



UNIVERSITÄT
DES
SAARLANDES



ZBI

ZENTRUM FÜR
BIOINFORMATIK

Error Tolerant Pattern Matching II

Algorithms for Sequence Analysis

Sven Rahmann

Summer 2021

Overview

Previous Lectures

- Distance and similarity measures between two sequences
- Error-tolerant pattern search (edit distance) in a text:
Algorithms: Basic DP, Ukkonen, Myers, NFA-Shift-And, NFA-FM, Four Russians.
- Alignments as visualization of edit process (global, semiglobal)

Overview

Previous Lectures

- Distance and similarity measures between two sequences
- Error-tolerant pattern search (edit distance) in a text:
Algorithms: Basic DP, Ukkonen, Myers, NFA-Shift-And, NFA-FM, Four Russians.
- Alignments as visualization of edit process (global, semiglobal)

Today

- From costs (distances) to scores (similarities)
- General scoring schemes
- More general introduction of alignments
- Four variants of alignments:
(1) global, (2) semiglobal (pattern search),
(3) free end gaps (overlap detection), (4) local (regions of similarity)

Scoring Schemes for Pairwise Sequence Comparison

Need for fine-grained similarity

- Comparison of biosequences (esp. protein sequences) needs a fine-grained notion of similarity instead of only “equal” vs. “not equal” amino acids.
- **Example:** Leucine (L) and Isoleucine (I) are physically and chemically similar. Tryptophan (Y) has very different properties than most other amino acids.

Scoring Schemes for Pairwise Sequence Comparison

Need for fine-grained similarity

- Comparison of biosequences (esp. protein sequences) needs a fine-grained notion of similarity instead of only “equal” vs. “not equal” amino acids.
- **Example:** Leucine (L) and Isoleucine (I) are physically and chemically similar. Tryptophan (Y) has very different properties than most other amino acids.

Change of paradigm: Zero-centered similarity

- Evaluate similarity (positive and negative) instead of distances (non-negative)
- Value of 0 means “neutral”, positive means “similar”, negative means “dissimilar”.

Scoring Schemes for Pairwise Sequence Comparison

Need for fine-grained similarity

- Comparison of biosequences (esp. protein sequences) needs a fine-grained notion of similarity instead of only “equal” vs. “not equal” amino acids.
- **Example:** Leucine (L) and Isoleucine (I) are physically and chemically similar. Tryptophan (Y) has very different properties than most other amino acids.

Change of paradigm: Zero-centered similarity

- Evaluate similarity (positive and negative) instead of distances (non-negative)
- Value of 0 means “neutral”, positive means “similar”, negative means “dissimilar”.
- **Therefore:** Use a general **score matrix** $s = s(a, b)$ for any $a, b \in \Sigma$, and (negative) similarity values (**gap scores**) for insertions and deletions.

