alignments can be found [here](#).

Global Alignments

Global alignments attempt to align every base (or amino acid) in each aligned sequence.

Local Alignments

Local alignments will align only similar regions between sequences, and leave regions with too many differences unaligned. Local alignments can be better suited for the alignment of very dissimilar sequences. In sequence assembly, the program Phrap demonstrate that local alignments can be used to reduce or eliminate the need to remove low-quality sequence (end clipping) before assembly.

Multiple Sequence Alignment

Multiple sequence alignment refers to the process of aligning three or more nucleotide or protein sequences to identify similarities between the sequences. Alignments that include many sequences can be computational intensive, and require more sophisticated algorithms than pairwise alignments.

Pairwise Alignment

In pairwise sequence alignment, exactly two nucleotide or protein sequences are aligned to each other to determine the similarity between the two sequences.

Word-based Alignment Methods

Word-based alignment methods are an optimization often used in sequence alignment and assemblies. Instead of examining every single nucleotide or amino acid, "words" of a fixed length are analyzed. This can lead to substantial reductions in memory use and alignment times. One common application is to use the number of shared words between two sequences to estimate the similarity in early phases of sequence alignments, or to identify sequences that share overlaps in sequence assembly.