

Recovering a Hidden Hamiltonian Cycle via Linear Programming

Yihong Wu

Department of Statistics and Data Science
Yale University

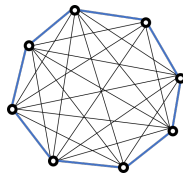
Joint work with
Vivek Bagaria (Stanford), Jian Ding (Penn), David Tse (Stanford) and Jiaming Xu
(Purdue \rightarrow Duke)

Princeton, May 13, 2018

Mathematical problem: Hidden Hamiltonian cycle model

- Observe: a weighted undirected complete graph on n vertices with weighted adjacency matrix W
- Latent: a Hamiltonian cycle C^*
- Edge weight

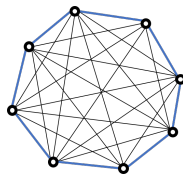
$$W_e \stackrel{\text{ind.}}{\sim} \begin{cases} P & e \in C^* \\ Q & e \notin C^* \end{cases}$$



Mathematical problem: Hidden Hamiltonian cycle model

- Observe: a weighted undirected complete graph on n vertices with weighted adjacency matrix W
- Latent: a Hamiltonian cycle C^*
- Edge weight

$$W_e \stackrel{\text{ind.}}{\sim} \begin{cases} P & e \in C^* \\ Q & e \notin C^* \end{cases}$$

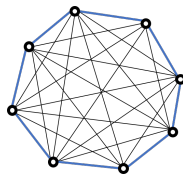


- Goal: observe W , recover C^* with high probability

Mathematical problem: Hidden Hamiltonian cycle model

- Observe: a weighted undirected complete graph on n vertices with weighted adjacency matrix W
- Latent: a Hamiltonian cycle C^*
- Edge weight

$$W_e \stackrel{\text{ind.}}{\sim} \begin{cases} P & e \in C^* \\ Q & e \notin C^* \end{cases}$$



- Goal: observe W , recover C^* with high probability

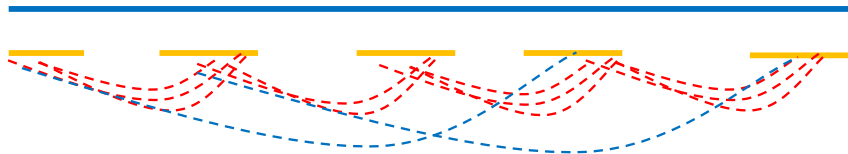
Remarks:

- P, Q depends on the graph size n
- For this talk, $Q = N(0, 1)$ and $P = N(\mu, 1)$, so that

$$W = \mu \cdot \underbrace{\text{adj matrix of } C^*}_{\text{"signal"}} + \text{noise}$$

Link information in Chicago datasets

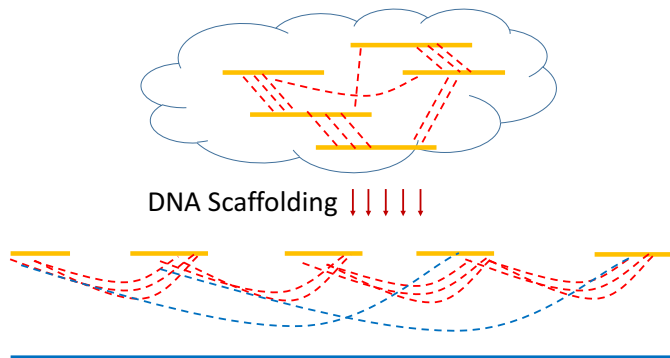
- 1 Reconstitute chromatin in vitro upon naked DNA
- 2 Produce cross-links by fixing chromatin with formaldehyde



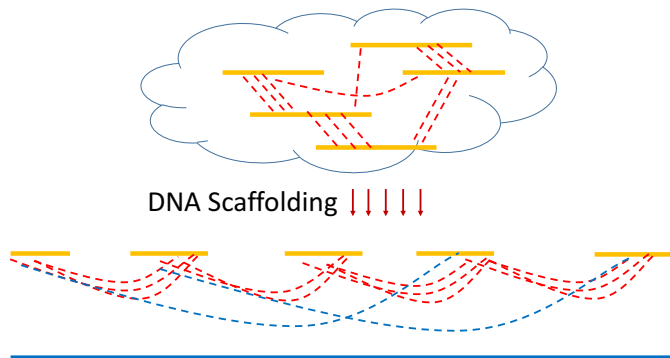
Chicago datasets generate cross-links among contigs [Putnam et al. '16]

On average **more** cross-links exist between **adjacent** contigs

Ordering DNA contigs with Chicago cross-links



Ordering DNA contigs with Chicago cross-links



Reduces to traveling salesman problem (TSP)

Find a path (tour) that visits every contig exactly once with the maximum number of cross-links

Key challenges for DNA scaffolding with Chicago data

- Computational: TSP is NP-hard in the **worst-case**
- Statistical: spurious cross-links between contigs that are far apart

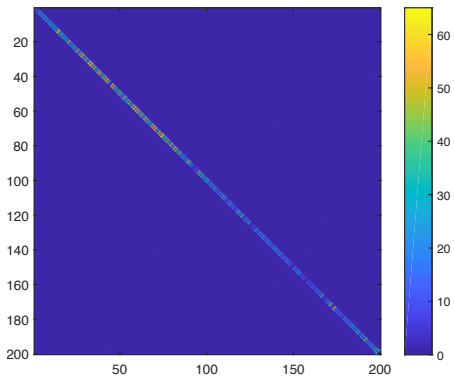
Key challenges for DNA scaffolding with Chicago data

- Computational: TSP is NP-hard in the **worst-case**
- Statistical: spurious cross-links between contigs that are far apart

Key questions:

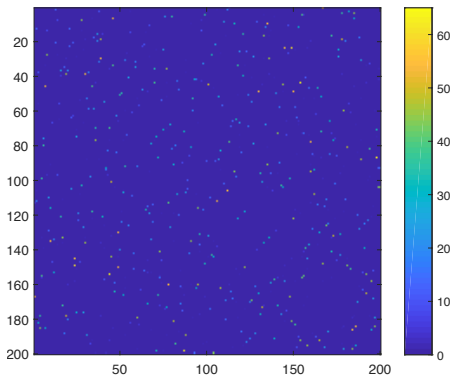
- How to **efficiently** order hundreds of thousands of contigs?
- How much **noise** can be tolerated for accurate DNA scaffolding?

