# Determinantal Point Processes

Jose Gallego-Posada

April 2021



0 0

sli.do -- #MilaDPP

## Today's agenda

- Why DPPs?
- Definition and properties
- Sampling
- Applications





### Bible for DPP in ML:

Foundations and Trends in Machine Learning Determinantal Point Processes for Machine Learning Alex Kulesza and Ben Taskar (2012) [<u>link</u>]

### Presentation based on slides by :

- Simon Barthelmé, Nicolas Tremblay, EUSIPCO19 [link]
- Alex Kulesza, Ben Taskar and Jennifer Gillenwater CVPR13 [link]









### 0 0

### **DPPy: Sampling Determinantal Point Processes with** Python

docs passing build passing coverage 80% RyPI package

"Anything that can go wrong, will go wrong" - Murphy's Law

### Introduction

Determinantal point processes (DPPs) are specific probability distributions over clouds of points that have been popular as models or computational tools across physics, probability, statistics, and more recently of booming interest in machine learning. Sampling from DPPs is a nontrivial matter, and many approaches have been proposed. DPPy is a Python library that puts together all exact and approximate sampling algorithms for DPPs.

Guillaume Gautier, Rémi Bardenet, Guillermo Polito, Michal Valko







### BioDiversity











### Variance reduction – Mean estimation

0

0



**IID Samples** 



### Variance reduction – Mean estimation

0

0





### Variance reduction – Mean estimation



0

 $\mathcal{P}(\mathbf{Y} = \mathbf{Y})$  depends on the determinant of a

matrix selected based on the elements of Y.

# Determinantal Point Process

Base set  $\mathcal{Y} = \{1, ..., n\}$  from which we sample a random subset Y.

**Y** is distributed according to a point process  $\mathcal{P}$  over  $2^{y}$ .



### Poisson Process

- Simplest point process... too simple!
- Element memberships are parameterized by independent Bernoulli rvs. •
- Special case of a DPP with marginal kernel  $\Re = D_p$ .

 $-p_i)$ 

## Representing repulsion

### Desiderata:

- Density is tractable; including normalization constant Ι.
- Inclusion probabilities (marginals) are tractable ii.
- iii. Sampling is tractable
- iv. Model is easy to understand

Contrary to most Gibbs processes (normalized, exponentiated potentials), DPPs tick all the boxes







### Loopy, negative interactions are hard

*(Inference becomes intractable; worst case)* 





### Global, negative interactions are easy



[KTG13]

### **L**-ensembles

0

0

- Model repulsion based on similarity between elements of  $\mathcal{Y}$ . •
- Similarity between elements *i* and *j* is stored in  $\mathfrak{L}_{ij}$ . •
- We assume  $\mathfrak{L}$  to be positive definite.
- £ is known as the likelihood kernel.

We say that **Y** is distributed according to a DPP if:  $\mathcal{P}(\boldsymbol{Y} = \boldsymbol{Y}) \propto \det \mathfrak{L}_{\boldsymbol{Y}}$ 





## Where did the repulsion go?

$$\mathfrak{L}_Y = [\mathfrak{B}^T \mathfrak{B}]_Y$$
$$\mathfrak{L}_{\{1,2,4\}} =$$

$$\mathcal{P}(\boldsymbol{Y}=\boldsymbol{Y}) \propto \det \mathfrak{L}_{\boldsymbol{Y}} = \det^2 \mathfrak{B}_{\boldsymbol{Y}}$$





Embedding of  ${\mathcal Y}$ 

## Where did the repulsion go?







 $\mathcal{P}(\{i,j\}) \propto \mathcal{P}(\{i\})\mathcal{P}(\{j\}) - \left(\frac{\mathfrak{L}_{ij}}{\det(\mathfrak{L}+\mathbb{I})}\right)^2$ 





## Where did the repulsion go?



Probability under a DPP grows with the spanned volume





[BT19, <u>dpp\_demo]</u>

 $I_{4}$ 

### Normalization

### Sum of squares of determinants of principal minors

Asked 9 years, 11 months ago Active 9 years, 11 months ago Viewed 2k times

I am interested in computing the sum of squares of determinants of principal minors. Let A be an  $n \times n$  positive semidefinite matrix and  $A_S$  be a principal minor of A indexed by the set  $S \subseteq \{1, \ldots, n\}$ . The classical result (without squares) is:

$$\sum_{S \subseteq \{1,\ldots,n\}} \det(A_S) = \det(A+I)$$

Are there any results on computing

$${\mathfrak O} \quad \sum_{S\subseteq\{1,\ldots,n\}} \det^2(A_S)$$

0

0

or any other powers?

