# Sequence Alignment

Motivation: assess similarity of sequences and learn about their evolutionary relationship Why do we want to know this? Example: Sequences Alignment ACCCGA ACCCGA  $\Rightarrow_{\text{align}}$ ACTA AC--TA TCCTA TCC-TA Homology: Alignment reasonable, if sequences homologous ACCGA ACCTA ACTA ACCCGA TCCTA Definition (Sequence Homology) Two or more sequences are *homologous* iff they evolved from a common ancestor. [Homology in anatomy

Will, 18.417, Fall 2011

# Plan (and Some Preliminaries)

- First: study only pairwise alignment.
   Fix alphabet Σ, such that ∉ Σ. is called the gap symbol.
   The elements of Σ\* are called sequences.
   Fix two sequences a, b ∈ Σ\*.
- For pairwise sequence comparison: define edit distance, define alignment distance, show equivalence of distances, define alignment problem and efficient algorithm gap penalties, local alignment
- Later: extend pairwise alignment to multiple alignment

### Definition (Alphabet, words)

An alphabet  $\Sigma$  is a finite set (of symbols/characters).  $\Sigma^+$  denotes the set of non-empty words of  $\Sigma$ , i.e.  $\Sigma^+ := \bigcup_{i>0} \Sigma^i$ . A word  $x \in \Sigma^n$  has length n, written |x|.  $\Sigma^* := \Sigma^+ \cup \{\epsilon\}$ , where  $\epsilon$  denotes the empty word of length 0.

# Levenshtein Distance

### Definition

The *Levenshtein Distance* between two words/sequences is the minimal number of substitutions, insertions and deletions to transform one into the other.

### Example

ACCCGA and ACTA have (at most) distance 3: ACCCGA  $\rightarrow$  ACCGA  $\rightarrow$  ACCTA  $\rightarrow$  ACTA

In biology, operations have different cost. (Why?)

# Edit Distance: Operations

# Definition (Edit Operations)

An edit operation is a pair  $(x, y) \in (\Sigma \cup \{-\} \neq (-, -)$ . We call (x,y)

- substitution iff  $x \neq -$  and  $y \neq -$
- deletion iff y = -
- *insertion* iff x = -

For sequences a, b, write  $a \rightarrow_{(x,y)} b$ , iff a is transformed to b by operation (x, y). Furthermore, write  $a \Rightarrow_S b$ , iff a is transformed to b by a sequence of edit operations S.

### Example

$$ACCCGA \rightarrow_{(C,-)} ACCGA \rightarrow_{(G,T)} ACCTA \rightarrow_{(-,T)} ATCCTA$$
  
 $ACCCGA \Rightarrow_{(C,-),(G,T),(-,T)} ATCCTA$ 

Recall:  $- \notin \Sigma$ , a, b are sequences in  $\Sigma^*$ 

Definition (Cost, Edit Distance) Let  $w : (\Sigma \cup \{-\})^2 \to \mathbb{R}$ , such that w(x, y) is the cost of an edit operation (x, y). The cost of a sequence of edit operations  $S = e_1, \ldots, e_n$  is  $\tilde{w}(S) = \sum_{i=1}^n w(e_i).$ 

The edit distance of sequences a and b is

$$d_w(a,b) = \min\{\tilde{w}(S) \mid a \Rightarrow_S b\}.$$



## Edit Distance: Cost and Problem Definition

Definition (Cost, Edit Distance) Let  $w : (\Sigma \cup \{-\})^2 \to \mathbb{R}$ , such that w(x, y) is the cost of an edit operation (x, y). The cost of a sequence of edit operations  $S = e_1, \ldots, e_n$  is  $\tilde{w}(S) = \sum_{i=1}^n w(e_i).$ 

The edit distance of sequences a and b is

$$d_w(a,b) = \min\{\tilde{w}(S) \mid a \Rightarrow_S b\}.$$

### Is the definition reasonable?

### Definition (Metric)

A function  $d: X^2 \to \mathbb{R}$  is called *metric* iff 1.) d(x, y) = 0 iff x = y2.) d(x, y) = d(y, x) 3.)  $d(x, y) \le d(x, z) + d(z, y)$ .

Remarks: 1.) for metric d,  $d(x, y) \ge 0$ , 2.)  $d_w$  is metric iff  $w(x, y) \ge 0$ , 3.) In the following, assume  $d_w$  is metric.

