

Convergence of an entropic Frank-Wolfe method

Nonlinear Approximation for High Dimensional Problems 2025

James Ashton Nichols (james.ashton.nichols@gmail.com)

July 3, 2025



Australian
National
University

Biological problems - IQ-Tree and COVID Variants



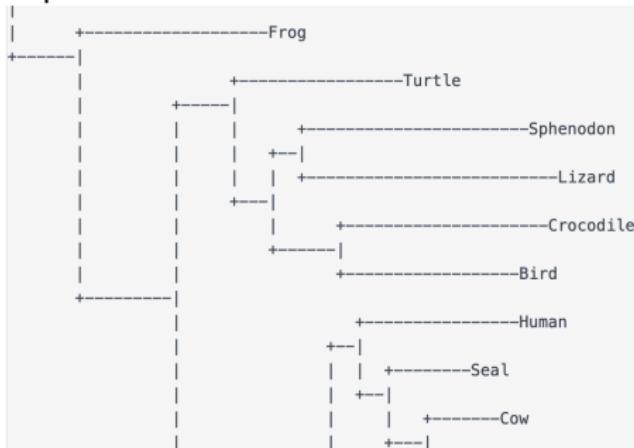
IQ-TREE

- Main “Phylogenetic inference” software.
- Used to distinguish COVID variants (Delta, Omicron...).
- Developed at ANU and Vienna University

Input:

```
7 28
Frog      AAATTGGTCCCTGTGATTCAAGCAGTGAT
Turtle     CTTCCACACCCCAGGGACTCAGCAGTGAT
Bird       CTACCACACCCCAGGGAAACAGCAGTGAT
Human      CTACCACACCCCAGGGAAACAGCAGTGAT
Cow        CTACCACACCCCAGGGAAACAGCAGTGAT
Whale      CTACCACGCCCCAGGACACAGCAGTGAT
Mouse      CTACCACACCCCAGGGACTCAGCAGTGAT
```

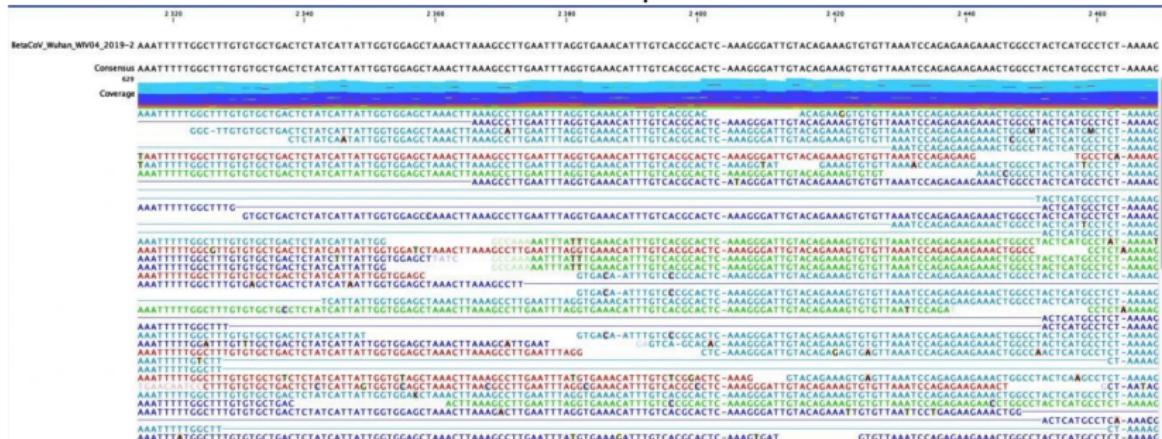
Output:



Biological problems - IQ-tree and COVID variants

- IQ-Tree ver1: 37 hours to infer a tree for 17,000 COVID sequences
 - IQ-Tree ver2 (published May 2020): reduced this to three minutes
 - At pandemic peak, 10,000 sequences were published per day

COVID genome as sequence at Institute L. Pasteur, January 24, 2020
~ 10^4 base pairs



The Frank-Wolfe Method

[FRANK & WOLFE, 1965]

To find $\min_{x \in K} f(x)$ where $K \subset \mathbb{R}^d$ is compact and convex, F.W. takes an iterative approach based on gradients

The Frank-Wolfe Method

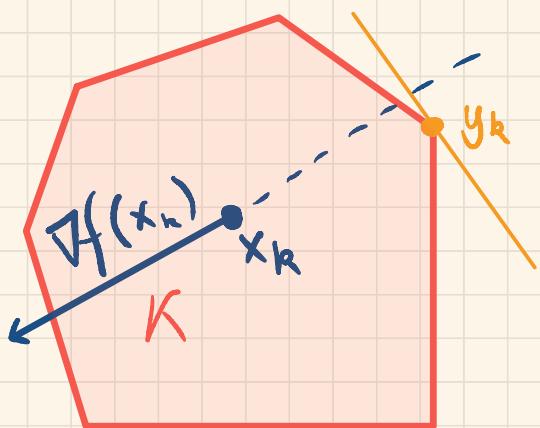
[FRANK & WOLFE, 1965]

To find $\min_{x \in K} f(x)$ where $K \subset \mathbb{R}^d$ is compact and convex, F.W. takes an iterative approach based on gradients

ALGORITHM (F.W.)

- Given $x_k \in K$, find

$$y_k = \operatorname{argmin}_{y \in K} \langle y, \nabla f(x_k) \rangle$$



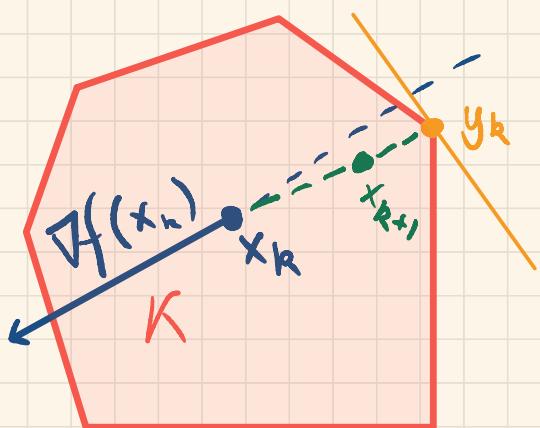
The Frank-Wolfe Method

[FRANK & WOLFE, 1965]

To find $\min_{x \in K} f(x)$ where $K \subset \mathbb{R}^d$ is compact and convex, F.W. takes an iterative approach based on gradients

ALGORITHM (F.W.)

- Given $x_k \in K$, find
 $y_k = \operatorname{argmin}_{y \in K} \langle y, \nabla f(x_k) \rangle$
- Set $x_{k+1} = x_k + \gamma_k (y_k - x_k)$



The Frank-Wolfe Method

[FRANK & WOLFE, 1965]

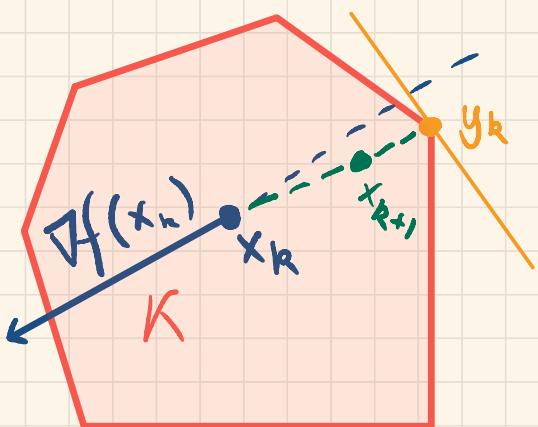
To find $\min_{x \in K} f(x)$ where $K \subset \mathbb{R}^d$ is compact and convex, F.W. takes an iterative approach based on gradients

ALGORITHM (F.W.)

