

Discrete Determinantal Point Processes

Stats/ML in Saclay

Jan 27, 2020

Victor-Emmanuel Brunel
Department of Statistics
ENSAE/CREST

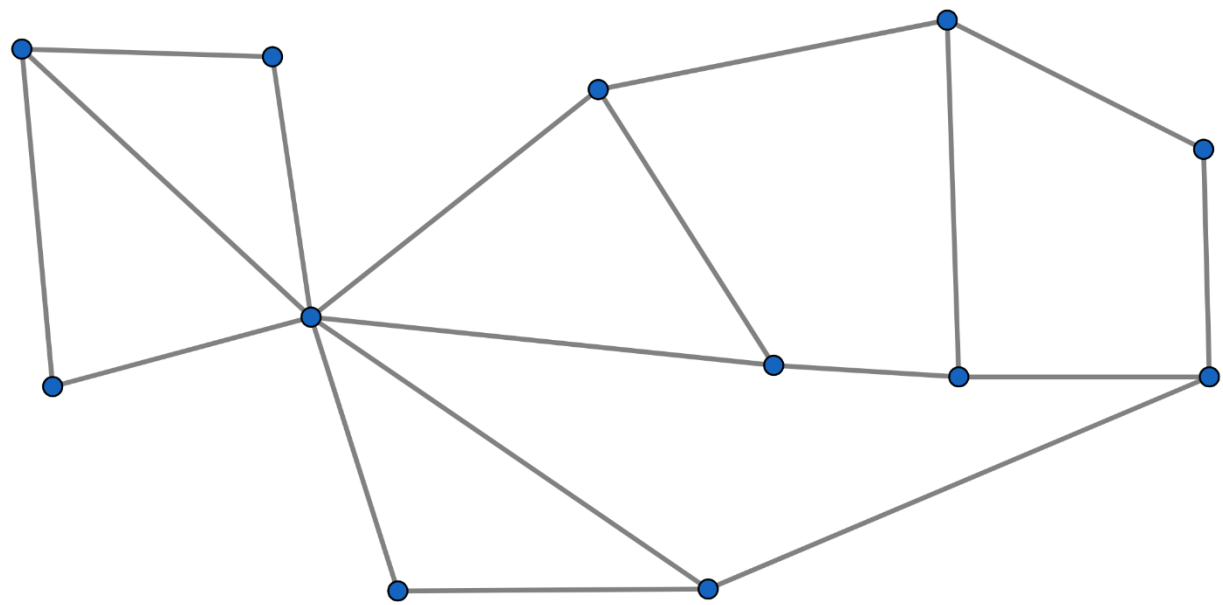


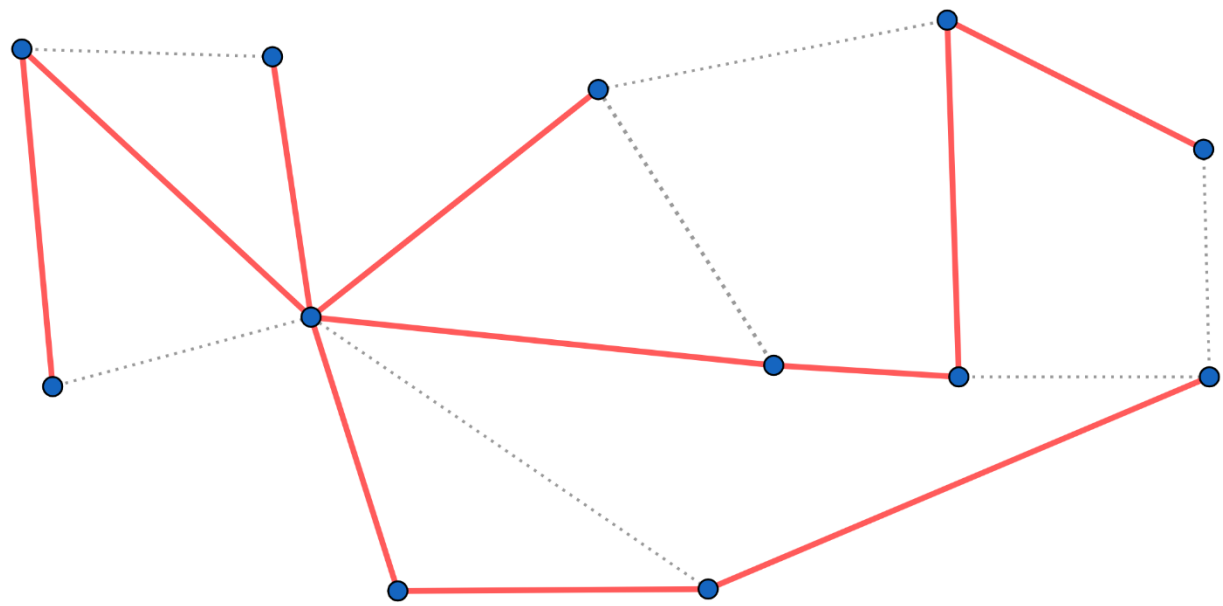
Random selection models