# Edit Distance: Cost and Problem Definition

Definition (Cost, Edit Distance) Let  $w : (\Sigma \cup \{-\})^2 \to \mathbb{R}$ , such that w(x, y) is the cost of an edit operation (x, y). The cost of a sequence of edit operations  $S = e_1, \ldots, e_n$  is  $\tilde{w}(S) = \sum_{i=1}^n w(e_i).$ 

The edit distance of sequences a and b is

$$d_w(a,b) = \min\{\tilde{w}(S) \mid a \Rightarrow_S b\}.$$

- Natural 'evolution-motivated' problem definition.
- Not obvious how to compute edit distance efficiently ⇒ define alignment distance

# Alignment Distance

# Definition (Alignment)

A pair of words  $a^{\diamond}, b^{\diamond} \in (\Sigma \cup \{-\})^*$  is called *alignment of* sequences a and b ( $a^{\diamond}$  and  $b^{\diamond}$  are called *alignment strings*), iff

- 1.  $|a^{\diamond}| = |b^{\diamond}|$
- 2. for all  $1 \leq i \leq |a^{\diamond}|$ :  $a_i^{\diamond} \neq -$  or  $b_i^{\diamond} \neq -$
- deleting all gap symbols from a<sup>◊</sup> yields a and deleting all – from b<sup>◊</sup> yields b

### Example

- a = ACGGAT
- b = CCGCTT

possible alignments are

 $\begin{array}{lll} a^\diamond &= \mbox{AC-GG-AT} & a^\diamond &= \mbox{ACGG---AT} \\ b^\diamond &= \mbox{-CCGCT-T} & b^\diamond &= \mbox{--CCGCT-T} \end{array} \mbox{ or } \dots \mbox{ (exponentially many)} \end{array}$ 

edit operations of first alignment: (A,-),(-,C),(G,C),(-,T),(A,-)

# Alignment Distance

Definition (Cost of Alignment, Alignment Distance) The *cost of the alignment*  $(a^{\diamond}, b^{\diamond})$ , given a cost function w on edit operations is

$$w(a^\diamond,b^\diamond)=\sum_{i=1}^{|a^\diamond|}w(a^\diamond_i,b^\diamond_i)$$

The alignment distance of a and b is

 $D_w(a,b) = \min\{w(a^\diamond,b^\diamond) \mid (a^\diamond,b^\diamond) \text{ is alignment of } a \text{ and } b\}.$ 

Theorem (Equivalence of Edit and Alignment Distance) For metric w,  $d_w(a, b) = D_w(a, b)$ .

Recall:

## Definition (Edit Distance)

The *edit distance of a and b* is  $d_w(a, b) = \min{\{\tilde{w}(S) \mid a \text{ transformed to } b \text{ by e.o.-sequence } S\}}.$ 

Definition (Alignment Distance)

The alignment distance of a and b is  $D_w(a, b) = \min\{w(a^\diamond, b^\diamond) \mid (a^\diamond, b^\diamond) \text{ is alignment of } a \text{ and } b\}.$  Theorem (Equivalence of Edit and Alignment Distance) For metric w,  $d_w(a, b) = D_w(a, b)$ .

Remarks

• Proof idea:

 $d_w(a, b) \leq D_w(a, b)$ : alignment yields sequence of edit ops  $D_w(a, b) \leq d_w(a, b)$ : sequence of edit ops yields equal or better alignment (needs triangle inequality)

- Reduces edit distance to alignment distance
- We will see: the alignment distance is computed efficiently by dynamic programming (using *Bellman's Principle of Optimality*).

# Principle of Optimality and Dynamic Programming

### Principle of Optimality:

'Optimal solutions consist of optimal partial solutions'

Example: Shortest Path



Idea of Dynamic Programming (DP):

- Solve partial problems first and materialize results
- (recursively) solve larger problems based on smaller ones

- The principle is valid for the alignment distance problem
- Principle of Optimality enables the programming method DP
- Dynamic programming is widely used in Computational Biology and you will meet it quite often in this class

# Alignment Matrix

*Idea:* choose alignment distances of prefixes  $a_{1..i}$  and  $b_{1..j}$  as partial solutions and define matrix of these partial solutions.

Let n := |a|, m := |b|.

## Definition (Alignment matrix)

The alignment matrix of a and b is the  $(n + 1) \times (m + 1)$ -matrix  $D := (D_{ij})_{0 \le i \le n, 0 \le j \le m}$  defined by

$$\begin{split} D_{ij} &:= D_w(a_{1..i}, b_{1..j}) \\ & \big( = \min\{w(a^\diamond, b^\diamond) \mid (a^\diamond, b^\diamond) \text{ is alignment of } a_{1..i} \text{ and } b_{1..j}\} \big). \end{split}$$

### Notational remarks

- *a<sub>i</sub>* is the i-th character of *a*
- $a_{x..y}$  is the sequence  $a_x a_{x+1} \dots a_y$  (subsequence of a).
- by convention  $a_{x..y} = \epsilon$  if x > y.

