Bioinformatics I

Lecture 07: Multiple Sequence Alignment II

Dr. Sebastian Will

 $Bioinformatics\ Group,\ University\ Freiburg$

Summer 2012



Multiple Alignment

Example: Sequences Alignment
$$a^{(1)} = ACCCGAG$$
 $ACCCGA-G-GACC$ $a^{(2)} = ACTACC$ $ACCCGA-G-ACC-CC$ $A^{(3)} = TCCTACGG$ $ACCCGA-G-ACC-CC$

Definition

A multiple alignment A of K sequences $a^{(1)}...a^{(K)}$ is a $K \times N$ -matrix $(A_{i,j})_{\substack{1 \le i \le K \\ 1 \le j \le N}}$ (N is the number of columns of A)

where

- 1. each entry $A_{i,j} \in (\Sigma \cup \{-\})$
- 2. for each row i: deleting all gaps from $(A_{i,1}...A_{i,N})$ yields $a^{(i)}$
- 3. no column j contains only gap symbols



How to Score Multiple Alignments

As for pairwise alignment:

- Assume columns are scored independently
- Score is sum over alignment columns

$$S(A) = \sum_{j=1}^{N} s(A_{1j}, \ldots, A_{Kj})$$

Example

$$S\begin{pmatrix} ACCCGA-G-\\ AC--TAC-C\\ TCC-TACGG \end{pmatrix} = s\begin{pmatrix} A\\A\\T \end{pmatrix} + s\begin{pmatrix} C\\C\\C \end{pmatrix} + s\begin{pmatrix} C\\-\\C \end{pmatrix} + s\begin{pmatrix} C\\-\\- \end{pmatrix} + \cdots + s\begin{pmatrix} -\\C\\G \end{pmatrix}$$

How can we know similarities?
See later in this class how to learn pairwise similarities. Not enough data to learn similarities for triples, 4-tuples, ...



How to Score Multiple Alignments

As for pairwise alignment:

- Assume columns are scored independently
- Score is sum over alignment columns

$$S(A) = \sum_{j=1}^{N} s(A_{1j}, \ldots, A_{Kj})$$

Example

$$S\begin{pmatrix} ACCCGA-G-\\ AC--TAC-C\\ TCC-TACGG \end{pmatrix} = s\begin{pmatrix} A\\A\\T \end{pmatrix} + s\begin{pmatrix} C\\C\\C \end{pmatrix} + s\begin{pmatrix} C\\-\\C \end{pmatrix} + s\begin{pmatrix} C\\-\\- \end{pmatrix} + \cdots + s\begin{pmatrix} -\\C\\G \end{pmatrix}$$

How can we know similarities?

See later in this class how to learn pairwise similarities. Not enough data to learn similarities for triples, 4-tuples, ...



How to Score Multiple Alignments

As for pairwise alignment:

- Assume columns are scored independently
- Score is sum over alignment columns

$$S(A) = \sum_{j=1}^{N} s(A_{1j}, \ldots, A_{Kj})$$

Example

$$S\begin{pmatrix} ACCCGA-G-\\ AC--TAC-C\\ TCC-TACGG \end{pmatrix} = s\begin{pmatrix} A\\A\\T \end{pmatrix} + s\begin{pmatrix} C\\C\\C \end{pmatrix} + s\begin{pmatrix} C\\-\\C \end{pmatrix} + s\begin{pmatrix} C\\-\\- \end{pmatrix} + \cdots + s\begin{pmatrix} -\\C\\G \end{pmatrix}$$

How can we know similarities?

See later in this class how to learn pairwise similarities. *Not enough data to learn similarities for triples, 4-tuples, . . .*



Idea: approximate column scores by pairwise scores

$$s\binom{x_1}{\ldots}_{x_K} = \sum_{1 \le k \le l \le K} s(x_k, x_l)$$

Sum-of-pairs is the most commonly used scoring scheme for multiple alignments.

Definition: The sum-of-pairs score of a multiple alignment A is the sum of the scores of all pairwise alignments induced by A.

Drawbacks: No theoretical justification as log-odd scores Overcounting.



Idea: approximate column scores by pairwise scores

$$s\binom{x_1}{\ldots}_{x_K} = \sum_{1 \le k < l \le K} s(x_k, x_l)$$

Sum-of-pairs is the most commonly used scoring scheme for multiple alignments.

Definition: The sum-of-pairs score of a multiple alignment \mathcal{A} is the sum of the scores of all pairwise alignments induced by \mathcal{A} .

Drawbacks: No theoretical justification as log-odd scores Overcounting.



