

Distribution of Patterns in Pairwise Sequence Alignments



PARIS
DESCARTES

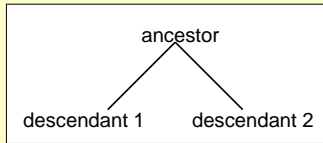
S. Wolfsheimer and G. Nuel

MAP5 – Mathématiques Appliquées à Paris 5, University Paris Descartes, France

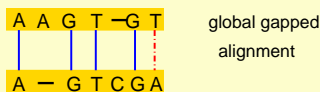
Sequence Alignment [1]

Detection of homological relationships in **DNA-** or **protein-sequences**.

Applications: Search tools for molecular databases, e.g. BLAST



- Given: Pair of sequences
 $\mathbf{a} = a_1 \dots a_M \in \Sigma^M$,
 $\mathbf{b} = b_1 \dots b_N \in \Sigma^N$.
- amino acids: $\Sigma = \{A, C, D, E, \dots\}$
nucleotides: $\Sigma = \{A, G, C, T\}$
- Alignment \mathcal{A}** : set of pairings
 $(i_1, j_1), \dots, (i_n, j_n)$
with $i_k < i_{k+1}, j_k < j_{k+1}$
- Gaps**: insertions and deletions of subsequences



- Alignment score:

$$S(\mathcal{A}; \mathbf{a}, \mathbf{b}) = \sum_{\text{pairs}} \sigma_{a_i, b_j} - \sum_{\text{gaps}} g(l_{\text{gap}})$$
- score matrix: $\sigma_{a,b} = \log \frac{P_{a,b}}{f_a f_b}$
gap penalty: $g(l) = \alpha + \beta(l - 1)$

Score based alignment: Similarity between \mathbf{a} and \mathbf{b} measured by **optimal alignment score** S_0 .

- $S_0 = \max_{\mathcal{A}} S(\mathcal{A}; \mathbf{a}, \mathbf{b})$
- $\mathcal{A}_0 = \operatorname{argmax}_{\mathcal{A}} S(\mathcal{A}; \mathbf{a}, \mathbf{b})$

Dynamic Programming

Global alignment:

Needleman-Wunsch algorithm

- Here: linear gap costs $g(l_{\text{gap}}) = \delta l_{\text{gap}}$
Affine gap-costs straight forward.
- $A_{i,j}$: optimal alignment of subproblem $a_1, \dots, a_i, b_1, \dots, b_j$

$$A_{i,j} = \max \begin{cases} A_{i-1, j-1} + \sigma_{a_i, b_j} \\ A_{i-1, j} - \delta \\ A_{i, j-1} - \delta \end{cases}$$

$$S_0 = A_{M,N}$$

- $\mathcal{O}(M \cdot N)$ time complexity

Alignment biases [3]

Close to gaps often many competitive alignments decrease the accuracy of the score based alignment. Typical effects are:

- Gap wander: a gap shifted by a few positions.



- Gap attraction: two close gaps merge into a single gap.



- Gap annihilation: two gaps of opposite signature (insertion / deletion) cancel each other.



Probabilistic alignment: more quantitative description. Here: distributions of typical patterns that might lead to weakly reliable alignment segments.

Probabilistic alignment [2]

- $P_T(\mathcal{A})$: **distribution of global alignments** of \mathbf{a}, \mathbf{b} .
- Canonical ensemble of alignments
 $P_T(\mathcal{A}; \mathbf{a}, \mathbf{b}) = \frac{1}{Z_T} \exp \left[\frac{1}{T} S(\mathcal{A}; \mathbf{a}, \mathbf{b}) \right]$
(Boltzmann distribution with "temperature" $T = 1$)
- Partition function**:
 $Z_T = \sum_{\mathcal{A}} \exp [S(\mathcal{A})/T]$
- Partition function calculation:
 $D_{i,j} \rightarrow Z_{i,j}, \max \rightarrow \sum$, and $+ \rightarrow \times$
- Forward algorithm**:
 $Z_{i,j}$: sum over all alignments of $a_1 \dots a_i$ and $b_1 \dots b_j$

