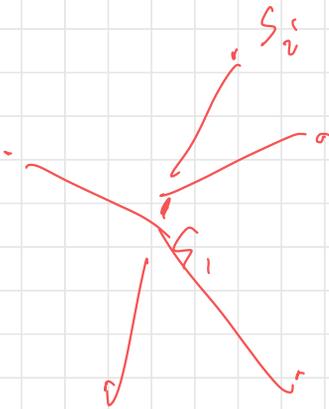


Review:

- "Multiple alignment shout loudly."
- Merge 2 alignments into 1.
- pairwise alignment induced by the multiple alignment.



# Exact Algorithm for Multiple Alignment

---

- When the optimal alignment is needed. There is an exact algorithm as well.

# Score Function

---

- For each column, assuming the letters are  $a_1, \dots, a_n$  for the  $n$  sequences. Define a column score  $S(a_1, \dots, a_n)$ .
- The alignment score is the total of the column score.
- There are more than one ways to define the column score. Let us examine the simplest one first.

# SP-score

---

- SP (Sum of Pair) score is the most widely used because its simplicity.
- We have a similarity matrix  $s(a,b)$  for any two given letters.
- For the  $k$ -th column with  $n$  letters  $x_1, \dots, x_n$ . Define

$$S_k = \sum_{i,j} s(x_i, x_j)$$

- The SP-score of the alignment is

$$\sum_{1 \leq k \leq m} S_k$$

		↓ Motif 1	↓	Motif 2
H.s. Wee1 409-457	i →	QVGRGLRYI	SMS-LV	MDIKPSNIFISRTSIPNAASEEGEDDWASNK----
H.s. Tik 614-659		NMLEAVHTI	QHG-IV	SDLKPAFLIVDG-----MLKLIQFQIANQMCPD--
S.c. Ste7 313-358		GVLNGLDHLRQYK	RII	HRDIKPSVVLINSK----GQIKLCPQFVSKKLI----
S.c. Mkk1 332-376		AVLRGLSYL	EKK-VI	HRDIKPSILLNEN----GQVKLCPQFVSGEAV----
S.p. Byr1 168-213	j →	SMVVKGLIYLYNV	HI	HRDLKPSVVVNSR----GEIKLCPQFVSGELV----
S.c. St20 722-767		ETLSGLEFL	SKG-VL	HRDISDILLSM-----GDIKLTPQFCAQINE----
S.c. Cc15 129-172		QTLLGLKYL	GEG-VI	HRDIKAANILLSAD----NTVKLADPQFVSTIV----
S.p. Byr2 505-553		QTLKGLEYL	SRG-IV	HRDIKGANILVDNK----GKIKISDFQISKKLELNST
S.c. Spk1 302-348		QILTAIKYI	SMG-IS	HRDLKPDILLIEQDD--PVLVKITDFQFLAKVQG----
S.p. Kin1 249-293		QIGSALSYL	QNS-VV	HRDLKIEILISKT----GDIKIIPQFCLSNLYR----
S.p. Cdr1 111-156		QILDVAVHC	RFR-FR	HRDLKLEILIKVN---EQQIKIADPQFMATVEP----
M.m. K6a1 507-556		TISKTV EYL	SGQ-VV	HRDLKPSHILYVDESGNPECLRICPQFQFAKQLRA----
R.n. Kpbh 136-180		SLLEAVNFL	VNN-IV	HRDLKPEILLDDN----MQIRLSDFQFSCHE----
H.s. Erk2 132-176		QILRGLKYI	SAM-VL	HRDLKPSLLLNTT---CLSKICPQFGLARVA----
S.c. Kss1 137-182		QILRALKSI	SAQ-VI	HRDIKPSMLLNSN-----CKVCPQFLARCLASS--

© 1999-2004 New Science Press

$$s(-, -) = 0$$

$$s(S_i, S_j) = \sum_k s(S_i[k], S_j[k])$$

$$\sum_{\substack{i, j \\ i < j}} s(S_i, S_j)$$

# SP-score

---

- Note that SP-score can be viewed from a different perspective as follows.
- For each pair of sequences  $S_i$  and  $S_j$  (including the gaps) in the multiple sequence alignment, define
- $S_{i,j} = \sum_k s(S_i[k], S_j[k])$
- The SP score is then  $\sum_{i,j} S_{i,j}$
- The two perspectives are equivalent if  $s(-,-)=0$ .