determinants sums-of-squares

Share Cite Improve this question Follow

asked Apr 11 '11 at 3:04



## Analytic normalization constant!





# Exploit linear-algebraic properties to make inference/sampling easy (or feasible in high-dims)



## Marginal kernels

- Consider a DPP with L-ensemble  $\mathfrak{L}$ .
- The inclusion (marginal) probability that **Y** contains a set S is given by: •

$$\mathcal{P}(S \subset Y) = \frac{1}{\Im} \sum_{S \subset Y} \det \mathfrak{L}_Y = \det \mathfrak{L}_Y$$

with  $\Re = \mathfrak{L}(\mathfrak{L} + \mathbb{I})^{-1}$ .

- R is known as the marginal kernel of the DPP.
- $\mathcal{P}(i \in \mathbf{Y}) = \mathfrak{K}_{ii}$ .

0

0

•  $\mathbb{E}[|Y|] = \mathbb{E}[\sum_{i} \mathbb{1}_{i \in Y}] = \sum_{i} \mathcal{P}(i \in Y) = \operatorname{tr} \mathfrak{K}.$ 





## Conditioning

0

0



$$\mathcal{P}(B \subset \mathbf{Y} \mid A \subset \mathbf{Y}) = \frac{\mathcal{P}(A \cup B \subset \mathbf{Y})}{\mathcal{P}(A \subset \mathbf{Y})} = \frac{\det \mathfrak{K}_{A \cup B}}{\det \mathfrak{K}_A} = \frac{\det \mathfrak{K}_{A \cup B}}{\det \mathfrak{K}_A}$$

### Schur complement

### $\det \mathfrak{K}_{A \cup B} = \det \mathfrak{K}_A \det (\mathfrak{K}_B - \mathfrak{K}_{BA} \mathfrak{K}_A^{-1} \mathfrak{K}_{AB})$

### DPPs are closed under conditioning!

### $= \det(\mathfrak{K}_B - \mathfrak{K}_{BA}\mathfrak{K}_A^{-1}\mathfrak{K}_{AB})$



## Complexity?

- Evaluation of  $\mathfrak{L}$   $\mathcal{O}(n^2)$
- Normalization constant  $O(n^3)$  [determinant]
- Marginal probabilities  $O(n^3)$  [matrix inversion]
- Conditional probabilities  $O(n^3)$  [Schur complement]









### Extensions





### DPPs



### k-DPPs

- In practical applications, often preferred to limit cardinality of the set
  - Search results
  - Minibatch selection
  - Summarization

 $\mathcal{P}(\boldsymbol{Y}=\boldsymbol{Y}) \propto \det \mathfrak{L}_{\boldsymbol{Y}} \, \mathbf{1}_{|\boldsymbol{Y}|=k}$ 

- Normalization constant  $\sum_{|Y|=k} \det \mathfrak{L}_Y = e_k(\lambda_1, ..., \lambda_N)$  [k-th elementary sym. polynomial]
- Special case: 1-DPP

0

0

Need not have a corresponding marginal kernel

### Elementary π-DPPs

- Special case: k-DPP with  $k = \operatorname{rank} \mathfrak{L}$  and  $\mathfrak{L} = V \Lambda V^T$ , has marginal kernel  $\mathfrak{K} = V V^T$
- A DPP is called **elementary** if the spectrum of its marginal kernel is {0, 1}.

$$\mathfrak{K}^V = \sum_{\boldsymbol{v} \in V} \boldsymbol{v} \boldsymbol{v}^T$$

• We denote this process as  $\mathcal{P}^V$ .