$$Z_{i,j} = \underbrace{Z_{i-1, j-1} e^{\sigma_{a_i, b_j}/T}}_{\text{match/mismatch}} + \underbrace{(Z_{i-1, j} + Z_{i, j-1}) e^{-\delta/T}}_{\text{gap}}$$

$$Z_T = Z_{M,N}$$

- Backward algorithm**:
 $Z'_{i,j}$: sum over all alignments of $a_{i+1} \dots a_M$ and $b_{j+1} \dots b_N$ given that $(i, j) \in \mathcal{A}$.
- Posterior probabilities**:
 $p_{ij} := P[(i, j) \in \mathcal{A}] = \frac{1}{Z_T} Z_{i,j} Z'_{i,j}$

$$P \begin{pmatrix} a_i \\ - \end{pmatrix} = 1 - \sum_j p_{ij}$$



Number of gaps

Distribution of the number of gaps (see [4] for a general HMM setup).

Insertion $a_{i-1} \dots a_{i-1} a_i \quad a_{i+1} \dots a_{i+l}$
 $b_j \quad - \quad - \quad - \quad - \quad b_{j+l}$

Deletion $a_i \quad - \quad - \quad - \quad - \quad a_{i+1}$
 $b_{j-1} \dots b_{j-1} b_j \quad b_{j+1} \dots b_{j+l}$

Insertions: auxiliary matrices $Z_{i,j}^{(I,k)}$ and $I_{i,j}^{(I,k)}$ for $k = 1, 2, \dots$, sum over alignments, k insertions seen so far.

$Z_{i,j}^{(I,k)}$: sum over $\dots \dots a_i$
 $\dots b_j \quad - \quad -$

$$Z_{i,j}^{(I,k)} = \left(I_{i-1, j-1}^{(I,k)} + Z_{i-1, j-1}^{(I,k)} \right) e^{\sigma_{a_i, b_j}/T} + \left(Z_{i, j-1}^{(I,k)} + I_{i, j-1}^{(I,k)} \right) e^{-\delta/T}$$

$$I_{i,j}^{(I,k)} = \left(Z_{i-1, j}^{(I, k-1)} + I_{i-1, j}^{(I, k-1)} \right) e^{-\delta/T}$$

$$P(k) = \left(Z_{M,N}^{(I,k)} + I_{M,N}^{(I,k)} \right) / Z_{M,N}$$

$$P \begin{pmatrix} a_i \\ - \end{pmatrix}, k\text{th insertion} = \sum_j \frac{I_{i,j}^{(I,k)} I'_{i,j+1}}{Z_{M,N}}$$



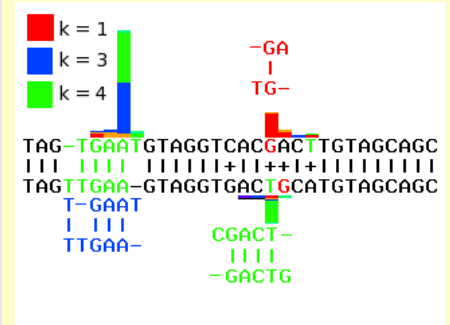
Gap Annihilation Pattern

Detecting candidates for gap annihilation. Search for the pattern:

Insertion $\dots - a_{i-k+1} \dots a_i a_{i+1}$
 $\dots b_{j-k} b_{j-k+1} \dots b_j -$

Deletion $\dots a_{i-k} a_{i-k+1} \dots a_i -$
 $\dots - \underbrace{b_{j-k+1} \dots b_j}_{k} b_{j+1}$

Decode alternative alignment segments:



Bibliography

- [1] Durbin et al., Biological sequence analysis (1998)
- [2] Miyazawa., Prot. Eng. 8 (1995)
- [3] Lunter et al., Genome Research. 18 (2008)
- [4] Aston and Martin, Ann. Appl. Stat. 1 (2007)