

Multiple alignments

Plan

- (1) Today: What is a multiple alignment of a set of sequences, how to compute an optimal multiple alignment using SP-score (and similar column-based score function) by dynamic programming.
- (2) Today and next week: How to approximate an optimal SP-score multiple alignment efficiently
- (3) In a couple of weeks: Speeding up the implementation of MSA using forward dynamic programming.

Multiple Alignment

A generalization of pairwise alignment, comparison of multiple sequence which makes it possible to identify weaker similarities

(A) p110 α cAMP-kinase

```
TFILGIGDRHNSNIMVKDDG-QLFHI DFGHFLDHKKKKFGYKRERVPFVLT--QDFLIVI 142
QIVLTFEYLSLDLIYRD LKPENLLIDQQGYIQVT DFGFAKRVKGRTWXLCGTPEYLAPE 179
```

(B) p110 β 136
p110 δ 136
p110 α 135
p110 γ 135
p110_dicti 135
cAMP-kinase 177

```
SYVLGIG-----DRHSDNINVKKTGQLFHI DFGHILGNFKSKFGIKRERVPFILT 136
TYVLGIG-----DRHSDNIMIRESGQLFHI DFGHFLGNFKTKFGINRERVPFILT 136
TFILGIG-----DRHNSNIMVKDDGQLFHI DFGHFLDHKKKKFGYKRERVPFVLT 135
TFVLGIG-----DRHNDNIMITETGNLFHI DFGHILGNYKSFLGINKERVPFVLT 135
TYVLGIG-----DRHNDNLMVTKGGRLFHI DFGHFLGNYSKKFGFKRERAPFVFT 135
QIVLTFEYLSLDLIYRD LKPENLLIDQQGYIQVT DFGFAKRVKGRTWXLCG--TPEYLA 177
```

Multiple Alignment

A generalization of pairwise alignment, comparison of multiple sequence which makes it possible to identify weaker similarities

A multiple alignment of $S_1, \dots, S_k \in \Sigma^*$ is a

$k \times l$ -matrix where

- $\mathcal{H}_{ij} \in \Sigma \cup \{-\}$

$$\mathcal{H} = \begin{bmatrix} - & R_1 & - \\ & \vdots & \\ - & R_k & - \end{bmatrix}$$

- row i , R_i , with "-"s removed, spells S_i

The starting point of many bioinformatics workflows

Multiple Alignment

A generalization of pairwise alignment, comparison of multiple sequence which makes it possible to identify weaker similarities

(A)	p1	Establish evolutionary relationship	I 142
	cA		E 179
(B)	p1	Inference of structural and functional properties:	
	p1	Sequence homology implies structural similarity,	T 136
	p1	and conserved residues implies conserved role	T 136
	p1		T 135
	p1		T 135
	p1		T 135
	cA	Identify and define family patterns and motifs	A 177
		Sequence assembly	

The starting point of many bioinformatics workflows

Multiple Alignment

A generalization of pairwise alignment, comparison of multiple sequence which makes it possible to identify weaker similarities

(A) p1
cA

Establish evolutionary relationship

I 142
E 179

To computationally find an optimal multiple alignment, we must:

Define the cost of an alignment

Define an optimal alignment

Construct an efficient algorithm

Sequence assembly

The starti

$$\mathcal{M}^* = \underset{\mathcal{M}}{\operatorname{argmax}} \operatorname{SCORE}(\mathcal{M})$$

workflows

The cost of an multiple alignment

Idea: Define cost such that ideas from pairwise alignment can be reused

The cost of an multiple alignment

Idea: Define cost such that ideas from pairwise alignment can be reused

Induced pairwise alignments



Sum-of-pairs score

The SP score of an multiple alignment is the **sum of the cost of the induced pairwise alignments**.

$$\textcircled{1} \begin{bmatrix} A & A & G & A & A & - & A \\ A & T & - & A & A & T & G \end{bmatrix}$$

$$\textcircled{2} \begin{bmatrix} A & T & - & A & A & T & G \\ C & T & G & - & G & - & G \end{bmatrix}$$

$$\textcircled{3} \begin{bmatrix} A & A & G & A & A & - & A \\ C & T & G & - & G & \diagup & G \end{bmatrix}$$