Example: BLOSUM62 Scoring Matrix for Amino Acids

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	4	-3	-1	-1	-3	-2	0	1	-3	-2	-3	-3	-2	-5	1	1	1	-7	-4	0
R	-3	7	-2	-4	-5	1	-3	-5	1	-3	-5	2	-1	-6	-1	-1	-3	1	-6	-4
N	-1	-2	5	3	-5	-1	1	-1	2	-3	-4	1	-4	-5	-2	1	0	-5	-2	-3
D	-1	-4	3	5	-7	0	4	-1	-1	-4	-6	-1	-5	-8	-3	-1	-2	-9	-6	-4
C	-3	-5	-5	-7	9	-8	-8	-5	-4	-3	-8	-8	-7	-7	-4	-1	-4	-9	-1	-3
Q	-2	1	-1	0	-8	6	2	-3	3	-4	-2	0	-2	-7	-1	-2	-2	-7	-6	-3
E	0	-3	1	4	-8	2	5	-1	-1	-3	-5	-1	-4	-8	-2	-1	-2	-9	-5	-3
G	1	-5	-1	-1	-5	-3	-1	5	-4	-5	-6	-3	-4	-6	-2	0	-2	-9	-7	-3
H	-3	1	2	-1	-4	3	-1	-4	7	-4	-3	-2	-4	-3	-1	-2	-3	-4	-1	-3
I	-2	-3	-3	-4	-3	-4	-3	-5	-4	6	1	-3	1	0	-4	-3	0	-7	-3	3
L	-3	-5	-4	-6	-8	-2	-5	-6	-3	1	6	-4	3	0	-4	-4	-3	-3	-3	0
K	-3	2	1	-1	-8	0	-1	-3	-2	-3	-4	5	0	-7	-3	-1	-1	-6	-6	-4
M	-2	-1	-4	-5	-7	-2	-4	-4	-4	1	3	0	9	-1	-4	-3	-1	-6	-5	1
F	-5	-6	-5	-8	-7	-7	-8	-6	-3	0	0	-7	-1	8	-6	-4	-5	-1	4	-3
P	1	-1	-2	-3	-4	-1	-2	-2	-1	-4	-4	-3	-4	-6	7	0	-1	-7	-7	-3
S	1	-1	1	-1	-1	-2	-1	0	-2	-3	-4	-1	-3	-4	0	4	2	-3	-4	-2
T	1	-3	0	-2	-4	-2	-2	-2	-3	0	-3	-1	-1	-5	-1	2	5	-7	-4	0
W	-7	1	-5	-9	-9	-7	-9	-9	-4	-7	-3	-6	-6	-1	-7	-3	-7	12	-2	-9
Y	-4	-6	-2	-6	-1	-6	-5	-7	-1	-3	-3	-6	-5	4	-7	-4	-4	-2	9	-4
V	0	-4	-3	-4	-3	-3	-3	-3	-3	3	0	-4	1	-3	-3	-2	0	-9	-4	5

Reminder: Alignments

Definition (Alignment, Projections π_1, π_2)

An **alignment** is a string A over the **alignment alphabet** $(\Sigma \cup \{-\})^2 \setminus \{(-, -)\}$ (pairs of characters, or one character paired with a gap).

The **first (second) projection** π_1 (π_2) reads the first (second) elements without gaps, so π_1 is the string homomorphism with $\pi_1((a, b)) := a$ and $\pi_1((- , b)) := \epsilon$, etc.

Reminder: Alignments

Definition (Alignment, Projections π_1, π_2)

An **alignment** is a string A over the **alignment alphabet** $(\Sigma \cup \{-\})^2 \setminus \{(-, -)\}$ (pairs of characters, or one character paired with a gap).

The **first (second) projection** π_1 (π_2) reads the first (second) elements without gaps, so π_1 is the string homomorphism with $\pi_1((a, b)) := a$ and $\pi_1((- , b)) := \epsilon$, etc.

Definition (Global alignment)

A **global alignment** between $s, t \in \Sigma^*$ is an alignment with $\pi_1(A) = s$, $\pi_2(A) = t$.

Reminder: Alignments

Definition (Alignment, Projections π_1, π_2)

An **alignment** is a string A over the **alignment alphabet** $(\Sigma \cup \{-\})^2 \setminus \{(-, -)\}$ (pairs of characters, or one character paired with a gap).

The **first (second) projection** π_1 (π_2) reads the first (second) elements without gaps, so π_1 is the string homomorphism with $\pi_1((a, b)) := a$ and $\pi_1((- , b)) := \epsilon$, etc.

Definition (Global alignment)

A **global alignment** between $s, t \in \Sigma^*$ is an alignment with $\pi_1(A) = s$, $\pi_2(A) = t$.