Mathematical model for DNA scaffolding



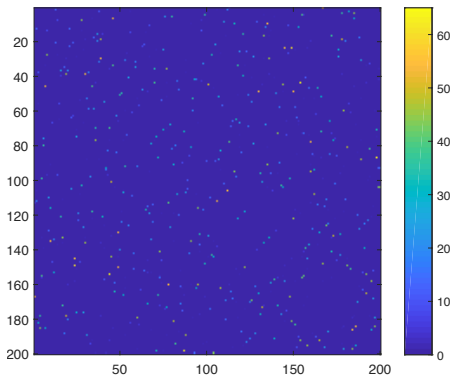
Chicago dataset [Putnam et al. '16]

Mathematical model for DNA scaffolding

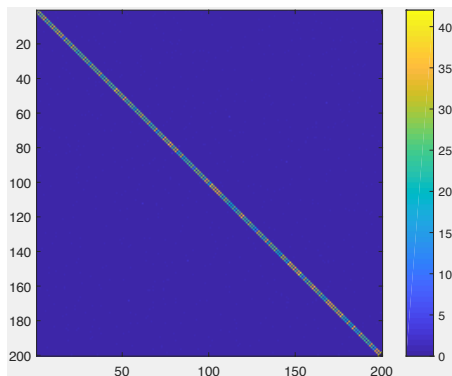


Chicago dataset [Putnam et al. '16]

Mathematical model for DNA scaffolding

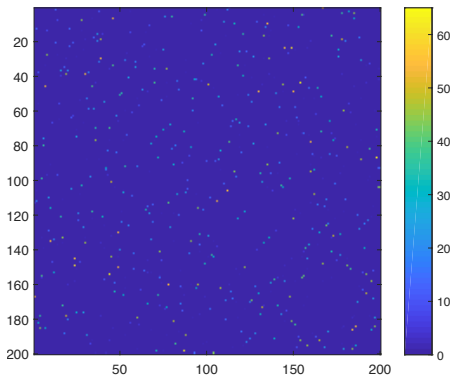


Chicago dataset [Putnam et al. '16]

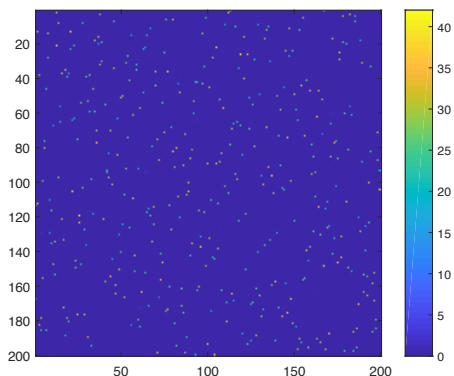


Simulated Poisson data

Mathematical model for DNA scaffolding



Chicago dataset [Putnam et al. '16]



Simulated Poisson data

What is known information-theoretically

Maximum likelihood estimator reduces to TSP

$$\hat{X}_{\text{TSP}} = \arg \max_X \langle L, X \rangle$$

s.t. X is the adjacency matrix of some Hamiltonian cycle

where L is the log likelihood ratio matrix $L_{ij} = \log \frac{dP}{dQ}(W_{ij})$. For Gaussian or Poisson, simply take $L = W$.

What is known information-theoretically

Maximum likelihood estimator reduces to TSP

$$\hat{X}_{\text{TSP}} = \arg \max_X \langle L, X \rangle$$

s.t. X is the adjacency matrix of some Hamiltonian cycle

where L is the log likelihood ratio matrix $L_{ij} = \log \frac{dP}{dQ}(W_{ij})$. For Gaussian or Poisson, simply take $L = W$.

Theorem (Sharp threshold)

If $\mu^2 < 4 \log n$, exact recovery is information-theoretically impossible

If $\mu^2 > 4 \log n$, MLE succeeds in exact recovery

What is known algorithmically

- Spectral methods fails miserably:
 - ▶ $\mu \gg n^{2.5}$ (spectral gap of cycle is too small)

What is known algorithmically

- **Spectral methods** fails miserably:
 - ▶ $\mu \gg n^{2.5}$ (spectral gap of cycle is too small)
- **Thresholding:**
 - ▶ $\mu > \sqrt{8 \log n}$

What is known algorithmically

- **Spectral methods** fails miserably:
 - ▶ $\mu \gg n^{2.5}$ (spectral gap of cycle is too small)
- **Thresholding**:
 - ▶ $\mu > \sqrt{8 \log n}$
- **Greedy merging** [Motahari-Bresler-Tse '13]:
 - ▶ $\mu > \sqrt{6 \log n}$

What is known algorithmically

- **Spectral methods** fails miserably:
 - ▶ $\mu \gg n^{2.5}$ (spectral gap of cycle is too small)
- **Thresholding**:
 - ▶ $\mu > \sqrt{8 \log n}$
- **Greedy merging** [Motahari-Bresler-Tse '13]:
 - ▶ $\mu > \sqrt{6 \log n}$
- This talk: **linear programming** achieves sharp threshold

$$\frac{\mu^2}{\log n} > 4 : \quad \text{LP succeeds}$$
$$\frac{\mu^2}{\log n} < 4 : \quad \text{Everything fails}$$

