

1. Poiščite globalno poravnavo zaporedij $s_1 = \text{ATTCG}$, $s_2 = \text{GGAGT}$, $s_3 = \text{TGACG}$ in $s_4 = \text{GAACT}$ z uporabo

- (a) algoritma za progresivno poravnavo
- (b) 2-aproksimacijskega algoritma

Za razdaljo med dvema zaporedjema uporabite Needleman-Wunschev algoritem z nagrado 0 za ujemanje in kaznijo za zamenjavo, vstavljanje in brisanje enako 1.

a) Progresivna poravnava

1. Matrica razdalj
2. Drevo UPGMA
3. MSA

- $\binom{k}{2}$ parov zaporedij
- k zaporedij

1. Matrica razdalj (več možnih razdalj)

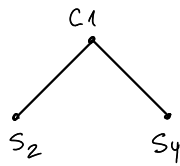
• Razdalja = ocena globalne poravnave (absolutne vrednosti)

	s_1	s_2	s_3	s_4
s_1	0	5	3	4
s_2		0	3	2
s_3			0	3
s_4				0

Poravnave so bile že ustvarjene, drugače moras vse narediti sam.

2. Drevo UPGMA (več metod: UPGMA, neighbor joining, iz aditivnih matrik)

	s_1	s_2	s_3	s_4
s_1	0	5	3	4
s_2		0	3	2
s_3			0	3
s_4				0

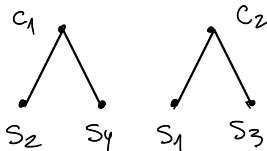


$$d(C_1, s_1) = \frac{(d(s_2, s_1) + d(s_1, s_4))}{|C_1|} = \frac{5+4}{2} = 4.5$$

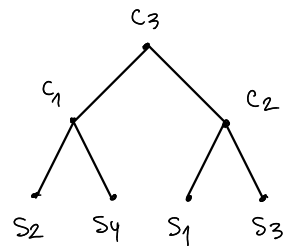
$d(C_1, s_3) = \text{analogno}$

	s_1	C_1	s_3
s_1	0	4.5	3
C_1		0	3
s_3			0

izbereš enega



= očitno bo $C_3 = (C_2, C_1) \Rightarrow$



3. MSA

- vsaki grupi določimo najboljšo poravnavo zaporedja ene skupine z zaporedjem iz druge skupine (najboljša = najmanjša):

- C_1 ima poravnavo (s_2, s_4)
- C_2 ima poravnavo (s_1, s_3)
- C_3 ima poravnavo $\min\{C_1, C_2\} = \min\{d(s_2, s_1), d(s_2, s_3), d(s_4, s_1), d(s_4, s_3)\} = \min\{5, 3, 4, 3\}$; izberemo (s_2, s_3) , saj se s_4 boljše poravnava z s_1 .

- Začnemo na vrhu s $C_3 = (s_2, s_3)$ in progresivno dodajamo poravnave.

- s_4 dodaš v (s_2, s_3) glede na to kako se poravnava z s_2
- s_1 dodaš v (s_4, s_2, s_3) glede na to kako se poravnava z s_3

s_2 | G G A - G T
 s_3 | T G A C G -

→ s_2 | G G A G T
 s_4 | G A A C T

s_4 | G A A - C T
 s_2 | G G A - G T
 s_3 | T G A C G -

→ s_1 | A T T C G
 s_3 | T G A C G

s_4 | G A A - C T
 s_2 | G G A - G T
 s_3 | T G A C G -
 s_1 | A T T C G -

b) 2-aproksimacijski algoritem

- 1. korak je enak kot pri (a); 2. korak ni potreben - nadomestimo ga z naslednjim algoritmom:
- Vrednost poravnave je enaka vsoti vrednosti poravnav vseh $\binom{k}{2}$ parov zaporedij, induciranih s poravnavo k zaporedij.
- Predp., da iščemo čim manjše vrednosti
- Vsota parov poravnav k-zaporedij bo kvečjemu za faktor $2^{-\frac{k}{2}}$ večja od optimalne vsote parov.

	S ₁	S ₂	S ₃	S ₄	
S ₁	0	5	3	4	= 12
S ₂	5	0	3	2	= 10
S ₃	3	3	0	3	= 9
S ₄	4	2	3	0	= 9

seštejemo vrstice ...

MSA:

S ₁	A	T	T	-	C	G	-
S ₂	G	G	A	-	-	G	T
S ₃	T	G	A	-	C	G	-
S ₄	-	G	A	A	C	T	-

Izberemo enega izmed teh dveh najmanjših. Izberimo npr. S₃.

- Poravnava po vrsti (S₁, S₃), (S₁, S₂, S₃), (S₁, S₂, S₄, S₃).
glede na (S₁, S₃), (S₂, S₃) in (S₄, S₃).