Sum-of-pairs score

The SP-score $SP(\mathcal{M})$ is the sum of the scores of all induced pairwise alignments ...

$$\mathcal{M} = \begin{bmatrix} \text{---} R_1 \text{---} \\ \vdots \\ \text{---} R_k \text{---} \end{bmatrix} \quad SP(\mathcal{M}) = \sum_{1 \leq i < j \leq k} \text{PAIRSCORE}(R_i, R_j)$$

Is SP-score **column based**, i.e. is the cost of an multiple alignment the sum of the cost of each column?

If PAIRSCORE is column-based (e.g. global alignment with linear gap cost) then ...

$$SP(\mathcal{H}) = \sum_{1 \leq i, j \leq k} \text{PAIRSCORE}(R_i, R_j)$$

$$= \sum_{1 \leq i, j \leq k} \sum_{1 \leq c \leq L} d(R_i[c], R_j[c])$$

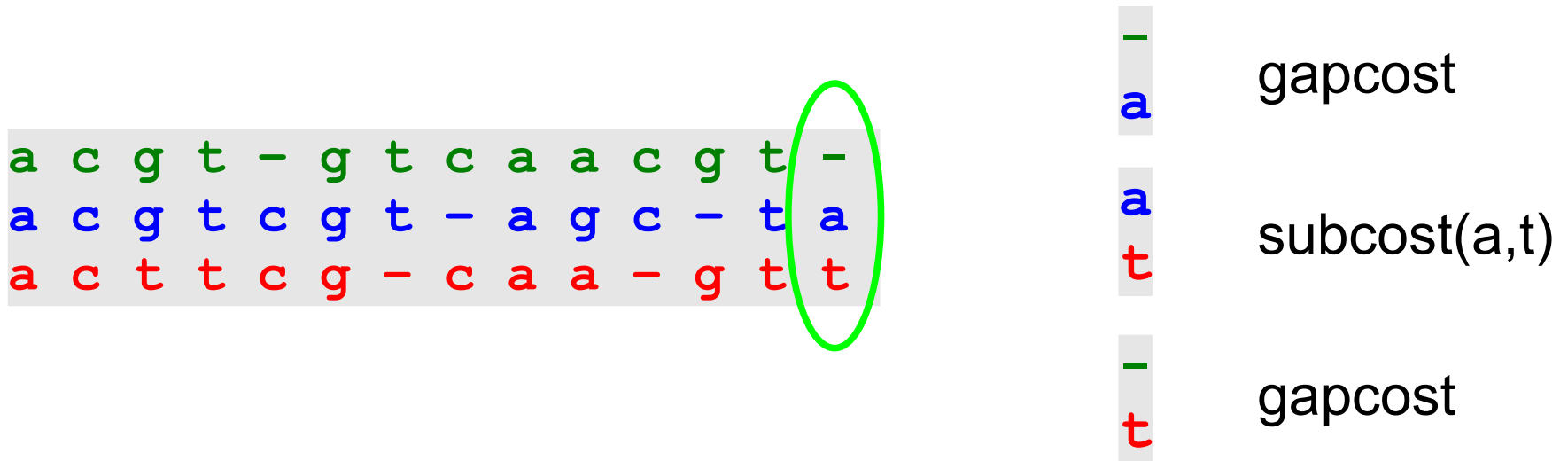
$$= \sum_{1 \leq c \leq L} \sum_{1 \leq i, j \leq k} d(\mathcal{H}_{ic}, \mathcal{H}_{jc})$$

$$= \sum_{1 \leq c \leq L} SP \left(\begin{bmatrix} \mathcal{H}_{1c} \\ \vdots \\ \mathcal{H}_{kc} \end{bmatrix} \right)$$

↑
column c in \mathcal{H}

... SP-score is column-based

Sum-of-pairs score



Cost of alignment = “sum of the cost of each column”

