Topics of today:

- 1. DCJ double distance
- 2. DCJ halving
- 3. DCJ median
- 4. Inversion distance (unichromosomal genomes)

DCJ double distance

The DCJ distance of **balanced genomes** \mathbb{A}^f and \mathbb{B}^f is:

$$\mathsf{d}_{\mathrm{DCJ}}(\mathbb{A}^{f},\mathbb{B}^{f})=\min_{f_{m}}\mathsf{d}_{\mathrm{DCJ}}(\mathbb{A}^{f_{m}}_{\rhd},\mathbb{B}^{f_{m}}_{\rhd})$$

where f_m is any function that produces a maximal matching of the families defined by f



DCJ halving

DCJ Halving Distance Problem:

Compute the minimum DCJ double distance for a (rearranged) duplicated genome \mathbb{D}^{f} :

$$\mathsf{h}_{\text{DCJ}}(\mathbb{D}^f) = \min_{\text{singular } \mathbb{H}^f} \mathsf{d}^2_{\text{DCJ}}(\mathbb{H}^f, \mathbb{D}^f)$$

Find a singular genome $\widehat{\mathbb{H}^{f}}$ and a perfectly duplicated genome $\widehat{\mathbb{P}^{f}} \in 2 \cdot \widehat{\mathbb{H}^{f}}$ such that $h_{DCJ}(\mathbb{D}^{f}) = d_{DCJ}^{2}(\widehat{\mathbb{H}^{f}}, \mathbb{D}^{f}) = d_{DCJ}(\widehat{\mathbb{P}^{f}}, \mathbb{D}^{f})$ **DCJ Halving Sorting Problem:** Give a sequence of $h_{DCJ}(\mathbb{D}^{f})$ DCJ operations that transform \mathbb{D}^{f} into $\widehat{\mathbb{P}^{f}}$ $\widehat{\mathbb{P}^{f}}$ **DCJ Halving Problem:**

Natural graph of a duplicated genome

Natural graph $NG(\mathbb{D}^f) = (V, E)$ of a duplicated genome \mathbb{D}^f :

First assign arbitrarily indices a and b to the two genes of each family in \mathbb{D}^{f} , obtaining \mathbb{D}^{f}

- 1. $V = \Gamma(\mathbb{D}'^f) \cup \Theta(\mathbb{D}'^f)$ (each adjacency or telomere of \mathbb{D}'^f is a vertex of $NG(\mathbb{D}^f)$)
- 2. For each family $X \in \mathcal{F}(\mathbb{D}^{f})$, each pair of paralogous extremities is connected by an edge in $NG(\mathbb{D}^{f})$, i.e.:
 - there is an edge connecting the vertex u that contain X_a^h and the vertex v that contain X_b^h
 - there is an edge connecting the vertex u' that contain X_a^t and the vertex v that contain X_b^t

Note that:

- ▶ There can be adjacencies/vertices of type $X_a^h X_b^h$ and/or $X_a^t X_b^t$ ($NG(\mathbb{D}^f)$ can contain 1-cycles)
- Let $n = |\mathcal{F}(\mathbb{D}^f)| = \frac{|\mathcal{G}(\mathbb{D}^f)|}{2}$. The number of edges in $NG(\mathbb{D}^f) = 2n$ (two edges per element of $\mathcal{F}(\mathbb{D}^f)$).