# Alignment Matrix Example

### Example

• 
$$a = AT$$
,  $b = AAGT$   
•  $w(x, y) = \begin{cases} 0 & \text{iff } x = y \\ 1 & \text{otherwise} \end{cases}$ 



Remark: The alignment matrix D contains the alignment distance (=edit distance) of a and b in  $D_{n,m}$ .



# Alignment Matrix Example

### Example

• 
$$a = AT$$
,  $b = AAGT$   
•  $w(x, y) = \begin{cases} 0 & \text{iff } x = y \\ 1 & \text{otherwise} \end{cases}$ 



Remark: The alignment matrix D contains the alignment distance (=edit distance) of a and b in  $D_{n,m}$ .



# Needleman-Wunsch Algorithm

# Claim For $(a^{\diamond}, b^{\diamond})$ alignment of a and b with length $r = |a^{\diamond}|$ , $w(a^{\diamond}, b^{\diamond}) = w(a_{1..r-1}^{\diamond}, b_{1..r-1}^{\diamond}) + w(a_r^{\diamond}, b_r^{\diamond}).$

### Theorem

For the alignment matrix D of a and b, holds that

• 
$$D_{0,0} = 0$$
  
• for all  $1 \le i \le n$ :  $D_{i,0} = \sum_{k=1}^{i} w(a_k, -) = D_{i-1,0} + w(a_i, -)$   
• for all  $1 \le j \le m$ :  $D_{0,j} = \sum_{k=1}^{j} w(-, b_k) = D_{0,j-1} + w(-, b_j)$   
•  $D_{ij} = \min \begin{cases} D_{i-1,j-1} + w(a_i, b_j) & (match) \\ D_{i-1,j} + w(a_i, -) & (deletion) \\ D_{i,j-1} + w(-, b_j) & (insertion) \end{cases}$ 

Remark: The theorem claims that each prefix alignment distance can be computed from a constant number of smaller ones. Proof ???

# Needleman-Wunsch Algorithm

# Claim For $(a^{\diamond}, b^{\diamond})$ alignment of a and b with length $r = |a^{\diamond}|$ , $w(a^{\diamond}, b^{\diamond}) = w(a_{1..r-1}^{\diamond}, b_{1..r-1}^{\diamond}) + w(a_r^{\diamond}, b_r^{\diamond}).$

### Theorem

For the alignment matrix D of a and b, holds that

• 
$$D_{0,0} = 0$$
  
• for all  $1 \le i \le n$ :  $D_{i,0} = \sum_{k=1}^{i} w(a_k, -) = D_{i-1,0} + w(a_i, -)$   
• for all  $1 \le j \le m$ :  $D_{0,j} = \sum_{k=1}^{j} w(-, b_k) = D_{0,j-1} + w(-, b_j)$   
•  $D_{ij} = \min \begin{cases} D_{i-1,j-1} + w(a_i, b_j) & (match) \\ D_{i-1,j} + w(a_i, -) & (deletion) \\ D_{i,j-1} + w(-, b_j) & (insertion) \end{cases}$ 

Remark: The theorem claims that each prefix alignment distance can be computed from a constant number of smaller ones. Proof: Induction over i+j

# Needleman-Wunsch Algorithm (Pseudocode)

```
D_{0.0} := 0
for i := 1 to n do
   D_{i,0} := D_{i-1,0} + w(a_i, -)
end for
for j := 1 to m do
   D_{0,i} := D_{0,i-1} + w(-, b_i)
end for
for i := 1 to n do
   for i := 1 to m do
      D_{i,j} := \min \begin{cases} D_{i-1,j-1} + w(a_i, b_j) \\ D_{i-1,j} + w(a_i, -) \\ D_{i,j-1} + w(-, b_j) \end{cases}
   end for
```

end to end for



# Back to Example

### Example

• 
$$a = AT$$
,  $b = AAGT$   
•  $w(x, y) = \begin{cases} 0 & \text{iff } x = y \\ 1 & \text{otherwise} \end{cases}$ 



Open: how to find best alignment?



# Back to Example

### Example

• 
$$a = AT$$
,  $b = AAGT$   
•  $w(x, y) = \begin{cases} 0 & \text{iff } x = y \\ 1 & \text{otherwise} \end{cases}$ 



Open: how to find best alignment?



# Traceback

$$w(x,y) = \begin{cases} 0 & \text{iff } x = y \\ 1 & \text{otherwise} \end{cases}$$



- Start in (n, m). For every (i, j) determine optimal case.
- Not necessarily unique.
- Sequence of trace arrows let's infer best alignment.