- Given $x_k \in K$, find
 $y_k = \operatorname{argmin}_{y \in K} \langle y, \nabla f(x_k) \rangle$

- Set $x_{k+1} = x_k + \gamma_k (y_k - x_k)$

where
OPTION 1: $\gamma_k = \frac{2}{k+2}$
OPTION 2 (LINE-SEARCH): γ_k minimises $f(x_k + \gamma_k (y_k - x_k))$



Linear Minimisation

One drawback of FW: we must solve

$$\underset{y \in K}{\operatorname{argmin}} \langle y, \nabla f(x_k) \rangle$$

- If K is a convex polytope, we use linear programming, $O(d^3)$ complexity in \mathbb{R}^d
- Otherwise we seek an approximate linear solver

Regularising with entropy

Assume now we're dealing with **probability vectors**, i.e.

$$K \subseteq \mathbb{S}_d := \left\{ x \in [0,1]^d : \sum_{i=1}^d x_i = 1 \right\}$$

We have an **approximate linear program**

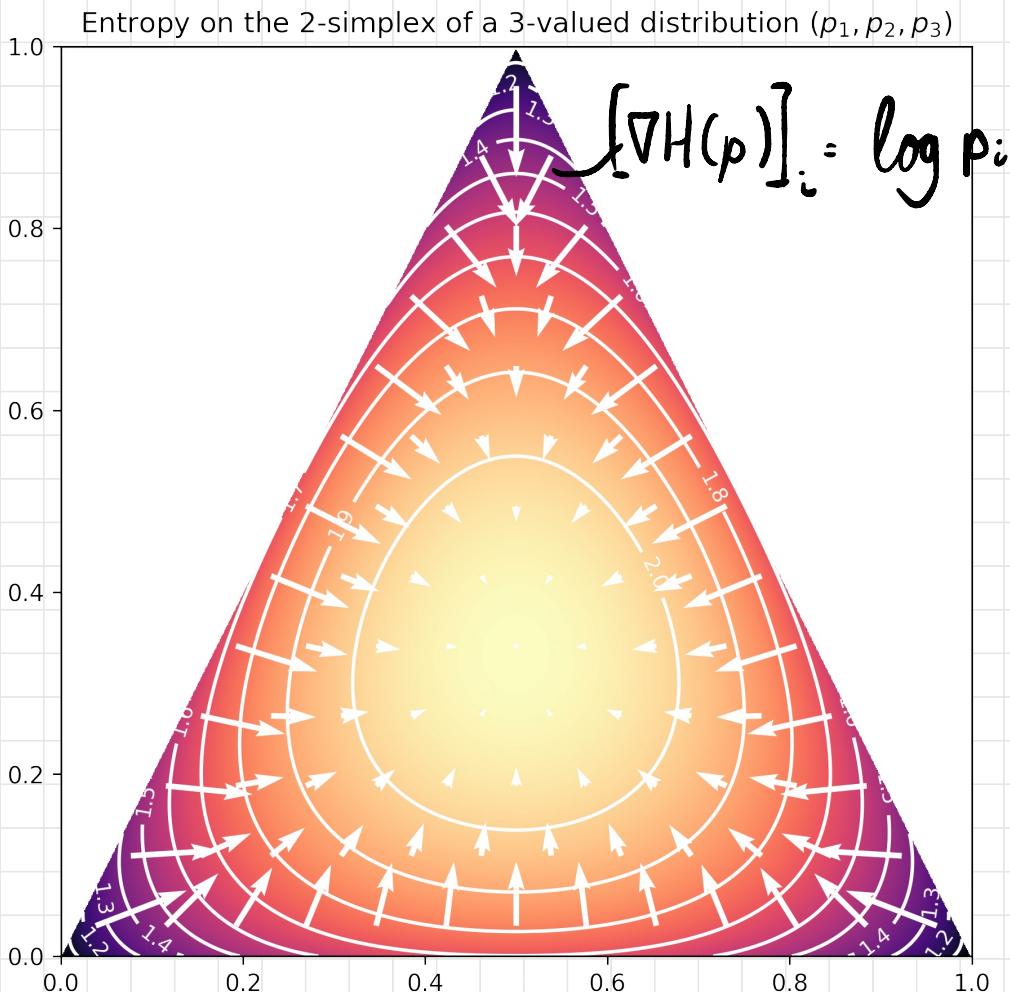
$$\min_{x \in K} \langle x, c \rangle \approx \min_{x \in K} \langle x, c \rangle - \varepsilon H(x)$$

where $\varepsilon > 0$ is small and **entropy** H is

$$H(x) := -\sum_{i=1}^d x_i \log x_i$$

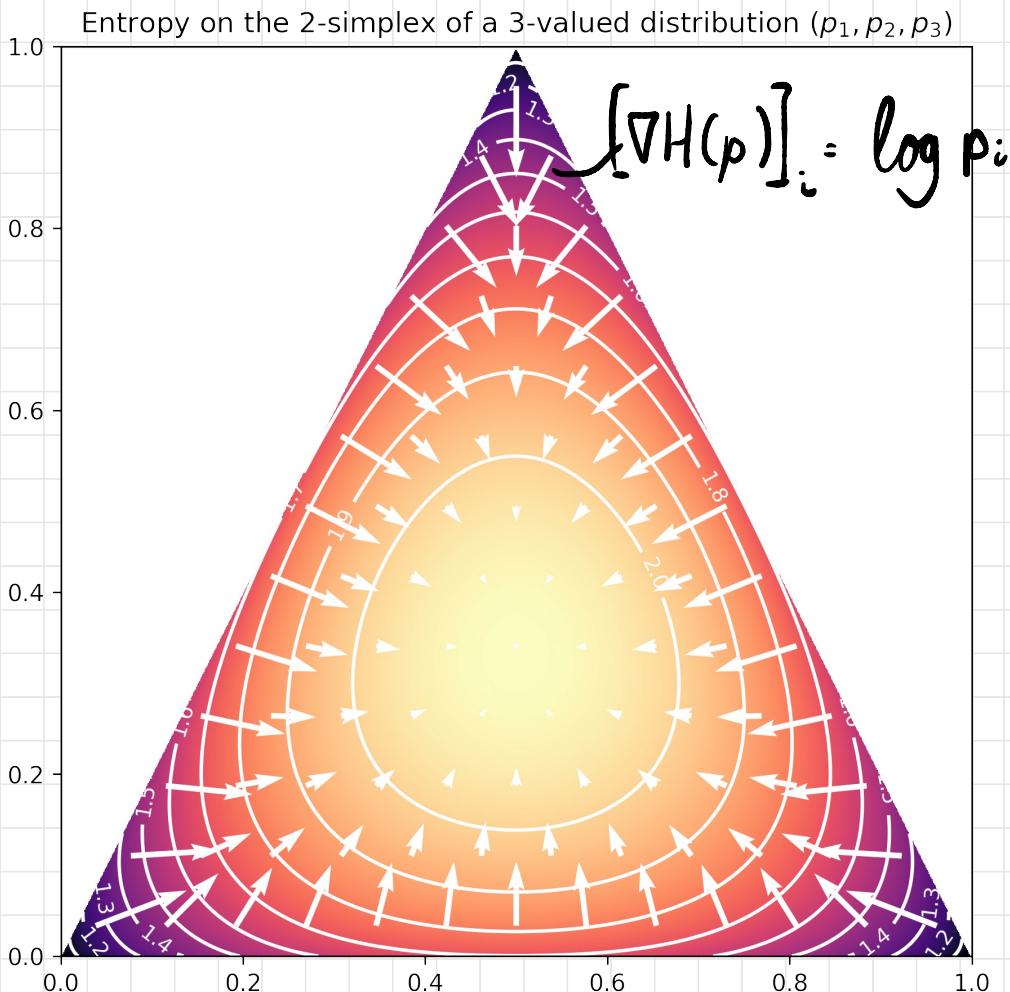
DH Illustrated

- $\nabla H: \sum_d \rightarrow \mathbb{R}^{d-1}$
is a bijection
- $\nabla H(x) = c/\varepsilon$
always has a solution
- $\|\nabla H(x)\| \rightarrow \infty$
on $\partial \sum_d$