Threshold are determined by **Rényi divergence** of order $\rho > 0$ from P to Q :

$$D_\rho(P\|Q) \triangleq \frac{1}{\rho - 1} \log \int (dP)^\rho (dQ)^{1-\rho}.$$

- LP works when

$$D_{1/2}(P\|Q) - \log n \rightarrow \infty$$

optimal under mild assumptions

Threshold are determined by **Rényi divergence** of order $\rho > 0$ from P to Q :

$$D_\rho(P\|Q) \triangleq \frac{1}{\rho - 1} \log \int (dP)^\rho (dQ)^{1-\rho}.$$

- LP works when

$$D_{1/2}(P\|Q) - \log n \rightarrow \infty$$

optimal under mild assumptions

- Thresholding works when

$$D_{1/2}(P\|Q) - 2 \log n \rightarrow \infty$$

- Greedy works when

$$D_{1/3}(Q\|P) - \log n \rightarrow \infty$$

Convex relaxations of TSP

$$\begin{aligned}\hat{X}_{\text{TSP}} &= \arg \max_X \langle W, X \rangle \\ \text{s.t.} \quad & \sum_j X_{ij} = 2, \quad \forall i \\ & X_{ij} \in \{0, 1\} \\ & \sum_{i \in I, j \notin I} X_{ij} \geq 2, \quad \forall \emptyset \neq I \subset [n]\end{aligned}$$

$$\begin{aligned}\hat{X}_{\text{TSP}} &= \arg \max_X \langle W, X \rangle \\ \text{s.t.} \quad & \sum_j X_{ij} = 2, \quad \forall i \\ & X_{ij} \in \{0, 1\} \\ & \sum_{i \in I, j \notin I} X_{ij} \geq 2, \quad \forall \emptyset \neq I \subset [n]\end{aligned}$$

- The last constraint: subtour elimination

$$\begin{aligned}\hat{X}_{\text{SUB}} &= \arg \max_X \langle W, X \rangle \\ \text{s.t.} \quad & \sum_j X_{ij} = 2, \quad \forall i \\ & X_{ij} \in [0, 1] \\ & \sum_{i \in I, j \notin I} X_{ij} \geq 2, \quad \forall \emptyset \neq I \subset [n]\end{aligned}$$

$$\begin{aligned} \hat{X}_{\text{SUB}} &= \arg \max_X \langle W, X \rangle \\ \text{s.t. } \sum_j X_{ij} &= 2, \quad \forall i \\ X_{ij} &\in [0, 1] \\ \sum_{i \in I, j \notin I} X_{ij} &\geq 2, \quad \forall \emptyset \neq I \subset [n] \end{aligned}$$

- Replacing the integrality constraint with box constraint: **SUBTOUR LP** relaxation [Dantzig-Fulkerson-Johnson '54, Held-Karp '70]
- Exponentially many linear constraints, nevertheless solvable using interior point method

$$\begin{aligned}\hat{X}_{\text{F2F}} &= \arg \max_X \langle W, X \rangle \\ \text{s.t.} \quad & \sum_j X_{ij} = 2, \quad \forall i \\ & X_{ij} \in [0, 1]\end{aligned}$$

- Further dropping subtour elimination constraints \implies **Fractional 2-factor (F2F) LP**

$$\begin{aligned} \hat{X}_{\text{F2F}} &= \arg \max_X \langle W, X \rangle \\ \text{s.t.} \quad & \sum_j X_{ij} = 2, \quad \forall i \\ & X_{ij} \in [0, 1] \end{aligned}$$

- Further dropping subtour elimination constraints \implies **Fractional 2-factor (F2F) LP**
- Extensively studied in worst case [Boyd-Carr '99, Schalekamp-Williamson-van Zuylen '14]
 - ▶ The integrality gap $\frac{2F}{\text{F2F}} \leq \frac{4}{3}$ for **metric TSP** (min formulation)

$$\begin{aligned} \hat{X}_{\text{F2F}} &= \arg \max_X \langle W, X \rangle \\ \text{s.t.} \quad &\sum_j X_{ij} = 2, \quad \forall i \\ &X_{ij} \in [0, 1] \end{aligned}$$

- Further dropping subtour elimination constraints \implies **Fractional 2-factor (F2F) LP**
- Extensively studied in worst case [Boyd-Carr '99, Schalekamp-Williamson-van Zuylen '14]
 - ▶ The integrality gap $\frac{2\text{F}}{\text{F2F}} \leq \frac{4}{3}$ for **metric TSP** (min formulation)
- What is the integrality gap whp in our random instance?

Theorem

If $\mu^2 - 4 \log n \rightarrow \infty$, then $\widehat{X}_{\text{F2F}} = X^$ with high probability.*

Theorem

If $\mu^2 - 4 \log n \rightarrow \infty$, then $\widehat{X}_{\text{F2F}} = X^$ with high probability.*

Remarks

- The integrality gap is 1 whp!
- Achieving the IT-limit $\mu^2 = 4 \log n$

Max-Product Belief Propagation

$$m_{i \rightarrow j}(t) = w_{ij} - 2\text{nd max}_{\ell \neq j} \{m_{\ell \rightarrow i}(t-1)\}$$

$$m_{i \rightarrow j}(0) = w_{ij}$$

After T iterations, for each vertex i , keep the two largest incoming messages $m_{\ell \rightarrow i}(T)$ and delete the rest.