S ₁	A	T	T	C	G
S ₃	T	G	A	C	G

S ₂	G	G	A	-	G	T
S ₃	T	G	A	C	G	-

S ₃	T	G	A	-	C	G
S ₄	-	G	A	A	C	T

Vrednost poravnave:

$$v(A) := \sum_{i < j} v_A(s_i, s_j)$$

Iz MSA dobimo Levenshteinovo razdaljo za vsak par:

- $v_A(s_1, s_2) = 5$
- $v_A(s_2, s_3) = 3$
- $v_A(s_3, s_4) = 3$
- $v_A(s_2, s_4) = 5$

S ₁	A	T	T	-	C	G	-
S ₂	G	G	A	-	-	G	T

- $v_A(s_1, s_3) = 3$
- $v_A(s_1, s_4) = 5$

$$v(A) = 5 + 3 + 5 + 3 + 5 + 3 = 24$$

Za optimalno poravnavo A* velja:

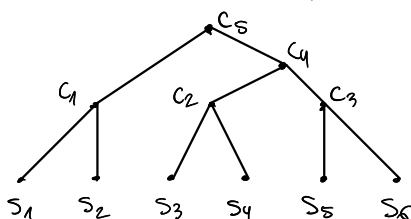
$$v(A_c) \leq (2 - 2/k) v(A^*)$$

$$v(A_c) / (2 - 2/k) \leq v(A^*)$$

Dobimo jo tako, da iz matrice razdalj odčitamo razdalje med zaporedji; vzamemo manjšo izmed te in one razdalje iz MSA:

$$v(A^*) = 5 + 3 + 4 + 3 + 2 + 3 = 20$$

*. Kako bi poravnali zaporedja glede na drevo:



- C₁ ima (S₁, S₂)
- C₂ ima (S₃, S₄)
- C₃ ima (S₅, S₆)

- C₄ ima $\min\{C_2, C_3\} = \min\{d(S_3, S_5), d(S_3, S_6), d(S_4, S_5), d(S_4, S_6)\}$
npr. (S₃, S₅)
- C₅ ima $\min\{C_1, C_4\} = \min\{d(S_1, S_3), d(S_1, S_5), d(S_2, S_3), d(S_2, S_5)\}$
npr. (S₁, S₅)

par, ki smo ga izbrali za C₄, ne vsa zaporedja spadaj...

Za poravnavo potrebujemo poravnave zaporedij:
(S₁, S₂), (S₃, S₄), (S₅, S₆), (S₃, S₅) in (S₁, S₅).

3. Izračunajte zaporedje s , ki vsebuje 4-terice podane v multimnožici

$$S = \{CGAA, ATAG, TAGT, GAAT, TAGA, AGTA, AATA, AGAT, GTAG, GATA, ATAG\}$$

tako, da bo $\text{Spekter}(s, 4) = S$, z uporabo

- (a) Hamiltonske poti
- (b) Eulerjeve sledi

Ali je rešitev enolično določena?