Idea: approximate column scores by pairwise scores

$$s\binom{x_1}{\ldots}_{X_K} = \sum_{1 \le k < l \le K} s(x_k, x_l)$$

Sum-of-pairs is the most commonly used scoring scheme for multiple alignments.

Definition: The sum-of-pairs score of a multiple alignment \mathcal{A} is the sum of the scores of all pairwise alignments induced by \mathcal{A} .

$$\begin{split} & \textit{Example:} \\ & \textit{sum-of-pairs-score} \left(\begin{array}{c} \textit{ACCCGA-G-} \\ \textit{AC--TAC-C} \\ \textit{TCC-TACGG} \end{array} \right) \\ & = \textit{sc} \left(\begin{array}{c} \textit{ACCCGA-G-} \\ \textit{AC--TAC-C} \end{array} \right) + \textit{sc} \left(\begin{array}{c} \textit{ACCCGA-G-} \\ \textit{TCC-TACGG} \end{array} \right) + \textit{sc} \left(\begin{array}{c} \textit{AC-TAC-C} \\ \textit{TCCTACGG} \end{array} \right) \end{split}$$

Drawbacks: No theoretical justification as log-odd scores Overcounting.



Idea: approximate column scores by pairwise scores

$$s\binom{x_1}{\ldots}_{X_K} = \sum_{1 \le k < l \le K} s(x_k, x_l)$$

Sum-of-pairs is the most commonly used scoring scheme for multiple alignments.

Definition: The sum-of-pairs score of a multiple alignment A is the sum of the scores of all pairwise alignments induced by A.

$$= \operatorname{sc} \left(\begin{array}{c} \operatorname{ACCCGA-G-} \\ \operatorname{AC--TAC-C} \end{array} \right) + \operatorname{sc} \left(\begin{array}{c} \operatorname{ACCCGA-G-} \\ \operatorname{TCC-TACGG} \end{array} \right) + \operatorname{sc} \left(\begin{array}{c} \operatorname{AC-TAC-C} \\ \operatorname{TCCTACGG} \end{array} \right)$$

Drawbacks: No theoretical justification as log-odd scores. Overcounting.



Optimal Multiple Alignment

For 3 sequences a, b, c, use 3-dimensional matrix $D_{i,j,k} := \text{minimal sum of pairwise distances in MSA of prefix sequences } a[1..i], <math>b[1..j]$, and c[1..k]

(after initialization)

$$D_{i,j,k} = \min \begin{cases} D_{i-1,j-1,k-1} & +d(a_i,b_j) + d(a_i,c_k) + d(b_j,c_k) \\ D_{i-1,j-1,k} & +d(a_i,b_j) + 2\gamma \\ D_{i-1,j,k-1} & +d(a_i,c_k) + 2\gamma \\ D_{i,j-1,k-1} & +d(b_j,c_k) + 2\gamma \\ D_{i-1,j,k} & +2\gamma \\ D_{i,j-1,k} & +2\gamma \\ D_{i,j,k-1} & +2\gamma \end{cases}$$

For K sequences use K-dimensional matrix. (complexity $\mathit{O}(\mathit{n}^K)$)

One can show: MSA is NP-complete



Optimal Multiple Alignment

For 3 sequences a, b, c, use 3-dimensional matrix $D_{i,j,k} := \text{minimal sum of pairwise distances in MSA of prefix sequences } a[1..i], <math>b[1..j]$, and c[1..k]

(after initialization)

$$D_{i,j,k} = \min \begin{cases} D_{i-1,j-1,k-1} & +d(a_i,b_j) + d(a_i,c_k) + d(b_j,c_k) \\ D_{i-1,j-1,k} & +d(a_i,b_j) + 2\gamma \\ D_{i-1,j,k-1} & +d(a_i,c_k) + 2\gamma \\ D_{i,j-1,k-1} & +d(b_j,c_k) + 2\gamma \\ D_{i-1,j,k} & +2\gamma \\ D_{i,j-1,k} & +2\gamma \\ D_{i,j,k-1} & +2\gamma \end{cases}$$

For K sequences use K-dimensional matrix. (complexity $O(n^K)$)



Optimal Multiple Alignment

For 3 sequences a, b, c, use 3-dimensional matrix $D_{i,j,k} := \text{minimal sum of pairwise distances in MSA of prefix sequences } a[1..i], <math>b[1..j]$, and c[1..k]

(after initialization)

$$D_{i,j,k} = \min \begin{cases} D_{i-1,j-1,k-1} & +d(a_i,b_j) + d(a_i,c_k) + d(b_j,c_k) \\ D_{i-1,j-1,k} & +d(a_i,b_j) + 2\gamma \\ D_{i-1,j,k-1} & +d(a_i,c_k) + 2\gamma \\ D_{i,j-1,k-1} & +d(b_j,c_k) + 2\gamma \\ D_{i-1,j,k} & +2\gamma \\ D_{i,j-1,k} & +2\gamma \\ D_{i,j,k-1} & +2\gamma \end{cases}$$

For K sequences use K-dimensional matrix. (complexity $O(n^K)$)

One can show: MSA is NP-complete.