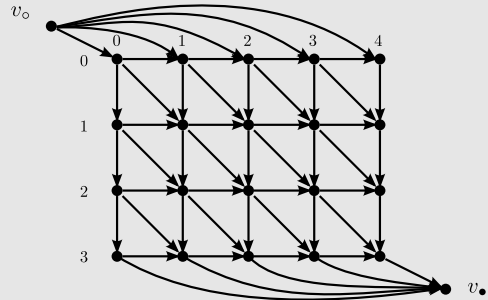
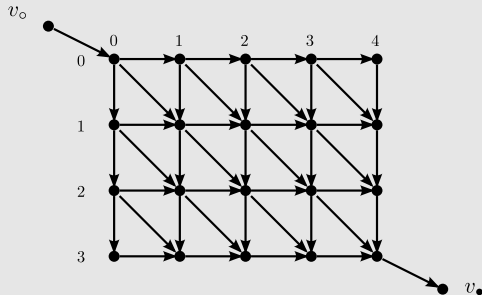
Definition (Semiglobal alignment)

A **semiglobal alignment** between $P, T \in \Sigma^*$ is an alignment with $\pi_1(A) = P$, $\pi_2(A) = T'$, where T' is any substring of T .

Universal Alignment Algorithm

Given

- Sequences s, t
- Scoring scheme
- Alignment graph topology (e.g., for global or semiglobal alignment)



Universal Alignment Algorithm

- **Given:** sequences s, t , scoring scheme, graph topology
- **Sought:**
 - Maximum score among all paths $v_o \rightarrow v_\bullet$ (**optimal alignment score**)
 - A path that maximizes the scores (**optimal alignment**)
- Let $S(v)$ be the maximal score of all paths $v_o \rightarrow v$, and $S(v_o) := 0$.
- Let $T(v)$ be the predecessor of v , from which the maximum $S(v)$ is obtained.

Universal Alignment Algorithm

- **Given:** sequences s, t , scoring scheme, graph topology
- **Sought:**
 - Maximum score among all paths $v_o \rightarrow v_\bullet$ (**optimal alignment score**)
 - A path that maximizes the scores (**optimal alignment**)
- Let $S(v)$ be the maximal score of all paths $v_o \rightarrow v$, and $S(v_o) := 0$.
- Let $T(v)$ be the predecessor of v , from which the maximum $S(v)$ is obtained.
- For $v \neq v_o$:
$$S(v) = \max_{w: w \rightarrow v \in E} \{S(w) + \text{score}(w \rightarrow v)\},$$
$$T(v) = \arg \max_{w: w \rightarrow v \in E} \{S(w) + \text{score}(w \rightarrow v)\}.$$
- Compute nodes in topological order (graph is acyclic!)
- The optimal score is obtained as $S(v_\bullet)$.
- The optimal path (alignment) is obtained by traceback from v_\bullet :
$$v_\bullet \rightarrow T(v_\bullet) \rightarrow T(T(v_\bullet)) \rightarrow \dots \rightarrow T^k(v_\bullet) \rightarrow \dots \rightarrow v_o.$$

Traceback

Traceback, also Backtracing

- Reconstruction of the optimal path by tracing back the predecessor nodes that lead to the optimal score value in each node
- Do not confuse with Backtracking!

Traceback

Traceback, also Backtracing

- Reconstruction of the optimal path by tracing back the predecessor nodes that lead to the optimal score value in each node
- Do not confuse with Backtracking!

Optimal path and alignment

- Optimal path is reconstructed backwards, can be flipped when done
- Optimal alignment:
Read off the edge labels along the optimal path

Traceback

Traceback, also Backtracing

- Reconstruction of the optimal path by tracing back the predecessor nodes that lead to the optimal score value in each node
- Do not confuse with Backtracking!

Optimal path and alignment

- Optimal path is reconstructed backwards, can be flipped when done
- Optimal alignment:
Read off the edge labels along the optimal path

Time and memory requirements

- **Running time:** $O(m + n)$ for an $m \times n$ matrix (maximum length of a path)
- **Memory:** $O(mn)$ because the full matrix T must be stored (improvement soon)

