

KODIRANJE

Stiskanje z enakim številom bitov: vsak ASCII znak ima 8 bitov oz. 1 bajt. Genom je velik cca 3 milijarde znakov, to je 3×10^9 in posledično 2.4×10^{10} bitov (3×10^9 bajtov). Primer, ko ima vsak nukleotid svoje bite, ki so njegova koda:

A: 1 0 0 0 C: 0 1 0 0 G: 0 0 1 0 T: 0 0 0 1 \Rightarrow pol manj kot prej (iz 8 v 4 bite), kar je 1.5×10^9 bajtov = 1.5 GB.

A: 0 0 C: 1 0 G: 0 1 T: 1 1

\Rightarrow 1 bajt na 4 znake $\Rightarrow 3 \cdot 1/4 \times 10^9 = 0.75$ GB

Primer zapisa: A G A G C C T = 00 01 00 01 10 10 11

stiskanje z različnim številom bitov: Npr. dajmo vsakemu znaku neko kodo naključno:

A: 0 C: 0 1 G: 0 0 1 T: 0 0 0

\Rightarrow Problem tu je da ne vemo kdaj preberemo A(0) ali C(01) ali T(000). Ne vemo kdaj nehati brati kodo za znak. Problem je, da je A pripoma vsem (vsi se začnejo z 0), ali pa npr. AC = 001 = G?! Problem pri dešifriranju.

Iščemo tako stiskanje, da bo predpansko prosto.

HUFFMANOVO KODIRANJE: efektivno zakodira podatke

glede na frekvence znakov, zato je odvisno od značilnosti besedila (ne zakodira vseh besedil enako). Znakom s manjšo frekvenco dodeli daljšo kodo in obratno.

PORAVNAVNA ŽAPOREDIJ

LEVENSTHEINOVA RAZDALJA: število razlik med dvema nizoma oz. najmanjše število sprememb potrebnih, da enega pretvorimo v drugega. Npr. niza s = ATATAT- in t = -TATATATA:

$$\begin{array}{c|cccccccc} s & A & T & A & T & A & T & A & T & - \\ t & - & T & A & T & A & T & A & T & A \end{array} \Rightarrow d(s,t) = 2$$

HAMMINGOVA RAZDALJA: število mest na katerih se niza razlikujeta. Za zgornji primer je to 2 (začetek in konec), za niza s' = ACT in t' = GTA pa je 3 (vse črke). Problem: Hammingova razdalja je lahko velika, čeprav sta si zaporedji zelo podobni:

$$\begin{array}{c|cccccc} s & A & T & A & T & A & T \\ t & T & A & T & A & T & A \end{array} \Rightarrow d(s,t) = 8$$

Hevristični algoritmi so hitrejši od optimalnih, toda ne izračunajo optimalne rešitve.

Poravnave ovrednotimo tako, da večje vrednosti ustrezajo boljšim poravnavam.

Problem poravnave lahko modeliramo kot iskanje najdaljše poti v acikličnem grafu (predstavljaj si matriko) \Rightarrow vsaka pot v grafu od (0,0) do (i,j) predstavlja neko poravnano podzaporedij $s[1,...,i]$ in $t[1,...,j]$ in obratno. Če povezane utežimo z ustreznimi nagradami in kaznimi, je največja vrednost poravnave zaporedij s in t (optimalna) enaka dolžini najdaljše poti.

Izračun poteka v treh korakih:

1. Inicializacija tabele
2. Izračun vrednosti v tabeli
3. Izračun ene (ali več) optimalnih poravnav zaporedij

Motivacija za lokalne poravnave: globalne ne odkrijejo podobnih podzaporedij. Algoritem: Smith-Waterman.

Vsaka pot od elementa 0 do elementa z največjo vrednostjo (i,j) predstavlja optimalno lokalno poravnano.

Primer:

	t	T	T	C	G	T	A	C
s	0	1	2	3	4	5	6	7
0	0	0	0	0	0	0	0	0
A	1	0	0	0	0	0	1	0
T	2	0	1	1	0	0	1	0
G	3	0	0	0	0	1	0	0
T	4	0	1	1	0	0	2	1
T	5	0	1	2	1	0	1	0
A	6	0	1	1	0	0	2	1
T	7	0	1	1	0	0	1	1

	t	A	T	C	G	T	A	C
s	0	1	2	3	4	5	6	7
0	0	0	0	0	0	0	0	0
A	1	1	0	0	0	0	0	0
T	2	0	1	1	0	0	1	0
G	3	0	0	0	0	1	0	0
T	4	0	1	1	0	0	2	1
T	5	0	1	2	1	0	1	0
A	6	0	1	1	0	0	2	1
T	7	0	1	1	0	0	1	1

Algoritem:
Needleman - Wunsch

Kako razdeliti množico eksperimentalno pridobljenih podatkov na skupine (grmče), tako da bodo podatki znotraj ene grmče čim bolj podobni drug drugemu in podatki v različnih grmčah čim bolj različni? Primeri uporabe: identifikacija družin genov s podobnimi funkcijami, rekonstrukcija filogenetskih dreves, poravnave več zaporedij...