Idea: Don't fill entries in K-dimensional matrix that are provably not on optimal trace. (Here only K=3).

Definition



Idea: Don't fill entries in K-dimensional matrix that are provably not on optimal trace. (Here only K=3).

Definition

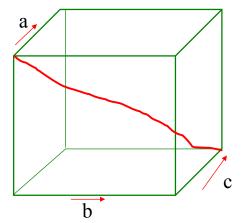
Let $d_{a,b}(i,j)$ be edit distance for suffixes a[i+1..|a|] and b[j+1..|b|]. Analogously define $d_{a,c}(i,k)$ and $d_{b,c}(j,k)$.

Remark: These *suffix distances* can be computed like prefix distances after reversing the sequences (cf. Hirschberg-algorithm).



Idea: Don't fill entries in K-dimensional matrix that are provably not on optimal trace. (Here only K=3).

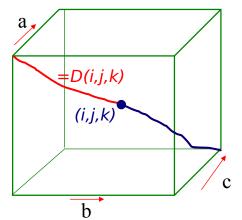
Definition





Idea: Don't fill entries in K-dimensional matrix that are provably not on optimal trace. (Here only K=3).

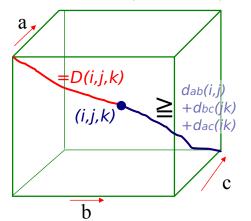
Definition





Idea: Don't fill entries in K-dimensional matrix that are provably not on optimal trace. (Here only K=3).

Definition





Key idea

Guess score z of some "good" alignment of a, b, and c. If $D_{i,j,k}+d_{a,b}(i,j)+d_{a,c}(i,k)+d_{b,c}(j,k)\geq z$, then $D_{i,j,k}$ does not contribute to optimum.

This saves time in *forward dynamic programming* (vs. backward DP, which we used before). If many results are not sent forward, many entries are never computed (\rightarrow speedup)

Implemented in tool MSA. Up to 6 sequences, n = 200.



Key idea

Guess score z of some "good" alignment of a, b, and c. If $D_{i,j,k}+d_{a,b}(i,j)+d_{a,c}(i,k)+d_{b,c}(j,k)\geq z$, then $D_{i,j,k}$ does not contribute to optimum.

This saves time in *forward dynamic programming* (vs. backward DP, which we used before). If many results are not sent forward, many entries are never computed (\rightarrow speedup)

Implemented in tool MSA. Up to 6 sequences, n = 200.



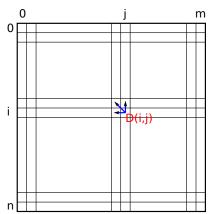
Forward Dynamic Programming

Case K=2:
$$D_{i,j} = \min \left\{ \begin{array}{lcl} D_{i-1,j-1} & + & d(a_i,b_j) \\ D_{i-1,j} & + & \gamma \\ D_{i,j-1} & + & \gamma \end{array} \right.$$

Forward Dynamic Programming

Case K=2:
$$D_{i,j} = \min \left\{ \begin{array}{lcl} D_{i-1,j-1} & + & d(a_i,b_j) \\ D_{i-1,j} & + & \gamma \\ D_{i,j-1} & + & \gamma \end{array} \right.$$

"Backward" DP: (classical evaluation)

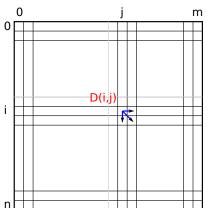




Forward Dynamic Programming

Case K=2:
$$D_{i,j} = \min \left\{ egin{array}{lll} D_{i-1,j-1} & + & d(a_i,b_j) \\ D_{i-1,j} & + & \gamma \\ D_{i,j-1} & + & \gamma \end{array} \right.$$

"Forward DP: Send scores in DP entries forward



Forward Dynamic Programming: Pseudocode

```
Pseudocode K=3
Init queue Q := empty
Push Q, (0, 0, 0) \mapsto 0
While Q not empty
  (i, j, k) \mapsto s := top(Q)
  send s+2\gamma to (i,j,k+1)
  send s+2\gamma to (i, j+1, k)
  send s+2\gamma to (i+1,i,k)
  send s + d(a_{i+1}, b_{i+1}, c_{k+1}) to (i+1, j+1, k+1)
```