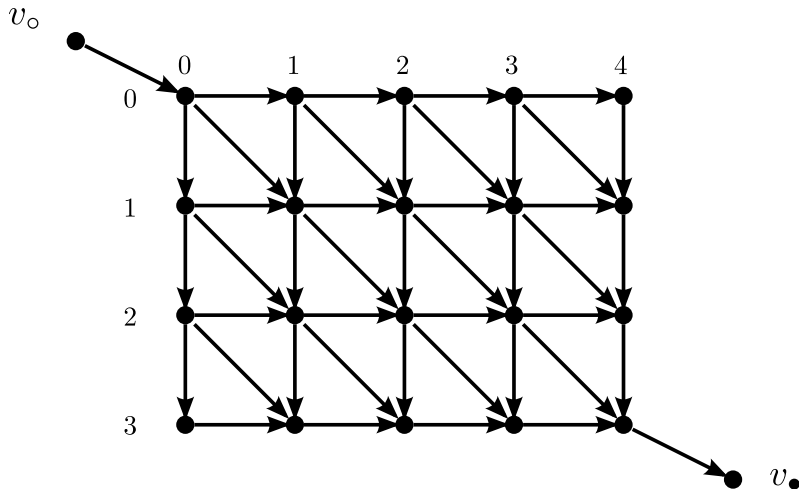
Variants of Alignments

Four Variants

- 1 Global alignment (similarity of full sequences)
- 2 Semiglobal alignment (pattern search)
- 3 Free end gaps alignment (good/optimal overlap)
- 4 local alignment (region[s] of high/optimal similarity)

In the following, we discuss the associated **graph topology** for each variant.
All variants can be handled uniformly with the **universal alignment algorithm**.

Global Alignment



Global Alignment

Definition (global alignment graph)

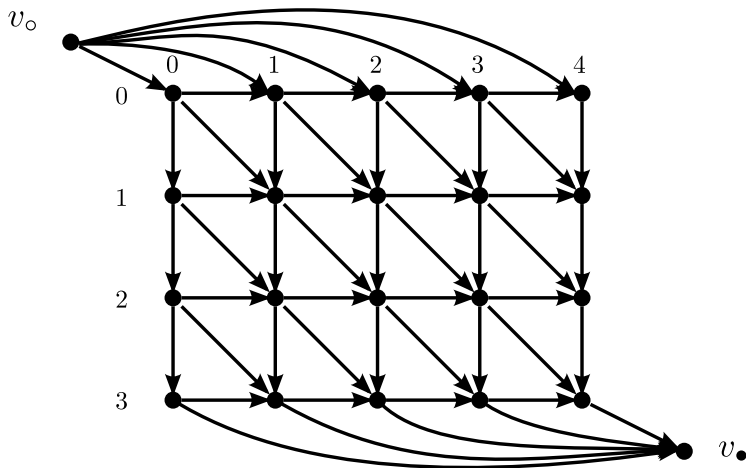
- Nodes $V := \{(i, j) : 0 \leq i \leq m, 0 \leq j \leq n\} \cup \{v_o, v_\bullet\}$
- Edges:

	Edge	label	score
horizontal	$(i, j) \rightarrow (i, j + 1)$	$\begin{bmatrix} - \\ t_j \end{bmatrix}$	$< 0 (*)$
vertical	$(i, j) \rightarrow (i + 1, j)$	$\begin{bmatrix} s_i \\ - \end{bmatrix}$	$< 0 (*)$
diagonal	$(i, j) \rightarrow (i + 1, j + 1)$	$\begin{bmatrix} s_i \\ t_j \end{bmatrix}$	beliebig (*)
Initialization	$v_o \rightarrow (0, 0)$	ϵ	0
Finalization	$(m, n) \rightarrow v_\bullet$	ϵ	0

(*): Meaningful scoring schemes have negative scores for gaps and most substitutions, and positive scores for identities.

Semiglobal Alignment (Pattern Search)

Additional initialization edges $v_o \rightarrow (0,j)$ and finalization edges $(m,j) \rightarrow v_\bullet$:



“Free End Gaps” Alignment (Overlap Detection)

Question

- (How) Do two sequences overlap?



- Gaps (overhangs) at either border of either sequence shall not be penalized.

“Free End Gaps” Alignment (Overlap Detection)

Question

- (How) Do two sequences overlap?



- Gaps (overhangs) at either border of either sequence shall not be penalized.

