Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

## Efektywne algorytmy do porównań sekwencji

Bartek Wilczyński

12. marca, 2018

・ロト ・ 同ト ・ ヨト ・ ヨト

э

Sac

### How sequences evolve?

・ロト ・ 同ト ・ ヨト ・ ヨト

= √Q (~

Efektywne algorytmy do porównań sekwencji

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach



image (c) BW

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Efektywne algorytmy do porównań sekwencji

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

- How far in evolution are sequences we can observe in different living species?
- More formally: Can we define a measure of sequence similarity

$$d: \Sigma^* \times \Sigma^* \to \mathcal{R}^+$$

approximating the true evolutionary distance?

• Hint: We should count the number of mutations leading to the observed divergence.

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

- Mutations occur on DNA level, but selection acts much higher: on the phenotype level.
- This makes the assumption of base independence invalid
- Long evolutionary times violate time-reversibility
- Multiplicative measure not too convenient in practice
- We can only account for substitutions, not for insertions or deletions

Suggested solutions:

- Use protein sequences for comparisons
- Define additive substitution matrices

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

- We are still assuming time-reversible Markov chain, but now in space of protein sequences.
- Matrix entries contain log-probabilities, leading to additive measures of similarity
- PAM (Point accepted mutations) matrices (Dayhoff, 1978) describe observed probabilities of occurence of point mutations for a given average divergence (PAM1 = one mutation/100 bases, mostly used PAM250)
- BLOSUM (BLOcks Substitution Matrix) (Henikoff, Henikoff 1992) were constructed using short protein alignments (Blocks) of given sequence identity.
   e.g.BLOSUM80 was derived from sequences of ≥ 80% identity

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

Υ v QEGHILKMFP S Т A R Ν D С -1 -1 0 -2 -1 -2 -1 -1 -3 -1 1 0 - 3 - 2 = 00 -3 0 -4 -3 3 -2 -3 -3 -1 -1 -3 -1 -3 5 - 2 - 1 - 2 - 1А R 1 -3 -4 0 -2 -4 -2 2 -1 -1 -4 -4 -1 -4 -5 -1 0 -1 -5 -3 -4 N 0 -3 -3 -3 -2 -2 -3 -2 -2 -4 -1 -1 -5 -3 -1 D -313 -3 -2 2 0 -4 -1 0 -1 -1 -1 -3 -4 -3 1 -2 -3 -1 -1 -1 -3 -2 -3 -3 0 Е -10 - 2 - 3 - 3 - 4-2 -4 -4 -2 -3 -4 -2 -28 0 - 3-310 -4 -3 0 -1 -1 -2 -1 -2 -3 2 - 40 - 20 Η 0 - 3 - 3 - 1 - 32 - 3 - 4 - 4 - 4 52 -32 Т 5 - 31 - 4 - 3 - 1 - 2 - 1-2 -2 -3 -4 -32 3  $0 -3 -3 \quad 6 -2 -4 -1$ 2 -2-3 - 2 - 3Κ 0 - 2 - 3 - 13 - 2-1 - 2 - 2-22 7 0 - 3 - 2 - 1-4 8 - 4 - 3 - 2-4 - 3 - 4 - 10 1 -1-1 -2 -2 -3 -4 -1 -3 -4 **10** -1 $0 - 1 - 3 - 3 \quad 0 - 2 - 3 - 1$ S 1 - 15 -2 -2-1 -1 -2 -2 -1 -1 -1 -1 -2 -1 т 0 - 10 - 15 - 3 - 20 -3 -3 -4 -5 -5 -1 -3 -3 -3 -3 -2 -3 -1 1 -4 -4 -3 15 W 2 - 3-2 - 1 - 2 - 3 - 3 - 1 - 2 - 3 - 2 - 1 - 1 - 2 - 0 - 4 - 3 - 2 - 2Y 2 8 - 1 $\overline{0}$  -3 -3 -4 -1 -3 -3 -4 -4 4 1 -3 1 -1 -3 -2  $\overline{0}$  -3 -1 77 5

Log-odds log  $\frac{P_{x,y}}{Q_x Q_y}$  instead of probabilities or substitution rates. <sup>image</sup> (c) Durbin et al.