Forward Dynamic Programming: Pseudocode

```
Pseudocode K=3
Init queue Q := emptv
Push Q, (0, 0, 0) \mapsto 0
While Q not empty
  (i, j, k) \mapsto s := top(Q)
  send s+2\gamma to (i,j,k+1)
  send s+2\gamma to (i, i+1, k)
  send s + 2\gamma to (i + 1, i, k)
  send s + d(a_{i+1}, b_{i+1}, c_{k+1}) to (i+1, j+1, k+1)
Procedure send s to (i, j, k):
If (i, j, k) \mapsto s' in Q:
  update entry to (i,j,k) \mapsto min(s,s')
Otherwise:
  Push Q, (i, j, k) \mapsto s
```

Forward DP Allows Lazy Evaluation

Recall: in Carillo-Lipman, we send forward only from entries that pass the criterion

$$D_{i,j,k} + d_{a,b}(i,j) + d_{a,c}(i,k) + d_{b,c}(j,k) < z$$

Only entries that received any input are processed!

In Carillo-Lipman, this saves the computation of many matrix entries.

Recall: MSA aligns 6 sequences, n = 200.

Otherwise very expensive: $200^6 = 6.4 \times 10^{13}$ entries



Most Common Heuristics: Iterative Alignment

IDEA: Iteratively merge alignments of subsets of sequences into multiple alignment

Methods differ in

- selection of subsets, i.e. order of merges
- merging procedure

Example: simplest iterative alignments strategy

- 1. align two sequences with minimum edit distance
- merge in sequence with minimum edit distance among remaining sequences
- 3. repeat (until all sequences are aligned)

More sophisticated: *progressive alignment*, e.g., Feng-Doolittle and CLUSTAL variants



Aligning Alignments in Progressive Alignment

Two (multiple) alignments A and B can be aligned by DP in the same way as two sequences.

Idea:

• An alignment is a sequence of alignment columns.

$$\begin{array}{ccc} & \text{ACCCGA-G-} \\ \text{Example: } & \text{AC--TAC-C} \\ & \text{TCC-TACGG} \end{array} & \equiv \begin{pmatrix} A \\ A \\ T \end{pmatrix} \begin{pmatrix} C \\ C \\ C \end{pmatrix} \begin{pmatrix} C \\ - \\ C \end{pmatrix} \dots \begin{pmatrix} - \\ C \\ G \end{pmatrix}.$$

• Assign similarity to two columns from resp. A and B, e.g.

$$s(\begin{pmatrix} - \\ C \\ G \end{pmatrix}, \begin{pmatrix} G \\ C \end{pmatrix})$$
 by sum-of-pairs.

We can use dynamic programming, which recurses over alignment scores of prefixes of alignments.

Consequences for progressive alignment scheme:

- Optimization only local.
- Commits to local decisions. "Once a gap, always a gap"



IN:
$$a^{(1)} = ACCG$$
, $a^{(2)} = TTGG$, $a^{(3)} = TCG$, $a^{(4)} = CTGG$

$$w(x,y) = \begin{cases} 0 \text{ iff } x = y \\ 2 \text{ iff } x = - \text{ or } y = - \\ 3 \text{ otherwise (for mismatch)} \end{cases}$$

Compute all against all edit distances and cluster

| Align ACCG and TTGG Align ACCG and TCG | | | | | | | | | | | | | |
|--|---|---|----|---|----|-----|--------------------|---|----|---|---|---|--|
| | | Т | Т | G | G | | _ | | Т | C | G | | |
| | 0 | 2 | 4 | 6 | 8 | | | 0 | 2 | 4 | 6 | | |
| Α | 2 | 3 | 5 | 7 | 9 | | A | 2 | 3 | 5 | 7 | | |
| C | 4 | 5 | 6 | 8 | 10 | | C | 4 | 5 | 3 | 6 | | |
| C | 6 | 7 | 8 | 9 | 11 | | C | 6 | 7 | 5 | 6 | | |
| G | 8 | 9 | 10 | 8 | 9 | | G | 8 | 9 | 8 | 5 | | |
| Align ACCG and CTGG | | | | | | | Align TTGG and TCG | | | | | | |
| , | | С | T | G | G | 7 | 8 | | T | C | G | | |
| | 0 | 2 | 4 | 6 | 8 | | | 0 | 2 | 4 | 6 | | |
| Α | 2 | 3 | 5 | 7 | 9 | - | Т | 2 | 0 | 3 | 6 | | |
| C | 4 | 2 | 5 | 8 | 10 | - | Т | 4 | 2 | 3 | 6 | | |
| C | 6 | 4 | 5 | 8 | 11 | | G | 6 | 5 | 5 | 3 | | |
| G | 8 | 7 | 7 | 5 | 8 | | G | 8 | 8 | 8 | 5 | | |
| Align TTGG and CTGG | | | | | | Ali | Align TCG and CTGG | | | | | | |
| | | C | Т | G | G | | o | | C. | Т | G | G | |
| | 0 | 2 | 4 | 6 | 8 | | | 0 | 2 | 4 | 6 | 8 | |
| Т | 2 | 3 | 2 | 5 | 8 | - | Т | 2 | 3 | 2 | 5 | 8 | |
| Т | 4 | 5 | 3 | 5 | 8 | | С | 4 | 2 | 5 | 5 | 8 | |
| G | 6 | 7 | 6 | 3 | 5 | | G | 6 | 5 | 5 | 5 | 5 | |
| G | Ω | O | ۵ | 6 | 3 | | | | | | | | |