Problem: Given a set of sequences S_1, S_2, S_3, \dots , a score matrix and a gap cost, find a multiple alignment of optimal sum-of-pairs cost

Computing an optimal MSA using a column based score

- * Let $D(i_1, \dots, i_k)$ be the cost of an optimal mult. align of $S_1[1..i_1], \dots, S_k[1..i_k]$
- * Compute $D(i_1, \dots, i_k)$ by optimizing over all possible last columns of such an alignment.

$$\begin{bmatrix} S_1[i_1] \\ S_2[i_2] \\ \vdots \\ S_k[i_k] \end{bmatrix}, \begin{bmatrix} - \\ S_2[i_2] \\ \vdots \\ S_k[i_k] \end{bmatrix}, \dots, \begin{bmatrix} S_1[i_1] \\ S_2[i_2] \\ \vdots \\ S_k[i_k] \end{bmatrix}, \dots$$

What is number of possible last columns?

$$D(i_1, \dots, i_k) = \min_{\text{last.col}} [D(i_1', \dots, i_k') + \text{cost}(\text{last.col})]$$

Computing an optimal MSA using a column based score

* Let $D(i_1, \dots, i_K)$ be the cost of an optimal mult. align of $S_1[1..i_1], \dots, S_K[1..i_K]$

* Compute $D(i_1, \dots, i_K)$ by optimizing over all possible last columns of such an alignment.

$$\begin{bmatrix} S_1[i_1] \\ S_2[i_2] \\ \vdots \\ S_K[i_K] \end{bmatrix}, \begin{bmatrix} - \\ S_2[i_2] \\ \vdots \\ S_K[i_K] \end{bmatrix}, \dots, \begin{bmatrix} S_1[i_1] \\ S_2[i_2] \\ \vdots \\ S_K[i_K] \end{bmatrix}, \dots$$

* possible last.columns = $2^K - 1$

$$D(i_1, \dots, i_K) = \min_{\text{last.col}} [D(i_1', \dots, i_K') + \text{cost}(\text{last.col})]$$

Computing an optimal MSA using a column based score

* Let $D(i_1, \dots, i_K)$ be the cost of an optimal mult. align of $S_1[1..i_1], \dots, S_K[1..i_K]$

* Compute $D(i_1, \dots, i_K)$ by optimizing over all

Time:

Space:

$$\begin{bmatrix} S_1[i_1] \\ S_2[i_2] \\ \vdots \\ S_K[i_K] \end{bmatrix}, \begin{bmatrix} S_2[i_2] \\ \vdots \\ S_K[i_K] \end{bmatrix}, \dots, \begin{bmatrix} S_1[i_1] \\ S_2[i_2] \\ \vdots \\ S_K[i_K] \end{bmatrix}, \dots$$

* possible last.columns = $2^K - 1$

$$D(i_1, \dots, i_K) = \min_{\text{last.col}} [D(i_1', \dots, i_K') + \text{cost}(\text{last.col})]$$

Computing an optimal MSA using a column based score

* Let $D(i_1, \dots, i_k)$ be the cost of an optimal mult. align of $S_1[1..i_1], \dots, S_k[1..i_k]$

* Compute $D(i_1, \dots, i_k)$ by optimizing over all

Time: $O(n^k \cdot 2^k \cdot \text{"time to compute cost of last column"})$

Space: $O(n^k)$

$$\begin{bmatrix} S_1[i_1] \\ S_2[i_2] \\ \vdots \\ S_k[i_k] \end{bmatrix}, \begin{bmatrix} S_2[i_2] \\ \vdots \\ S_k[i_k] \end{bmatrix}, \dots, \begin{bmatrix} S_1[i_1] \\ S_2[i_2] \\ \vdots \\ S_k[i_k] \end{bmatrix}, \dots$$

* possible last.columns = $2^k - 1$

$$D(i_1, \dots, i_k) = \min_{\text{last.col}} [D(i_1^*, \dots, i_k^*) + \text{cost}(\text{last.col})]$$

Computing an optimal MSA using a column based score

* Let $D(i_1, \dots, i_k)$ be the cost of an optimal mult. align of $S_1[1..i_1], \dots, S_k[1..i_k]$

