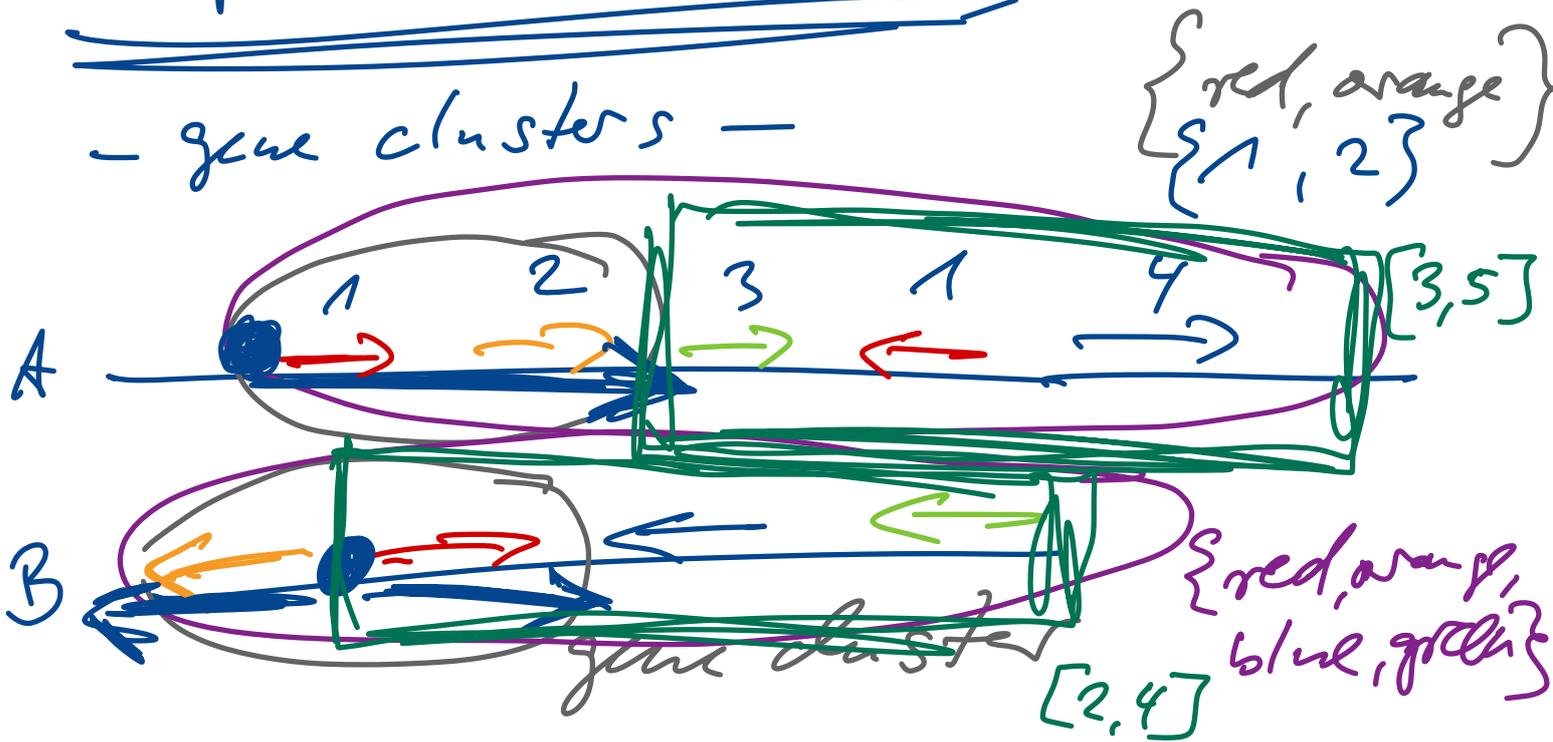


Comparative Genomics II

5 Feb. 2021

- gene clusters -



formally: common interval

S : sequence of "genes" $N = \{1, \dots, n\}$
(unsigned)

$S[i]$: the i -th element of S

$[x, y]$: the set of indices $\{x, x+1, x+2, \dots, y\}$

$S[x, y]$: the substring $S[x] S[x+1] \dots S[y]$

$\mathcal{Z}(S) = \{S[i] \mid 1 \leq i \leq |S|\}$ "character set of S "