IN:
$$a^{(1)} = ACCG$$
, $a^{(2)} = TTGG$, $a^{(3)} = TCG$, $a^{(4)} = CTGG$

$$w(x,y) = \begin{cases} 0 \text{ iff } x = y \\ 2 \text{ iff } x = - \text{ or } y = - \\ 3 \text{ otherwise (for mismatch)} \end{cases}$$

Compute all against all edit distances and cluster
 ⇒ distance matrix

 \Rightarrow Cluster (e.g. UPGMA) $a^{(2)}$ and $a^{(4)}$ are closest, Then: $a^{(1)}$ and $a^{(3)}$



IN:
$$a^{(1)} = ACCG$$
, $a^{(2)} = TTGG$, $a^{(3)} = TCG$, $a^{(4)} = CTGG$

$$w(x,y) = \begin{cases} 0 \text{ iff } x = y \\ 2 \text{ iff } x = - \text{ or } y = - \\ 3 \text{ otherwise (for mismatch)} \end{cases}$$

- Compute all against all edit distances and cluster \Rightarrow guide tree $((a^{(2)},a^{(4)}),(a^{(1)},a^{(3)}))$
- Align $a^{(2)}$ and $a^{(4)}$: ${}^{\rm TTGG}_{\rm CTGG}$, Align $a^{(1)}$ and $a^{(3)}$: ${}^{\rm ACCG}_{\rm -TCG}$
- Align the alignments!

| Align | TTGG CTGG | and | ACC -TC | :G :G | |
|-------|--------------|--------|------------|----------|----|
| | | A - | <u>C</u> | C | G |
| | 0 | 4 | 12 | 20 | 28 |
| TC | 8 | 10 | | | |
| TT | 16 | | ٠. | | |
| GG | 24 | | | | |
| GG | 32 | | | | |



IN:
$$a^{(1)} = ACCG$$
, $a^{(2)} = TTGG$, $a^{(3)} = TCG$, $a^{(4)} = CTGG$
 $w(x,y) = \begin{cases} 0 \text{ iff } x = y \\ 2 \text{ iff } x = - \text{ or } y = - \\ 3 \text{ otherwise (for mismatch)} \end{cases}$

- Compute all against all edit distances and cluster
 - \Rightarrow guide tree $((a^{(2)}, a^{(4)}), (a^{(1)}, a^{(3)}))$
- ACCG • Align $a^{(2)}$ and $a^{(4)}$: $\frac{\text{TTGG}}{\text{grad}}$, Align $a^{(1)}$ and $a^{(3)}$: -TCG
- Align the alignments!

$$w(T, -) + w(C, -) + w(T, -) + w(C, -) = 8$$

G
$$w(--,A-) = W(-,A) + w(-,-) + w(-,A) + w(-,-) = 4$$

$$w(T,A) + w(C,A) + w(T,-) + w(C,-) = 10$$

•
$$w(TC, CT) =$$

 $w(T, C) + w(C, C) + w(T, T) + w(C, T) = 6$



w(TC, −−) =



IN:
$$a^{(1)} = ACCG$$
, $a^{(2)} = TTGG$, $a^{(3)} = TCG$, $a^{(4)} = CTGG$

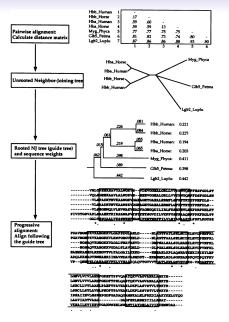
$$w(x,y) = \begin{cases} 0 \text{ iff } x = y \\ 2 \text{ iff } x = - \text{ or } y = - \\ 3 \text{ otherwise (for mismatch)} \end{cases}$$