DH Illustrated

- $\nabla H: K \rightarrow \mathbb{R}^{\dim(K)}$
is a bijection
- $\nabla H(x) = c/\varepsilon$
always has a solution
- $\|\nabla H(x)\| \rightarrow \infty$
on ∂K



Soft-max & Duality

Log-sum-exp

"Soft-max"

$$h(c) := \log \sum_{i=1}^d e^{c_i}$$
$$[\nabla h(c)]_i = \frac{e^{c_i}}{\sum_{i=1}^d e^{c_i}} = e^{c_i - h(c)}$$

```
c = np.array([0.5, 1, 14, 2])
log_sum_exp = np.log(np.exp(c).sum())
softmax = np.exp(c - log_sum_exp)
print(f'c = {c}'); print(f'h(c) = {log_sum_exp:.5f}'); print(f'nab_h(c) = {softmax}')
✓ 0.0s

c = [ 0.5  1.  14.  2. ]
h(c) = 14.00001
nab_h(c) = [1.3709e-06 2.2603e-06 9.9999e-01 6.1442e-06]
```

Note: $\nabla h(c) \in \mathbb{R}^d$ for any $c \in \mathbb{R}^d$

Soft-max & Duality

Log-sum-exp

"Soft-max"

$$h(c) := \log \sum_{i=1}^d e^{c_i}$$
$$[\nabla h(c)]_i = \frac{e^{c_i}}{\sum_{i=1}^d e^{c_i}} = e^{c_i} - h(c)$$

Log-sum-exp is convex dual to entropy
(FENCHEL-LEGENDRE)

$$h^*(x) := \sup_{c \in \mathbb{R}^d} \langle x, c \rangle - h(c) = -H(x)$$

[BOYD & VANDENBERGHE, 2004]

Soft-max & Duality

Log-sum-exp

"Soft-max"

$$h(c) := \log \sum_{i=1}^d e^{c_i}$$
$$[\nabla h(c)]_i = \frac{e^{c_i}}{\sum_{i=1}^d e^{c_i}} = e^{c_i - h(c)}$$

Log-sum-exp is convex dual to entropy
(FENCHEL-LEGENDRE)

$$h^*(x) := \sup_{c \in \mathbb{R}^d} \langle x, c \rangle - h(c) = -H(x)$$

[BOYD & VANDENBERGHE, 2004]

And gradients invert!

$$\nabla h(\nabla h^*(x)) = x \quad \nabla h^*(\nabla h(c)) = "c"$$
$$\nabla h: \mathbb{R}^d \rightarrow \Sigma_d \quad \nabla h^*: \Sigma^d \rightarrow \mathbb{R}^d$$

Solving entropic LP

[ERLANDER, 1981]
[CLASEN, 1965]

$$x_{c/\epsilon}^* = \underset{x \in \Sigma_d}{\operatorname{arg\,min}} \quad \langle x, c \rangle - \epsilon H(x)$$

$$\Rightarrow \nabla H(x_{c/\epsilon}^*) = c/\epsilon$$

$$\Rightarrow x_{c/\epsilon}^* = \nabla h(-c/\epsilon) = e^{-c/\epsilon} - h(-c/\epsilon)$$

Solving entropic LP

[ERLANDER, 1981]
[CLASEN, 1965]

$$x_{c/\epsilon}^* = \underset{x \in \Sigma_d}{\operatorname{arg\,min}} \quad \langle x, c \rangle - \epsilon H(x)$$

$$\Rightarrow \nabla H(x_{c/\epsilon}^*) = c/\epsilon$$

$$\Rightarrow x_{c/\epsilon}^* = \nabla h(-c/\epsilon) = e^{-c/\epsilon} - h(-c/\epsilon)$$

When satisfying constraints it's a little trickier.

Lagrangian: $\lambda(x, \lambda) = \langle x, c \rangle - \epsilon H(x) + \langle Ax - b, \lambda \rangle$

$$\nabla_\lambda \lambda = 0 \Rightarrow x_{c/\epsilon}^* = e^{-c/\epsilon + A^T \lambda} - h(-c/\epsilon + A^T \lambda)$$

$$x \in K := \{x \in \Sigma_d : Ax = b\}$$

Solving entropic LP

[ERLANDER, 1981]
[CLASEN, 1965]

$$x_{c/\epsilon}^* = \underset{x \in \Sigma_d}{\operatorname{arg\,min}} \quad \langle x, c \rangle - \epsilon H(x)$$

$$\Rightarrow \nabla H(x_{c/\epsilon}^*) = c/\epsilon$$

$$\Rightarrow x_{c/\epsilon}^* = \nabla h(-c/\epsilon) = e^{-c/\epsilon} - h(-c/\epsilon)$$

When satisfying constraints it's a little trickier.

Lagrangian: $\lambda(x, \lambda) = \langle x, c \rangle - \epsilon H(x) + \langle Ax - b, \lambda \rangle$

$$\nabla_\lambda \lambda = 0 \Rightarrow x_{c/\epsilon}^* = e^{-c/\epsilon + A^T \lambda} - h(-c/\epsilon + A^T \lambda)$$

$$K := \{x \in \Sigma_d : Ax = b\}$$

How do we find λ ?

Satisfying constraints

$$K := \{x \in \Sigma_d : Ax = b\}$$

$$x_{c/\epsilon}^* = \underset{x \in K}{\operatorname{argmin}} \langle x, c \rangle - \epsilon H(x) = e^{-c/\epsilon + A^T \lambda} - h(-c/\epsilon + A^T \lambda)$$

- "Newton-Kantorovich": Solve $A e^{-c/\epsilon + A^T \lambda} - h(-c/\epsilon + A^T \lambda) = b$ for λ with Newton's method
- Sinkhorn's method: If $A_{ij} \in \{0, 1\}$ we can scale $x_{c/\epsilon}^*$ in marginals, iteratively setting each λ_i in a loop.
[SINKHORN, 1965] [CUTURI, 2013]

A RELATIONSHIP BETWEEN ARBITRARY POSITIVE MATRICES AND DOUBLY STOCHASTIC MATRICES

BY RICHARD SINKHORN

University of Houston

1. Introduction. Suppose one observes n transitions of a Markov chain with N states and stochastic matrix $P = (p_{ij})$. The usual estimate of p_{ij} is $t_{ij} = a_{ij}/\lambda_i$ where a_{ij} is the number of transitions from i to j which are observed, and $\lambda_i = \sum_j a_{ij}$. (Cf. [1].) This amounts to a normalization of the rows of $A = (a_{ij})$, and can be expressed as a matrix equation $T = D_1 A$ where $T = (t_{ij})$ and $D_1 = \text{diag}[\lambda_1^{-1}, \dots, \lambda_N^{-1}]$.

Solving O.T.

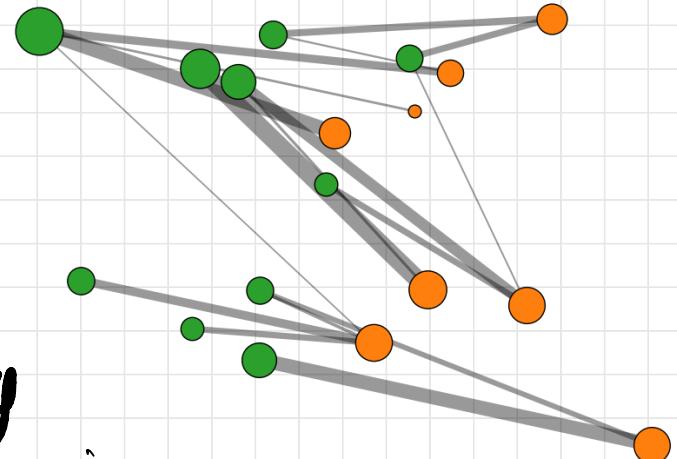
$$\min_{P \in V(\alpha, \beta)} \langle P, C \rangle_F$$

is a linear programming problem:

- Objective is linear
- Constraint set $V(\alpha, \beta)$ is a linear polytope

[KANTOROVIC, 1939]

But LP solvers will be $\sim O((mn)^3)$!