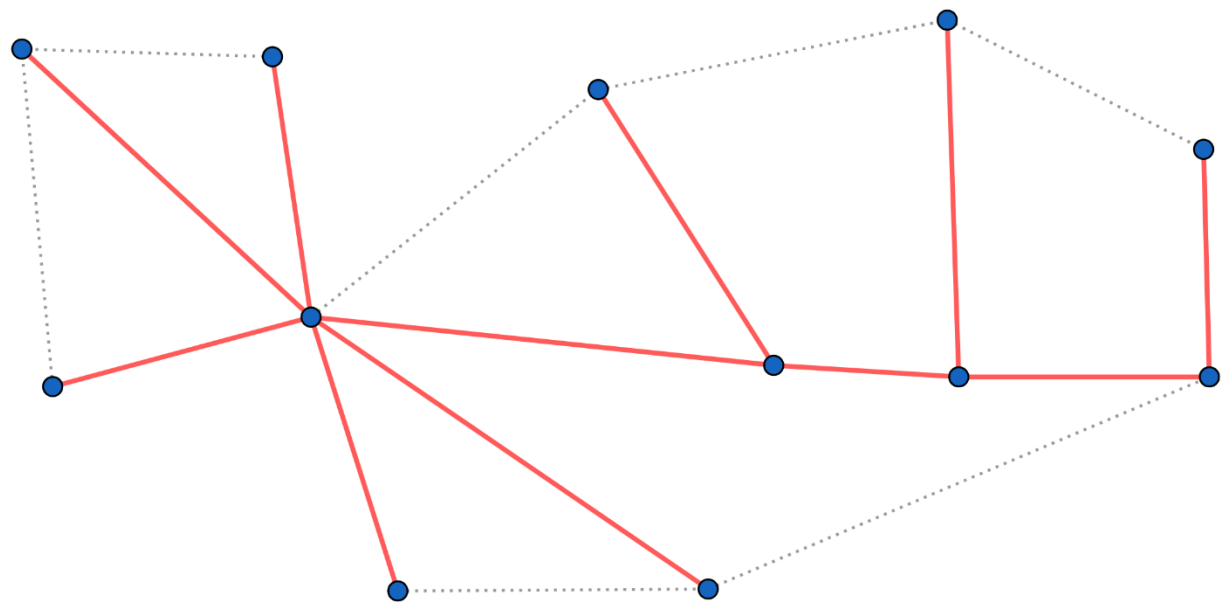
- Sampling
 - Independent Bernoulli sampling
 - with/without replacement
- Ising models
 - Ferromagnetic (attractive interactions)
 - Antiferromagnetic (repulsive interactions)
- DPPs