Ugotavljanje funkcij genov poteka dandanes s pomočjo mikromrež (DNA microarrays), ki analizirajo ravni izražanja genov (količino proizvedene mRNA) pod različnimi pogoji in ob različnih časih.

Na ta način ugotovimo, kateri geni so ob katerem času aktivni v celici.

Podatki takega eksperimenta so zbrani v t.i.

matriki izražanja, I ,

ki ima n vrstic, po eno za vsak gen, in m stolpcev, po enega za vsako meritev.

Element matrike

$I(i,j)$

predstavlja raven izražanja i -tega gena v j -tem poskusu.

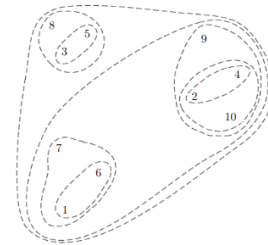
Celotna i -ta vrstica predstavlja t.i. **vzorec izražanja** i -tega gena.

Ideja: če imata dva gena podobna vzorca izražanja, potem sta zelo verjetno tudi biološko povezana (bodisi opravljata isto funkcijo ali sta prisotna pri istem biološkem procesu).

Grnčenje je dobro, ko zagotavlja **homogenost** (vrednosti D_{ij} morajo biti majhne za vsak par (i,j) iz iste grmče) ter **separacijo** (vrednosti D_{ij} morajo biti velike za vsak par (i,j) iz različnih grmč).

hierarhično grnčenje: elemente organizira v drevo. Ne predstavlja ene razdelitve množice na grmče, temveč družino takih razdelitev.

- Geni so v listih
- Vsaka povezava ima neko dolžino
- Razdalje med listi so sorazmerne z elementi matrike razdalj



- Metode:**
- UPGMA (Unweighted Pair Group Method with Arithmetic Averages)
 - NJ (Neighbor Joining) – to metodo uporabljajo programi **Clustal** (vrsta programov za poravnavo več zaporedij)
 - Rekonstrukcija dreves iz aditivnih matrik
- (ta metoda je direktno uporabna za konstrukcijo filogenetskih dreves, ob t.i. predpostavki molekularne ure)

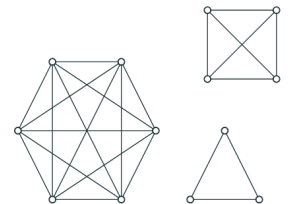
Nehierarhično grnčenje: npr. metoda voditeljev, ki zagotavlja homogenost vendar ne dobre separacije. Za te probleme ne obstajajo učinkoviti algoritmi, le heuristike – metoda poškodovanih klik, CA- σ .

(1)

- Klik v grafu je poln podgraf (poln graf je tak v katerem imata vsaki dve točki povezava).
- Klični graf ima za povezane komponente same polne grafe.

ti ustrezajo grmčam

- V idealnem primeru dobimo klični graf iz katerega lahko razberemo grnčenje, toda zaradi eksperimentalnih napak pogosto dobimo graf z dodatnimi povezavami ali z manj povezavami, zaradi česar ne dobimo kličnega grafa, therefore problem poškodovanih klik.



↳ Določi najmanjše število povezav, ki jih moramo odnesti/dodati grafu, da dobimo klični graf. NP-težek problem, le heuristike.

(2)

- CA- σ izračuna grmče glede na razdaljni graf ter konstanto σ . Iterativno izračuna particijo množice genov tako da na vsakem koraku najde grmčo C , za katero velja:
 - (i) $\exists i \notin C$ tako da je $d(i,C) < \sigma$ (bližnji gen izven grmče ne obstaja)
 - (ii) $\exists i \in C$ tako da je $d(i,C) > \sigma$ (noben gen v grmči ni oddaljen)
- Nima zagotovljene kvalitete rešitve ali zagotovljenega ustavitvenega pogoja, vendar se v praksi kar dobro obnese.