- Compute all against all edit distances and cluster \Rightarrow guide tree $((a^{(2)},a^{(4)}),(a^{(1)},a^{(3)}))$
- Align $a^{(2)}$ and $a^{(4)}$: $_{\rm CTGG}^{\rm TTGG}$, Align $a^{(1)}$ and $a^{(3)}$: $_{\rm -TCG}^{\rm ACCG}$
- Align the alignments!

| Align | TTGG CTGG | and | ACCG -TCG | | | | |
|----------------|--------------------|--------------|--------------|----|--------|--------------------------------|------------------------------|
| | | A | <u>C</u> | C | G G | \Longrightarrow | TTGG CTGG ACCG -TCG |
| TC TT GG | 0 8 16 24 | 4 10 : | 12 | 20 | 28 | after filling and traceback | |



A Classical Approach: CLUSTAL W



- prototypical progressive alignment
- similarity score with affine gap cost
- neighbor joining for tree construction
- sequence weighting
- 'tricks' for gap handling

Details: see last week's slides



Improving Progressive Alignment

- Limitations of progressive alignment: once a gap – always a gap: decisions are local (cf. consistency, Bioinformatics II)
- Iterative refinement: iteratively select one sequence or alignment of some sequences and realign to rest
- Performance bottleneck: guide tree construction requires $O(K^2)$ all-2-all comparisons (vs. O(K) comparisons construct MSA)
 - How to align hundreds or thousands of sequences?
 - use faster distance measure
 - construct draft alignment and reestimate distance



Improving Progressive Alignment

- Limitations of progressive alignment: once a gap – always a gap: decisions are local (cf. consistency, Bioinformatics II)
- Iterative refinement: iteratively select one sequence or alignment of some sequences and realign to rest
- Performance bottleneck: guide tree construction requires $O(K^2)$ all-2-all comparisons (vs. O(K) comparisons construct MSA)
 - How to align hundreds or thousands of sequences?
 - use faster distance measure
 - construct draft alignment and reestimate distance



Improving Progressive Alignment

- Limitations of progressive alignment: once a gap – always a gap: decisions are local (cf. consistency, Bioinformatics II)
- Iterative refinement: iteratively select one sequence or alignment of some sequences and realign to rest

• Performance bottleneck: guide tree construction requires $O(K^2)$ all-2-all comparisons (vs. O(K) comparisons construct MSA)

How to align hundreds or thousands of sequences?

- use faster distance measure
- construct draft alignment and reestimate distances



Improving Progressive Alignment

- Limitations of progressive alignment: once a gap – always a gap: decisions are local (cf. consistency, Bioinformatics II)
- Iterative refinement: iteratively select one sequence or alignment of some sequences and realign to rest
- Performance bottleneck: guide tree construction requires $O(K^2)$ all-2-all comparisons (vs. O(K) comparisons construct MSA)

How to align hundreds or thousands of sequences?

- use faster distance measure
- construct draft alignment and reestimate distances



Improving Progressive Alignment

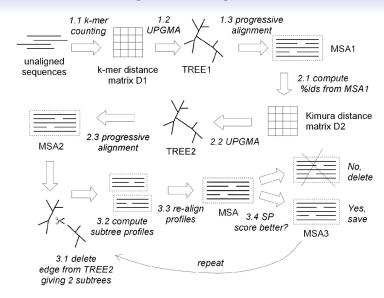
- Limitations of progressive alignment: once a gap – always a gap: decisions are local (cf. consistency, Bioinformatics II)
- Iterative refinement: iteratively select one sequence or alignment of some sequences and realign to rest
- Performance bottleneck: guide tree construction requires
 O(K²) all-2-all comparisons
 (vs. O(K) comparisons construct MSA)

How to align hundreds or thousands of sequences?

- use faster distance measure
- construct draft alignment and reestimate distances



Advanced Progressive Alignment in MUSCLE



1.) alignment draft and 2.) reestimation 3.) iterative refinement



How to compare MSA tools?

Criteria

- speed
- degree of optimization (of SP-score) by heuristics
- quality of computed alignments (in terms of true alignments)
- Benchmark
 - measuring alignment quality requires comparison to reference
 - benchmark instances: sets of protein sequences with hand-curated alignments
 - measure average alignment quality across instances
- Example: BAliBASE [Thompson et al., Bioinformatics. 1999]
 - reference alignments: 3D structural superpositions



How to compare MSA tools?