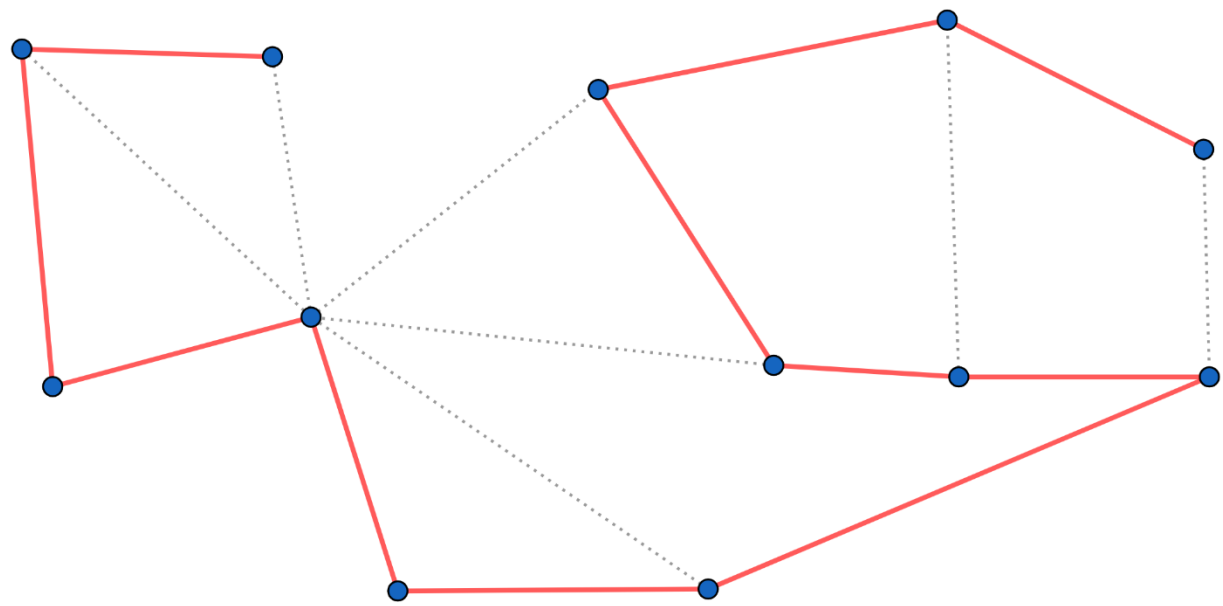
DPPs in probability

- Random spanning trees
- Increases in i.i.d. random sequences
- Non-intersecting random walks
- Eigenvalues of random matrices










DPPs in probability

- Random spanning trees
- Increases in i.i.d. random sequences
- Non-intersecting random walks
- Eigenvalues of random matrices

• $X_1, X_2, \dots, X_{N+1} \stackrel{\text{i.i.d.}}{\sim} P$

• $\{i = 1, \dots, n : X_i < X_{i+1}\}$ is a DPP

3 2 7 3 6 1 1 0 8 5 9 6 3 8 5 5 3 1



$$Y = \{2, 4, 8, 10, 13\} \subseteq [17]$$


DPPs in probability

- Random spanning trees
- Increases in i.i.d. random sequences
- Non-intersecting random walks
- Eigenvalues of random matrices

DPPs in probability

- Random spanning trees
- Increases in i.i.d. random sequences
- Non-intersecting random walks
- Eigenvalues of random matrices

DPPs in ML

- Model for random selections of items within a catalog/dictionary
-  Tractable models: marginalizing, conditioning, are computationally simple
- In some cases, sampling is simple and fast (e.g., via spectral decomposition or MCMC)

Applications of DPPs

(Symmetric) DPPs have become popular in various applications:

- Quantum physics (*fermionic processes*) [Macchi '74]
- Experimental design [Derezinski et al. '19]
- Document and timeline summarization [Lin, Bilmes '12; Yao *et al.* '16]
- Image search [Kulesza, Taskar '11; Affandi *et al.* '14]
- Bioinformatics [Batmanghelich *et al.* '14]
- Neuroscience [Snoek *et al.* '13]
- Wireless or cellular network modelisation [Li *et al.* '15; Deng *et al.* '15]
- Recommendation systems [Gartrell et al '16, '17]

And they remain an elegant and important tool in probability theory [Borodin '11]

Set representation of binary vectors

A binary vector of size N can be represented as a subset of $[N]$:

1 0 0 1 1 0 1 0 1 1 0 1 0 0 1 0 0 0 1 0 \leftrightarrow {1,4,5,7,9,10,12,15,19}

0 0 1 1 0 1 0 1 1 0 0 1 0 0 1 0 0 0 1 0 \leftrightarrow {3,4,6,8,9, 12,15,19}

$(X_1, \dots, X_N) \in \{0,1\}^N$ \leftrightarrow $Y \subseteq [N]$

$$X_i = 1 \Leftrightarrow i \in Y$$

DPPs: Definition

DPP: Random subset Y of $[N]$

- For all $J \subseteq [N]$,

$$\mathbb{P}[Y = J] \propto \det \mathbf{L}_J$$

- $\mathbf{L} \in \mathbb{R}^{N \times N}$: parameter of the DPP
- $L_J = (L_{i,j})_{i,j \in J}$
- Normalization constant: $\det(I + L)^{-1}$.
- What are the admissible matrices L ? [B. '18]

P_0 -matrices

- L must be a **P_0 -matrix** (all its principal minors are ≥ 0).
- Examples:
 - Any PSD matrix : **Symmetric DPPs**
 - If $L + L^T$ is PSD
- More about P_0 -matrices: *Convex sets of nonsingular and P -matrices*, C. R. Johnson, M. J. Tsatsomeros (1995)

DPPs: Alternative representation (1)

- For all $J \subseteq [N]$, $\mathbb{P}[J \subseteq Y] = \det \mathbf{K}_J$

$\mathbf{K} = L(I + L)^{-1}$: *kernel* of the DPP

- $\mathbb{P}[Y = J] = (-1)^{|\bar{J}|} \det(\mathbf{K} - I_{\bar{J}})$

DPPs: Alternative representation (2)

- $X_i \sim \text{Ber}(K_{i,i})$

- $\mathbb{P}[1,2 \in Y] = \begin{vmatrix} K_{1,1} & K_{1,2} \\ K_{2,1} & K_{2,2} \end{vmatrix} = K_{1,1}K_{2,2} - K_{1,2}K_{2,1}$

$$\Rightarrow \text{cov}(X_1, X_2) = -K_{1,2}K_{2,1}$$

- $\mathbb{E}[|Y|] = \mathbb{E}[\sum X_i] = \text{Tr} K$

- $\text{Var}[|Y|] = \text{Tr}(K(I - K))$

DPPs: Alternative representation (2)

- $X_i \sim \text{Ber}(K_{i,i})$

- $\mathbb{P}[1,2 \in Y] = \begin{vmatrix} K_{1,1} & K_{1,2} \\ K_{2,1} & K_{2,2} \end{vmatrix} = K_{1,1}K_{2,2} - K_{1,2}K_{2,1}$

$$\Rightarrow \text{cov}(X_1, X_2) = -K_{1,2}K_{2,1}$$

- $\mathbb{E}[|Y|] = \mathbb{E}[\sum X_i] = \text{Tr} K$

- $\text{Var}[|Y|] = \text{Tr}(K(I - K))$

DPPs: Alternative representation (2)

- $X_i \sim \text{Ber}(K_{i,i})$

- $\mathbb{P}[1,2 \in Y] = \begin{vmatrix} K_{1,1} & K_{1,2} \\ K_{2,1} & K_{2,2} \end{vmatrix} = K_{1,1}K_{2,2} - K_{1,2}K_{2,1}$

$$\Rightarrow \text{cov}(X_1, X_2) = -K_{1,2}K_{2,1}$$

- $\mathbb{E}[|Y|] = \mathbb{E}[\sum X_i] = \text{Tr } K$

- $\text{Var}[|Y|] = \text{Tr}(K(I - K))$

DPPs: Alternative representation (2)

- $X_i \sim \text{Ber}(K_{i,i})$

- $\mathbb{P}[1,2 \in Y] = \begin{vmatrix} K_{1,1} & K_{1,2} \\ K_{2,1} & K_{2,2} \end{vmatrix} = K_{1,1}K_{2,2} - K_{1,2}K_{2,1}$

$$\Rightarrow \text{cov}(X_1, X_2) = -K_{1,2}K_{2,1}$$

- $\mathbb{E}[|Y|] = \mathbb{E}[\sum X_i] = \text{Tr } K$

- $\text{Var}[|Y|] = \text{Tr}(K(I - K))$

DPPs: Symmetric case (1)