# Review

- Before we study the algorithm to do multiple sequence alignment, we review the pairwise alignment algorithm.

$$DP[i,j] = S_1[1..i] \text{ v.s. } S_2[1..j]$$

Case 1:

S1[1..i]	-
S2[1..j-1]	s2[j]

Case 2:

S1[1..i-1]	s1[i]
S2[1..j]	-

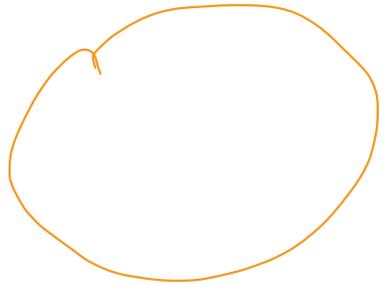
Case 3:

S1[1..i-1]	s1[i]
S2[1..j]	s2[j]

$$DP[i,j] = \max \begin{cases} DP[i-1, j-1] + f(s[i], t[j]); \\ DP[i-1, j] + f(s[i], -); \\ DP[i, j-1] + f(-, t[j]); \end{cases}$$

# Alignment of three sequences

- Use X to indicate a letter. The last column has 7 cases:



f one of

X	-	X	-	X	-	X
X	X	-	-	X	X	-
X	X	X	X	-	-	-

$DP[i_1, i_2, i_3]$

- Let  $f$  be the column-wise score function.
- $DP[i_1, i_2, i_3]$  is the optimal score for  $s_1[1..i_1]$ ,  $s_2[1..i_2]$ ,  $s_3[1..i_3]$ .
- $DP[i_1, i_2, i_3] = \text{maximum of the 7 cases}$

case XXX:  $DP[i_1, i_2, i_3] = DP[i_1-1, i_2-1, i_3-1] + S(s_1[i_1], s_2[i_2], s_3[i_3])$

case XX-:  $DP[i_1, i_2, i_3] = DP[i_1-1, i_2-1, i_3] + S(s_1[i_1], s_2[i_2], -)$

.....

- It is necessary to introduce some notation for n sequences ...

# Alignment of many sequences

---

X	-	X	-	X	-	X
X	X	-	-	X	X	-
X	X	X	X	-	-	-
1	0	0	1	1	0	1

- Here X indicates a letter and a dash indicates an indel.
- If you regard X as 1 and - as 0. Then these cases are all the 3-bit binary numbers but 000.
- In general, for n sequences, there are  $2^n - 1$  cases.
- Let us use  $\delta_j$  to indicate whether at sequence j the last column is a letter.  $\delta_j = 1$  if it is a letter.

# Complexity

---

- Time complexity  $O(m^n 2^n T)$ , where  $T$  is the time needed for computing a column score.
- Space complexity  $O(m^n)$ .
- This problem is NP-hard. So, no polynomial time algorithm exists (unless  $P=NP$ ).
- Computing the optimal alignment with large  $m$  and  $n$  is hopeless. But at least this works for small  $m$  and  $n$ .

# Better Scoring Function

---

- SP score (either maximization or minimization) is one of the simplest scoring functions in use.
- There are better proposals to the scoring function.
- Let's examine a more sophisticated score called relative entropy. We want to maximize.

## Score 2: Relative Entropy

---

- There are  $n$  letters  $x_1, \dots, x_n$  in a column. For letter  $a$  in alphabet, let  $n_a$  be the number of  $a$  in the  $n$  letters.
- Our two different models assume two different distributions of letters in that column.
- Model 1:  $a$  occurs with probability  $p_a = n_a / n$ .
- Model 2:  $a$  occurs with “background” probability  $q_a$ , which is learned from a database. *random*
- We want to compute the likelihood ratio.

# Relative Entropy

---

- Log likelihood ratio:

$$\Pr(x_1, \dots, x_n \mid \text{model1}) = \prod_{1 \leq i \leq n} p_{x_i}$$

$$\Pr(x_1, \dots, x_n \mid \text{model2}) = \prod_{1 \leq i \leq n} q_{x_i}$$

$$\log \prod_{1 \leq i \leq n} \frac{p_{x_i}}{q_{x_i}} = \sum_{1 \leq i \leq n} \log \frac{p_{x_i}}{q_{x_i}} = \sum_{a \in \Sigma} n_a \log \frac{p_a}{q_a}$$

- This is called the relative entropy at a column.
- Let  $E_j$  be the relative entropy of column  $j$ .
- The alignment score is equal to  $\sum_{1 \leq j \leq m} E_j$
- We want to maximize.

# Relative Entropy

---

- A few assumptions of relative entropy score:
  - Sequences are independent to each other.
  - Columns are independent to each other.
- These seem to be too strong but at least when these assumptions are true the score has a theory foundation. Other scores are more empirical.
- Note: If all  $q_a$  are equal to each other, then this is equivalent to the “minimum entropy” introduced in the reference book (Durbin page 138).