0

- If  $Y \sim \mathcal{P}^V$ , then |Y| = |V| with probability one. (|Y| is a sum of Bernoulli rvs.)
- R is a projection matrix also called projection DPPs



## Hierarchy of DPPs





### Cauchy-Binet Lemma

- Consider matrices A of size  $r \times s$  and B of size  $s \times r$
- For each r-subset  $Y \subset [1, ..., r]$ , construct square matrices  $A_{Y}$  and  $B_{Y}$ :

$$\det AB = \sum_{|Y|=r} \det A_{:Y} \det$$

### $F_n = \frac{\varphi^n - \psi^n}{\omega - \psi}$

 $B_{Y:}$ 



JPM Binet



### DPPs as mixture models

0

0

 $\mathcal{P}(Y = Y) \propto \det \mathfrak{L}_Y = \det [V \Lambda V]_Y$ 

 $= \det V_{Y:} \sqrt{\Lambda} \sqrt{\Lambda} V_{:Y}$ 

$$= \sum_{|Z|=|Y|} \det V_{YZ} \sqrt{\Lambda_{ZZ}} \quad \det \sqrt{\Lambda_{ZZ}}$$
$$= \sum_{|Z|=|Y|} \det V_{YZ} V_{YZ}^T \quad \det \Lambda_Z$$
$$Elementary \quad Diago$$
$$DPP \quad \pounds-enser$$

### $V_{ZY}$

### ZZ

### nal

mble



## Sampling

0

0

- Consider a DPP with L-ensemble  $\mathfrak{L} = \sum_n \lambda_n \boldsymbol{v}_n \boldsymbol{v}_n^T$ .
- For each subset  $J \subset \mathcal{Y}$ , let  $V_J$  denote the set  $\{\boldsymbol{v}_n\}_{n \in J}$  and the elementary DPP  $\mathcal{P}^{V_J}$ .

$$\mathcal{P} \propto \sum_{J \subset \mathcal{Y}} \mathcal{P}^{V_J} \prod_{n \in J} \lambda_n = \sum_{J \subset \mathcal{Y}} \mathcal{P}^{V_J}$$

Factorize the original DPP as a mixture of *elementary* DPPs





## Sampling via spectral decomposition



0

0



- Draw a sample from  $\mathcal{P}^J$
- by sequential exploiting closure
  - of DPPs under conditioning



[KT12 – p.145]

## Sampling in action



0 0 

Step 7

Step 8





## Advanced sampling

- Spectral method for sampling has cost  $\mathcal{O}(n^2 + n^3 + nk^2)$ 
  - Dual sampling: instead of using  $\mathfrak{L} = \mathfrak{B}^T \mathfrak{B}$  with  $\mathfrak{B} d \times n$  use  $\mathfrak{C} = \mathfrak{B} \mathfrak{B}^T$  [KT12§3.3]
  - Random projections
  - Nyström approximations: Low rank approximation [Li, Jegelka, Sra 16a]
- MCMC sampling [LJS16b] ullet
  - Add, remove, swap

0

- Prove fast mixing for chains in terms of total variation
- Distortion-free intermediate sampling [Derezinski 18; CDV20]
  - Suitably construct an intermediate subset  $\sigma$  and then subsample from it



## Learning DPPs

- Basic setting: Maximum Likelihood •
  - Given  $\{Y_t\}_{t=1}^T$  subsets of  $\mathcal{Y}$ . Parameterize  $\mathfrak{L}$ -ensemble as  $\mathfrak{L}(\theta)$

$$\underset{\theta}{\operatorname{argmax}} \log \prod_{t} \mathcal{P}_{\theta}(Y_{t}) = \sum_{t} \log \det \mathfrak{L}_{Y_{t}}(\theta)$$

- Can use gradient-based methods for optimizing  $\theta$ •
- Can be extended to conditioning on a covariate X:  $\mathfrak{L}(\theta, X)$ •
  - For each X we have a DPP
  - X may be a query during search on which we want to condition the distribution over results
- See [KT12§4] for more details

0

0

### $-\log \det(\mathfrak{L}(\theta) + \mathbb{I})$



## Applications

### Image search

### {Relevance vs Diversity}



Mila Kunis - Wikipedia en.wikipedia.org



Mila Kunis Facts | POPSUGAR Cel.. popsugar.com



Mila Kunis | POPSUGAR Celebrity popsugar.com



Mila Kunis - Wikipedia, la ... es.wikipedia.org



Mila Ximénez confiesa en 'Sálvame' .. lecturas.com



Mila - Quebec Artificial Intelli... sv.linkedin.com



NASA and the Russian Space Agency have agreed to set aside a last-minute Russian request to launch an international space station into an orbit closer to Mir, officials announced Friday....