### Blosum50 matrix

▲ロト ▲御 ト ▲ 臣 ト ▲ 臣 ト ○ 臣 - の Q ()

#### Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

## Quiz - using silent mutations

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

We know two types of mutations in DNA silent and coding

- Which of them are more interesting for calculating divergence between species?
- And which are more interesting for paternity testing?

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

- Hamming distance: a metric originating from Information theory
- Given two vectors of the same length, it returns the number of positions where they differ.

 $D_H(s_1, s_2) = \sum_{i=1}^n \{1 : s_1[i] \neq s_2[i]; 0 : otherwise\}$ 

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

• A proper distance (satisfies triangle inequality)

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

- DNA polymerase can (rarely) slide over nucleotides
- especially over stretches of low complexity
- this leads to short deletions of DNA after replication
- Transposable elements lead to insertions of larger segments
- Chromosome recombination leads to duplications and deletions on different chromosomes at the same time

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

・ロト ・ 同ト ・ ヨト ・ ヨト

3

Sar

Number of mutations needed to *evolve* two sequences from a common ancestor is the same (under parsimony assumption) as the number of mutations needed to *evolve* one into the other



Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

◆□ ▶ ◆□ ▶ ◆三 ▶ ◆□ ▶ ◆□ ◆ ●

- Classically, genes are the **basic units of heritability**
- Gregor Mendel (1822-1884) laid foundations of genetics with his experiments on peas
- He also introduced the term allele and formulated laws of inheritance (segregation and independence)
- He knew nothing about DNA!



### Genes - modern view

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Efektywne algorytmy do porównań sekwencji

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

- $\bullet\,$  Currently, we know that genes are carried by DNA
- Current definition of a gene is substantially more complex:
- a locatable region of genomic sequence, corresponding to a unit of inheritance, which is associated with regulatory regions, transcribed regions, and or other functional sequence regions (Pearson, Nature, 2006)
- This is overly complex for our purposes, so
- We will be most concerned with *protein coding genes*, i.e. DNA sequences encoding proteins

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

- We can introduce *edit distance*: the number of editing operations needed to transform one sequence into the other. These operations are:
  - Substitutions
  - Insertions
  - Deletions
- The *procedural* definition of the distance makes it difficult to work with
- Does it matter in what order I make the operations (*If i delete a character, I cannot substitute it anymore...*)
- It turns out the *optimal* edit distances are simpler and can be described in a formal way as sequence *alignments*

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ うへつ

For a given sequences s, t over an alphabet  $\Sigma$ , their alignment is a pair of words s', t' over the extended alphabet  $\Sigma' = \Sigma \cup \{-\}$ . Sequences s', t' need to satisfy the following:

• 
$$|s'| = |t'|$$

1

• 
$$s'_{|\Sigma} = s$$
 and  $t'_{|\Sigma} = t$ 

• for no position i, s'[i] = t'[i] = -

For example, one of the words HEAGAWGHEE and PAWHEAE is HEAGAWGHE-E --P-AW-HEAE

Number of possible alignments for words of length n

$$\binom{2n}{n} = \frac{(2n)!}{(n!)^2} \simeq \frac{2^{2n}}{\sqrt{\pi n}}$$

## Scoring alignments: binary dotplots

Efektywne algorytmy do porównań sekwencji

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach



#### Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

# Scoring alignments: BLOSUM score matrix

-

1

590



Bartek Wilczyński

Reminding sequence evolution

From evolution t distance

Sequence alignment

Dynamic programming approach

## Recursive equation for sequence alignment

A B > A B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B

Sac

Э

	Н	Е	A	G	A	W	G	Н	Е	Ε	
P	-2	-1	-1	-2	-1	-4	-2	-2	-1	-1	
A	-2	-1	5	0	5	-3	0	-2	-1	-1	
W	-3	-3	-3	-3	-3	15	-3	-3	-3	-3	
Η	10	0	$^{-2}$	-2	-2	-3	-2	10	0	0	
Ε	0	6	-1	-3	-1	-3	-3	0	6	6	
A	-2	-1	5	0	5	-3	0	-2	-1	-1	
Ε	0	6	-1	-3	-1	-3	-3	0	6	6	

