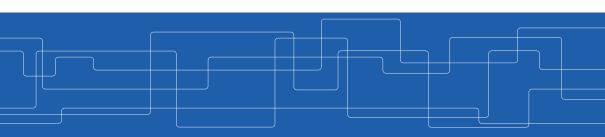


Aligning proteins

Lukas Käll lukask@kth.se





Why aligning protein rather than DNA sequences?

► Amino acid sequences are often more conserved than their underlying DNA. That is synonymous mutations are more more common than non-synonymous ones.



Why aligning protein rather than DNA sequences?

- Amino acid sequences are often more conserved than their underlying DNA. That is synonymous mutations are more more common than non-synonymous ones.
- ► Even non-synonymous mutations are more frequently causing shifts to amino acids with similar properties (polarity, size) than different properties.



Scoring functions for amino acid sequences.

In principle one could use the same type of score functions as for DNA sequences. However, we can create better scoring systems by using *score matrices*, i.e. score functions that are dependent on which amino acids that are evaluated.



There are two major types of scoring matrices:

- ► PAM = Percentage Accepted Mutations (Margeret Dayhoff)
- ▶ BLOSUM = Blocks Substitution Matrix (Henikoff & Henikoff)

PAM

- Created from global alignments, from tertiary structures.
- Better for similar sequences.
- Higher numbers indicates suitability for more diverse sequences.

BLOSUM45 \sim PAM250

BLOSUM

- Created from local alignments, from blocks of similar sequences (the BLOCKS DB)
- ▶ Better for distant sequences.
- Lower numbers indicates suitability for more diverse sequences.



PAM62

```
Ala
Arg
Asn
Asp −2
        -2
Cys
        -3
            -3
                -3
Gln
Glu
                     -4
Gly
His
                     -1
                         -3
Leu
                     -1
Lys
Met
Phe
            -3
                 -3
                    -2
                         -3
                            -3 -3 -1
                                        0
                                               -3
Pro
Ser
Thr
                                -2 -2 -1
                     -2
                         -2
                            -3
                                -2 -2 -3
                                           -2
Trp
                            -2 -3
                                     2 -1
                                           -1
Val
                -3 -1 -2 -2 -3 -3
                                       3
                                               -2
                                                       -1 -2 -2
    Ala Arg Asn Asp Cys Gln Glu Gly His IIe Leu Lys Met Phe Pro Ser Thr Trp Tyr Val
```



What is the probability that one amino acid is replaced by another?

When scoring a position in an alignment containing the amino acid a and b, we take interest in the ratio between the probability that they appear together if they stem from homologue sequences and if they do not stem from homologues.

$$\frac{\Pr(a, b | \text{homologues})}{\Pr(a, b | \text{not homologues})} = \frac{\Pr(a, b)}{\Pr(a) \Pr(b)}.$$



Substitution scores in score matrices

In scoring matrices this property is used in the following form

$$d(a,b) = \frac{1}{\lambda} \log \frac{\Pr(a,b)}{\Pr(a) \Pr(b)}.$$

Here λ is selected in a manner that the d(a,b)'s can be rounded to integer value with as little rounding errors as possible.



The approximate reasoning behind the scores

For the full length sequences we are interested in evaluating

 $\frac{\text{Pr(Sequence alignment given the sequences are homologues)}}{\text{Pr(Sequence alignment given the sequences are not homologues)}} \approx$

$$\approx \frac{\prod_{i} \Pr(\text{align pos } i | \text{homologues})}{\prod_{i} \Pr(\text{align pos } i | \text{not homologues})} \approx \prod_{i} \frac{\Pr(a_{i}, b_{i})}{\Pr(a_{i}) \Pr(b_{i})} =$$

$$= \exp\left(\log\left(\prod_{i} \frac{\Pr(a_{i}, b_{i})}{\Pr(a_{i}) \Pr(b_{i})}\right)\right) = \exp\left(\sum_{i} \log\left(\frac{\Pr(a_{i}, b_{i})}{\Pr(a_{i}) \Pr(b_{i})}\right)\right)$$

This resembles $\exp(\sum_i d(a_i, b_i))$, and hence it makes sense to score alignments based on the sums of $d(a_i, b_i)$.



Thanks!