# Traceback

$$w(x,y) = \begin{cases} 0 & \text{iff } x = y \\ 1 & \text{otherwise} \end{cases}$$



- Start in (n, m). For every (i, j) determine optimal case.
- Not necessarily unique.
- Sequence of trace arrows let's infer best alignment.



# Complexity

- compute one entry: three cases, i.e. constant time
- nm entries  $\Rightarrow$  fill matrix in O(nm) time
- traceback: O(n + m) time
- TOTAL:  $O(n^2)$  time and space (assuming  $m \le n$ )

- assuming  $m \le n$  is w.l.o.g. since we can exchange a and b
- space complexity can be improved to O(n) for computation of distance (simple, "store only current and last row") and traceback (more involved; Hirschberg-algorithm uses "Divide and Conquer" for computing trace)

• We have seen how to compute the pairwise edit distance and the corresponding optimal alignment.

18.417, Fall 2011

- Before going multiple, we will look at two further special topics for pairwise alignment:
  - more realistic, non-linear gap cost and
  - similarity scores and local alignment

# Alignment Cost Revisited

Motivation:

- The alignments  $\begin{array}{c} GA--T\\ GAAGT \end{array}$  and  $\begin{array}{c} G-A-T\\ GAAGT \end{array}$  have the same edit distance.
- The first one is biologically more reasonable: it is more likely that evolution introduces one large gap than two small ones.
- This means: gap cost should be non-linear, sub-additive!

# Gap Penalty

### Definition (Gap Penalty)

A gap penalty is a function  $g:\mathbb{N}\to\mathbb{R}$  that is sub-additive, i.e.

$$g(k+l) \leq g(k) + g(l).$$

A gap in an alignment string  $a^{\diamond}$  is a substring of  $a^{\diamond}$  that consists of only gap symbols – and is maximally extended.  $\Delta^{a^{\diamond}}$  is the multi-set of gaps in  $a^{\diamond}$ .

The alignment cost with gap penalty g of  $(a^\diamond, b^\diamond)$  is

$$\begin{split} w_{g}(a^{\diamond}, b^{\diamond}) &= \sum_{\substack{1 \leq r \leq |a^{\diamond}|, \\ \text{where } a_{r}^{\diamond} \neq -, b_{r}^{\diamond} \neq -}} w(a_{r}^{\diamond}, b_{r}^{\diamond}) \quad (\textit{cost of mismatchs}) \\ &+ \sum_{x \in \Delta^{a^{\diamond}} \uplus \Delta^{b^{\diamond}}} g(|x|) \quad (\textit{cost of gaps}) \end{split}$$
Example:

$$a^{\diamond} = \text{ATG} - -\text{CGAC} - \text{GC} \Rightarrow \Delta^{a^{\diamond}} = \{---, --\}, \ \Delta^{b^{\diamond}} = \{-, -\}$$
  
 $b^{\diamond} = -\text{TGCGGCG} - \text{CTTTC}$ 

# General sub-additive gap penalty

### Theorem

Let D be the alignment matrix of a and b with cost w and gap penalty g, such that  $D_{ij} = w_g(a_{1..i}, b_{1..j})$ . Then:

• 
$$D_{0,0} = 0$$
  
• for all  $1 \le i \le n$ :  $D_{i,0} = g(i)$   
• for all  $1 \le j \le m$ :  $D_{0,j} = g(j)$   
•  $D_{ij} = \min \begin{cases} D_{i-1,j-1} + w(a_i, b_j) & (match) \\ \min_{1 \le k \le j} D_{i-k,j} + g(k) & (deletion of length k) \\ \min_{1 \le k \le j} D_{i,j-k} + g(k) & (insertion of length k) \end{cases}$ 

- Complexity  $O(n^3)$  time,  $O(n^2)$  space
- pseudocode, correctness, traceback left as exercise
- much more realistic, but significantly more expensive than Needleman-Wunsch ⇒ can we improve it?

# Affine gap cost

### Definition

A gap penalty is affine, iff there are real constants  $\alpha$  and  $\beta$ , such that for all  $k \in \mathbb{N}$ :  $g(k) = \alpha + \beta k$ .

- Affine gap penalties almost as good as general ones: Distinguishing gap opening  $(\alpha)$  and gap extension cost  $(\beta)$  is "biologically reasonable".
- The minimal alignment cost with affine gap penalty can be computed in  $O(n^2)$  time! (Gotoh algorithm)

## Gotoh algorithm: sketch only

In addition to the alignment matrix D, define two further matrices/states.