* Compute $D(i_1, \dots, i_k)$ by optimizing over all

Time: $O(n^k \cdot 2^k \cdot \text{"time to compute cost of last column"})$

Space: $O(n^k)$

$S_1[i_1]$	$S_2[i_2]$	$S_3[i_3]$	$S_4[i_4]$
$S_1[i_1]$	$S_2[i_2]$	$S_3[i_3]$	$S_4[i_4]$

Comment: An optimal alignment can be found by backtracking (in time?), and space can be reduced by a factor n using "Hirschberg's trick"

* possible last.columns = $k - 1$

$$D(i_1, \dots, i_k) = \min_{\text{last.col}} [D(i_1, \dots, i_k) + \text{cost}(\text{last.col})]$$

Computing an optimal MSA using a column based score

* Let $D(i_1, \dots, i_k)$ be the cost of an optimal mult. align of $S_1[1..i_1], \dots, S_k[1..i_k]$

* Compute $D(i_1, \dots, i_k)$ by optimizing over all

Time: $O(n^k \cdot 2^k \cdot \text{"time to compute cost of last column"})$

Space: $O(n^k)$

Is it useful?

Aligning 10 sequences of length 100 takes:

$100^{10} \text{ op} / 100.000.000 \text{ op/sec} = 32.000 \text{ years !!!}$

A huge need for fast and useful alternatives.

$$D(i_1, \dots, i_k) = \min_{\text{last.col}} [D(i_1, \dots, i_k) + \text{cost}(\text{last.col})]$$

Recursive formulation for 3 seqs

Input: Three sequences $A[1..n]$, $B[1..n']$, $C[1..n'']$

Recursion:

$$D(0,0,0) = 0$$

$$D(i,j,k) = \min \left\{ \begin{array}{ll} D(i-1, j-1, k-1) + SP(A[i], B[j], C[k]) & \text{if } i>0, j>0, k>0, \\ D(i-1, j-1, k) + SP(A[i], B[j], -) & \text{if } i>0, j>0, k\geq 0, \\ D(i-1, j, k-1) + SP(A[i], -, C[k]) & \text{if } i>0, j\geq 0, k>0, \\ D(i, j-1, k-1) + SP(-, B[j], C[k]) & \text{if } i\geq 0, j>0, k>0. \\ D(i-1, j, k) + SP(A[i], -, -) & \text{if } i>0, j\geq 0, k\geq 0, \\ D(i, j-1, k) + SP(-, B[j], -) & \text{if } i\geq 0, j>0, k\geq 0, \\ D(i, j, k-1) + SP(-, -, C[k]) & \text{if } i\geq 0, j\geq 0, k>0 \end{array} \right.$$

$D(n, n', n'')$ can be computed in time $O(n \cdot n' \cdot n'')$ using memorization

Iterative formulation for 3 seqs

```
T[0..n][0..n'][0..n''] = undef
for i = 0 to n:
  for j = 0 to n':
    for k = 0 to n'':
      v0 = v1 = v2 = v3 = v4 = v5 = v6 = v7 = undef
      if i=0 and j=0 and k=0 then
        v0 = 0
      if i>0 and j>0 and k>0 then
        v1 = T[i-1][j-1][k-1] + SP(A[i], B[j], C[k])
      if i>0 and j>0 and k≥0 then
        v2 = T[i-1][j-1][k] + SP(A[i], B[j], - )
      if i>0 and j≥0 and k>0 then
        v3 = T[i-1][j][k-1] + SP(A[i], -, C[k])
      if i≥0 and j>0 and k>0 then
        v4 = T[i][j-1][k-1] + SP(-, B[j], C[k])
      if i>0 and j≥0 and k≥0 then
        v5 = T[i-1][j][k] + SP(A[i], -, - )
      if i≥0 and j>0 and k≥0 then
        v6 = T[i][j-1][k] + SP(-, B[j], - )
      if i≥0 and j≥0 and k>0 then
        v7 = T[i][j][k-1] + SP(-, -, C[k])
      T[i][j][k] = min(v0, v1, v2, v3, v4, v5, v6, v7)
print T[n][n'][n'']
```