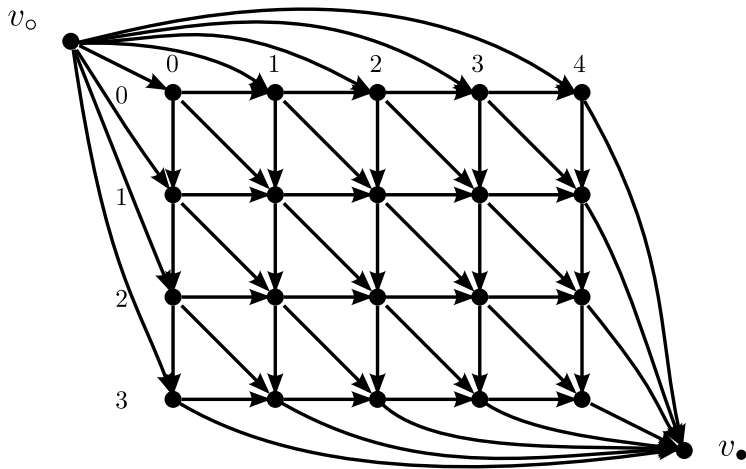
Graph construction

Additional initialization edges $v_o \rightarrow (i, 0)$ and $v_o \rightarrow (0, j)$,

and finalization edges $(i, n) \rightarrow v_\bullet$ and $(m, j) \rightarrow v_\bullet$.

(All such edges have empty labels and contribute score 0.)

“Free End Gaps” Alignment (Overlap Detection)



Local Alignment

Question

- (Where) Are there regions (substrings) of high similarity between two sequences?
- Where are the most similar substrings (maximal score) ?



- Formally: Find alignment with maximal score among all substrings s' of s and t' of t .

Local Alignment

Question

- (Where) Are there regions (substrings) of high similarity between two sequences?
- Where are the most similar substrings (maximal score) ?



- Formally: Find alignment with maximal score among all substrings s' of s and t' of t .

Graph construction

- Initialization edges $v_o \rightarrow (i, j)$ for all $0 \leq i \leq m, 0 \leq j \leq n$,
- Finalization edges $(i, j) \rightarrow v_\bullet$ for all $0 \leq i \leq m, 0 \leq j \leq n$.
- Visualization is not helpful (far too many edges)

Variants and Distances vs. Scores

Meaningful combinations

Variant	Distances	Scores
Global alignment (similarity of full sequences)	✓	✓
Semiglobal alignment (pattern search)	✓	✓
Free end gap alignment (good/optimal overlap)		✓
Local alignment (region[s] of high/optimal similarity)		✓

Variants and Distances vs. Scores

Meaningful combinations

Variant	Distances	Scores
Global alignment (similarity of full sequences)	✓	✓
Semiglobal alignment (pattern search)	✓	✓
Free end gap alignment (good/optimal overlap)		✓
Local alignment (region[s] of high/optimal similarity)		✓

Why?

Optimal distance is always zero ($d \geq 0$).

Free end gap and local alignments allow trivial “empty” alignments, which always have distance zero. No incentive for non-trivial alignments.

Specialization of Algorithms

For each alignment variant (graph topology):

- What is the interpretation of the score $S(v)$ for any $v = (i, j)$?
("Score of an optimal alignment of ...")
- How does the universal algorithm specialize to matrix form ?
First row and column ?
 $S[i, j] = \max\{\dots\}$?
Collection of interesting results or optimal result ?
- How do running time and memory requirements change vs. global alignment?
(They don't.)

Specialization of Algorithms

For each alignment variant (graph topology):

- What is the interpretation of the score $S(v)$ for any $v = (i, j)$?
("Score of an optimal alignment of ...")
- How does the universal algorithm specialize to matrix form ?
First row and column ?
 $S[i, j] = \max\{ \dots \}$?
Collection of interesting results or optimal result ?
- How do running time and memory requirements change vs. global alignment?
(They don't.)