• 
$$B_{i,j} := \text{cost of best alignment of } a_{1..i}, b_{1..j},$$

Recursions:  

$$A_{i,j} = \min \begin{cases} A_{i-1,j} + \beta & (deletion \ extension) \\ D_{i-1,j} + g(1) & (deletion \ opening) \end{cases}$$

$$B_{i,j} = \min \begin{cases} B_{i,j-1} + \beta & (insertion \ extension) \\ D_{i,j-1} + g(1) & (insertion \ opening) \end{cases}$$

$$D_{ij} = \min \begin{cases} D_{i-1,j-1} + w(a_i, b_j) & (match) \\ A_{i,j} & (deletion \ closing) \\ B_{i,j} & (insertion \ closing) \end{cases}$$

b:

Remark:  $O(n^2)$  time and space

# Similarity

### Definition (Similarity)

The similarity of an alignment  $(a^\diamond, b^\diamond)$  is

$$s(a^\diamond,b^\diamond)=\sum_{i=1}^{|a^\diamond|}s(a^\diamond_i,b^\diamond_i),$$

where  $s: (\Sigma \cup \{-\})^2 \to \mathbb{R}$  is a *similarity function*, where for  $x \in \Sigma$ : s(x, x) > 0, s(x, -) < 0, s(-, x) < 0.

Observation. Instead of minimizing alignment cost, one can maximize similarity:  $S_{ii} = \max \begin{cases} S_{i-1,j-1} + s(a_i, b_j) \\ S_{i-1,j} + s(a_{i,j} - 1) \end{cases}$ 

$$\begin{cases} S_{i-1,j} + s(a_i, -) \\ S_{i,j-1} + s(-, b_j) \end{cases}$$

Motivation:

- defining similarity of 'building blocks' could be more natural, e.g. similarity of amino acids.
- similarity is useful for *local alignment*



# Local Alignment Motivation

Local alignment asks for the best alignment of any two subsequences of *a* and *b*. Important Application: Search! (e.g. BLAST combines heuristics and local alignment)

## Example

- a =AWGVIACAILAGRS
- b = VIVT<u>AIAVAG</u>YY

In contrast, all previous methods compute "global alignments". Why is distance not useful?

### Example

a)  $\begin{array}{c} XXXAAXXXX \\ YYAAYY \end{array}$  b)  $\begin{array}{c} XXAAAAAXXXX \\ YYYAAAAAYY \end{array}$ 

Where is the stronger local motif? Only similarity can distinguish.

# Local Alignment

### Definition (Local Alignment Problem)

Let s be a similarity on alignments.

$$egin{alignment} S_{ ext{global}}(a,b) &:= \max_{\substack{(a^\diamond,b^\diamond) \ alignment ext{ of } a ext{ and } b}} s(a^\diamond,b^\diamond) & (global ext{ similarity}) \ S_{ ext{local}}(a,b) &:= \max_{\substack{1 \leq i' < i \leq n \ 1 \leq j' < j \leq m}} S_{ ext{global}}(a_{i'..i},b_{j'..j}) & (local ext{ similarity}) \end{split}$$

The local alignment problem is to compute  $S_{local}(a, b)$ .

- That is, local alignment asks for the subsequences of *a* and *b* that have the best alignment.
- How would we define the local alignment matrix for DP?
- For example, why does " $H_{i,j} := S_{\text{local}}(a_{1..i}, b_{1..j})$ " not work?

# Local Alignment Matrix

# Definition The *local alignment matrix* H of a and b is $(H_{i,j})_{0 \le i \le n, 0 \le j \le m}$ defined by

$$H_{i,j} := \max_{0 \le i' \le i, 0 \le j' \le j} S_{global}(a_{i'+1..i}, b_{j'+1..j}).$$

• 
$$S_{\text{local}}(a,b) = \max_{i,j} H_{i,j}$$
 (!)

- all entries  $H_{i,j} \ge 0$ , since  $S_{global}(\epsilon, \epsilon) = 0$ .
- $H_{i,j} = 0$  implies no subsequences of *a* and *b* that end in respective *i* and *j* are similar.
- Allows case distinction / Principle of optimality holds!



# Local Alignment Algorithm — Case Distinction

Cases for 
$$H_{i,j}$$
  
1.)  $\dots \begin{vmatrix} a_i \\ b_i \end{vmatrix}$  2.)  $\dots \begin{vmatrix} a_i \\ - \end{vmatrix}$  3.)  $\dots \begin{vmatrix} - \\ b_j \end{vmatrix}$ 

4.) 0, since if each of the above cases is dissimilar (i.e. negative), there is still  $(\epsilon, \epsilon)$ .