SP-cost of a column of 3 symbols

Let **sub**(x,y) be the substitution cost and **gap** be the gapcost

$$SP(A[i], B[j], C[k]) = \text{sub}(A[i], B[j]) + \text{sub}(B[j], C[k]) + \text{sub}(A[i], C[k])$$

$$SP(A[i], B[j], -) = \text{sub}(A[i], B[j]) + \text{gap} + \text{gap}$$

$$SP(A[i], -, C[k]) = \text{gap} + \text{sub}(A[i], C[k]) + \text{gap}$$

$$SP(-, B[j], C[k]) = \text{gap} + \text{gap} + \text{sub}(B[j], C[k])$$

$$SP(A[i], -, -) = \text{gap} + \text{gap}$$

$$SP(-, B[j], -) = \text{gap} + \text{gap}$$

$$SP(-, -, C[k]) = \text{gap} + \text{gap}$$

Other column based score functions?

Consensus score

Tree score

...

Consensus score

A consensus seq., or row, of \mathcal{M} is a seq. that summarize \mathcal{M} .

The consensus score $CS(\cdot)$ measures to what extend this is possible ...

$$\mathcal{M} = \begin{bmatrix} \text{---} R_1 \text{---} \\ \vdots \\ \text{---} R_K \text{---} \end{bmatrix}$$

--- R_{ref} ---

$$CS(\mathcal{M}) = \sum_{1 \leq i < j \leq K} \text{PairScore}(R_i, R_j)$$

CS alignment problem

Find \mathcal{M}^* such that $CS(\mathcal{M}^*)$ is minimized (or maximized)

.... complexity depends on definition of R_{ref} .

Consensus score

A typical definition of R_{sc} is:

$$R_{sc} = c_1 c_2 \dots c_n, \quad c_j \in \Sigma \cup \{-\}$$

where

$$c_j = \operatorname{argmin}_{c \in \Sigma \cup \{-\}} \sum_{i=1}^n d(c, \mathcal{K}_{ij})$$

Using this definition and assuming PairScore is col-based yields:

$$CS(\mathcal{K}) = \sum_{1 \leq i < k \leq n} \underbrace{\sum_{1 \leq j \leq n} d(c_j, \mathcal{K}_{ij})}_{\text{Pair Score}(R_i, R_k)}$$

$$= \sum_{1 \leq j \leq n} \underbrace{\sum_{1 \leq i < k \leq n} d(c_j, \mathcal{K}_{ij})}_{\text{Score of col. } j}$$

Consensus score

A typical definition of R_{dc} is:

$CS(\cdot)$ is turned into col.-based score function, i.e. opt. score $CS(c^*)$

can be computed as:

$$D(i_1, i_2, \dots, i_k) = \min_{c, \text{histed.}} \left[D(i_1^*, \dots, i_k^*) + \sum_{1 \leq i \leq k} d(c, \text{lastcol}[i]) \right]$$

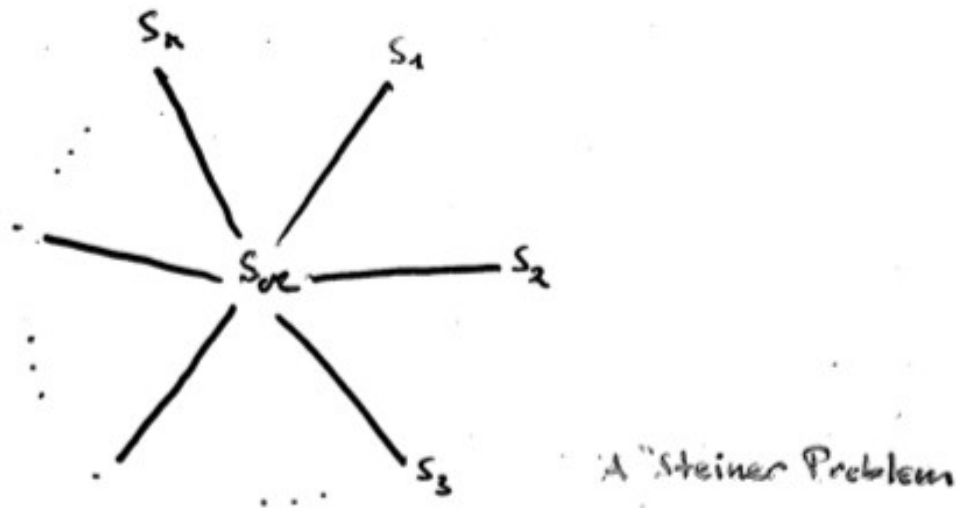
$$\text{Time: } O(|\Sigma| \cdot 2^K \cdot k \cdot n^K), \quad \text{space: } O(n^K)$$