- Criteria
 - speed
 - degree of optimization (of SP-score) by heuristics
 - quality of computed alignments (in terms of true alignments)
- Benchmark
 - measuring alignment quality requires comparison to reference
 - benchmark instances: sets of protein sequences with hand-curated alignments
 - measure average alignment quality across instances
- Example: BAliBASE [Thompson et al., Bioinformatics. 1999]
 - reference alignments: 3D structural superpositions



How to compare MSA tools?

- Criteria
 - speed
 - degree of optimization (of SP-score) by heuristics
 - quality of computed alignments (in terms of true alignments)
- Benchmark
 - measuring alignment quality requires comparison to reference
 - benchmark instances: sets of protein sequences with hand-curated alignments
 - measure average alignment quality across instances
- Example: BAliBASE [Thompson et al., Bioinformatics. 1999]
 - reference alignments: 3D structural superpositions



BAliBASE 3.0: 6255 sequences

Defines core regions of reliably aligned residues

Measure accuracy of computed (K-way) alignment A by comparison to reference alignment R (of same sequences): SPS(A|R) = average over pairwise alignments of rows k and l:

valid columns in rows k and l of A that identically occur in R total valid columns in rows k and l of A

(valid: ≥ 1 non-gap; identical: same sequence position)

For example:
$$A = \left(\begin{array}{c} \mathbf{UVW-XZ} \\ \mathbf{UWX--Z} \\ \mathbf{UV-XYZ} \end{array} \right) \text{ vs. } R = \left(\begin{array}{c} \mathbf{UVWX-Z} \\ \mathbf{U-W-XZ} \\ \mathbf{UV-XYZ} \end{array} \right)$$



BAliBASE 3.0: 6255 sequences

Defines core regions of reliably aligned residues

Measure accuracy of computed (K-way) alignment A by comparison to reference alignment R (of same sequences): SPS(A|R) = average over pairwise alignments of rows k and l:

valid columns in rows k and l of A that identically occur in R total valid columns in rows k and l of A

(valid: ≥ 1 non-gap; identical: same sequence position)

For example:
$$A = \left(\begin{array}{c} \tt UVW-XZ \\ \tt UWX--Z \\ \tt UV-XYZ \end{array} \right) \text{ vs. } R = \left(\begin{array}{c} \tt UVWX-Z \\ \tt U-W-XZ \\ \tt UV-XYZ \end{array} \right)$$

Consider k = 1, l = 2: total=5, identical=3, ratio = 0.6



BAliBASE 3.0: 6255 sequences

Defines core regions of reliably aligned residues

Measure accuracy of computed (K-way) alignment A by comparison to reference alignment R (of same sequences): SPS(A|R) = average over pairwise alignments of rows k and l:

valid columns in rows k and l of A that identically occur in R total valid columns in rows k and l of A

(valid: ≥ 1 non-gap; identical: same sequence position)

For example:
$$A = \left(\begin{array}{c} \mathbf{UVW-XZ} \\ \mathbf{UWX--Z} \\ \mathbf{UV-XYZ} \end{array}\right) \text{ vs. } R = \left(\begin{array}{c} \mathbf{UVWX-Z} \\ \mathbf{U-W-XZ} \\ \mathbf{UV-XYZ} \end{array}\right)$$

Consider k = 1, l = 2: total=5, identical=3, ratio = 0.6



BAliBASE 3.0: 6255 sequences

Defines core regions of reliably aligned residues

Measure accuracy of computed (K-way) alignment A by comparison to reference alignment R (of same sequences): SPS(A|R) = average over pairwise alignments of rows k and l:

valid columns in rows k and l of A that identically occur in R total valid columns in rows k and l of A

(valid: ≥ 1 non-gap; identical: same sequence position)

For example:
$$A = \left(\begin{array}{c} \mathbf{UVW-XZ} \\ \mathbf{UWX--Z} \\ \mathbf{UV-XYZ} \end{array} \right) \text{ vs. } R = \left(\begin{array}{c} \mathbf{UVWX-Z} \\ \mathbf{U-W-XZ} \\ \mathbf{UV-XYZ} \end{array} \right)$$

Consider k = 1, l = 2: total=5, identical=3, ratio = 0.6



BAliBASE 3.0: 6255 sequences

Defines core regions of reliably aligned residues

Measure accuracy of computed (K-way) alignment A by comparison to reference alignment R (of same sequences): SPS(A|R) = average over pairwise alignments of rows k and l:

valid columns in rows k and l of A that identically occur in R total valid columns in rows k and l of A

(valid: ≥ 1 non-gap; identical: same sequence position)

For example:
$$A = \left(\begin{array}{c} \tt UVW-XZ \\ \tt UWX--Z \\ \tt UV-XYZ \end{array} \right) \text{ vs. } R = \left(\begin{array}{c} \tt UVWX-Z \\ \tt U-W-XZ \\ \tt UV-XYZ \end{array} \right)$$