Entropic optimal transport

$$\min_{P \in V(\alpha, \beta)} \langle P, C \rangle_F - \varepsilon H(P)$$

Sinkhorn Distances:
Lightspeed Computation of Optimal Transport

[CUTURI, 2013]

Marco Cuturi

Entropic optimal transport

Sinkhorn Distances:
Lightspeed Computation of Optimal Transport

[CUTURI, 2013]

Marco Cuturi

$$\max_{\lambda \in \mathbb{R}^n} \min_{P \in \mathbb{R}^{n \times m}} \langle P, C \rangle_F - \varepsilon H(P) - \langle P \mathbf{1}_m - \alpha, \lambda \rangle - \langle P^T \mathbf{1}_n - \beta, \theta \rangle$$

$= \mathcal{L}(P, \lambda, \theta)$

Entropic optimal transport

Sinkhorn Distances:
Lightspeed Computation of Optimal Transport

[CUTURI, 2013]

Marco Cuturi

$$\max_{\lambda \in \mathbb{R}^n, \theta \in \mathbb{R}^m} \min_{P \in \mathbb{R}^{n \times m}} \langle P, C \rangle_F - \varepsilon H(P) - \langle P \mathbf{1}_m - \alpha, \lambda \rangle - \langle P^T \mathbf{1}_n - \beta, \theta \rangle$$

$= h(P, \lambda, \theta)$

• Critical point $\nabla_{\lambda, \theta} h(P, \lambda, \theta) = 0 \Rightarrow P_{ij} = e^{\frac{\lambda_i}{\varepsilon}} e^{-\frac{C_{ij}}{\varepsilon}} e^{\frac{\theta_j}{\varepsilon}}$

Entropic optimal transport

Sinkhorn Distances:
Lightspeed Computation of Optimal Transport

$$\max_{\lambda \in \mathbb{R}^n, \theta \in \mathbb{R}^m} \min_{P \in \mathbb{R}^{n \times m}} \underbrace{\langle P, C \rangle_F - \varepsilon H(P) - \langle P \mathbf{1}_m - \alpha, \lambda \rangle - \langle P^T \mathbf{1}_n - \beta, \theta \rangle}_{= h(P, \lambda, \theta)}$$

- Critical point $\nabla_{\lambda, \theta} h(P, \lambda, \theta) = 0 \Rightarrow P_{ij} = e^{\frac{\lambda_i}{\varepsilon}} e^{\frac{C_{ij}}{\varepsilon}} e^{\frac{\theta_j}{\varepsilon}}$
- To satisfy the marginals, $P \mathbf{1}_m = \alpha$ & $P^T \mathbf{1}_n = \beta$, we use Sinkhorn's method.

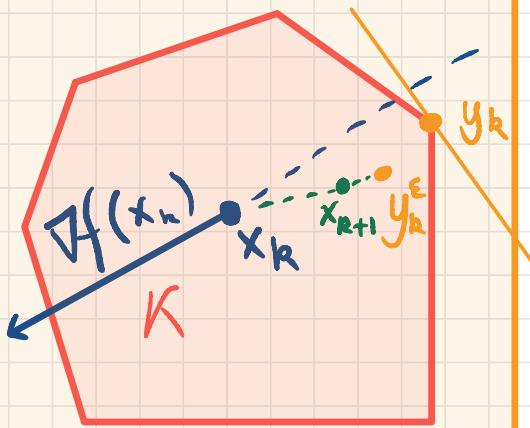
- ① Set u s.t. $P \mathbf{1}_m = \alpha$
- ② Set v s.t. $P^T \mathbf{1}_n = \beta$
- ③ Repeat. $P \rightarrow P_\varepsilon^*$

$$\begin{array}{c} e^{\theta_i/\varepsilon} = v \\ \odot \\ \begin{array}{c} e^{\lambda_i/\varepsilon} \\ \odot \\ u \end{array} \end{array} \xrightarrow{e^{-C_{ij}/\varepsilon} = K_{ij}} \begin{array}{c} \oplus \\ \beta \end{array} \Rightarrow \alpha$$

Entropic Frank-Wolfe

ALGORITHM [N. Q STONE, 2025]

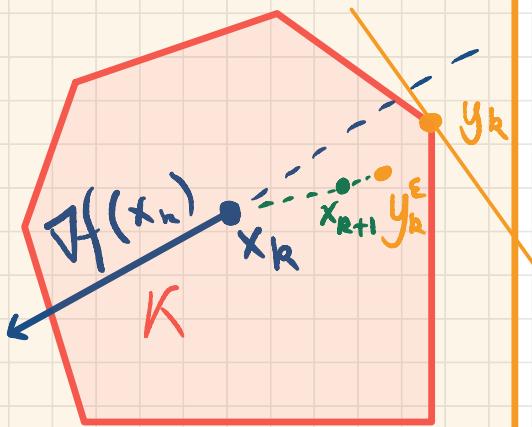
1. Solve $y_k^\epsilon = \operatorname{argmin}_{y \in K} \langle y, \nabla f(x_k) \rangle - \epsilon H(y)$



Entropic Frank-Wolfe

ALGORITHM [N. Q STONE, 2025]

1. Solve $y_k^\epsilon = \operatorname{argmin}_{y \in K} \langle y, \nabla f(x_k) \rangle - \epsilon H(y)$
2. If $\langle y_k^\epsilon - x_k, \nabla f(x_k) \rangle \geq 0$, halve ϵ ,
repeat step 1.



Entropic Frank-Wolfe

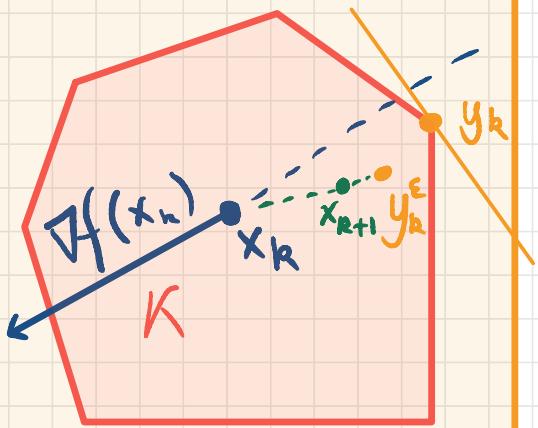
ALGORITHM [N. Q STONE, 2025]

1. Solve $y_k^\epsilon = \operatorname{argmin}_{y \in K} \langle y, \nabla f(x_k) \rangle - \epsilon H(y)$

2. If $\langle y_k^\epsilon - x_k, \nabla f(x_k) \rangle \geq 0$, halve ϵ ,
repeat step 1.

3. Set $x_{k+1} = x_k + \gamma_k (y_k^\epsilon - x_k)$

where γ_k minimises $f(x_k + \gamma_k (y_k^\epsilon - x_k))$



Entropic Frank-Wolfe

ALGORITHM [N. Q STONE, 2025]