# Local Alignment Algorithm (Smith-Waterman Algorithm)

### Theorem

For the local alignment matrix H of a and b,

- $H_{0,0} = 0$
- for all  $1 \le i \le n$ :  $H_{i,0} = 0$

• for all 
$$1 \le j \le m$$
:  $H_{0,j} = 0$   
•  $H_{ij} = \max \begin{cases} 0 & (empty \ alignment) \\ H_{i-1,j-1} + s(a_i, b_j) \\ H_{i-1,j} + s(a_i, -) \\ H_{i,j-1} + s(-, b_j) \end{cases}$ 



# Local Alignment Remarks

- Complexity  $O(n^2)$  time and space, again space complexity can be improved
- Requires that similarity function is centered around zero, i.e. positive = similar, negative = dissimilar.
- Extension to affine gap cost works
- Traceback?



# Local Alignment Example

### Example

• 
$$a = AAC, b = ACAA$$
  
•  $s(x, y) = \begin{cases} 2 & \text{iff } x = y \\ -3 & \text{otherwise} \end{cases}$ 



Traceback: start at maximum entry, trace back to first 0 entry

# Substitution/Similarity Matrices

- In practice: use similarity matrices learned from closely related sequences or multiple alignments
- PAM (Percent Accepted Mutations) for proteins
- BLOSUM (BLOcks of Amino Acid SUbstitution) for proteins
- RIBOSUM for RNA
- Scores are (scaled) log odd scores: log <u>Pr[x,y|Related]</u>

### For example, BLOSUM62:

Ala	4																				
Arg	- 1	5																			
Asn	- 2	0	6																		
Asp	- 2	- 2	1	б																	
Cys	0	- 3	- 3	- 3	9																
Gln	- 1	1	0	0	- 3	5															
Glu	- 1	0	0	2	- 4	2	5														
Gly	0	- 2	0	- 1	- 3	- 2	- 2	б													
His	- 2	0	1	- 1	- 3	0	0	- 2	8												
lle	- 1	- 3	- 3	- 3	- 1	- 3	- 3	- 4	- 3	4											
Leu	- 1	- 2	- 3	- 4	- 1	- 2	- 3	- 4	- 3	2	4										
Lys	- 1	2	0	- 1	- 3	1	1	- 2	- 1	- 3	- 2	5									
Met	- 1	- 1	- 2	- 3	- 1	0	- 2	- 3	- 2	1	2	- 1	5								
Phe	- 2	- 3	- 3	- 3	- 2	- 3	- 3	- 3	- 1	0	0	- 3	0	б							
Pro	- 1	- 2	- 2	- 1	- 3	- 1	- 1	- 2	- 2	- 3	- 3	- 1	- 2	- 4	7						
Ser	1	- 1	1	0	- 1	0	0	0	- 1	- 2	- 2	0	- 1	- 2	- 1	4					
Thr	0	- 1	0	- 1	- 1	- 1	- 1	- 2	- 2	- 1	- 1	- 1	- 1	- 2	- 1	1	5				
Trp	- 3	- 3	- 4	- 4	- 2	- 2	- 3	- 2	- 2	- 3	- 2	- 3	- 1	1	- 4	- 3	- 2	11			
Tyr	- 2	- 2	- 2	- 3	- 2	- 1	- 2	- 3	2	- 1	- 1	- 2	- 1	3	- 3	- 2	- 2	2	7		
Val	0	- 3	- 3	- 3	- 1	- 2	- 2	- 3	- 3	3	1	- 2	1	- 1	- 2	- 2	0	- 3	- 1	4	
	Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	lle	Leu	Lys	Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val	



# Multiple Alignment



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

$$S(A) = s \begin{pmatrix} A \\ A \\ T \end{pmatrix} + s \begin{pmatrix} C \\ C \\ C \end{pmatrix} + s \begin{pmatrix} C \\ - \\ C \end{pmatrix} + s \begin{pmatrix} C \\ - \\ - \end{pmatrix} + \dots + s \begin{pmatrix} - \\ C \\ G \end{pmatrix}$$

How do we know similarities?



# How to Score Multiple Alignments

As for pairwise alignment:

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

$$S(A) = \sum_{j=1}^{N} s inom{A_{1j}}{ \cdots \ A_{\kappa_j}}$$

Example  

$$S(A) = s \begin{pmatrix} A \\ A \\ T \end{pmatrix} + s \begin{pmatrix} C \\ C \\ C \end{pmatrix} + s \begin{pmatrix} C \\ - \\ C \end{pmatrix} + s \begin{pmatrix} C \\ - \\ - \end{pmatrix} + \dots + s \begin{pmatrix} - \\ C \\ G \end{pmatrix}$$

How to define 
$$s \begin{pmatrix} x \\ y \\ z \end{pmatrix}$$
? as log odds  $s \begin{pmatrix} x \\ y \\ z \end{pmatrix} = log \frac{Pr[x,y,z| \text{ Related}]}{Pr[x,y,z| \text{ Background}]}$ ?  
Problems? Can we learn similarities for triples, 4-tuples, ...?