- BP is exact provided the solution is integral [[Bayati-Borgs-Chayes-Zecchina '11](#)]
- It can be shown that $T = O(n^2 \log n)$ whp

Add more constraints to F2F LP

- SDP1 [Cvetković et al '99]: PSD constraint based on second largest eigenvalue of cycle

$$X \preceq \frac{2}{n}J + 2 \cos \frac{2\pi}{n} \left(I - \frac{1}{n}J \right)$$

Add more constraints to F2F LP

- SDP1 [Cvetković et al '99]: PSD constraint based on second largest eigenvalue of cycle

$$X \preceq \frac{2}{n}J + 2 \cos \frac{2\pi}{n} \left(I - \frac{1}{n}J \right)$$

- ▶ provably weaker than Subtour LP [Goemans-Rendl '00]

Add more constraints to F2F LP

- SDP1 [Cvetković et al '99]: PSD constraint based on second largest eigenvalue of cycle

$$X \preceq \frac{2}{n}J + 2 \cos \frac{2\pi}{n} \left(I - \frac{1}{n}J \right)$$

- ▶ provably weaker than Subtour LP [Goemans-Rendl '00]
- SDP2 [Zhao et al '98]: Quadratic Assignment Problem

$$\langle W, X \rangle = \langle W, \Pi \underbrace{X_0}_{\text{fixed cycle}} \Pi^T \rangle = \left\langle W \otimes X_0, \underbrace{\text{vec}(\Pi)\text{vec}(\Pi)^T}_{\text{relax..}} \right\rangle$$

Add more constraints to F2F LP

- SDP1 [Cvetković et al '99]: PSD constraint based on second largest eigenvalue of cycle

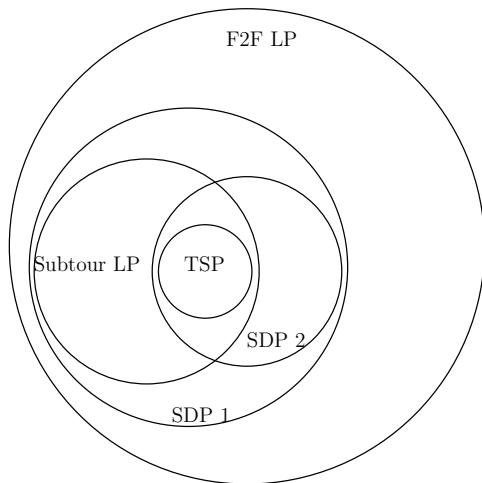
$$X \preceq \frac{2}{n}J + 2 \cos \frac{2\pi}{n} \left(I - \frac{1}{n}J \right)$$

- ▶ provably weaker than Subtour LP [Goemans-Rendl '00]
- SDP2 [Zhao et al '98]: Quadratic Assignment Problem

$$\langle W, X \rangle = \langle W, \Pi \underbrace{X_0}_{\text{fixed cycle}} \Pi^\top \rangle = \left\langle W \otimes X_0, \underbrace{\text{vec}(\Pi)\text{vec}(\Pi)^\top}_{\text{relax..}} \right\rangle$$

- ▶ decision variable: $n^2 \times n^2$ matrix
- ▶ provably stronger than SDP1 [de Klerk et al '08]

Different relaxations



F2F LP succeeds \implies all other relaxations succeed.

Theoretical analysis of convex relaxation

Primal approach vs Dual approach: high level

- Dual argument:
 - ▶ Construct **dual witness** that certify the ground truth whp (KKT conditions)

- Dual argument:
 - ▶ Construct **dual witness** that certify the ground truth whp (KKT conditions)
 - ▶ Successful in proving SDP relaxation attaining sharp threshold for graph partitions: community detection, densest subgraph, etc
[Abbe-Bandeira-Hall '14,Hajek-W-Xu '14,'15,Bandeira '15,Perry-Wein '15]

- Dual argument:
 - ▶ Construct **dual witness** that certify the ground truth whp (KKT conditions)
 - ▶ Successful in proving SDP relaxation attaining sharp threshold for graph partitions: community detection, densest subgraph, etc [Abbe-Bandeira-Hall '14,Hajek-W-Xu '14,'15,Bandeira '15,Perry-Wein '15]
 - ▶ Limitations: construction is **ad hoc**

Primal approach vs Dual approach: high level

- Dual argument:
 - ▶ Construct **dual witness** that certify the ground truth whp (KKT conditions)
 - ▶ Successful in proving SDP relaxation attaining sharp threshold for graph partitions: community detection, densest subgraph, etc [Abbe-Bandeira-Hall '14,Hajek-W-Xu '14,'15,Bandeira '15,Perry-Wein '15]
 - ▶ Limitations: construction is **ad hoc**
- Primal argument:
 - ▶ No feasible solution other than the ground truth has a better objective value whp

Primal approach vs Dual approach: high level

- Dual argument:
 - ▶ Construct **dual witness** that certify the ground truth whp (KKT conditions)
 - ▶ Successful in proving SDP relaxation attaining sharp threshold for graph partitions: community detection, densest subgraph, etc [Abbe-Bandeira-Hall '14,Hajek-W-Xu '14,'15,Bandeira '15,Perry-Wein '15]
 - ▶ Limitations: construction is **ad hoc**
- Primal argument:
 - ▶ No feasible solution other than the ground truth has a better objective value whp
 - ▶ Key: for LP, can restrict to **extremal points** (vertices of the feasible polytope)

- KKT conditions (Farkas' lemma): $\widehat{X}_{\text{F2F}} = X^* \iff \exists u \in \mathbb{R}^n$ (dual certificate):

$$u_i + u_j \leq W_{ij}, \quad \text{for } i \sim j \text{ in } C^*$$

$$u_i + u_j \geq W_{ij}, \quad \text{for } i \not\sim j \text{ in } C^*$$

- KKT conditions (Farkas' lemma): $\widehat{X}_{\text{F2F}} = X^* \iff \exists u \in \mathbb{R}^n$ (dual certificate):

$$u_i + u_j \leq W_{ij}, \quad \text{for } i \sim j \text{ in } C^*$$

$$u_i + u_j \geq W_{ij}, \quad \text{for } i \not\sim j \text{ in } C^*$$

- One feasible choice of dual:

$$u_i = \frac{1}{2} \min\{W_{ij} : j \sim i\}$$

- KKT conditions (Farkas' lemma): $\widehat{X}_{\text{F2F}} = X^* \iff \exists u \in \mathbb{R}^n$ (dual certificate):