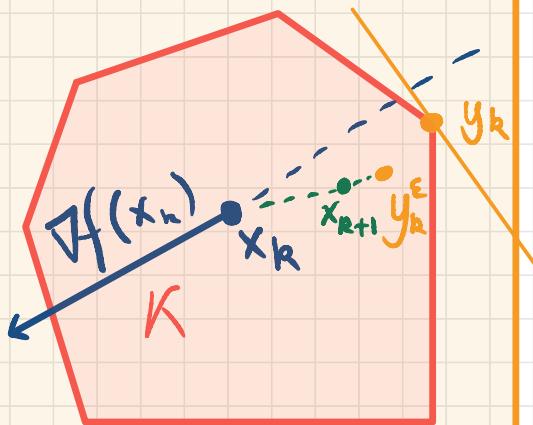
1. Solve $y_k^\epsilon = \operatorname{argmin}_{y \in K} \langle y, \nabla f(x_k) \rangle - \epsilon H(y)$

2. If $\langle y_k^\epsilon - x_k, \nabla f(x_k) \rangle \geq 0$, halve ϵ ,
repeat step 1.

3. Set $x_{k+1} = x_k + \gamma_k (y_k^\epsilon - x_k)$

where γ_k minimises $f(x_k + \gamma(y_k^\epsilon - x_k))$

4. Stop when $\langle \nabla f(x_k), y_k^\epsilon - x_k \rangle \leq \eta$ ($\eta > 0$) STOPPING CRITERIA



Error Bounds

LP:

$$x^* = \operatorname{argmin}_{x \in K} \langle x, c \rangle$$

ENTROPIC LP

$$x_{\epsilon/\epsilon}^* = \operatorname{argmin}_{x \in K} \langle x, c \rangle - \epsilon H(x)$$

We have

$$\bullet 0 \leq \langle x^*, c \rangle - \langle x_{\epsilon/\epsilon}^*, c \rangle \leq \epsilon (H(x_{\epsilon/\epsilon}^*) - H(x^*))$$

[PEYRÉ & CUTURI, 2018]

$$\bullet \| x^* - x_{\epsilon/\epsilon}^* \| \leq e^{-\frac{1}{\epsilon}}$$

[COMINETTI & SAN MARTIN, 1994]

Convergence

① FW: f is convex then $f(x_k) - \min_{x \in K} f(x) \leq C_f/k$

[FRANK & WOLFE 1965] [CLARKSON 2010] [JAGGI 2013]

② FW: f is non-convex then

$$\min_{0 \leq e < k} f(x_e) - f(x_{e+1}) \leq C_f/\sqrt{k}$$

[LACOSTE-JULIEN 2013]

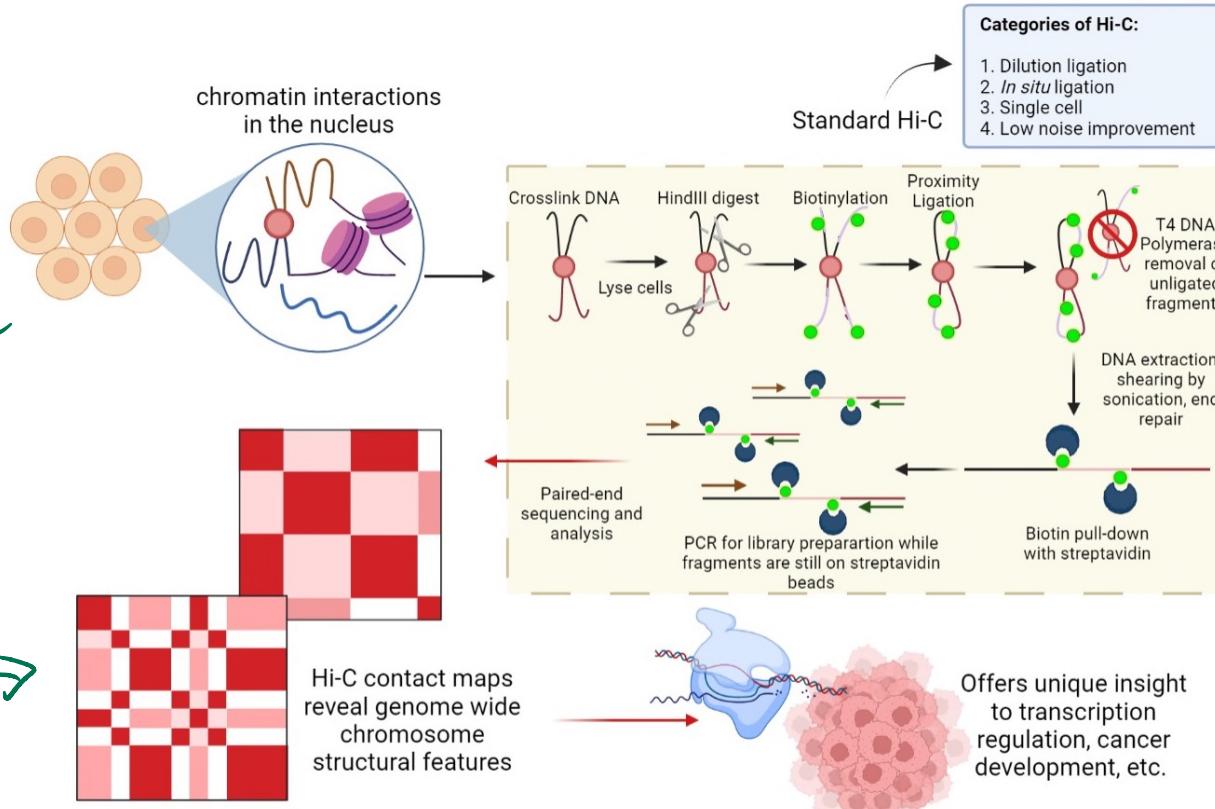
③ ENTROPIC FW: f is non-convex then

$$\min_{0 \leq e < k} f(x_e) - f(x_{e+1}) \leq C_f/\sqrt{k}$$

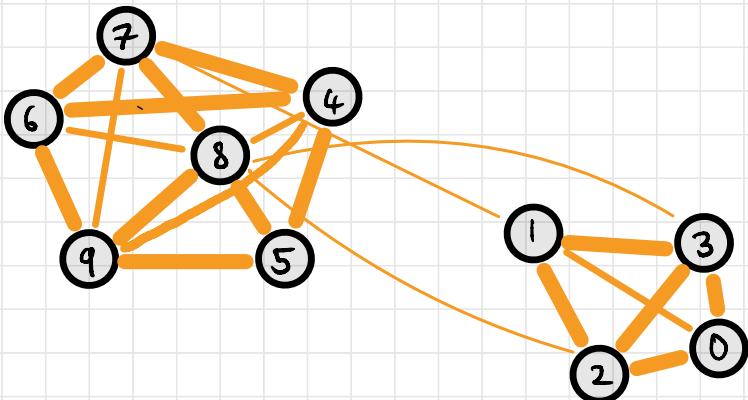
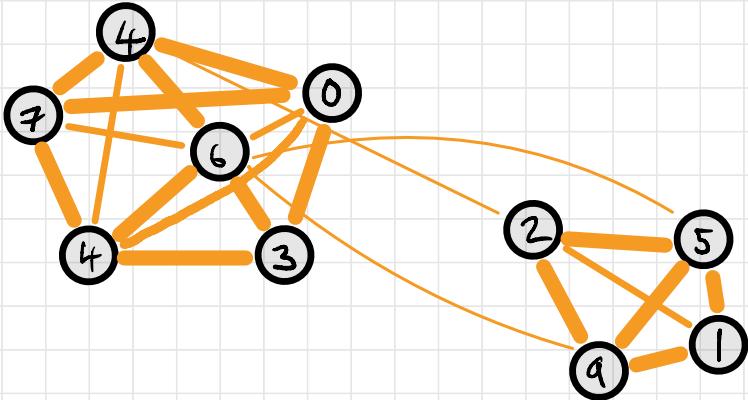
[N. & STONE, 2025]

