

Exact sampling of determinantal point processes with sublinear time preprocessing

Michał Dereziński^{†*}, Daniele Calandriello^{‡*}, Michal Valko[§]

[†]UC Berkeley, [‡]LCSL - Istituto Italiano di Tecnologia, [§]DeepMind Paris ^{*}Equal contribution

Berkeley
UNIVERSITY OF CALIFORNIA

iit
ISTITUTO ITALIANO DI
TECNOLOGIA

MalGa

DeepMind

In a nutshell

Determinantal point processes (DPP) are used in machine learning for **randomized selection of diverse subsets** to capture **negative dependencies** between samples.

Existing algorithms for DPP sampling are **expensive** and require the eigendecomposition of an $n \times n$ similarity matrix at **$O(n^3)$** cost.

Our contribution: first **exact** DPP sampler with **$\tilde{O}(n)$** cost (the cost is **sublinear** in the n^2 size of the similarity matrix!)

Determinantal point processes

Goal: Given an $n \times n$ p.s.d. matrix \mathbf{L} , sample $S \subseteq \{1, \dots, n\}$ from:

$$\text{(variant 1)} \quad \text{DPP}(\mathbf{L}) : \quad \Pr(S) = \frac{\det(\mathbf{L}_S)}{\det(\mathbf{I} + \mathbf{L})},$$

$$\text{(variant 2)} \quad k\text{-DPP}(\mathbf{L}) : \quad \text{above, restricted to } |S| = k.$$

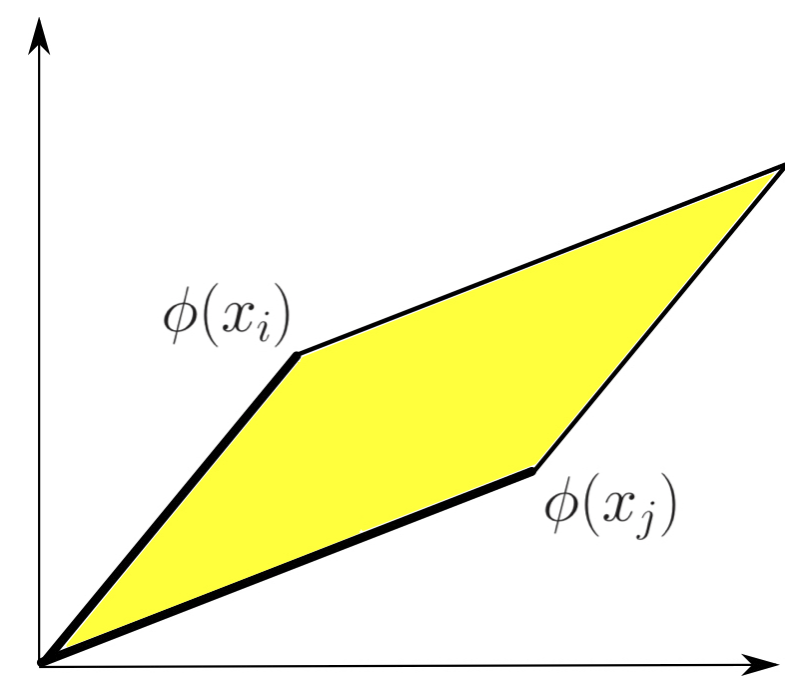
Similarity/Kernel interpretation

$\mathbf{L} = [\phi(x_i)^\top \phi(x_j)]_{ij}$ for a mapping $\phi : \mathcal{X} \rightarrow \mathbb{R}^m$

Determinant is *volume squared* in kernel space:

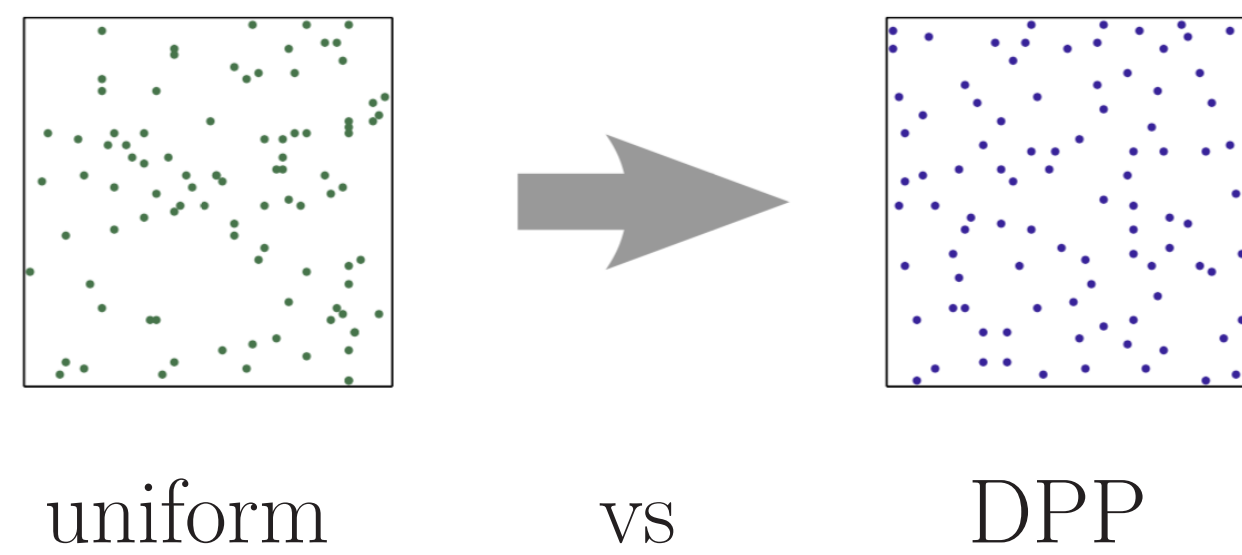
$$\det(\mathbf{L}_S) = \text{Vol}^2(\{\phi(x_i) : i \in S\})$$

Expected size $k = \mathbb{E}[|S|]$ is effective dimension of \mathbf{L} .



Applications: to learn more about DPPs, see [1]

- Recommender systems,
- Data summarization,
- Stochastic optimization,
- Gaussian processes.



Existing DPP samplers

Exact sampler [2]

Eigendecompose \mathbf{L} : $O(n^3)$
performed only once

Volume sampling: $O(nk^2)$
performed for every sample

Cost of $S_1 \sim \text{DPP}(\mathbf{X})$: $O(n^3)$

Cost of $S_2 \sim \text{DPP}(\mathbf{X})$: $O(nk^2)$

MCMC sampler [3]

1. Start from $S \subseteq [n]$

2. Sample $i \in S$ and $j \notin S$

3. Swap w.p. $\frac{1}{2} \min \left\{ 1, \frac{\Pr(S-i+j)}{\Pr(S)} \right\}$

$\tilde{\epsilon}$: ϵ -close in total variation distance

Cost of $S_1 \stackrel{\tilde{\epsilon}}{\sim} k\text{-DPP}(\mathbf{X})$: $n \cdot \text{poly}(k)$

Cost of $S_2 \stackrel{\tilde{\epsilon}}{\sim} k\text{-DPP}(\mathbf{X})$: $n \cdot \text{poly}(k)$

Distortion-free intermediate sampling

💡 Sample intermediate set σ out of $\{1..n\}$, then downsample S out of σ

🤖 How to ensure that intermediate σ *always* contains DPP sample S ?

1. Use marginal inclusion probabilities of $S \sim \text{DPP}(\mathbf{L})$

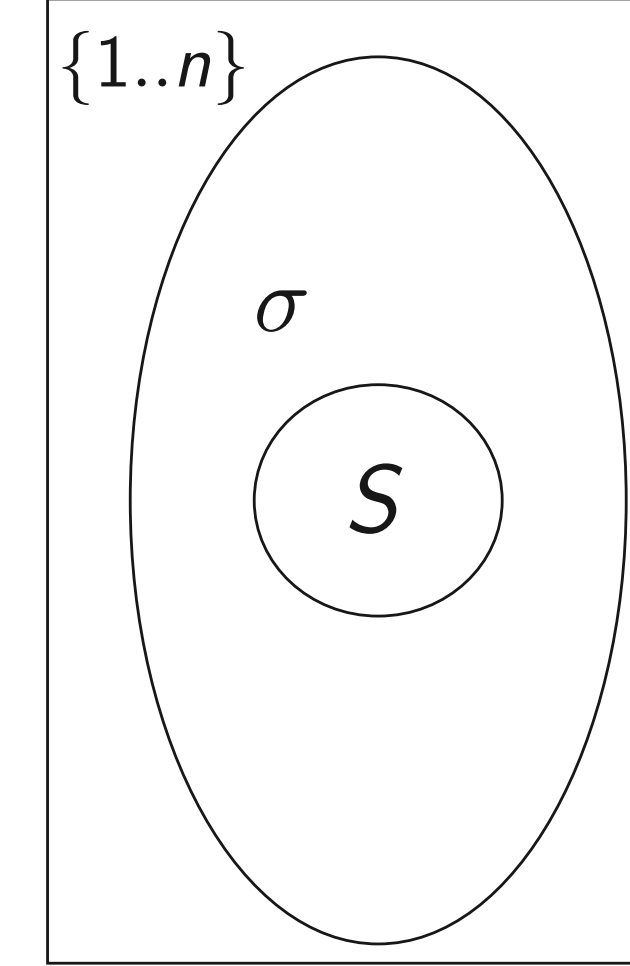
$$\Pr(i \in S) = \ell_i = [\mathbf{L}(\mathbf{I} + \mathbf{L})^{-1}]_{ii}$$