# Approximation Algorithm

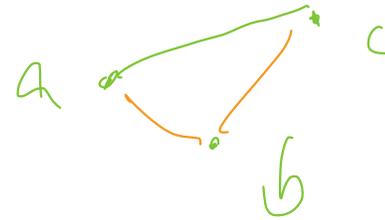
---

- In the heuristic algorithm, it would be good if we can get some guarantee on the quality of the suboptimal result.
- Suppose we want to maximize a score, and the optimal value is  $OPT$ , how about computing a solution with score at least  $OPT/c$  for a small constant  $c > 1$ ?
  - If  $c = 1.001$  this is not so bad.
- If this can be done, this is called a ratio- $c$  approximation algorithm for the problem.
- For minimization problem, ratio  $c$  means the algorithm gives score at most  $c * OPT$ .

## The Minimization Version of Multiple Alignment

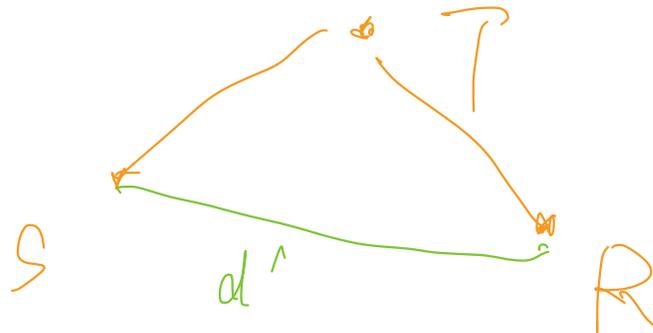
---

- The known approximation algorithm for multiple alignment only works for the minimization version of the problem.
- Let  $M[a,b]$  be a distance metric between two letters (possibly –)  $a$  and  $b$  satisfying **triangular inequality**.
  - $M[a,a] = 0$
  - $M[a,b] \geq 0$
  - $M[a,c] \leq M[a,b] + M[b,c]$
- The distance between two sequences  $d'(S_i, S_j)$  is the total of their column scores using  $M$  as the score scheme.
- The SP score is  $\sum_{i,j} d'(S_i, S_j)$ .
- We want to **minimize**.
- Note that this is almost the same as before, except that now we want to minimize the difference.



# Triangular Inequality

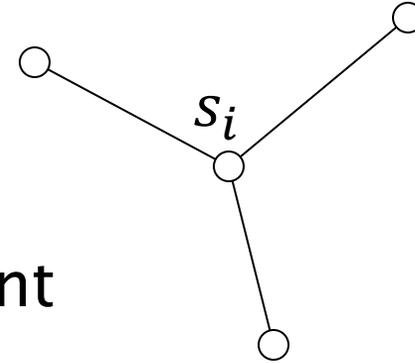
- For sequences  $S, T, R$  taken from the same multiple alignment
  - $d'(S, R) \leq d'(S, T) + d'(T, R)$
- Note that  $S, T, R$  may have the ‘-’ symbols in the sequence.
- Notations:
  - When context is clear (the multiple alignment is given), we also use  $\rightarrow d'(s, t)$  to denote the distance between the original sequences  $s$  and  $t$  (without ‘-’). It is equal to  $d'(S, T)$ .
  - When there are more than one multiple alignments  $A_1, A_2, \dots, A_n$ , we use  $d_i'(s, t)$  to indicate the distance induced by  $A_i$ .
  - We use  $d(s, t)$  to denote the optimal pairwise alignment between  $s$  and  $t$ . Note that  $d(s, t) \leq d'(s, t)$



# Approximation Algorithm

---

- Consider the following algorithm:
- Input:  $s_1, s_2, \dots, s_n$
- 1. Select an  $s_i$ . Build optimal pairwise alignment between  $s_i$  and  $s_j$  for each  $j=1,2,\dots,n$ .
- 2. Now use these pairwise alignment to build a multiple alignment  $A_i$ .
- 3. Repeat 1 and 2 for every  $i$ , and output the  $A_i$  with the minimum SP score.



We want to prove that this algorithm has ratio 2.

# Approximation Ratio

---

- Fact 1:  $d'_i(s_i, s_j) = d(s_i, s_j)$  for every  $j$ .