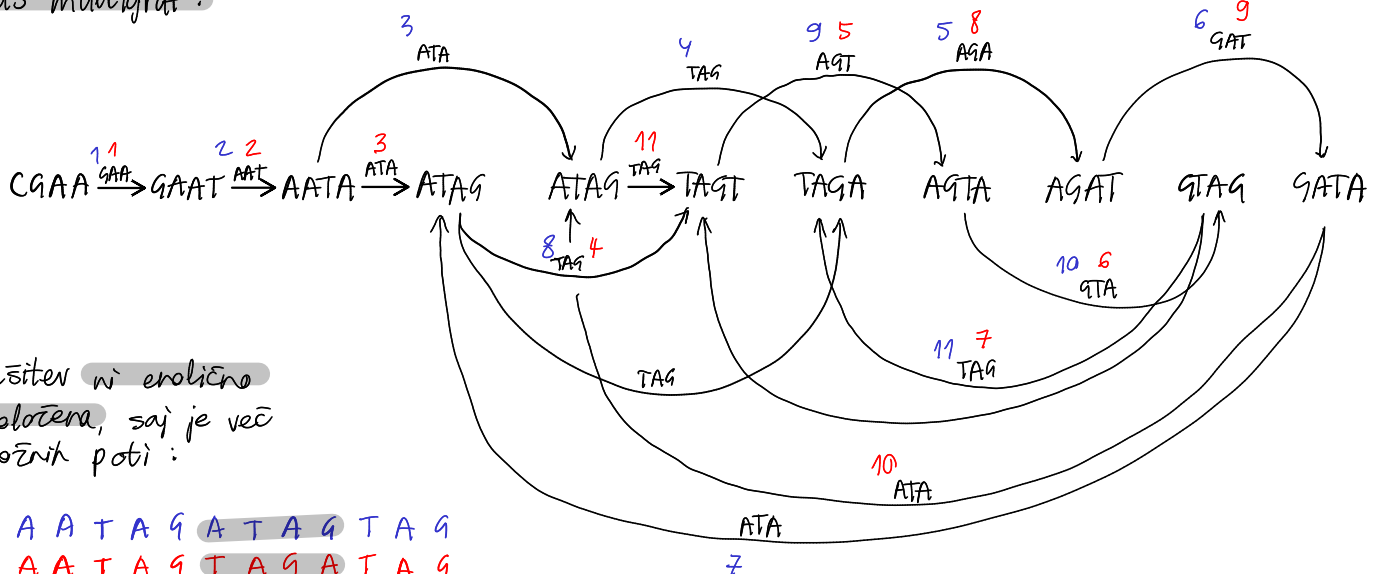
a) Hamiltonske poti

- Hamiltonska pot obišče vsako točko natanko enkrat.
- Hamiltonski cikel je cikel, ki gre skozi vse točke enkrat.
- Ne moremo podati preprostega pogoja potrebnega za obstoj hamiltonskega cikla, saj je problem NP-poln (enako za ham. poti).
- Metode sekvenciranja najprej deložijo množico podzaporedij neznanega zaporedja, ki se med seboj prekrivajo. Kako ugotoviti celotno zaporedje?
 - Dano je zap. s nad abecedo Σ in $l \in \mathbb{N}$. l -spekter zaporedja je multimnožica vseh podzap. zap. s dolžine l . Označimo ga $\text{Spekter}(s, l)$.
 - Problem izračuna zaporedja s , za katerega velja $\text{Spekter}(s, l) = S$, je učinkovito rešljiv.

Naj bo $V=S$. $p, q \in V$ natanko tedaj, ko se p in q prekrivata \Leftrightarrow ko je zadnjih $l-1$ znakov p enakih prvih $l-1$ znakov q .

$p: A \begin{matrix} C & C & T & G \\ C & C & T & G \end{matrix} \quad q: \begin{matrix} C & C & T & G & G \end{matrix}$

Naš multigraf:



Namesto, da pišes nad povezane $l-1$ besede, smo naredi pikice na točkah:



3. Izračunajte zaporedje s , ki vsebuje 4-terice podane v multimnožici

$$S = \{CGAA, ATAG, TAGT, GAAT, TAGA, AGTA, AATA, AGAT, GTAG, GATA, ATAG\}$$

tako, da bo $\text{Spekter}(s, 4) = S$, z uporabo

- (a) Hamiltonske poti
- (b) Eulerjeve sledi

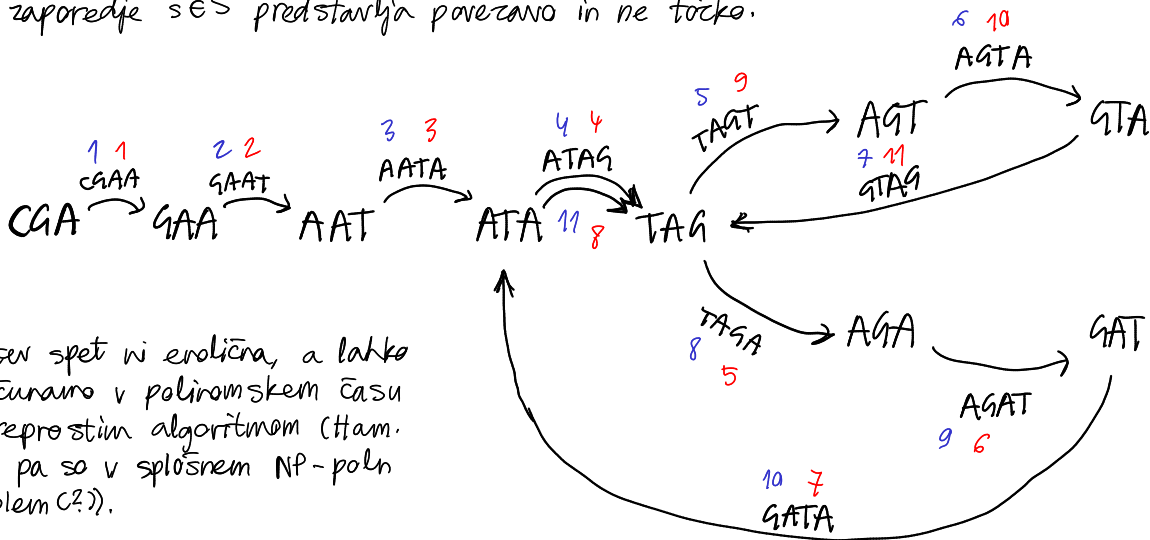
Ali je rešitev enolično določena?

b) Eulerjeve sledi

Za vhodne podatke (l, S) konstruiramo digraf $D = (V, E)$ (v katerem so možne večkratne povezave), kjer je:

- V množica vseh zaporedij dolžine $l-1$, ki nastopajo kot podzaporedja v S .
- $(p, q) \in E$ natanko tedaj, ko \exists zaporedje $s \in S$, da je prvih $l-1$ znakov s enakih p , zadnjih $l-1$ s pa q . Povezavo označimo z l -terico s .

Vsako zaporedje $s \in S$ predstavlja povezano in ne točko.



Rešitev spet ni enolična, a lahko izračunamo v polinomskem času s preprostim algoritmom (Ham. poti pa so v splošnem NP-poln problem (?)).

Rešitvi:

CGAATA TAGATA
CGAATA TAGATA

(enako kot prej :D)