# Sum-Of-Pairs Score

Idea: approximate column scores by pairwise scores

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

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

(Extensible to gap penalties, in particular affine gap cost)

Drawbacks?

# **Optimal Multiple Alignment**

Idea: use dynamic programming

Example

For 3 sequences a, b, c, use 3-dimensional matrix (after initialization:)

$$S_{i,j,k} = \max \begin{cases} S_{i-1,j-1,k-1} + s(a_i, b_j, c_k) \\ S_{i-1,j-1,k} + s(a_i, b_j, -) \\ S_{i-1,j,k-1} + s(a_i, -, c_k) \\ S_{i,j-1,k-1} + s(-, b_j, c_k) \\ S_{i-1,j,k} + s(a_i, -, -) \\ S_{i,j-1,k} + s(-, b_j, -) \\ S_{i,j,k-1} + s(-, -, c_k) \end{cases}$$

For K sequences use K-dimensional matrix. *Complexity*?

S.Will, 18.417, Fall 2011

# Heuristic Multiple Alignment: Progressive Alignment

Idea: compute optimal alignments only pairwise

### Example

4 sequences  $a^{(1)}, a^{(2)}, a^{(3)}, a^{(4)}$ 

- 1. determine how they are related  $\Rightarrow$  tree, e.g.  $((a^{(1)}, a^{(2)}), (a^{(3)}, a^{(4)}))$
- 2. align most closely related sequences first  $\Rightarrow$  (optimally) align  $a^{(1)}$  and  $a^{(2)}$  by DP
- 3. go on  $\Rightarrow$  (optimally) align  $a^{(3)}$  and  $a^{(4)}$  by DP
- go on?! ⇒ (optimally) align the two alignments How can we do that?
- 5. Done. We produced a multiple alignment of  $a^{(1)}, a^{(2)}, a^{(3)}, a^{(4)}$ .

Remarks: - Optimality is not guaranteed. Why?

- The tree is known as guide tree. How can we get it?



# Guide tree

The guide tree determines the order of pairwise alignments in the progressive alignment scheme.

The order of the progressive alignment steps is crucial for quality!

Heuristics:

- 1. Compute pairwise distances between all input sequences
  - align all against all
  - in case, transform similarities to distances (e.g. Feng-Doolittle)
- 2. Cluster sequences by their distances, e.g. by
  - Unweighted Pair Group Method (UPGMA)
  - Neighbor Joining (NJ)



# Aligning Alignments

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.

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

• 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"

# Progressive Alignment — Example 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 7	TTGG								
		Т	Т	G	G						
	0	2	4	6	8						
Α	2	3	5	7	9						
С	4	5	6	8	10						
С	6	7	8	9	11						
G	8	9	10	8	9						
Align	Align ACCG and CTGG										
		С	Т	G	G						
	0	2	4	6	8						
Α	2	3	5	7	9						
С	4	2	5	8	10						
С	6	4	5	8	11						
G	8	7	7	5	8						
Align	TTGG	and	стбб								
		С	Т	G	G						
	0	2	4	6	8						
Т	2	3	2	5	8						
т	4	5	3	5	8						
G	6	7	6	3	5						
G	8	9	9	6	3						

A	ugn	ACCG	and	ICG		
			Т	С	G	
		0	2	4	6	
	А	2	3	5	7	
	С	4	5	3	6	
	С	6	7	5	6	
	G	8	9	8	5	
А	lign	TTGG	and	TCG		
	0		т	С	G	
		0	2	4	6	
	Т	2	0	3	6	
	Т	4	2	3	6	
	G	6	5	5	3	
	G	8	8	8	5	
A	lign	TCG a	ind C	TGG		
			С	т	G	G
		0	2	4	6	8
	Т	2	3	2	5	8
	С	4	2	5	5	8
	G	6	5	5	5	5



# Progressive Alignment — Example

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  $\Rightarrow$  distance matrix

$$\begin{array}{ccccccc} a^{(1)} & a^{(2)} & a^{(3)} & a^{(4)} \\ a^{(1)} & 0 & 9 & 5 & 8 \\ a^{(2)} & 0 & 5 & 3 \\ a^{(3)} & & 0 & 5 \\ a^{(4)} & & & 0 \end{array}$$

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

### Progressive Alignment — Example

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  $a^{(4)}$ :  $\operatorname{TTGG}_{\operatorname{CTGG}}$ , Align  $a^{(1)}$  and  $a^{(3)}$ :  $\operatorname{ACCG}_{\operatorname{-TCG}}$
- Align the alignments!