$$u_i + u_j \leq W_{ij}, \quad \text{for } i \sim j \text{ in } C^*$$

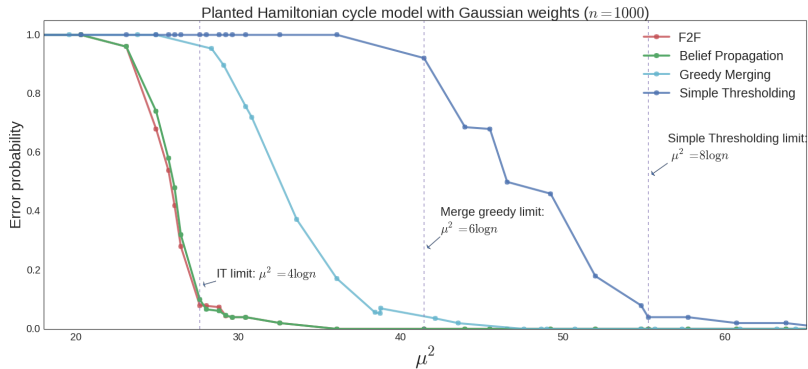
$$u_i + u_j \geq W_{ij}, \quad \text{for } i \not\sim j \text{ in } C^*$$

- One feasible choice of dual:

$$u_i = \frac{1}{2} \min\{W_{ij} : j \sim i\}$$

- This certificate shows correctness if $\mu^2 > 6 \log n$ (same as greedy merging)

Synthetic data experiment



- Show whp for all extremal points $X \neq X^*$:

$$\langle W, X \rangle < \langle W, X^* \rangle$$

- F2F polytope:

$$\left\{ X \in [0, 1]^{n \times n} : \sum_{j=1}^n X_{ij} = 2 \right\}$$

- The proof heavily exploits the characterization of extremal points

- Show whp for all extremal points $X \neq X^*$:

$$\langle W, X \rangle < \langle W, X^* \rangle$$

- F2F polytope:

$$\left\{ X \in [0, 1]^{n \times n} : \sum_{j=1}^n X_{ij} = 2 \right\}$$

- The proof heavily exploits the characterization of extremal points
 - ▶ F2F polytope is not integral: fractional vertices exist

- Show whp for all extremal points $X \neq X^*$:

$$\langle W, X \rangle < \langle W, X^* \rangle$$

- F2F polytope:

$$\left\{ X \in [0, 1]^{n \times n} : \sum_{j=1}^n X_{ij} = 2 \right\}$$

- The proof heavily exploits the characterization of extremal points
 - ▶ F2F polytope is not integral: fractional vertices exist
 - ▶ Characterization [Balinski '65]: for any vertex X of F2F polytope
 - Half integrality

$$X_{ij} \in \{0, 1/2, 1\}$$

- Show whp for all extremal points $X \neq X^*$:

$$\langle W, X \rangle < \langle W, X^* \rangle$$

- F2F polytope:

$$\left\{ X \in [0, 1]^{n \times n} : \sum_{j=1}^n X_{ij} = 2 \right\}$$

- The proof heavily exploits the characterization of extremal points
 - ▶ F2F polytope is not integral: fractional vertices exist
 - ▶ Characterization [Balinski '65]: for any vertex X of F2F polytope
 - Half integrality

$$X_{ij} \in \{0, 1/2, 1\}$$

- 1/2's form disjoint odd cycle connected by path of 1's.

- Show whp for all extremal points $X \neq X^*$:

$$\langle W, X \rangle < \langle W, X^* \rangle$$

- F2F polytope:

$$\left\{ X \in [0, 1]^{n \times n} : \sum_{j=1}^n X_{ij} = 2 \right\}$$

- The proof heavily exploits the characterization of extremal points
 - ▶ F2F polytope is not integral: fractional vertices exist
 - ▶ Characterization [Balinski '65]: for any vertex X of F2F polytope
 - Half integrality
$$X_{ij} \in \{0, 1/2, 1\}$$
 - 1/2's form disjoint odd cycle connected by path of 1's.

Why half integral?

Usual proofs:

- combinatorial proof [Lovasz-Plummer '86, Schrijver '04]
- linear-algebraic proof
 - ▶ F2F polytope (in adjacency vector):

$$\{x \in \mathbb{R}^{\binom{n}{2}} : Ax = 2\mathbf{1}\}$$

- ▶ A is $n \times \binom{n}{2}$ zero-one matrix: $A_{ie} = \mathbf{1}_{\{i \in e\}}$
- ▶ Each column of A has exactly two 1's

Why half integral?

Extremal feasible solution x is of the following form

$$x = \left(\underbrace{x_S}_{\text{fractional}}, \underbrace{x_{S^c}}_{\text{integral}} \right)$$

for some $S \subset \binom{[n]}{[2]}$ of size n , where

- x_S is the solution to the following linear system:

$$A_S x_S = b'$$

Why half integral?

Extremal feasible solution x is of the following form

$$x = \left(\underbrace{x_S}_{\text{fractional}}, \underbrace{x_{S^c}}_{\text{integral}} \right)$$

for some $S \subset \binom{[n]}{2}$ of size n , where

- x_S is the solution to the following linear system:

$$A_S x_S = b'$$

- Cramer's rule:

$$(x_S)_i = \frac{\det(A_S^{(i)})}{\det(A_S)}$$

- ▶ $A_S^{(i)}$ is obtained by substituting the i th column by b' , hence $\det(A_S^{(i)}) \in \mathbb{Z}$.
- ▶ Each column of A_S has two 1's $\implies \det(A_S) \in \{0, \pm 1, \pm 2\}$ [Balinski '65]