- L is symmetric and PSD ($\Leftrightarrow K$ is symmetric with $0 \preceq K \preceq I$)
- Write $K = V^T V$, $V \in \mathbb{R}^{r \times N}$
- Columns of V : v_1, \dots, v_N (v_i : vector of r features of item i)
- $K_{i,i} = \|v_i\|^2$: *popularity of item i*
- $\det K_J = (\text{Volume spanned by } v_i, i \in J)^2$

DPPs: Symmetric case (2)

- $\mathbb{P}[i, j \in Y] = \det K_{\{i,j\}} = \begin{vmatrix} K_{i,i} & K_{i,j} \\ K_{i,j} & K_{j,j} \end{vmatrix} = K_{i,i}K_{j,j} - K_{i,j}^2$
- $\text{cov}(X_i, X_j) = \mathbb{P}[i, j \in Y] - \mathbb{P}[i \in Y]\mathbb{P}[j \in Y] = -K_{i,j}^2 \leq 0$
- More generally, if $S \cap T = \emptyset$:

$$\begin{aligned} \text{cov}\left(\prod_{i \in S} X_i, \prod_{j \in T} X_j\right) &= \mathbb{P}[S, T \subseteq Y] - \mathbb{P}[S \subseteq Y]\mathbb{P}[T \subseteq Y] \\ &= \det(K_{S \cup T}) - \det(K_S) \det(K_T) \\ &\leq 0 \end{aligned}$$

DPPs: Symmetric case (2)

- $\mathbb{P}[i, j \in Y] = \det K_{\{i,j\}} = \begin{vmatrix} K_{i,i} & K_{i,j} \\ K_{i,j} & K_{j,j} \end{vmatrix} = K_{i,i}K_{j,j} - K_{i,j}^2$
- $\text{cov}(X_i, X_j) = \mathbb{P}[i, j \in Y] - \mathbb{P}[i \in Y]\mathbb{P}[j \in Y] = -K_{i,j}^2 \leq 0$
- More generally, if $S \cap T = \emptyset$:

$$\begin{aligned} \text{cov}\left(\prod_{i \in S} X_i, \prod_{j \in T} X_j\right) &= \mathbb{P}[S, T \subseteq Y] - \mathbb{P}[S \subseteq Y]\mathbb{P}[T \subseteq Y] \\ &= \det(K_{S \cup T}) - \det(K_S) \det(K_T) \\ &\leq 0 \end{aligned}$$

DPPs: Symmetric case (2)

- $\mathbb{P}[i, j \in Y] = \det K_{\{i,j\}} = \begin{vmatrix} K_{i,i} & K_{i,j} \\ K_{i,j} & K_{j,j} \end{vmatrix} = K_{i,i}K_{j,j} - K_{i,j}^2$
- $\text{cov}(X_i, X_j) = \mathbb{P}[i, j \in Y] - \mathbb{P}[i \in Y]\mathbb{P}[j \in Y] = -K_{i,j}^2 \leq 0$
- More generally, if $S \cap T = \emptyset$:

$$\begin{aligned} \text{cov}\left(\prod_{i \in S} X_i, \prod_{j \in T} X_j\right) &= \mathbb{P}[S, T \subseteq Y] - \mathbb{P}[S \subseteq Y]\mathbb{P}[T \subseteq Y] \\ &= \det(K_{S \cup T}) - \det(K_S) \det(K_T) \\ &\leq 0 \end{aligned}$$

Negative Association (1)

Definition: Boolean rv's are *negatively associated** iff

$$\text{cov} \left(f(X_i, i \in S), g(X_j, j \in T) \right) \leq 0$$

for all $S, T \subseteq [N]$ and nondecreasing functions f, g .

*Remark: Hard to prove in general...

Negative Association (2)

Definition: Boolean rv's are *negatively associated** iff

$$\text{cov} \left(f(X_i, i \in S), g(X_j, j \in T) \right) \leq 0$$

for all $S, T \subseteq [N]$ and nondecreasing functions f, g .

