# Algoritmi v bioinformatiki - 2. Domača naloga

Jan Panjan

May 7, 2025

# 1. Dano imamo naslednje zaporedje izidov metov kovanca

$$V = CCCGCGGCCGC$$

pri čemer C označuje, da je bil izid meta cifra, G pa da je bil izid meta grb. Za mete imamo na voljo 3 kovance, A, B in C, veljajo naslednje verjetnosti:

Prehod:

%	A	В	C
A	40	30	30
В	30	40	30
$\overline{C}$	30	30	40

Izpis:

_%_	С	G
A	75	25
В	80	20
$\overline{C}$	20	80

Katera od možnosti je najbolj verjetna?

- (a) za vse mete smo uporabili kovanec A
- (b) za vse mete smo uporabili kovanec C
- (c) za vse mete smo uporabili kovanec B
- (d)  $\Pi = AAACBCCBBCA$

Odgovor ustrezno utemeljite.

# Za vse mete smo uporabili kovanec A

	7-krat vržemo C z verjetnostjo 0.75
\ /	7-krat vržemo G z verjetnostjo 0.75
$(0.4)^{10}$	10-krat ne zamenjamo kovanca A z verjetnostjo 0.4

$$p(A) = (0.75)^7 \cdot (0.25)^4 \cdot (0.4)^{10} = 0.00209$$

# Za vse mete smo uporabili kovan<u>ec</u> B

\ /	7-krat vržemo C z verjetnostjo 0.8
` /	7-krat vržemo G z verjetnostjo 0.2
$(0.4)^{10}$	10-krat ne zamenjamo kovanca $B$ z verjetnostjo $0.4$

$$p(B) = (0.8)^7 \cdot (0.2)^4 \cdot (0.4)^{10} = 0.00134$$

# Za vse mete smo uporabili kovanec C

$(0.2)^{\gamma}$	7-krat vržemo C z verjetnostjo 0.8
( /	7-krat vržemo G z verjetnostjo 0.2
$(0.4)^{10}$	10-krat ne zamenjamo kovanca $C$ z verjetnostjo $0.4$

$$p(C) = (0.2)^7 \cdot (0.8)^4 \cdot (0.4)^{10} = 0.0000209$$

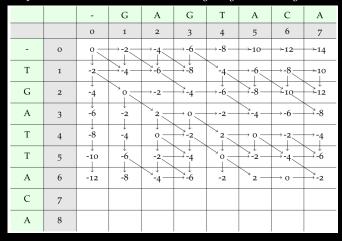
#### $\Pi = AAACBCCBBCA$

\ /	6-krat ostanemo v istem kovancu (vsi kovanci imajo enake verjetnosti)
$(0.4)^4$	4-krat zamenjamo kovanec (tudi tu imajo enako verjetnosti)
$(0.75)^4$	4-krat vržemo kovanec $A$ , vsakič vržemo cifro z verjetnostjo $0.75$
$(0.8)^3$	3-krat vržemo kovanec $B$ , vsakič vržemo cifro z verjetnostjo $0.8$
$(0.8)^4$	4-krat vržemo kovanec C, vsakič vržemo grb z verjetnostjo 0.8

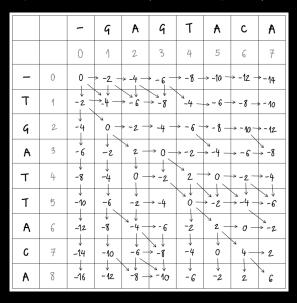
$$p(\Pi) = (0.3)^6 \cdot (0.4)^4 \cdot (0.75)^4 \cdot (0.8)^3 \cdot (0.8)^4 = 0.00000124$$

**Rešitev:** Najbolj verjetna je možnost z največjo verjetnostjo. To je možnost (a) z verjetnostjo 0.00209.

- 2. Dani imamo zaporedji s = GAGTACA in t = TGATTACA ter vrednostno funkcijo s parametroma  $\mu = 4, \sigma = 2$  in nagrado za ujemanje 2.
  - (a) Zuporabo Needleman-Wunsch-evega algoritma za globalno poravnavo smo dobili naslednjo tabelo:

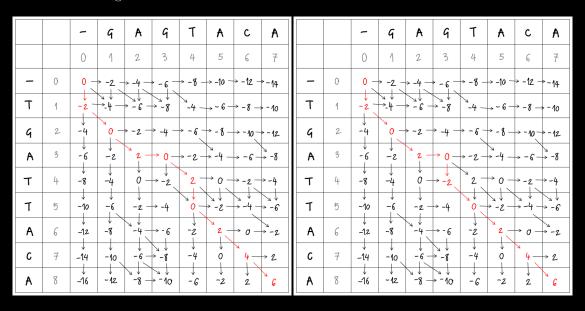


Dopolnite tabelo tako, da poračunate vrednosti (in ustrezne puščice) za zadnji dve vrstici.