A last-minute alarm forced NASA to halt Thursday's launching of the space shuttle Endeavour, on a mission to start assembling the international space station. This was the first time in three years . . .

The planet's most daring construction job began Friday as the shuttle Endeavour carried into orbit six astronauts and the first U.S.-built part of an international space station that is expected to cost more than \$100 billion....

Following a series of intricate maneuvers and the skillful use of the space shuttle Endeavour's robot arm, astronauts on Sunday joined the first two of many segments that will form the space station . . .

....

### document cluster



### Extractive summarization

On Friday the shuttle Endeavor carried six astronauts into orbit to start building an international space station. The launch occurred after Russia and U.S. officials agreed not to delay the flight in order to orbit closer to MIR, and after a last-minute alarm forced a postponement. On Sunday astronauts joining the Russian-made Zarya control module cylinder with the American-made module to form a 70,000 pounds mass 77 feet long....

human summary

- NASA and the Russian Space Agency have agreed to set aside . . .
- A last-minute alarm forced NASA to halt Thursday's launching . . .
- This was the first time in three years, and 19 flights . . .
- After a last-minute alarm, the launch went off flawlessly Friday . . .
- Following a series of intricate maneuvers and the skillful . . .
- It looked to be a perfect and, hopefully, long-lasting fit. . . .

### extractive summary



IKT121

## Applications

0

0

(Quasi) Monte-Carlo integration (Gautier et al., On two ways to use DPPs for Monte Carlo integration, 2019) •

$$\int f(x)\mu(dx) \simeq \sum_{n=1}^N \omega_n f(x_n)$$

Mini-batch sampling for SGD (Zhang et al., DPPs for Mini-Batch Diversification, 2017) 

$$\theta \leftarrow \theta - \eta \sum_{i \in \mathcal{X}} \nabla L_i(\theta)$$

**Coresets** (Tremblay et al., DPPs for Coresets, 2018) 

$$\hat{\mathcal{L}}(\mathcal{S}, heta) = \sum_{oldsymbol{s}\in\mathcal{S}} \omega_{oldsymbol{s}} f(oldsymbol{s}, heta) \hspace{1.5cm} orall heta \in \Theta \hspace{1.5cm} (1-\epsilon)\mathcal{L}(\mathcal{X}, heta) \leqslant \epsilon$$

### $\hat{L}(\mathcal{S}, \theta) \leqslant (1 + \epsilon) L(\mathcal{X}, \theta)$



## DPPs in Randomized LinAlg

• Consider a linear regression problem with a tall, full-rank matrix  $X \in \mathbb{R}^{n \times d}$  with  $n \gg d$ 

$$w^* = \underset{w}{\operatorname{argmin}} \|Xw - y\|^2 = X^{\dagger}y$$

- Sketching: approximating matrix  $\widetilde{X}$  (subset of rows, low-rank)
- Usual bounds have  $(\varepsilon, \delta)$ -PAC flavour

0

0

• If  $S \sim d$ -DPP( $XX^T$ ), then  $\mathbb{E}[X_{S}^{-1}y] = w^*$  [leverage scores]

• If 
$$S \sim \text{DPP}\left(\frac{1}{\lambda}XX^T\right)$$
, then  $\mathbb{E}[X_{S:}^{\dagger}y] = \underset{w}{\operatorname{argmin}} \|Xw - y\|^2$ 

 $+\lambda ||w||^2$  [ridge l.s.]





## Minibatch sampling for LinReg

- Previously we related sampling with properties of analytic solution ullet
- What is the influence of non-iid sampling **during stochastic optimization**?  $\bullet$ 
  - Previous work by [<u>Zhang, Kjellström, Mandt 17</u>] for variance reduction
- Toy example: linear model ullet

0

- Gradients are 'constant' and correspond to points
- Redundant points lead to redundant sampled gradients  $\bullet$
- Sample minibatches  $S \sim d$ -DPP( $XX^T$ ) and run SGD with momentum



## Minibatch sampling for LinReg







## Overparameterized regime









## Determinantal Point Processes

## are elegant, efficient and useful

## models of repulsion