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach  $F(i, j) = \max \begin{cases} F(i - 1, j - 1) + s(x_i, y_j), \\ F(i - 1, j) - d, \\ F(i, j - 1) - d. \end{cases}$ 



image (c) Durbin et al.

### Filling in the alignment matrix

・ロト ・ 同ト ・ ヨト ・ ヨト

Э

Sac

### Tracing back alignments

イロト イポト イヨト イヨト

E

Efektywne algorytmy do porównań sekwencji

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

Η Ε Α G Α W G Η Ε Ε -8 -16 --24 - $-32 \leftarrow -40 \leftarrow -48 \leftarrow -56 \leftarrow -64 \leftarrow -72 \leftarrow$ -80 0 • ĸ Ρ -17 🗲 -8-2-9 -25 -33 -10 --57-65-73 ٠ . κ. ٠ K Α -16-10-20 --3-4 4 -12-28 --36 --44 --52 -60 ٠ ۰ ĸ W -24-18-7 -15-5 -13 < -21 -29 -37 -6 π. ۰ . Η -14-18-13-8 -9 -13-3-11 -19 -7 4 ŧ ÷ κ. ۰ Е -40-8-16-16\_9 -12-15-5 + 3 ٠ ŧ Α -48 -30-12-16-3 + -11-12-152 ŧ ŧ ŧ Ε -56 -38 -24 -12-14-15 -12-11-6 \_9 1

HEAGAWGHE-E --P-AW-HEAE

Bartek Wilczyński

Reminding sequence evolution

From evolution t distance

Sequence alignment

Dynamic programming approach Finding local alignments - Smith, Waterman '82

$$F(i, j) = \max \begin{cases} 0, \\ F(i-1, j-1) + s(x_i, y_j), \\ F(i-1, j) - d, \\ F(i, j-1) - d. \end{cases}$$

		Н	Ε	А	G	А	W	G	Η	Ε	Ε
	0	0	0	0	0	0	0	0	0	0	0
Р	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	5	0 ~	5	0	0	0	0	0
W	0	0	0	0	2	0	20 ←	12 🗲	4	0	0
н	0	10 ←	2	0	0	0	12	18	22 🕂	14 🔶	6
Е	0	2	16 🔶	8	0	0	4	10	18	28	20
A	0	0	₹ ~ 8	21 🔶	13	5	0	4	10	20	27
Е	0	0	6	13	18	12 🔶	4	0	4	16	26
AWGHE											
AW-HE											

Bartek Wilczyński

Reminding sequence evolution

From evolution to distance

Sequence alignment

Dynamic programming approach

## Scoring alignments: general gap penalty

・ロト ・ 同ト ・ ヨト ・ ヨト

Sac

3

### General gap penalty

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j), \\ F(k, j) + \gamma(i-k), & k = 0, \dots, i-1, \\ F(i,k) + \gamma(j-k), & k = 0, \dots, j-1. \end{cases}$$

Affine gap penalty (caching)

$$M(i, j) = \max \begin{cases} M(i-1, j-1) + s(x_i, y_j), \\ I_x(i-1, j-1) + s(x_i, y_j), \\ I_y(i-1, j-1) + s(x_i, y_j); \end{cases}$$

$$I_x(i, j) = \max \begin{cases} M(i-1, j) - d, \\ I_x(i-1, j) - e; \end{cases}$$

$$I_y(i, j) = \max \begin{cases} M(i, j-1) - d, \\ I_y(i, j-1) - e. \end{cases}$$