(b) Koliko optimalnih globalnih poravnav dobite? Izpišite vse rešitve. Dobim dve optimalni globalni poravnavi, in sicer:

Z matrikami to izgleda tako:



Spremeni se mesto vrzeli in sicer iz mesta (4,4) na mesto (4,3).

3. Dano imamo naslednjo matriko izražanja:

	$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$
$g_1$	2	2	6	2	3	4
$g_2$	3	7	3	1	9	3
$g_3$	2	2	7	2	6	3
$g_4$	3	2	3	2	1	3
$g_5$	2	1	5	1	0	4
$g_6$	3	5	5	8	2	3
$g_7$	1	3	1	5	4	2
$g_8$	5	4	2	4	7	5

4. Določite gruče z uporabo metode voditeljev, če je začetna množica voditeljev enaka  $X = \{g_1, g_5, g_6\}$ . Vsak gen lahko obravnavamo kot vektor  $g_i = (T_1, \dots, T_6)$ .

#### Prva iteracija

**Prvi korak** Najprej je potrebno izračunati razdalje med geni. (Evklidska) razdalja med vsakim genom je definirana kot:

$$d(g_i, g_j) = \sqrt{\left(T_i^{(1)} - T_j^{(1)}\right)^2 + \dots + \left(T_i^{(6)} - T_j^{(6)}\right)^2} \quad ; \quad 1 \le i, j \le 8$$

Primer za prvi gen:

$$d(g_1, g_1) = \sqrt{(2-2)^2 + (2-2)^2 + (6-6)^2 + (2-2)^2 + (3-3)^2 + (4-4)^2} = 0$$

Očitno je razdalja med istim genom 0, kar pravi tudi prva lastnost metrike:  $d(x,y) = 0 \iff x = y$ . Ko poračunamo vse, dobimo matriko razdalj. Potrebujemo razdalje samo do voditeljev:

3

	$d(g_1,g_i)$	$d(g_5,g_i)$	$d(g_6,g_i)$
$g_1$	0	18.166	17.493
$g_2$	8.544	11.091	10.296
$g_3$	3.317	6.557	8.124
$g_4$	3.873	3.000	7.071
$g_5$	18.166	0	10.198
$g_6$	17.493	10.198	0
$g_7$	5.657	7.071	5.568
$g_8$	7.071	9.274	7.681

Drugi korak Za vsak gen izberemo najkrajšo razdaljo med njim in voditeljem.

	$d(g_1,g_i)$	$d(g_5,g_i)$	$d(g_6,g_i)$
$g_1$	0	18.166	17.493
$g_2$	8.544	11.091	10.296
$g_3$	3.317	6.557	8.124
$g_4$	3.873	3.000	7.071
$g_5$	18.166	0	10.198
$g_6$	17.493	10.198	0
$g_7$	5.657	7.071	5.568
$g_8$	7.071	9.274	7.681

Tretji korak Iz vsakega stolpca odčitamo nove voditelje (vrednosti označene z rdečo), katere označimo s  $C_i, i \in \mathbb{N}$ , in sicer:

$$C_1 = \{g_1, g_2, g_3, g_8\}$$

$$C_2 = \{g_4, g_5\}$$

$$C_3 = \{g_6, g_7\}$$

**Četrti korak** Za gručo  $C_i$  z n geni  $\{g_1, \ldots, g_k \mid 1 \le k \le 8\}$ , izračunamo nov vektor vrednosti  $v_i = (v_{i1}, \ldots, v_{in})$  z enačbo:

$$v_i = \frac{1}{n} \sum_{i=1}^n g_k$$

Nove vrednosti so torej aritmetična sredina vseh genov v gruči:

$$v_1 = (3, 3.75, 4.5, 2.25, 6.25, 3.75)$$
  

$$v_2 = (2.5, 1.5, 4, 1.5, 0.5, 3.5)$$
  

$$v_3 = (2, 4, 3.5, 6.5, 3, 2.5)$$

Zdaj ponovimo korake dokler ne dosežemo konvergence:

- ko se gruče med iteracijama ne spremenijo
- ko postanejo razlike med radaljami gruč manjše od neke vnaprej določene vrednosti.