*Remark: Hard to prove in general...

Strongly Rayleigh distributions (1)

- Let μ be the joint distribution of a Boolean vector (X_1, \dots, X_N)
- μ : Distribution on the cube $\{0,1\}^N$
- Generating polynomial: $g_\mu(z) = \mathbb{E}[z_1^{X_1} \dots z_N^{X_N}]$, $z = (z_1, \dots, z_N) \in \mathbb{C}^N$
- Ex:
 - Independent Bernoulli: $g_\mu(z) = \prod_{j=1}^N (p_j z_j + 1 - p_j)$
 - DPP: $g_\mu(z) = \sum_{J \subseteq [N]} \frac{\det L_J}{\det(I + L)} z^J = \frac{\det(I + LZ)}{\det(I + L)} = \det(I - K + KZ)$

Strongly Rayleigh distributions (1)

- Let μ be the joint distribution of a Boolean vector (X_1, \dots, X_N)
- μ : Distribution on the cube $\{0,1\}^N$
- Generating polynomial: $g_\mu(z) = \mathbb{E}[z_1^{X_1} \dots z_N^{X_N}]$, $z = (z_1, \dots, z_N) \in \mathbb{C}^N$

• Ex:

- Independent Bernoulli:
$$g_\mu(z) = \prod_{j=1}^N (p_j z_j + 1 - p_j)$$

- DPP:
$$g_\mu(z) = \sum_{J \subseteq [N]} \frac{\det L_J}{\det(I + L)} z^J = \frac{\det(I + LZ)}{\det(I + L)} = \det(I - K + KZ)$$

Strongly Rayleigh distributions (2)

Definition: μ is *Strongly Rayleigh* iff g_μ is real stable, i.e.,

$$\operatorname{Im}(z_j) > 0, \forall j \in [N] \Rightarrow g_\mu(z) \neq 0$$

Property: SR \Rightarrow NA

See [Borcea et al. '08]

Strongly Rayleigh distributions (2)

Definition: μ is *Strongly Rayleigh* iff g_μ is real stable, i.e.,

$$\operatorname{Im}(z_j) > 0, \forall j \in [N] \Rightarrow g_\mu(z) \neq 0$$

Property: SR \Rightarrow NA

See [Borcea et al. '08]

Strongly Rayleigh distributions (2)

Examples:

- Independent Bernoulli
- **Symmetric** DPPs
- **Symmetric** k -DPPs
- Sampling without replacement

Useful stability properties: Marginalization and conditioning.

Challenges in Stats and ML

Learning symmetric DPPs (1)

- Given $Y_1, Y_2, \dots, Y_n \stackrel{\text{i.i.d.}}{\sim} \text{DPP}(L)$, estimate L
- Problem 1: Identifiability

Learning symmetric DPPs (1)

- Given $Y_1, Y_2, \dots, Y_n \stackrel{\text{i.i.d.}}{\sim} \text{DPP}(L^*)$, estimate L^*
- Problem 1: Identifiability

Identification (symmetric case): \mathcal{D} -similarity

- $\text{DPP}(L') = \text{DPP}(L) \Leftrightarrow \det(L'_J) = \det(L_J), \forall J \subseteq [N]$

Identification (symmetric case): \mathcal{D} -similarity

- $\text{DPP}(L') = \text{DPP}(L) \Leftrightarrow \det(L'_J) = \det(L_J), \forall J \subseteq [N]$

[Oeding '11]

$$\Leftrightarrow L' = DLD \quad \text{for some } D = \begin{pmatrix} \pm 1 & & & \mathbf{0} \\ & \pm 1 & & \\ & \mathbf{0} & \ddots & \\ & & & \pm 1 \end{pmatrix}.$$

Identification (symmetric case): \mathcal{D} -similarity

- $\text{DPP}(L') = \text{DPP}(L) \Leftrightarrow \det(L'_J) = \det(L_J), \forall J \subseteq [N]$