Def: Given a family of k sequences $\mathcal{S} = (S_1, \dots, S_k)$,

then a common interval of \mathcal{S} is a k -tuple

$\mathcal{C} = ([l_1, u_1], [l_2, u_2], \dots, [l_k, u_k])$

if and only if $\mathcal{Z}(S_1[l_1, u_1]) = \dots = \mathcal{Z}(S_k[l_k, u_k])$.

genome models

① genome = permutation of $\{1, \dots, n\}$

w.l.o.g. assume that genome $A = \text{id}_n$
 $l=1, u=3, f=0$

A: 1 2 3 4 5 6 7 8
 B: 4 6 5 1 3 2 8 7
 $l=1, u=5, f=1$
 $l=1, u=5, f=5-1-(5-3)=2$

test functions:
 $l(x, y) = \min B[x, y]$
 $u(x, y) = \max B[x, y]$
 $f(x, y) = u(x, y) - l(x, y) - (y - x)$

observation: $f(x, y) = \text{number of "missing" genes in } B[x, y]$.

Lemma: $f(x, y) = 0 \iff A[x, y]$ is a common interval of A and B.

Algorithms

1. straight forward — $O(n^3)$ time
 → build all $\frac{n(n-1)}{2}$ intervals and compute f .
2. a bit more clever, using running min and running max — $O(n^2)$ time
3. Uno & Yagiura (2000) [no details] — $O(n+k)$ time
 where $k = \text{out put size}$

$k=2$
 genomes

4. Extension to Luo & Yagiura
for $k \geq 2$ genomes - $O(kn+k)$ time
Aeber, Mayr, Stojke (2011)

(2) genome = sequence over $N = \{1, \dots, n\}$
(Schmidt & Stojke, 2004/07)
[$O(n^2)$ time and space (for $k=2$ sequences)
[$O(kn^2)$ time for k sequences

Preprocessing

$POS[c] =$ list of all occurrences of c in A

$NUM[x, y] = |\mathcal{Z}(A[x, y])|$, $x < y$

Algorithm: enumerate all intervals (i, j) of B ,
while marking maximal intervals of A .
→ whenever such an interval $[x, y]$ fulfills
the test that $|\mathcal{Z}(B[i, j])| = NUM[x, y]$
then output a common interval.

$A = \begin{matrix} \overbrace{3} & \overbrace{1} & \overbrace{2} & \overbrace{3} & \overbrace{1} & \overbrace{5} & \overbrace{2} & \overbrace{6} \\ \underline{1} & \underline{2} & \underline{3} & \underline{4} & \underline{5} & \underline{6} & \underline{7} & \underline{8} \end{matrix} \quad (x,y)$
 $B = \begin{matrix} 4 & 3 & 5 & 5 & 5 & 1 & 4 & 2 & 2 \\ \underline{1} & \underline{2} & \underline{3} & \underline{4} & \underline{5} & \underline{6} & \underline{7} & \underline{8} & \underline{9} \end{matrix} \quad (i,j)$

POS [1] = 2, 5
 POS [2] = 3, 7
 POS [3] = 1, 4
 POS [4] = /
 POS [5] = 6
 POS [6] = 8

NUM

x \ y	1	2	3	4	5	6	7	8
1	1	2	3	3	3	4	4	5
2	1	1	2	3	3	4	4	5
3	1	2	1	2	3	4	4	5
4	1	2	3	1	2	3	4	5
5	1	2	3	4	1	2	3	4
6	1	2	3	4	5	1	2	3
7	1	2	3	4	5	6	1	2
8	1	2	3	4	5	6	7	1

output:
 ([4, 6], [2, 6])
 ...

variants of the gene cluster model

- (A) above: common interval of A and B.
- (B) more general: approximate common interval

→ add some error tolerance to the common interval model

common interval: $Z(A[x,y]) = Z(B[i,j])$

approx. c. i.: $|Z(A[x,y]) - Z(B[i,j])| \leq t$

where

\triangle is the symmetric set difference:

$$A \triangle B = (A \cup B) \setminus (A \cap B)$$

and t is a threshold value ≥ 0

(c) Weak common intervals in indeterminate strings

example: A: 1 2 ~~(3)~~ 2 ~~(3)~~ ~~(2)~~ ~~(4)~~ 3...

B: 2 3 ~~(4)~~ 3 ~~(2)~~ ~~(5)~~ 3...

def w.c.i. if by removing some alternatives of indeterminate positions, we get a c.i.