Natural graph of a duplicated genome [5a 5b] 2-yelds are sorted Ex: $[\bar{4}_{a} \ 1_{a} \ \bar{4}_{b} \ \bar{3}_{a} \ 2_{a}]$ $[\bar{2}_{b} \ 3_{b} \ 1_{b}]$ $\Gamma(\mathbb{D}) \cup \Theta(\mathbb{D}) = \{ 4_{a}^{h}, 4_{a}^{t} 1_{a}^{t}, 1_{a}^{h} 4_{b}^{h}, 4_{b}^{t} 3_{a}^{h}, 3_{a}^{t} 2_{a}^{t}, 2_{a}^{h}, 2_{b}^{h}, 2_{b}^{t} 3_{b}^{t} \}, 3_{b}^{h} 1_{b}^{t}, 1_{b}^{h}, 5_{a}^{t}, 5_{a}^{h} 5_{b}^{h}, 5_{b}^{t} \}$ $n = |\mathcal{F}(\mathbb{D}^f)| = 5$ and $\kappa(\mathbb{D}^f) = 3$ Every vertex has degree one or two: 2-yde $NG(\mathbb{D}^{f})$ is a collection of paths and cycles 2-p.h 3 2 2 a cycle with k edges: k-cycle or c_k 1-cycle path with k edges: k-path or p_k 3 t 2t $\mathcal{C}_e = \{c_k : k \text{ is even}\}$: set of even cycles $\mathcal{P}_e = \{p_k : k \text{ is even}\}$: set of even paths 4t 1t 3, 1, t $C_o = \{c_k : k \text{ is odd}\}$: set of odd cycles 2, $\mathcal{P}_{o} = \{p_{k} : k \text{ is odd}\} : \text{ set of odd paths}$ 2-putts 4 3 h $|\mathcal{C}_{o}| + |\mathcal{P}_{o}|$ is even (*NG* has 2*n* edges) S.t 5* $|\mathcal{P}_e| + |\mathcal{P}_o| = \kappa(\mathbb{D}^f)$ 3-cycle Otherwise, if a duplicated genome \mathbb{D}^{f} For a perfectly duplicated genome \mathbb{P}^{f} , $NG(\mathbb{P}^{f})$ has only 2-cycles and 1-paths: is not perfectly duplicated: $2n = 2|\mathcal{C}_e| + |\mathcal{P}_o| \Rightarrow n = |\mathcal{C}_e| + \frac{|\mathcal{P}_o|}{2}$ $n > |\mathcal{C}_e| + \left| \frac{|\mathcal{P}_o|}{2} \right|$

Types of DCJ operation

Goal: increase the number of even cycles $(|\mathcal{C}_e|)$ and/or the number of odd paths $(|\mathcal{P}_o|)$ in NG



Types of DCJ operation

Goal: increase the number of even cycles ($|C_e|$) and/or odd paths ($|P_o|$) in NG





DCJ Halving: Sorting & Distance

Recall that, if the genome is perfectly duplicated, we have $n = |C_e| + \frac{|\mathcal{P}_o|}{2}$, otherwise $n > |C_e| + \left\lfloor \frac{|\mathcal{P}_o|}{2} \right\rfloor$

A DCJ operation ρ is called **optimal** if $\begin{cases} \rho \text{ increases the number of even cycles by one, or} \\ \rho \text{ increases the number of odd paths by two, or} \\ \text{the number of odd paths is odd and} \\ \rho \text{ increases the number of odd paths by one} \\ \text{(can occur at most once)} \end{cases}$

Given a duplicated genome \mathbb{D}^{f} , it is possible to find an optimal DCJ operation at each sorting step. Therefore:

$$\mathsf{h}_{\mathrm{DCJ}}(\mathbb{D}^f) = n - |\mathcal{C}_e| - \left\lfloor \frac{|\mathcal{P}_o|}{2} \right\rfloor$$

DCJ halving: obtaining an optimal perfectly duplicated genome

Given a duplicated genome \mathbb{D}^{f} , with natural graph $NG(\mathbb{D}^{f})$, and DCJ halving distance $h = h_{DCJ}(\mathbb{D}^{f}) = n - |\mathcal{C}_{e}| - \lfloor \frac{|\mathcal{P}_{o}|}{2} \rfloor$:

1. $NG_0 \leftarrow NG(\mathbb{D}^f)$

2. For i = 1 to h:

- Find and apply one optimal DCJ operation, transforming NG_{i-1} into NG_i .
- NG_h is a simple collection of 2-cycles and 1-paths: reconstruct the perfectly duplicated genome P^f ∈ 2·Ⅲ from NG_h.



1 Which of the following statements about the Natural Graph are true?

A Merging two odd cycles is always optimal.
 A Breaking an odd cycle into an odd path cannot be optimal.
 C Breaking an even path into two odd paths is always optimal.
 A Breaking an even cycle into two cycles is always optimal.
 C Breaking an even cycle into two cycles is always optimal.
 A Breaking an even cycle into two cycles is always optimal.
 C Breaking an even cycle into two cycles is always optimal.