[Oeding '11]

$$\Leftrightarrow L' = DLD \quad \text{for some } D = \begin{pmatrix} \pm 1 & & & \mathbf{0} \\ & \pm 1 & & \\ & \mathbf{0} & \ddots & \\ & & & \pm 1 \end{pmatrix}.$$

- E.g.: $L = \begin{pmatrix} + & + & + & + \\ + & + & + & + \\ + & + & + & + \\ + & + & + & + \end{pmatrix}$

The diagram shows a 4x4 matrix L with all entries '+'. Red boxes highlight the first and fourth columns and the first and fourth rows. Red arrows point downwards from the top of the first and fourth columns, and red arrows point leftwards from the right side of the first and fourth rows, illustrating the DLD transformation.

$$\rightsquigarrow DLD = \begin{pmatrix} + & - & - & + \\ - & + & + & - \\ - & + & + & - \\ + & - & - & + \end{pmatrix}$$

- L and DLD are called **\mathcal{D} -similar**.



Learn L up to \mathcal{D} -similarity.

Learning symmetric DPPs (2)

- Error of an estimator \hat{L} : $\min_{D \in \mathcal{D}} \|\hat{L} - DLD\|$
- Generic methods:
 - Method of moments (Urschel et al, 2017): Noisy principal minor assignment problem
 - Maximum likelihood (B. et al, 2017)
 - Penalised likelihood (only empirical results)

Maximum likelihood estimation

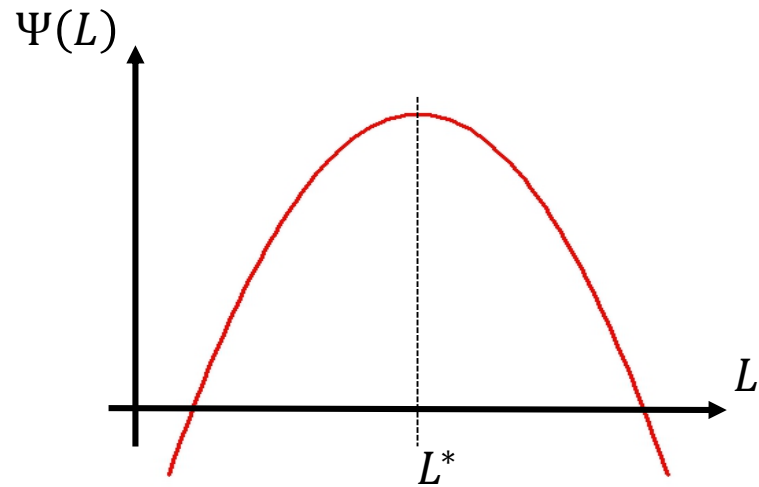
Samples: $Y_1, Y_2, \dots, Y_n \sim \text{DPP}(L^*), \quad L^* \succ 0$

- **Log-likelihood:** $\hat{\Psi}(L) = \sum_{J \subseteq [N]} \hat{p}_J \log \det(L_J) - \log \det(I + L)$
- **MLE:** $\hat{L} \in \operatorname{argmax} \hat{\Psi}(L)$

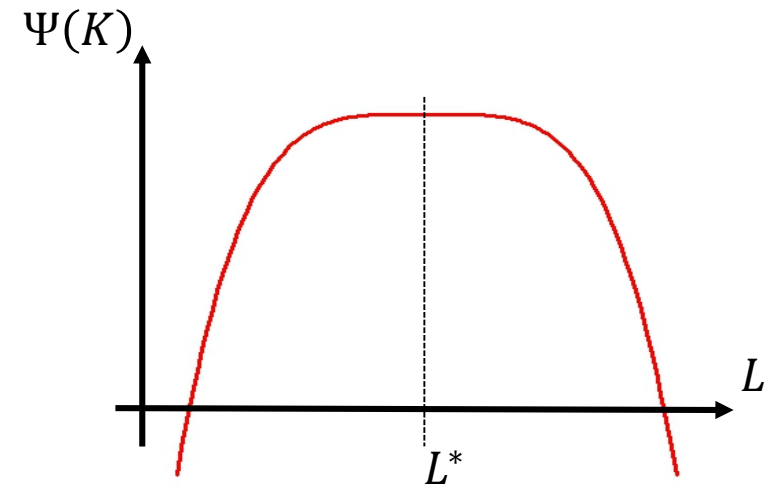
$$\begin{aligned} \Psi(L) &\triangleq \mathbb{E}[\hat{\Psi}(L)] = \sum_{J \subseteq [N]} p_J^* \log \det(L_J) - \log \det(I + L) \\ &= \Psi(L^*) - \text{KL}(\text{DPP}(L^*), \text{DPP}(L)) \end{aligned}$$