An application

This come in when



Matrix Reordering

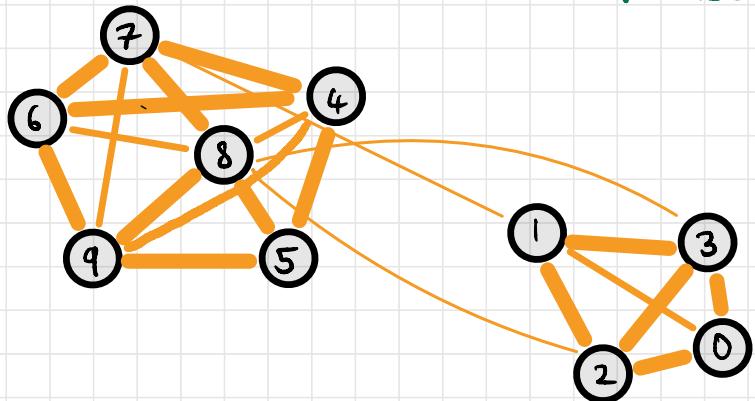
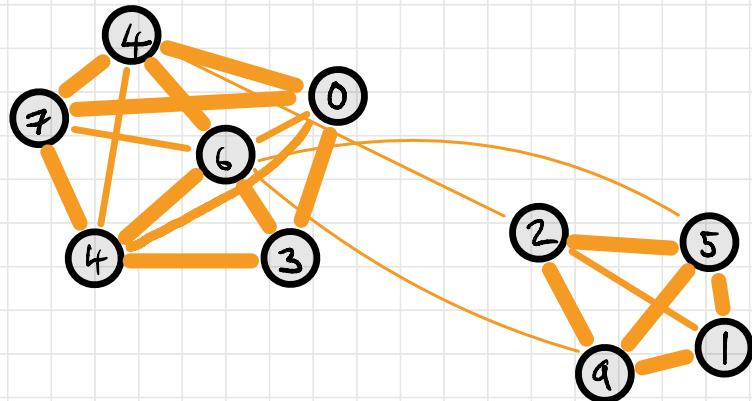


| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|-----|-----|------|-----|-----|------|------|-----|-----|------|
| 0 | 0 | 0 | 0 | 1.4 | 0 | 0 | 1.6 | 2.5 | 1.9 | 1.1 |
| 1 | 0 | 0 | 1.4 | 0 | 3.1 | 3 | 0 | 0 | 0 | 0 |
| 2 | 0 | 1.4 | 0 | 0 | 2.5 | 2.8 | 0 | 0 | 0 | 0.41 |
| 3 | 1.4 | 0 | 0 | 0 | 0 | 0 | 1.8 | 1.2 | 2.3 | 1.4 |
| 4 | 0 | 3.1 | 2.5 | 0 | 0 | 3.2 | 0.6 | 0 | 0 | 0 |
| 5 | 0 | 3 | 2.8 | 0 | 3.2 | 0 | 0.56 | 0 | 0 | 0 |
| 6 | 1.6 | 0 | 0 | 1.8 | 0.6 | 0.56 | 0 | 2.7 | 1.6 | 3 |
| 7 | 2.5 | 0 | 0 | 1.2 | 0 | 0 | 2.7 | 0 | 1.6 | 3.1 |
| 8 | 1.9 | 0 | 0 | 2.3 | 0 | 0 | 1.6 | 1.6 | 0 | 1.2 |
| 9 | 1.1 | 0 | 0.41 | 1.4 | 0 | 0 | 3 | 3.1 | 1.2 | 0 |

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|-----|------|-----|------|-----|-----|-----|------|------|-----|
| 0 | 0 | 1.4 | 3.1 | 3 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 1.4 | 0 | 2.5 | 2.8 | 0 | 0 | 0 | 0.41 | 0 | 0 |
| 2 | 3.1 | 2.5 | 0 | 3.2 | 0 | 0 | 0 | 0 | 0.6 | 0 |
| 3 | 3 | 2.8 | 3.2 | 0 | 0 | 0 | 0 | 0 | 0.56 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 1.4 | 2.5 | 1.1 | 1.6 | 1.9 |
| 5 | 0 | 0 | 0 | 0 | 1.4 | 0 | 1.2 | 1.4 | 1.8 | 2.3 |
| 6 | 0 | 0 | 0 | 0 | 2.5 | 1.2 | 0 | 3.1 | 2.7 | 1.6 |
| 7 | 0 | 0.41 | 0 | 0 | 1.1 | 1.4 | 3.1 | 0 | 3 | 1.2 |
| 8 | 0 | 0 | 0.6 | 0.56 | 1.6 | 1.8 | 2.7 | 3 | 0 | 1.6 |
| 9 | 0 | 0 | 0 | 0 | 1.9 | 2.3 | 1.6 | 1.2 | 1.6 | 0 |

Graph Sorting

Sorting a graph such that heavily connected nodes are enumerated to be maximally adjacent is the same as permuting rows and columns of the adjacency matrix to make it maximally diagonal. (A little like Cuthill-McKee)



Off-diagonal penalties

Consider a penalty matrix $C_{ij} = c(|i-j|)$,

e.g. $C_{ij} = |i-j|$

$C \in \mathbb{R}^{5 \times 5}$ looks like

$$\begin{bmatrix} 0 & 1 & 2 & 3 & 4 \\ 1 & 0 & 1 & 2 & 3 \\ 2 & 1 & 0 & 1 & 2 \\ 3 & 2 & 1 & 0 & 1 \\ 4 & 3 & 2 & 1 & 0 \end{bmatrix}$$

Then, if $\langle A, C \rangle_F \leq \langle B, C \rangle_F$ we say
A is more diagonally dominant than B.

Permutations

Let \mathcal{P}_n denote the set of permutation matrices:

$$\mathcal{P}_n = \{ P \in \{0, 1\}^{n \times n} : \sum_i P_{ij} = \sum_j P_{ij} = 1 \}$$

Given a permutation σ & P , $[PGP^T]_{ij} = G_{\sigma(i)\sigma(j)}$

So, our objective is to solve

$$\min_{P \in \mathcal{P}_n} \langle PGP^T, C \rangle_F = \min_{P \in \mathcal{P}_n} \sum_{i,j=1}^n [PGP^T]_{ij} C_{ij}$$

This is NP-hard!

Equivalent to

- Quadratic assignment (QAP)
- Travelling salesman (TSP)

[Koopman & Beckmans '57] [Yusile & Kosowsky '94]

Relaxation

We propose the following relaxation:

Solve $\min_{P \in B_n} \langle PGP^T, C \rangle_F$

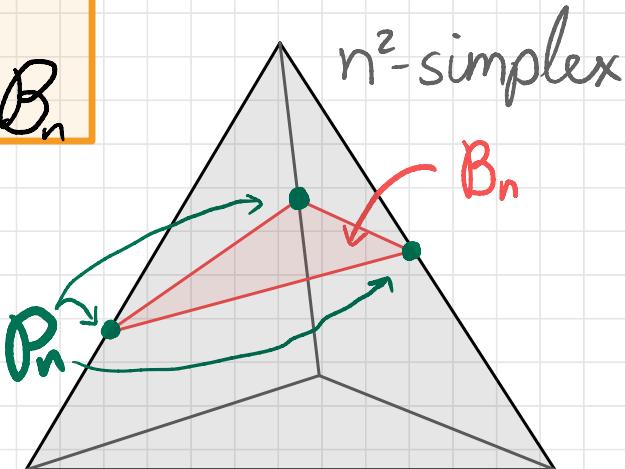
where $B_n = \{ P \in [0,1]^{n \times n} : \sum_i P_{ij} = \sum_j P_{ij} = 1 \}$
is the set of *bistochastic matrices*.