Proof of correctness for F2F LP

- 1 Encode the solution: for any extremal point X , represent $2(X - X^*)$ as a **bicolored multigraph** G_X

$$w(G_X) = \langle W, 2(X - X^*) \rangle$$

- 1 Encode the solution: for any extremal point X , represent $2(X - X^*)$ as a **bicolored multigraph** G_X

$$w(G_X) = \langle W, 2(X - X^*) \rangle$$

- 2 Divide and conquer: decompose G_X as edge-disjoint union of graphs in some family \mathcal{F}

$$w(G_X) = \sum_i w(F_i), \quad F_i \in \mathcal{F}$$

- 1 Encode the solution: for any extremal point X , represent $2(X - X^*)$ as a **bicolored multigraph** G_X

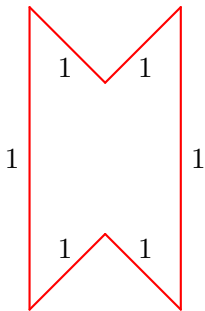
$$w(G_X) = \langle W, 2(X - X^*) \rangle$$

- 2 Divide and conquer: decompose G_X as edge-disjoint union of graphs in some family \mathcal{F}

$$w(G_X) = \sum_i w(F_i), \quad F_i \in \mathcal{F}$$

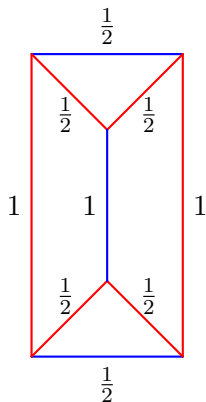
- 3 Counting: Show that whp $w(F) < 0$ for all $F \in \mathcal{F}$

Step 1: Bicolored multigraph representation



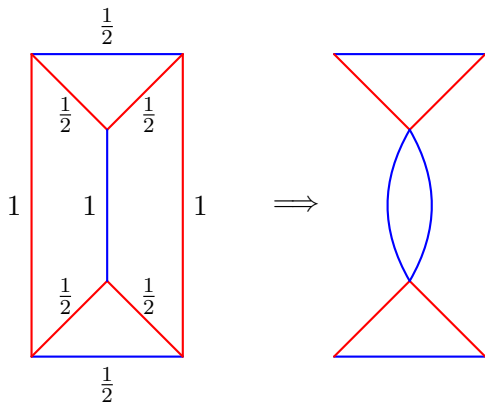
X^* : true cycle

Step 1: Bicolored multigraph representation



X : extremal solution

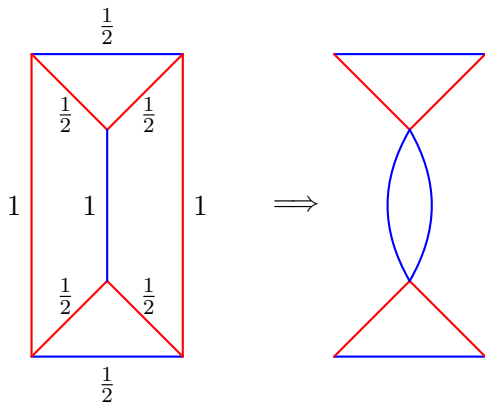
Step 1: Bicolored multigraph representation



X : extremal solution

G_X

Step 1: Bicolored multigraph representation

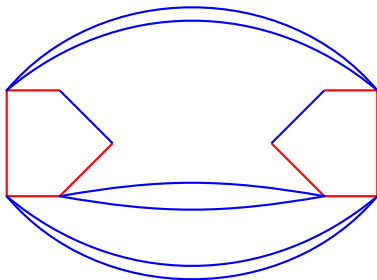
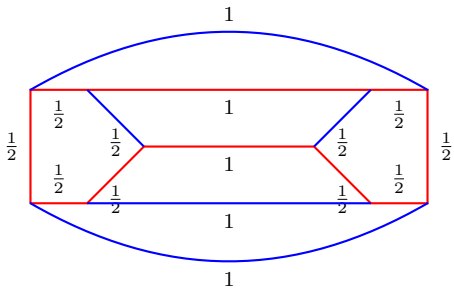


X : extremal solution

G_X

key observation

G_X is always balanced: red degree = blue degree



Step 2: Edge decomposition

Theorem (Kotzig '68)

*Every connected balanced bicolored multigraph has an **alternating Eulerian circuit**.*

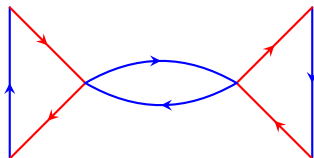
Step 2: Edge decomposition

Theorem (Kotzig '68)

Every connected balanced bicolored multigraph has an *alternating Eulerian circuit*.

Remarks

- An Eulerian circuit may traverse a double edge twice

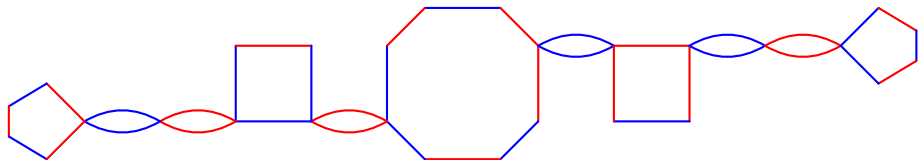


"Dumbbell" structure

Step 2: Edge decomposition

\mathcal{U} : collection of graphs recursively constructed

- ① Start with an even cycle in alternating colors
- ② **Blossoming procedure**: At each step, contract an edge in any cycle and attach a **flower** (path of double edges followed by an alternating odd cycle)

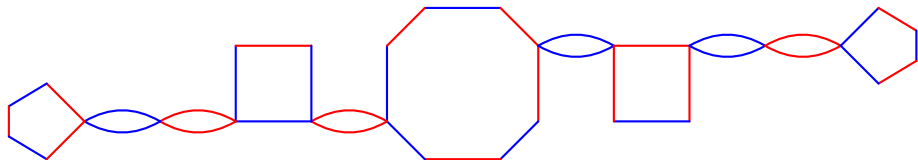


Obtained by starting with an 10-cycle and blossoming 4 times

Step 2: Edge decomposition

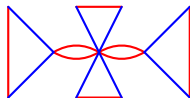
\mathcal{U} : collection of graphs recursively constructed