Solving the DCJ double distance

Let \mathbb{S}^f be a singular and \mathbb{D}^f be a duplicated genome.

We want to compute the double distance $d^2_{\scriptscriptstyle \mathrm{DCJ}}(\mathbb{S}^f,\mathbb{D}^f)$

Assign arbitrarily indices a and b to the two genes of each family in $\mathbb{D},$ obtaining \mathbb{D}'

All possible adjacencies in $2 \cdot S$:

For each $uv \in \Gamma(\mathbb{S})$ $\begin{cases}
\text{the paralogous adjacencies are} \begin{cases}
\text{either} & P(uv) = \{u_a v_a, u_b v_b\} \\
\text{or} & \widetilde{P}(uv) = \{u_a v_b, u_b v_a\} \\
\text{and the square of } uv \text{ is defined as } Q(uv) = P(uv) \cup \widetilde{P}(uv)
\end{cases}$

DCJ double distance: ambiguous breakpoint graph

The ambiguous breakpoint graph $ABD(\mathbb{D}', 2 \cdot \mathbb{S}) = (V, E)$:

1.
$$V = \bigcup_{\mathbf{X}\in\mathcal{G}_{\star}} \{ \mathbf{X}_{\mathbf{a}}^{h}, \mathbf{X}_{\mathbf{a}}^{t}, \mathbf{X}_{\mathbf{b}}^{h}, \mathbf{X}_{\mathbf{b}}^{t} \} \quad \Rightarrow V = \xi(\mathbb{D}'); \quad |V| = 4n$$

there are two vertices for each extremity of each gene in \mathcal{G}_{\star} each vertex v has a label $\ell(v)$, that corresponds to the extremity of \mathbb{D}' it represents

2.
$$E = E_{\Gamma}(\mathbb{D}') \cup E_Q(2 \cdot \mathbb{S})$$
, where:

- ▶ D-adjacency edges: $E_{\Gamma}(\mathbb{D}') = \{uv : u, v \in V(\xi(\mathbb{D}')) \text{ and } \ell(u)\ell(v) \in \Gamma(\mathbb{D}')\}$
- ► Ambiguous S-adjacency edges: $E_Q(2 \cdot S) = \bigcup_{uv \in \Gamma(S)} \{uv : u, v \in V(\xi(\mathbb{D}')) \text{ and } \ell(u)\ell(v) \in Q(uv)\}$

The number of edges is
$$|E| = |E_{\Gamma}(\mathbb{D})| + 4|E_{\Gamma}(\mathbb{S})| \quad \begin{cases} |E_{\Gamma}(\mathbb{D})| \leq 2n \\ |E_{\Gamma}(\mathbb{S})| \leq n \end{cases}$$

Ambiguous breakpoint graph



Solution: for each square Q(uv), fix either P(uv) or $\tilde{P}(uv)$ so that the number of cycles is maximized.

Bicolored graph of two unsigned canonical chromosomes

Each vertex of a bicolored graph has degree 0, 2 or 4:



Idea:

 \Rightarrow

Entirely decompose a bicolored graph

into edge-disjoint alternating even cycles

One possible decomposition:



 $A_{1} = (15\bar{3}\bar{2}\bar{4}6)$

Bicolored graph of two unsigned canonical chromosomes

Another possible decomposition:



Bicolored graph decomposition (BGDEC)

Each vertex of a bicolored graph has degree 0, 2 or 4 The number of red and of blue edges inciding in each vertex is identical



Problem:

Entirely decompose a bicolored graph into the maximum number of edge-disjoint alternating even cycles

> ↓ NP-hard

Reducing BGDec to the DCJ double distance

 \Rightarrow







2.5

DCJ median of three canonical genomes

Given three canonical genomes \mathbb{A} , \mathbb{B} , \mathbb{C} , find another canonical genome \mathbb{M} that minimizes the sum

 $\mathsf{d}_{\text{DCJ}}(\mathbb{M},\mathbb{A}) + \mathsf{d}_{\text{DCJ}}(\mathbb{M},\mathbb{B}) + \mathsf{d}_{\text{DCJ}}(\mathbb{M},\mathbb{C})$



Reducing BGDEC to the DCJ median of three canonical genomes $\left\{ F_{2} \left\langle 1, 2, 2, 3, 4, 5, 5' \right\rangle \right\}$ Breakpoint graph of A, B and C



Quiz 2

1 Which of the following statements are true?

The multi mixed/circular DCJ double distance is NP-hard, therefore the multi mixed/circular DCJ halving is also NP-hard.