$$= \sum_{1 \leq j \leq K} \underbrace{\sum_{1 \leq i \leq k} d(c_j, d_{i,j})}_{\text{Score of col. } j}$$

Consensus score – Alternative form

Find S_{ac} such that $\sum_{i=1}^k D(S_i, S_{ac})$ is minimized,

where $D(\cdot, \cdot)$ is the score of an opt. pairwise alignment



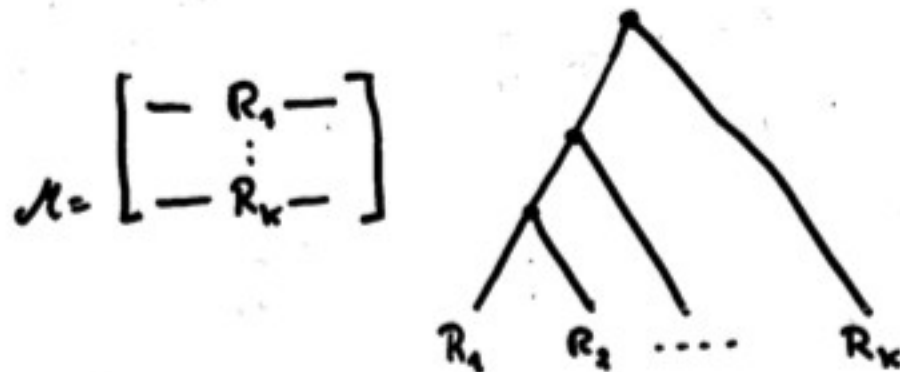
The optimal alignment \mathcal{A}^* is a combination of the

k pairwise alg. s.t. $\mathcal{A}^*(R_i, R_j) = \mathcal{A}(S_i, S_j)$.

This is possible for any tree cf. Thm 14.6.1. in [Gusfield].

Tree score

Generalized consensus score. Say seq's S_1, S_2, \dots, S_K are related by a tree:



Given an assignment of rows to internal nodes, then:

$$TS(A) = \sum_{\text{edge}(X,Y)} \text{PairScore}(X,Y)$$

Tree score

Alignment problem

Find \mathcal{M}^* , i.e. R_1, \dots, R_K , and assignments to internal nodes s.t.

$TS(\mathcal{M}^*)$ is minimized

If Pair Score (\cdot, \cdot) is col. based, then $TS(\cdot)$ becomes col. based similar to $CS(\cdot)$, we can compute an optimal TS-alignment in

Time: $O(|S|^K \cdot 2^K \cdot k \cdot n^K)$, Space: $O(n^K)$

↑
all possible assignments.

Tree score – Alternative form

Say S_1, \dots, S_k are related by a tree:



Find assignment of seq's $\in \Sigma^*$ to internal nodes s.t.

$$\sum_{\text{edge}(X,Y)} D(X,Y)$$

is minimized, the optimal alignment all^* is a combination of the optimal pairwise alignments.

Tree score – Alternative form

How do we get T ?

The "big" tree alignment problem:

Find tree T^* relating S_1, \dots, S_k and mult. align. \mathcal{A}^* of S_1, \dots, S_k

such $TS(\mathcal{A}^*)$ given T^* is minimized over all possible

choices of \mathcal{A}^* and T^*