# Druga iteracija

# Prvi + drugi korak

	$d(v_1,g_i)$	$d(v_2,g_i)$	$d(v_3,g_i)$
$g_1$	3.808	3.354	5.723
$g_2$	4.743	10.209	8.761
$g_3$	3.317	6.344	6.764
$g_4$	5.788	1.500	5.454
$g_5$	7.036	1.500	7.263
$g_6$	7.314	7.632	2.784
$g_7$	5.148	5.958	2.784
$g_8$	3.937	8.185	6.305

# Tretji korak

$$C_4 = \{g_1, g_2, g_3, g_8\}$$

$$C_5 = \{g_4, g_5\}$$

$$C_6 = \{g_6, g_7\}$$

Ker so gruče enake kot v prejšnji iteraciji, lahko postopek tu končamo...

 ${\bf Re \check{s}itev:}\,$ gruče določene z metodo voditeljev s k=3 za dano matriko izražanja so

$$\{g_1, g_2, g_3, g_8\}$$
  
 $\{g_4, g_5\}$   
 $\{g_6, g_7\}$ 

5. Izračunajte drevo hierarhičnega gručenja z uporabo algoritma UPGMA.

Osnova za algoritem UPGMA je matrika razdalj genov (matrika iz prve iteracije prejšnje naloge), za katero uporabimo sledečo enačbo

$$d_{\text{avg}}(C, C^*) = \frac{1}{|C||C^*|} \sum_{x \in C, y \in C^*} d(x, y)$$

kjer staC in  $C^{\ast}$  dve gruči (na začetku so to geni).

	$g_1$	$g_2$	$g_3$	$g_4$	$g_5$	$g_6$	$g_7$	$g_8$
$g_1$	0	8.544	3.317	3.873	3.464	7.000	5.657	7.071
$g_2$		0	8.124	10.198	11.091	10.296	10.677	8.602
$g_3$			0	6.325	6.557	8.124	6.403	7.000
$g_4$				0	3.000	7.071	5.099	6.403
$g_5$					0	10.198	7.071	9.274
$g_6$						0	5.568	7.681
$g_7$							0	5.916
$g_8$								0

Gručenje deluje tako, da vsako iteracijo izberemo najbližja gena in ju združimo v gručo. Na začetku je vsak gen v svoji gruči, do konca postopka pa ustvarimo eno celovito gručo, ki bo vsebovala vse gene.

Začetni dendrogram:

 $g_1$   $g_2$   $g_3$   $g_4$   $g_5$   $g_6$   $g_7$   $g_8$ 

Prvi korak Izberemo najmanjšo vrednost med razdaljami.

	$g_1$	$g_2$	$g_3$	$g_4$	$g_5$	$g_6$	$g_7$	$g_8$
$g_1$	0	8.544	3.317	3.873	3.464	7.000	5.657	7.071
$g_2$		0	8.124	10.198	11.091	10.296	10.677	8.602
$g_3$			0	6.325	6.557	8.124	6.403	7.000
$g_4$				0	3.000	7.071	5.099	6.403
$g_5$					0	10.198	7.071	9.274
$g_6$						0	5.568	7.681
$g_7$							0	5.916
$g_8$								0

**Drugi korak** Dodamo gena v gručo  $C_1$ , torej  $C_1 = \{g_4, g_5\}$  in poračunamo novi vektor  $v_1$ , ki bo predstavljal novo gručo v matriki razdalj (tako kot prej uporabimo aritmetično sredino komponent):

$$v_1 = (2.5, 1.5, 4, 1.5, 0.5, 3.5)$$

**Tretji korak** Izračunamo razdaljo gruče (C) do vseh ostalih gruč  $(C^*)$  s pomočjo prejšnje enačbe. V matriki razdalj odstranimo gena, ki smo ju združili v gručo ter dodamo novo gručo:

	$g_1$	$g_2$	$g_3$	$C_1$	$g_6$	$g_7$	$g_8$
$g_1$	0	8.544	3.317	3.669	7.000	5.657	7.071
$g_2$		0	8.124	10.645	10.296	10.677	8.602
$g_3$			0	6.441	8.124	6.403	7.000
$C_1$				0	8.635	6.085	7.839
$g_6$					0	5.568	7.681
$g_7$						0	5.916
$g_8$							0

<u>Četrti korak</u> Gena povežemo na dendrogramu na izračunani višini:

$$g_1$$
  $g_2$   $g_3$   $g_4$   $g_5$   $g_6$   $g_7$   $g_8$