The multi linear breakpoint double distance is polynomial, therefore the multi linear breakpoint halving is also polynomial.

- 2 We prove that DCJ median is NP-hard...
 - A ... by reducing it to the bicolored graph decomposition.
 - B. by reducing the bicolored graph decomposition to it.

Canonical inversion model - circular chromosomes

(Unichromosomal genomes \equiv chromosomes)

Given two canonical circular chromosomes $\mathbb A$ and $\mathbb B,\ldots$

Canonical Inversion Distance Problem:	Compute the minimum number of inversions required to transform \mathbb{A} into \mathbb{B} .
	Denote by $d_{\rm INV}(\mathbb{A},\mathbb{B})$ the inversion distance of \mathbb{A} and $\mathbb{B}.$
Canonical Inversion Sorting Problem:	Find a sequence of $d_{INV}(\mathbb{A}, \mathbb{B})$ inversions that transform \mathbb{A} into \mathbb{B} .

Breakpoint diagram of canonical circular chromosomes

Let \mathbb{A} and \mathbb{B} be canonical circular chromosomes, with $n = |\mathcal{G}_{\star}|$. The **breakpoint diagram** $BD(\mathbb{A}, \mathbb{B}) = (V, E)$ is described as follows:

1.
$$V = \bigcup_{\mathbf{X} \in \mathcal{G}_{\star}} {\mathbf{X}^{h}, \mathbf{X}^{t}} \Rightarrow V = \xi(\mathbb{A}) = \xi(\mathbb{B}) ; |V| = 2n$$

there is a vertex for each extremity of each gene in \mathcal{G}_{\star}

each vertex v has a label $\ell(v)$, that corresponds to the extremity it represents

The vertices are drawn in one line, next to each other.

The vertices must follow the same (circular) order of the corresponding extremities in chromosome \mathbb{A} , according to one of the two reading directions.

2. $E = E_{\Gamma}(\mathbb{A}) \cup E_{\Gamma}(\mathbb{B})$, where:

► Adjacency edges:
$$\begin{cases} E_{\Gamma}(\mathbb{A}) = \{uv : u, v \in V(\xi(\mathbb{A})) \text{ and } \ell(u)\ell(v) \in \Gamma(\mathbb{A})\} \\ E_{\Gamma}(\mathbb{B}) = \{uv : u, v \in V(\xi(\mathbb{B})) \text{ and } \ell(u)\ell(v) \in \Gamma(\mathbb{B})\} \end{cases}$$

The number of edges is |E| = 2n (*n* adjacency edges per chromosome)

Two equivalent breakpoint diagrams

 $BD(\mathbb{A},\mathbb{B}) \cong BD(\mathbb{B},\mathbb{A})$



Properties of the breakpoint diagram

 $\mathbb{A} = (1\,\bar{7}\,4\,5\,3\,\bar{6}\,\bar{2})$



Every vertex has degree two:

 $BD(\mathbb{A}, \mathbb{B})$ is a collection of (even) cycles (alternating edes in $E_{\Gamma}(\mathbb{A})$ and in $E_{\Gamma}(\mathbb{B})$) cycle with *k* edges: *k*-cycle (always even)

 $\mathcal{C} = \mathsf{set} \mathsf{ of cycles in } BD(\mathbb{A},\mathbb{B})$

 $n = |\mathcal{G}_{\star}| = 7$

If $\mathbb{A} = \mathbb{B}$, $RG(\mathbb{A}, \mathbb{B})$ has only 2-cycles: $2n = 2|\mathcal{C}| \implies n = |\mathcal{C}|$

Otherwise, if $\mathbb{A} \neq \mathbb{B}$:

n > |C|

References

Genome Halving under DCJ Revisited

(Julia Mixtacki)

LNCS, volume 5092, pages 276-286 (2008)

Multichromosomal median and halving problems under different genomic distances

(Eric Tannier, Chunfang Zheng and David Sankoff)