Likelihood geometry

Fisher information: $-\nabla^2\Psi(L^*)$



$$\nabla^2\Psi(L^*) < 0$$



$$\nabla^2\Psi(L^*) = 0$$

What is the order of the first non degenerate derivative of Ψ at $L = L^*$?

Asymptotic Results

- If L^* is irreducible:

$$\min_D \|\hat{L} - DL^*D\|_\infty = O_{\mathbb{P}}\left(n^{-\frac{1}{2}}\right)$$

- If L^* is block diagonal:

- $\min_D \|\hat{L} - DL^*D\|_\infty = O_{\mathbb{P}}\left(n^{-\frac{1}{4}}\right)$

- $\min_D \|\hat{L}_S - (DL^*D)_S\|_\infty = O_{\mathbb{P}}\left(n^{-\frac{1}{2}}\right)$ for all diagonal blocks S of L^* .

Non-asymptotic results in high dimension?

- Low rank assumption:

$$\text{rk}(L^*) \leq r \Leftrightarrow |Y| \leq r \text{ a.s.}$$

- If $\text{rk}(L^*) \leq r$, then $\hat{p}_J = 0$, $\forall J \subseteq [N]$ s.t. $|J| > r$
- However, $\text{rk}(\hat{L})$ can still be very large.
- Minimize a penalized version of the log-likelihood?
- How to compute \hat{L} (or \hat{L}_{pen}) efficiently?

Non-asymptotic results in high dimension?

- Low rank assumption:

$$\text{rk}(L^*) \leq r \Leftrightarrow |Y| \leq r \text{ a.s.}$$

- If $\text{rk}(L^*) \leq r$, then $\hat{p}_J = 0$, $\forall J \subseteq [N]$ s.t. $|J| > r$

- However, $\text{rk}(\hat{L})$ can still be very large.

- Minimize a penalized version of the log-likelihood?

- How to compute \hat{L} (or \hat{L}_{pen}) efficiently?

Non-asymptotic results in high dimension?

- Low rank assumption:

$$\text{rk}(L^*) \leq r \Leftrightarrow |Y| \leq r \text{ a.s.}$$

- If $\text{rk}(L^*) \leq r$, then $\hat{p}_J = 0$, $\forall J \subseteq [N]$ s.t. $|J| > r$

- However, $\text{rk}(\hat{L})$ can still be very large.

- Minimize a penalized version of the log-likelihood?

- How to compute \hat{L} (or \hat{L}_{pen}) efficiently?

Non-asymptotic results in high dimension?

- Low rank assumption:

$$\text{rk}(L^*) \leq r \Leftrightarrow |Y| \leq r \text{ a.s.}$$

- If $\text{rk}(L^*) \leq r$, then $\hat{p}_J = 0$, $\forall J \subseteq [N] \text{ s.t. } |J| > r$
- However, $\text{rk}(\hat{L})$ can still be very large.
- Minimize a penalized version of the log-likelihood?
- How to compute \hat{L} (or \hat{L}_{pen}) efficiently?

Non-asymptotic results in high dimension?

- Low rank assumption:

$$\text{rk}(L^*) \leq r \Leftrightarrow |Y| \leq r \text{ a.s.}$$

- If $\text{rk}(L^*) \leq r$, then $\hat{p}_J = 0$, $\forall J \subseteq [N]$ s.t. $|J| > r$
- However, $\text{rk}(\hat{L})$ can still be very large.
- Minimize a penalized version of the log-likelihood?
- How to compute \hat{L} (or \hat{L}_{pen}) efficiently?

Thank you for your attention!