Consider k = 1, l = 2: total=5, identical=3, ratio = 0.6 ratio for k = 2, l = 3: total=6, identical= , ratio=



some results from [Edgar, MUSCLE 2004]

| | BAliBASE score | BAliBASE time (s) |
|-------------|----------------|-------------------|
| MUSCLE | 0.896 | 81 |
| MAFFT | 0.844 | 16 |
| MUSCLE-fast | 0.849 | 11 |
| CLUSTALW | 0.860 | 160 |



- Uses of Multiple Sequence Alignments (MSA)
- Sum-Of-Pairs Scoring (SP)
- Methods for solving the MSA problem: exact/heuristic
- Heuristics: no guarantee, but effective in practice
- Iterative/Progressive Alignment
- ClustalW: classic progressive alignment
 - distance matrix
 - guide tree
 - optimized gap penalties
- MUSCLE: progressive alignment with various improvements
 - guide tree re-estimation
 - iterative refinement
- Benchmarks (BAliBASE)



- Uses of Multiple Sequence Alignments (MSA)
- Sum-Of-Pairs Scoring (SP)
- Methods for solving the MSA problem: exact/heuristic
- Heuristics: no guarantee, but effective in practice
- Iterative/Progressive Alignment
- ClustalW: classic progressive alignment
 - distance matrix
 - guide tree
 - optimized gap penalties
- MUSCLE: progressive alignment with various improvements
 - guide tree re-estimation
 - iterative refinement
- Benchmarks (BAliBASE)



- Uses of Multiple Sequence Alignments (MSA)
- Sum-Of-Pairs Scoring (SP)
- Methods for solving the MSA problem: exact/heuristic
- Heuristics: no guarantee, but effective in practice
- Iterative/Progressive Alignment
- ClustalW: classic progressive alignment
 - distance matrix
 - guide tree
 - optimized gap penalties
- MUSCLE: progressive alignment with various improvements
 - guide tree re-estimation
 - iterative refinement
- Benchmarks (BAliBASE)



- Uses of Multiple Sequence Alignments (MSA)
- Sum-Of-Pairs Scoring (SP)
- Methods for solving the MSA problem: exact/heuristic
- Heuristics: no guarantee, but effective in practice
- Iterative/Progressive Alignment
- ClustalW: classic progressive alignment
 - distance matrix
 - guide tree
 - optimized gap penalties
- MUSCLE: progressive alignment with various improvements
 - guide tree re-estimation
 - iterative refinement
- Benchmarks (BAliBASE)



- Uses of Multiple Sequence Alignments (MSA)
- Sum-Of-Pairs Scoring (SP)
- Methods for solving the MSA problem: exact/heuristic
- Heuristics: no guarantee, but effective in practice
- Iterative/Progressive Alignment
- ClustalW: classic progressive alignment
 - distance matrix
 - guide tree
 - optimized gap penalties
- MUSCLE: progressive alignment with various improvements
 - guide tree re-estimation
 - iterative refinement
- Benchmarks (BAliBASE)



- Uses of Multiple Sequence Alignments (MSA)
- Sum-Of-Pairs Scoring (SP)
- Methods for solving the MSA problem: exact/heuristic
- Heuristics: no guarantee, but effective in practice
- Iterative/Progressive Alignment
- ClustalW: classic progressive alignment
 - distance matrix
 - guide tree
 - optimized gap penalties
- MUSCLE: progressive alignment with various improvements
 - guide tree re-estimation
 - iterative refinement
- Benchmarks (BAliBASE)



- Uses of Multiple Sequence Alignments (MSA)
- Sum-Of-Pairs Scoring (SP)
- Methods for solving the MSA problem: exact/heuristic
- Heuristics: no guarantee, but effective in practice
- Iterative/Progressive Alignment
- ClustalW: classic progressive alignment
 - distance matrix
 - guide tree
 - optimized gap penalties
- MUSCLE: progressive alignment with various improvements
 - guide tree re-estimation
 - iterative refinement
- Benchmarks (BAliBASE)



- Uses of Multiple Sequence Alignments (MSA)
- Sum-Of-Pairs Scoring (SP)
- Methods for solving the MSA problem: exact/heuristic
- Heuristics: no guarantee, but effective in practice
- Iterative/Progressive Alignment
- ClustalW: classic progressive alignment
 - distance matrix
 - guide tree
 - optimized gap penalties
- MUSCLE: progressive alignment with various improvements
 - guide tree re-estimation
 - iterative refinement
- Benchmarks (BAliBASE)