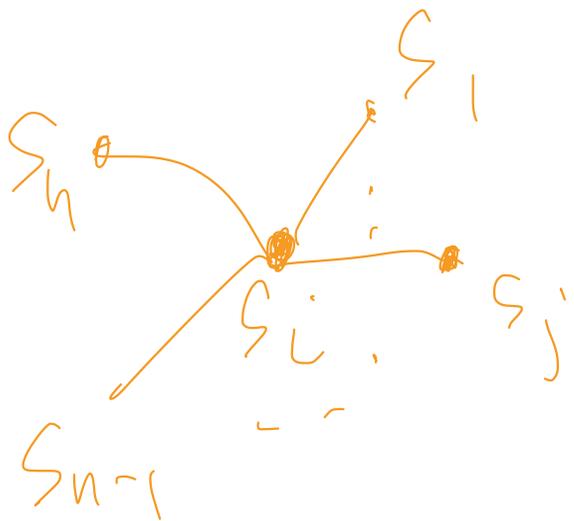


$A_i$

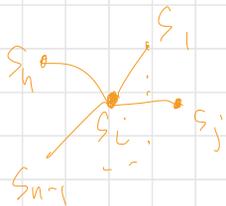
$$d'_i(s_k, s_e) \geq d(s_k, s_e)$$



$$\leq d'_i(s_i, s_k) + d'_i(s_i, s_e)$$



$A_i$



$$d'_i(s_k, s_l) \geq d(s_k, s_l)$$

↓

$$\leq d'_i(s_i, s_k) + d'_i(s_i, s_l)$$

$$\sum_{i=1}^n SP(A_i) = \sum_i \sum_{k,l} d'_i(s_k, s_l)$$

$A^*$  is optimal  
 $d^*$  is induced by  $A^*$

$$\leq \sum_i \sum_{k,l} (d'_i(s_i, s_k) + d'_i(s_i, s_l))$$

$$= \sum_i \sum_{k,l} d'_i(s_i, s_k) + \sum_{k,l} d'_i(s_i, s_l)$$

$$= \sum_i n \cdot \sum_k d'_i(s_i, s_k) + n \cdot \sum_l d'_i(s_i, s_l)$$

$$= \sum_i 2n \cdot \sum_k d'_i(s_i, s_k)$$

$$= \sum_i 2n \cdot \sum_k d(s_i, s_k)$$

$$\leq \sum_i 2n \cdot \sum_k d^*(s_i, s_k)$$

$$= 2n \cdot \left( \sum_i \sum_k d^*(s_i, s_k) \right)$$

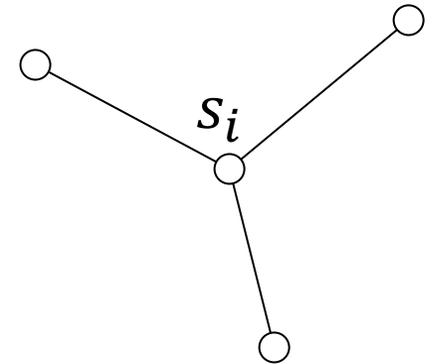
$$= 2n \cdot SP(A^*) = 2n \cdot OPT.$$

# Approximation Ratio

- Let  $A^*$  be the optimal multiple alignment and  $d^*$  be pairwise alignment score induced by optimal multiple alignment

$$\begin{aligned}
 \sum_{1 \leq i \leq n} \text{score}(A_i) &= \sum_{1 \leq i \leq n} \sum_{1 \leq k \leq l \leq n} d'_i(s_k, s_l) \leq \sum_{1 \leq i \leq n} \sum_{1 \leq k \leq l \leq n} (d'_i(s_i, s_k) + d'_i(s_i, s_l)) \\
 &\leq 2n \sum_{1 \leq i \leq n} \sum_{1 \leq k \leq l \leq n} d'_i(s_i, s_k) \\
 &= 2n \sum_{1 \leq i \leq n} \sum_{1 \leq k \leq l \leq n} d(s_i, s_k) \\
 &\leq 2n \sum_{1 \leq i \leq n} \sum_{1 \leq k \leq l \leq n} d^*(s_i, s_k) \\
 &= 2n \times \text{OPT}
 \end{aligned}$$

- Recall that
  - $d'_i$ : pairwise alignment score induced by  $A_i$
  - $d$ : optimal pairwise alignment score
- So, at least one of the  $A_j$  has score no more than  $2\text{OPT}$ .



# Approximation Algorithm

---

- Theorem: The above algorithm is a ratio-2 approximation to multiple alignment under SP score.
- Note: For general score scheme (general score matrix) this ratio-2 may not hold.
  - For example, when the score scheme does not satisfy triangular inequality. Or when the score scheme has both positive and negative scores.

# Applications

---

- Motif finding: conserved regions of a protein family may be important motifs.
- Structure prediction: after multiple alignment of a protein family, their structures can be predicted together. *Alpha Fold 2.*
- Phylogeny: the pairwise alignment induced by the multiple alignment can be more accurate than the optimal pairwise alignment in reality. Build multiple alignment first then do phylogeny by examining the evolution on each column.

*clustal w*  
*Ω.*

# What we have learned?

---

- 1. Multiple sequence alignments “shout out”.
- 2. extend the DP algorithm to multiple alignment
- 3. first NP–hard problem in this course
- 4. two ways to handle NP–hard problems: heuristics & approximation
- 5. apply our likelihood ratio score to multiple alignment