THEOREM (BIRKHOFF)

P_n is the set of extreme points of B_n

B_n is a *convex polytope*

and a subset of the n^2 -simplex P_n
of "probability" vectors.



Entropic Frank-Wolfe!

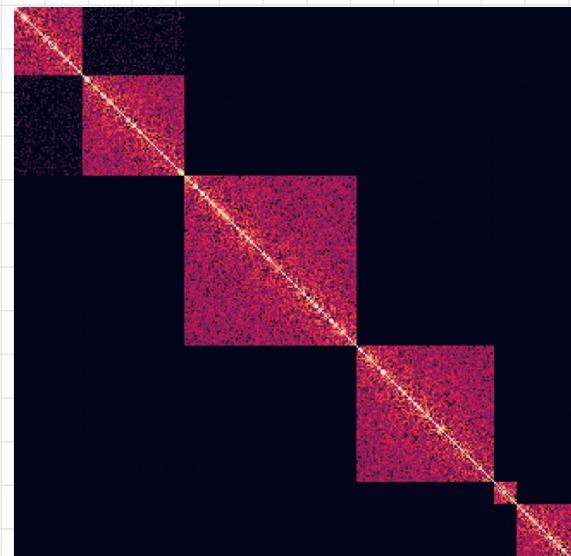
We solve our problem w/ entropic Frank-Wolfe!

- $f(P) = \langle PG^T, C \rangle_F$
- $\nabla_P (\langle PG^T, C \rangle_F) = CPG^T + C^T PG = 2CPG$
(As $C = C^T$ & $G = G^T$)
- $H(P) = - \sum_{i,j=1}^n P_{ij} (\log P_{ij} - 1)$
- $K = B_n$

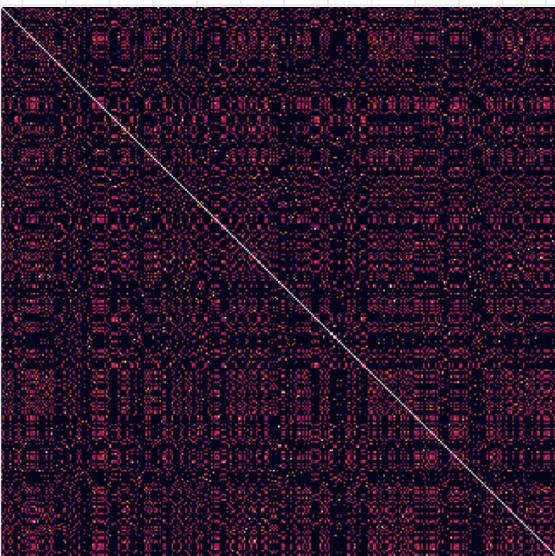
NOTE Our objective is quadratic and non-convex,
we have no guarantee that $P^* \in P_n \subset B_n$

Some results

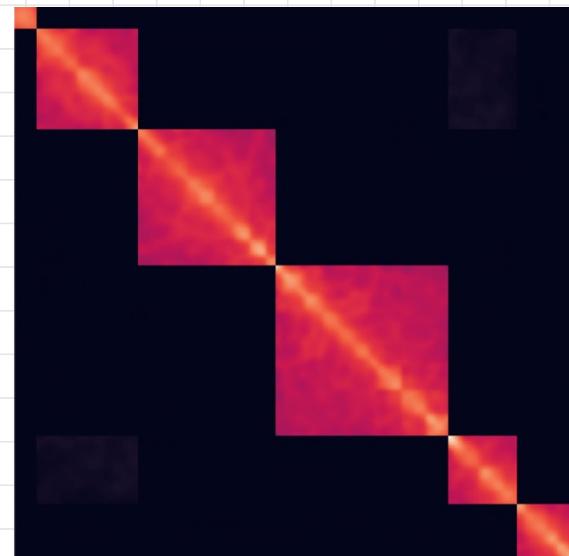
ORIGINAL RANDOM
BLOCK MATRIX



SCRAMBLED
MATRIX
(600×600)