- 1 Start with an even cycle in alternating colors
- 2 **Blossoming procedure**: At each step, contract an edge in any cycle and attach a **flower** (path of double edges followed by an alternating odd cycle)

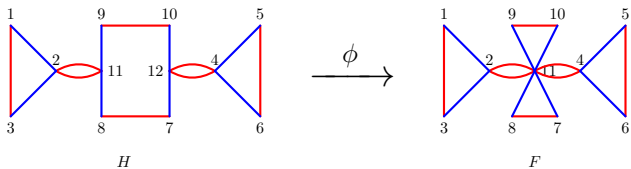


Obtained by starting with an 10-cycle and blossoming 4 times

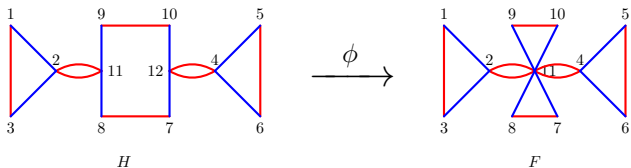
However, not every G_X is of this form...



- **Graph homomorphism** $\phi : H \rightarrow F$ is a vertex map that preserves edges and edge multiplicity



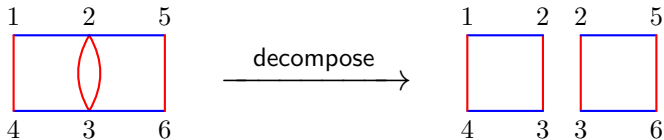
- **Graph homomorphism** $\phi : H \rightarrow F$ is a vertex map that preserves edges and edge multiplicity



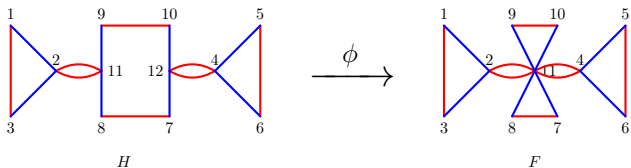
Lemma (Decomposition)

Every balanced bicolored multigraph G with edge multiplicity at most 2 can be decomposed as a union of elements in

$$\mathcal{F} = \{F : V(F) \subset [n], H \rightarrow F \text{ for some } H \in \mathcal{U}\}$$



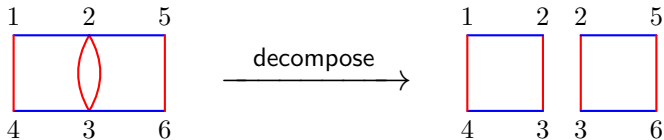
- **Graph homomorphism** $\phi : H \rightarrow F$ is a vertex map that preserves edges and edge multiplicity



Lemma (Decomposition)

Every balanced bicolored multigraph G with edge multiplicity at most 2 can be decomposed as a union of elements in

$$\mathcal{F} = \{F : V(F) \subset [n], H \rightarrow F \text{ for some } H \in \mathcal{U}\}$$



- It remains to show $\min_{F \in \mathcal{F}} w(F) < 0$ whp

Step 3: Counting

$\mathcal{F}_{k,\ell} = \{F \in \mathcal{F} : E(F) \text{ consists of } k \text{ double edges and } \ell \text{ single edges} \}$

Lemma (Counting isomorphism classes)

The number of distinct $H \in \mathcal{U}$ with k double edges and ℓ single edges is at most $C^{k+\ell}$ for universal constant C .

Lemma (Counting homomorphisms)

For each $H \in \mathcal{U}$, there exists $0 \leq r \leq \ell/2$

- Number of labelings for double edges:

$$\leq (Cn)^{k/2+r/2}$$

- Number of labelings for single edges conditioned on double edges

$$\leq (Cn)^{\ell/2-r}$$

Step 4: Probabilistic arguments

$$\mathcal{F}_{k,\ell} = \{F \in \mathcal{F} : E(F) \text{ consists of } k \text{ double edges and } \ell \text{ single edges} \}$$

Lemma

For any $k \geq 0$ and $\ell \geq 3$. With probability at least $1 - n^{-\Theta(k+\ell)}$,

$$\max_{F \in \mathcal{F}_{k,\ell}} (w(F) - \mathbb{E}[w(F)]) \leq (1 + \epsilon) (2k + \ell) \sqrt{\log n}$$

Step 4: Probabilistic arguments

$$\mathcal{F}_{k,\ell} = \{F \in \mathcal{F} : E(F) \text{ consists of } k \text{ double edges and } \ell \text{ single edges} \}$$

Lemma

For any $k \geq 0$ and $\ell \geq 3$. With probability at least $1 - n^{-\Theta(k+\ell)}$,

$$\max_{F \in \mathcal{F}_{k,\ell}} (w(F) - \mathbb{E}[w(F)]) \leq (1 + \epsilon) (2k + \ell) \sqrt{\log n}$$

Remarks

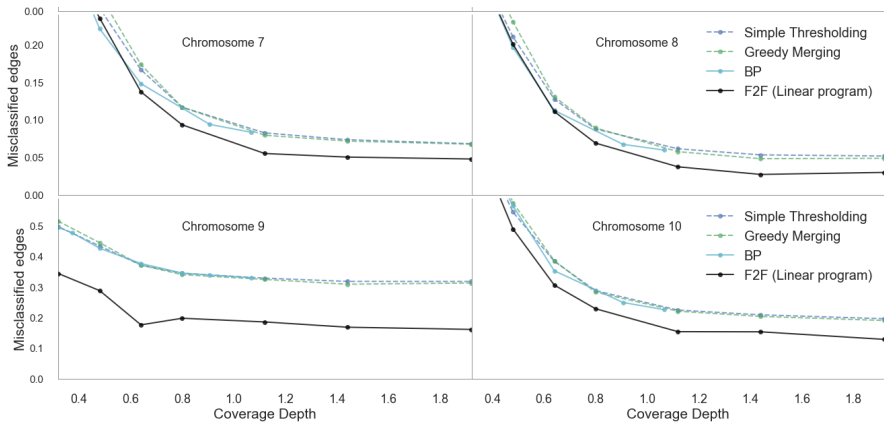
- Total: $2k + \ell$ edges, half red half blue. Weights on red edges $\sim N(\mu, 1)$. Weights on blue edges $\sim N(0, 1)$.