Algorithm Names

- Global alignment: Needleman-Wunsch algorithm (NW)
- Local alignment: Smith-Waterman alignment (SW)

Needleman-Wunsch Algorithm (Score-Based, Global, Full Matrix)

```
1 def needleman_wunsch(s, t, score):
2     m, n, gapscore = len(s), len(t), score(None)
3     S = np.zeros((m+1, n+1), dtype=np.int32) # scores
4     T = np.zeros((m+1, n+1), dtype=np.uint8) # traceback
5     S[0,:] = np.arange(n+1, dtype=S.dtype) * gapscore
6     S[:,0] = np.arange(m+1, dtype=S.dtype) * gapscore
7     T[0,0] = HOME; T[0,1:] = HORIZONTAL; T[1:,0] = VERTICAL
8     for i, si in zip(count(1), s): # (row, character in s)
9         for j, tj in zip(count(1), t): # (col, character in t)
10             d = S[i-1, j-1] + score(si, tj)
11             h = S[i, j-1] + gapscore
12             v = S[i-1, j] + gapscore
13             S[i,j] = opt = max(d, h, v)
14             T[i,j] = (d==opt)*DIAGONAL + (h==opt)*HORIZONTAL \
15                     + (v==opt)*VERTICAL
16     return S[m,n], traceback(m, n, T, s, t)
```

Smith-Waterman Algorithm (Score-Based, Local, Full Matrix)

```
1 def smith_waterman(s, t, score):
2     m, n, gapscore = len(s), len(t), score(None)
3     S = np.zeros((m+1, n+1), dtype=np.int32) # scores
4     T = np.zeros((m+1, n+1), dtype=np.uint8) # traceback
5     T[0,:] = HOME; T[:,0] = HOME # alignments end at border
6     best = (-1, -1, -1) # best (S, i, j)
7     for i, si in zip(count(1), s): # (row, character in s)
8         for j, tj in zip(count(1), t): # (col, character in t)
9             d = S[i-1, j-1] + score(si, tj)
10            h = S[i, j-1] + gapscore
11            v = S[i-1, j] + gapscore
12            S[i,j] = opt = max(0, d, h, v) # note additional 0
13            T[i,j] = (d==opt)*DIAGONAL + (h==opt)*HORIZONTAL \
14                    + (v==opt)*VERTICAL # can be HOME otherwise
15            if S[i,j] > best[0]: best = (S[i,j], i, j)
16     result, i, j = best
17     return result, traceback(i, j, T, s, t)
```

Implementation of Traceback

```
1 HOME, DIAGONAL, HORIZONTAL, VERTICAL = 0, 1, 2, 4
2 def traceback(i, j, T, s, t, *, GAP='-'):
3     # We reconstruct the alignment by traceback (T) from i, j
4     As, At = [], [] # rows of alignment: As (for s), At (for t)
5     while T[i,j] != HOME:
6         trace = T[i,j]
7         if (trace & DIAGONAL):
8             i -= 1; As.append(s[i])
9             j -= 1; At.append(t[j])
10        elif (trace & HORIZONTAL):
11            As.append(GAP)
12            j -= 1; At.append(t[j])
13        elif (trace & VERTICAL):
14            i -= 1; As.append(s[i])
15            At.append(GAP)
16    # create the final alignment (pair of strings)
17    return ("".join(As[::-1]), "".join(At[::-1]))
```


Summary

- Motivation of scoring schemes vs. cost functions
- Definition of pairwise alignments in general
- Definition of four pairwise sequence alignment variants
- Alignment graphs and four topology variants
- Universal alignment algorithm on graphs
- Universal traceback
- Specialization: Needleman-Wunsch (global)
- Specialization: Smith-Waterman (local)
- Other specializations: Homework

Possible Exam Questions

- Define alignment (in general).
- Define global / semiglobal / etc. alignment of two strings s, t .
- Explain four variants of alignments and their applications / use cases.
- What is the difference between score and cost function and why is it important?
- Why can't we use costs for free end gap and local alignment?
- How can sequence alignment be formulated as a graph problem?
- Show the alignment graph topology for each variant.
- Explain the universal alignment algorithm on the alignment graph.
- Give the DP formulation for computing an alignment score (any variant).
- Compute an optimal alignment (any variant) for two given strings.
- Explain traceback.