

# 1 Scoring Function Overview

The scoring function that I have designed builds upon the suggestions in the assignment brief, and to the best of my knowledge has strong biological intuition. The score function is made up of three key parts.

## 1.1 Convex Gap Scores

It is well-known that single mutation events - such as translocation and duplication - may create alignment gaps of different sizes. Therefore, there is a biological need to treat such gaps as a single entity, as opposed to individually penalising each successive *indel*. Two protein sequences may be relatively similar but differ only at a certain interval, and we don't want to excessively penalise such an alignment. Affine gap scores have been used extensively in industry, and are often the gap score method of choice when partnered with the *BLOSUM* score matrices. An affine gap score combines a constant 'gap-opening' penalty (favouring shorter gaps) with an additional penalty that is linear in the number of further residues in the gap (favouring fewer, larger gaps). Typically, the 'gap-opening' penalty is an order of magnitude larger than the additional penalty ((10, 1) for BLOSUM-62), and thus such a gap-score scheme goes some way to reducing the penalty given to large gaps. However, it has been shown empirically that an affine gap length is too rigid for use in a biological context<sup>1</sup>. Moreover, other studies have shown that the distribution of *indels* typically follows a power law (logarithmic) distribution<sup>2</sup>. Therefore, I have opted to use a convex gap penalty model, in which each additional space in a gap contributes less to the gap weight than the previous space. The model I propose for (internal) gap scores is:

$$P_g \cdot \ln(n)$$

where  $P_g$  is the gap opening penalty, and  $n$  is the length of the gap.

## 1.2 Trailing, Internal and Terminating Gap Score Differentiation

The second key feature of my scoring function is the implementation of different gap-scoring parameters for trailing, internal, and terminating gaps. It has been shown empirically, through optimization techniques, that the optimal model for matches often has opening and terminal gap-open penalties that are approximately half of the gap-open penalty used for internal gaps<sup>3</sup>. Intuitively, this makes sense - the query sequence could merely be a translation of the database sequence, and should not be excessively penalised. Therefore, the model I propose for trailing and terminating gap sequences is:

$$\frac{P_g}{2} \cdot \ln(n)$$

where  $P_g$  is the gap opening penalty, and  $n$  is the length of the gap. It is worth noting that this feature of the scoring function only applies to the global alignment case, as all (biologically-imitating) score functions assign negative scores to gaps, and so leading and trailing gaps would never be included in a local alignment.

## 1.3 Codon Match Reward

A codon is a set of three bases (technically three nucleotides) which codes for a certain amino acid. This sequence of contiguous triplets defines a protein's functionality, and the codons hold the key to the translation of genetic information for the synthesis of proteins. Therefore, I believe matches of (multiples of) three contiguous bases should be rewarded with a score that is higher than the summation of the scores of the individual matches. I propose a multiplicative scoring system, where:

$$\text{Adjusted Score} = C(t) \cdot \text{Additive Score},$$

where  $t$  is the number of codons matched consecutively. This codon scoring scheme could be extended further and the entire alignment scoring could be implemented on the codon level, using the empirical scoring matrices found here<sup>4</sup>.

## 1.4 Further Features

Further features that could be included to more accurately mimic biological reality include:

- Take into account the position within the current codon. Point mutation frequencies are not evenly distributed over the three positions within a codon, so different scoring matrices could be used for each of the three positions
- Have gap-specific indel scores. There is evidence to suggest that specific residue types are preferred in gap regions<sup>5</sup>, and so once a gap has been opened a secondary scoring matrix could be used for the following residues in the gap.
- Score specific mutation events individually. For example, in the case of a duplicated amino acid (a triplet/codon of three bases), the penalty induced for the insertion of three *indels* could be reduced if the three residues are identical to previous triplet of residues: as in the case of AACACGTCG and AACACGACGTCG, for example.

<sup>1</sup>Sung, Wing-Kin (2011). *Algorithms in Bioinformatics : A Practical Introduction*. CRC Press. pp. 42-47

<sup>2</sup>[http://elbo.gs.washington.edu/courses/GS\\_559\\_11\\_wi/slides/4A-Sequence\\_comparison-ScoreMatrices.pdf](http://elbo.gs.washington.edu/courses/GS_559_11_wi/slides/4A-Sequence_comparison-ScoreMatrices.pdf)

<sup>3</sup><https://www.ncbi.nlm.nih.gov/pmc/articles/PMC548345/>

<sup>4</sup><https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-6-134>

<sup>5</sup>Wrabl JO, Grishin NV (2004). "Gaps in structurally similar proteins: towards improvement of multiple sequence alignment"

## 2 Scoring Function Implementation

### 2.1 Parameters

As I have described an array of different features above, I shall implement only the convex gap score penalty, and the codon match reward function. The parameters required are therefore:

- The gap-opening penalty:  $P_g = -10$  (as is often used with some of the BLOSUM matrices)
- The codon match reward function:  $C(t) = \begin{cases} 1 + (0.1)t & 0 \leq t \leq 10 \\ 2 & t > 10 \end{cases}$  where  $t = \lfloor m/3 \rfloor$ , and  $m$  is the number of contiguously matched residues. This function increasingly rewards longer contiguous matches of codons.
- The scoring matrix:  $a = \begin{pmatrix} 1 & -1 & -2 \\ -1 & 2 & -4 \\ -2 & -4 & 3 \end{pmatrix}$ , indexed in the usual way for the alphabet  $\Sigma = \{A, B, C\}$ . In the absence of a real scoring matrix such as PAM or BLOSUM, I shall use the score matrix given in the assignment brief. In reality, the score matrix would be alphabet-specific, and would be found empirically using the probabilistic model discussed in lectures.

### 2.2 Algorithm

---

**Algorithm 1:** Local Sequence Alignment with Convex Gap Penalties and Codon Rewards

---

**Input** : A query sequence  $Q$ , and a database sequence  $D$ , both formed from the alphabet  $\Sigma = \{A, B, C\}$

**Output:** An array of [Alignment Score, The indices of  $Q$ , The indices of  $D$ ]

**Data:**

$\text{score\_matrix} \leftarrow [[1, -1, -2], [-1, 2, -4], [-2, -4, 3]]$

$P_g \leftarrow -4$

$C(t) \leftarrow 1 + (0.1 \cdot t)$  if  $0 \leq t \leq 10$  else 2

```
1 align_matrix ← zero_array(m + 1, n + 1);
2 pointers ← zero_array(m + 1, n + 1);
3 for i, l in enumerate(D) do
4     for j, k in enumerate(Q) do
5         left_gap, a, b ← backtrack(i, j, left);
6         up_gap, c, d ← backtrack(i, j, up);
7         diag_gap ← backtrack(i, j, diag);
8         left_score ← align_matrix[a, b] + (P_g · ln(left_gap + 2));
9         up_score ← align_matrix[c, d] + (P_g · ln(up_gap + 2));
10        diag_score ← align_matrix[i, j] + (score_matrix[l, k] · C( $\lfloor \frac{\text{diag\_gap}}{3} \rfloor$ ));
11        scores ← array(0, left_score, diag_score, up_score);
12        align_matrix[i + 1, j + 1] ← max(scores);
13        pointers[i + 1, j + 1] ← argmax(scores);
14    end
15 end
16 return [max(scores), generate_indices(pointers, argmax(scores))]
```

---