$$w(F) \sim N(-(k + \ell/2)\mu, 4k + \ell)$$

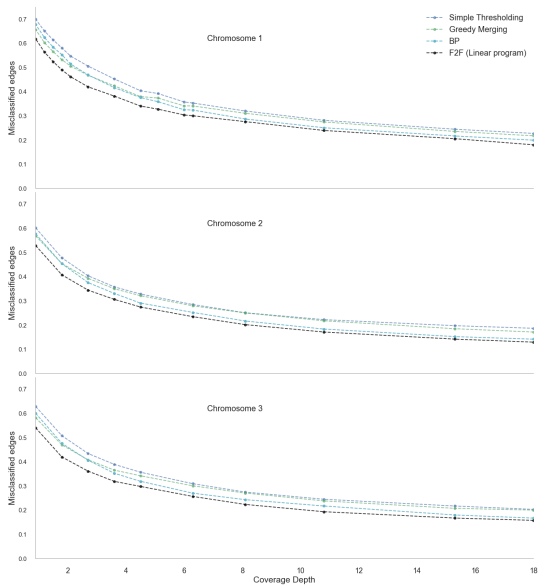
- Proof: Counting $\mathcal{F}_{k,\ell}$ and large deviation bounds

- 1000 DNA contigs of size 100 kbps
- 0.45 million Chicago cross-links
- Subsample each cross-link with probability p

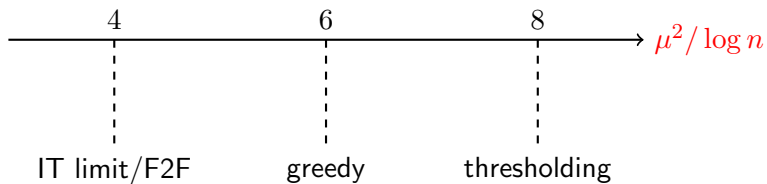
Homosapiens [Putnam et al 16, Genome Research]



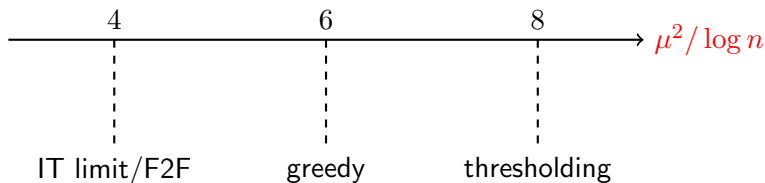
Aedes Aegypti (zika mosquito) [Dudchenko et al '16, Science]



Conclusion and remarks



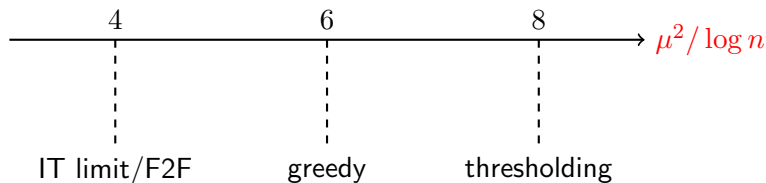
Conclusion and remarks



Future work

- More realistic models
 - ▶ 2-NN graph: IT limit becomes $\sqrt{2 \log n}$ not achieved by LP.

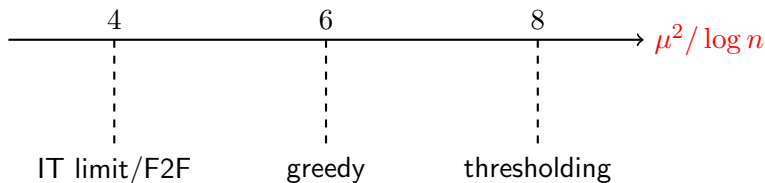
Conclusion and remarks



Future work

- More realistic models
 - ▶ 2-NN graph: IT limit becomes $\sqrt{2 \log n}$ not achieved by LP.
 - ▶ small-world graphs

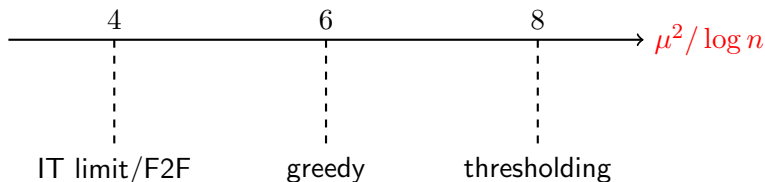
Conclusion and remarks



Future work

- More realistic models
 - ▶ 2-NN graph: IT limit becomes $\sqrt{2 \log n}$ not achieved by LP.
 - ▶ small-world graphs
- Smarter rounding algorithm in practice

Conclusion and remarks



Future work

- More realistic models
 - ▶ 2-NN graph: IT limit becomes $\sqrt{2 \log n}$ not achieved by LP.
 - ▶ small-world graphs
- Smarter rounding algorithm in practice
- Reduction from/to Hamiltonian cycle and path more elegantly

References

- Vivek Bagaria, Jian Ding, David Tse, W. & Jiaming Xu (2018). *Hidden Hamiltonian Cycle Recovery via Linear Programming*, <https://arxiv.org/abs/1804.05436>