TTGG CTGG	and	ACC -TC	G G			
	A	Ę	C	G G		
0	4	12	20	28		
8	10					
16		·				
24						
32						
	TTGG CTGG 0 8 16 24 32	TTGG CTGG         and           A         -           0         4           8         10           16         :           24         32	TTGG CTGG         and         ACC -TC           A         C           0         4         12           8         10            16             24	$\begin{array}{c} TTGG\\ CTGG \end{array}  and  \begin{array}{c} ACCG\\ -TCG \end{array} \\ \begin{array}{c} A \\ - \end{array}  \begin{array}{c} C \\ C \\ - \end{array}  \begin{array}{c} C \\ C \\ C \\ - \end{array} \\ \begin{array}{c} C \\ C $		



Progressive Alignment — Example 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  $a^{(4)}$ :  ${}^{\text{TTGG}}_{\text{CTGG}}$ , Align  $a^{(1)}$  and  $a^{(3)}$ :  ${}^{\text{ACCG}}_{\text{-TCG}}$
- Align the alignments!

Align	TTGG CTGG	and	ACC -TC	CG CG		
		A	C T	C C	G	
	0	4	12	20	28	
ΤС	8	10				
ΤT	16		·			
GG	24					
GG	32					

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

• 
$$w(--, A) =$$
  
 $w(-, A) + w(-, -) + w(-, A) + w(-, -) = 4$ 

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

• 
$$w(TC, CT) =$$
  
 $w(T, C) + w(C, C) + w(T, T) + w(C, T) = 6$ 

•

### Progressive Alignment — Example

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  $a^{(4)}$ :  $\operatorname{TTGG}_{\operatorname{CTGG}}$ , Align  $a^{(1)}$  and  $a^{(3)}$ :  $\operatorname{ACCG}_{\operatorname{-TCG}}$
- Align the alignments!

A

Align	CTGG	and	-TC	Ğ			
		A	С Т	C	G G	$\implies$	TTGG
тс	0	4 10	12	20	28	after filling	ACCG -TCG
TT	16	:	·			and traceback	
GG	24						
GG	32						



# A Classical Approach: CLUSTAL W



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



# Advanced Progressive Alignment in MUSCLE



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

# Consistency-based scoring in T-Coffee



- Progressive alignment + Consistency heuristic
- Avoid mistakes when optimizing locally by modifying the scores *"Library extension"*
- Modified scores reflect global consistency
- Details of consistency transformation: next slide
- Merges local and global alignments

# Consistency-based scoring in T-Coffee

#### Misalignment by standard procedure



#### Correct alignment after library extension



### Consistency Transformation

- For each sequence triplet: strengthen compatible edges
- This moves global information into scores
- Consistency-based scores guide pairwise alignments towards (global) consistency

### All-2-all alignments for weighting

SeqA SeqB	GARFIELD GARFIELD	THE THE	LAST FAST	FAT CAT	CAT	Prim. Weight = 88	SeqB SeqC	GARFIELD GARFIELD	THE THE	VERY	FAST FAST	CAT	Prim Weight = 100
SeqA SeqC	GARFIELD GARFIELD	THE THE	LAST VERY	FA-1 FAS1	r cat Cat	Prim. Weight = 77	SeqB SeqD	GARFIELD	THE THE	FAST FA-T	CAT CAT		Prim. Weight = 100
SegA SegD	GARFIELD	THE THE	LAST	FAT FAT	CAT CAT	Prim. Weight =100	SeqC SeqD	GARFIELD	THE THE	VERY	FAST FA-T	CAT CAT	Prim. Weight = 100

# **Alignment Profiles**



ACCG-

- A *profile* of a multiple alignment consists of character frequency vectors for each column.
- The profile describes sequences of the alignment in a rigid way.
- Modeling insertions/deletions requires profile HMMs.

# Hidden Markov Models (HMMs)



[The frog climbs the ladder more likely when the sun shines. Assume that the weather is hidden, but we can observe the frog.]

- *Idea:* the probability of an observation depends on a hidden state, where there are specific probabilities to change states.
- Hidden Markov Models generate observation sequences (e.g. TBTTT) according to an (encoded) probability distribution.
- One can compute things like "most probable path given an observation sequence", ... (no details here)

# Profile HMMs

- Profile HMMs describe (probability distribution of) sequences in a multiple alignment (observation ≡ sequence).
- hidden states = insertion (I<sub>i</sub>), match (M<sub>i</sub>), deletion (D<sub>i</sub>) in relation to consensus (state sequence ≡ alignment string)



- Profile HMMs are used to search for sequences that are similar to sequences of a given alignment (Pfam, HMMer)
- Profile HMMs can be used to construct multiple alignments
- We come back to HMMs when we discuss SCFGs.