RECOVERED /
ENTROPICALLY SORTED

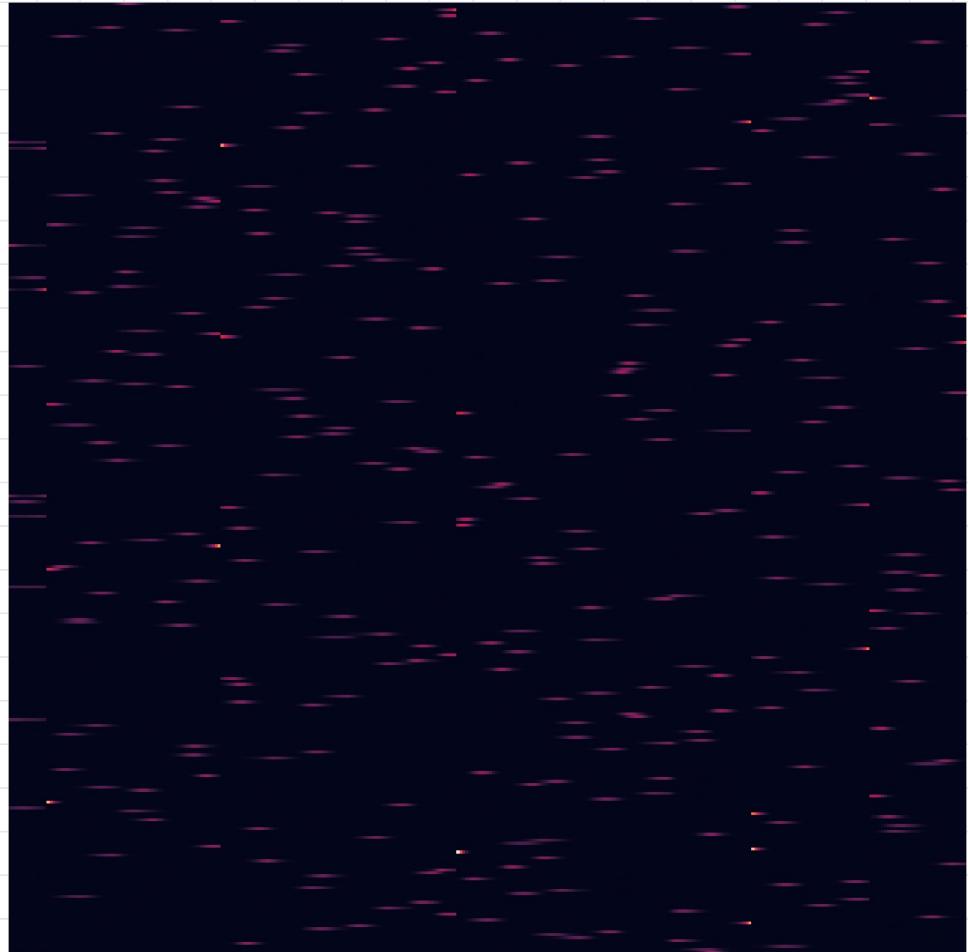


40 iterations
 $\epsilon_0 = 2 \times 10^{-3}$

Some results

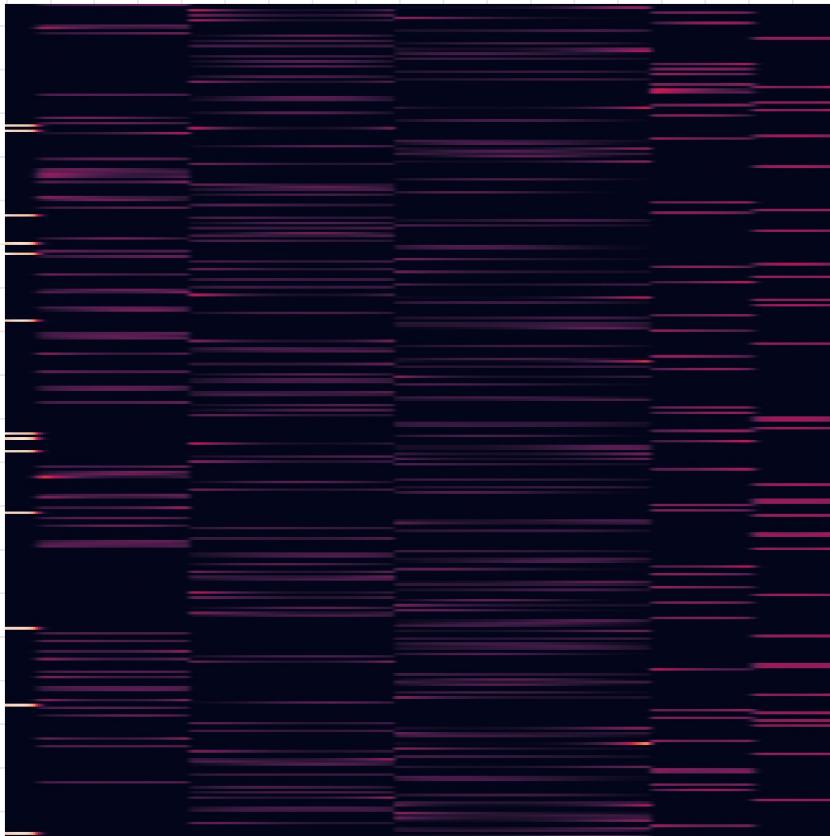
APPROXIMATE
PERMUTATION
(IN \mathfrak{S}_n)

(can recover
a real permutation
from this)

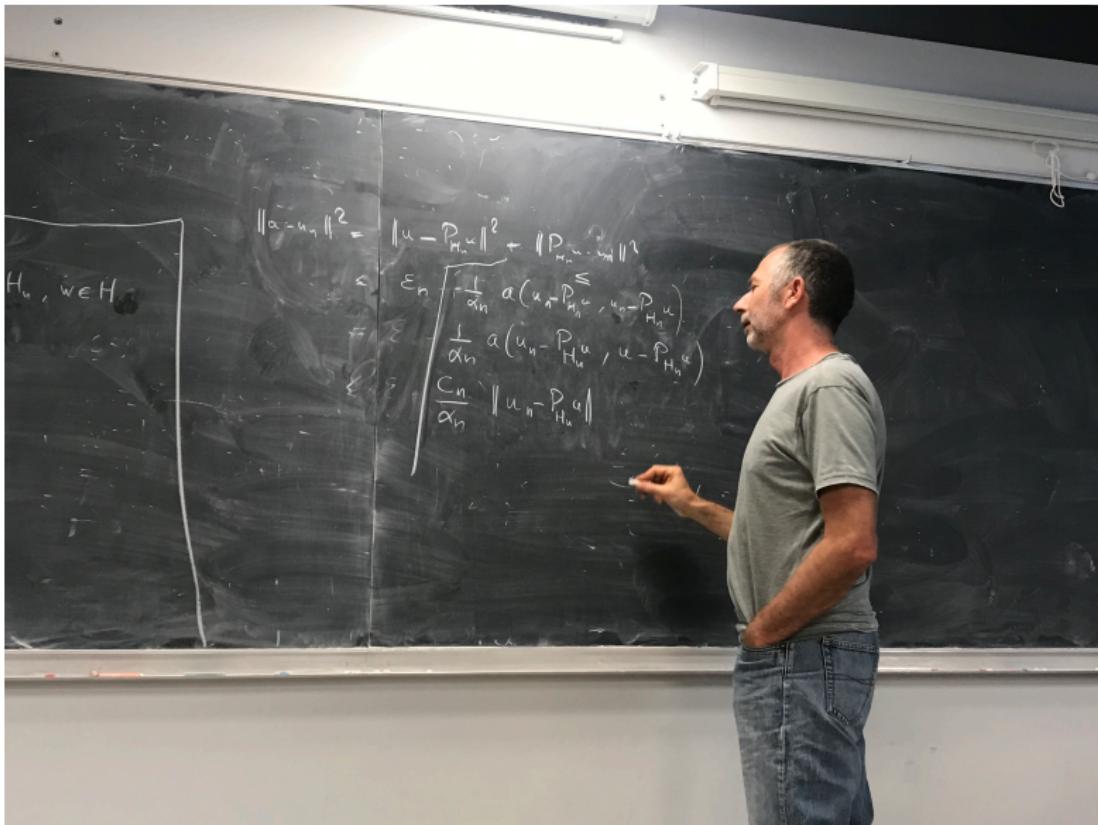


Only 3 iterations

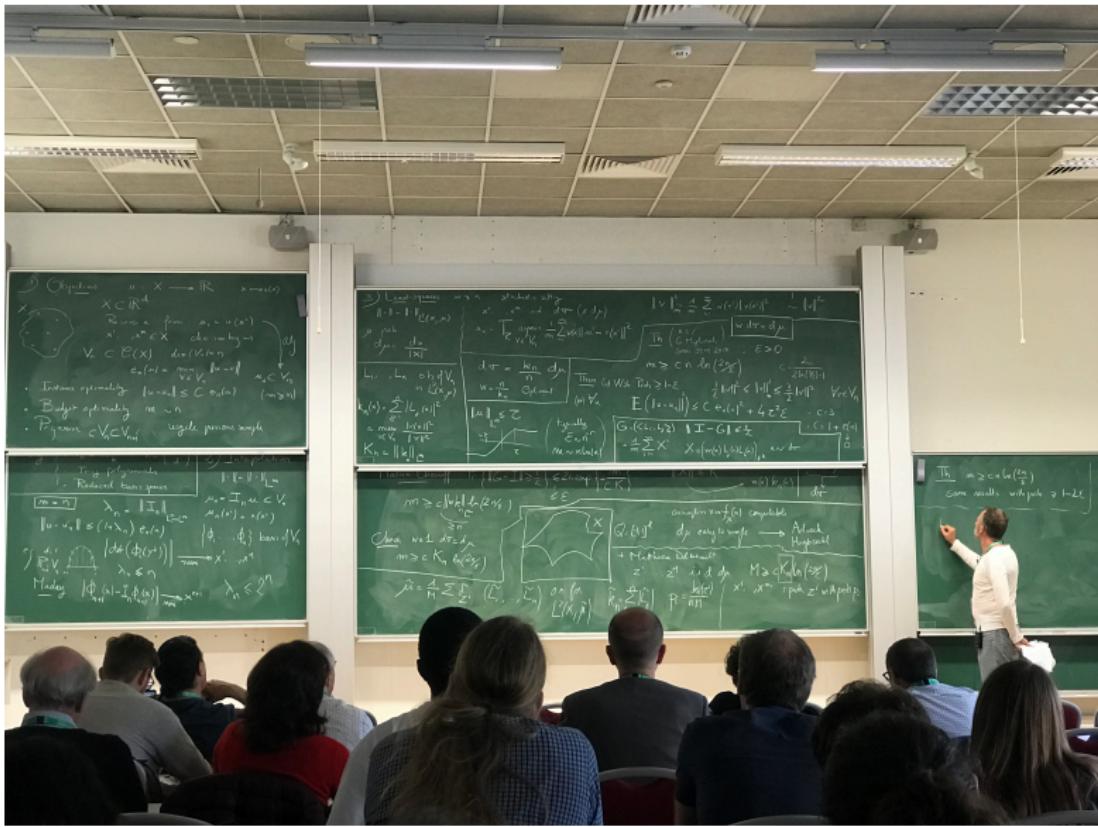
Same 600×600 problem w/ 3 iterations



POV - You are about to have a great day



POV - You are witnessing a most elegant talk



Playing music



POV - You had a great three years in Paris with this guy