a.k.a. 1-ridge leverage scores, $\sum_{i=1}^n \ell_i = \mathbb{E}[|S|] = k$

2. Sample i.i.d. $\sigma_1, \dots, \sigma_t \sim [\ell_1/k, \dots, \ell_n/k]$

😞 Ignores negative dependence: $S \not\subseteq \sigma$ possible

💡 Reject σ if $S \not\subseteq \sigma$! Sampling becomes *exact*!



🤖 **Key question:** how many rejections occur before we get a good σ ?

Trade-offs: 1. How large should the intermediate sample σ be?
2. How accurate should the leverage score estimates be?

DPP-VFX: first sub-linear time exact DPP sampler

Thm. For any $\text{DPP}(\mathbf{L})$ or $k\text{-DPP}(\mathbf{L})$, we can sample

- the first subset S_1 in: $n \cdot \text{poly}(k)$ polylog(n) time,
- each successive S_i in: $\text{poly}(k)$ time.

😊 **Our answer:** we only suffer a constant number of rejections!

- Size of the intermediate sample is $t = O(k^2)$, independent of n !
- $\ell_i \approx (1 \pm O(1/k)) \ell_i$ leverage scores estimates are accurate enough
↳ efficient to compute with off-the-shelf Nyström approximations

😊 Easier in practice: $\ell_i \approx (1 \pm 1/2) \ell_i$ enough for constant rejections

Comparison with existing DPP samplers

	exact	DPP	$k\text{-DPP}$	first sample	next sample
Hough et al. [2]	✓	✓	✗	n^3	nk^2
Kulesza and Taskar [1]	✓	✗	✓	n^3	nk^2
Anari et al. [3]	✗	✗	✓	$n \cdot \text{poly}(k)$	$n \cdot \text{poly}(k)$
Li et al. [4]	✗	✓	✗	$n^2 \cdot \text{poly}(k)$	$n^2 \cdot \text{poly}(k)$
DPP-VFX	✓	✓	✓	$n \cdot \text{poly}(k)$	$\text{poly}(k)$

Very Fast and eXact DPP sampler (DPP-VFX)

Input: $\mathbf{L} \in \mathbb{R}^{n \times n}$, its Nyström approximation $\hat{\mathbf{L}}$

$\ell_i = [(\mathbf{L} - \hat{\mathbf{L}}) + \hat{\mathbf{L}}(\mathbf{I} + \hat{\mathbf{L}})^{-1}]_{ii} \approx \ell_i,$ $s = \sum_i \ell_i, z = \text{tr}(\hat{\mathbf{L}}(\mathbf{I} + \hat{\mathbf{L}})^{-1}), [\tilde{\mathbf{L}}]_{ij} = [\mathbf{L}]_{ij} / (s\sqrt{\ell_i \ell_j})$	} preprocessing
1: repeat 2: sample $t \sim \text{Poisson}(s^2 e^{1/s})$ 3: sample $\sigma_1, \dots, \sigma_t \sim (\ell_1/s, \dots, \ell_n/s),$ 4: sample $\text{Acc} \sim \text{Bernoulli}\left(e^{z-t/s} \cdot \frac{\det(\mathbf{I} + \tilde{\mathbf{L}}_\sigma)}{\det(\mathbf{I} + \tilde{\mathbf{L}})}\right)$ 5: until $\text{Acc} = \text{true}$	
6: sample $\tilde{S} \sim \text{DPP}(\tilde{\mathbf{L}}_\sigma)$ 7: return $S = \{\sigma_i : i \in \tilde{S}\}$	} downsample

General reduction from $k\text{-DPP}$ to DPP sampling

Folklore heuristic: $k\text{-DPP}$ rejection sampler

repeat $S_\alpha \sim \text{DPP}(\alpha \mathbf{L})$ **until** $|S_\alpha| = k$

🤖 Rejecting is expensive. How often should we do it?

The folklore heuristic was right (with the right α^*)

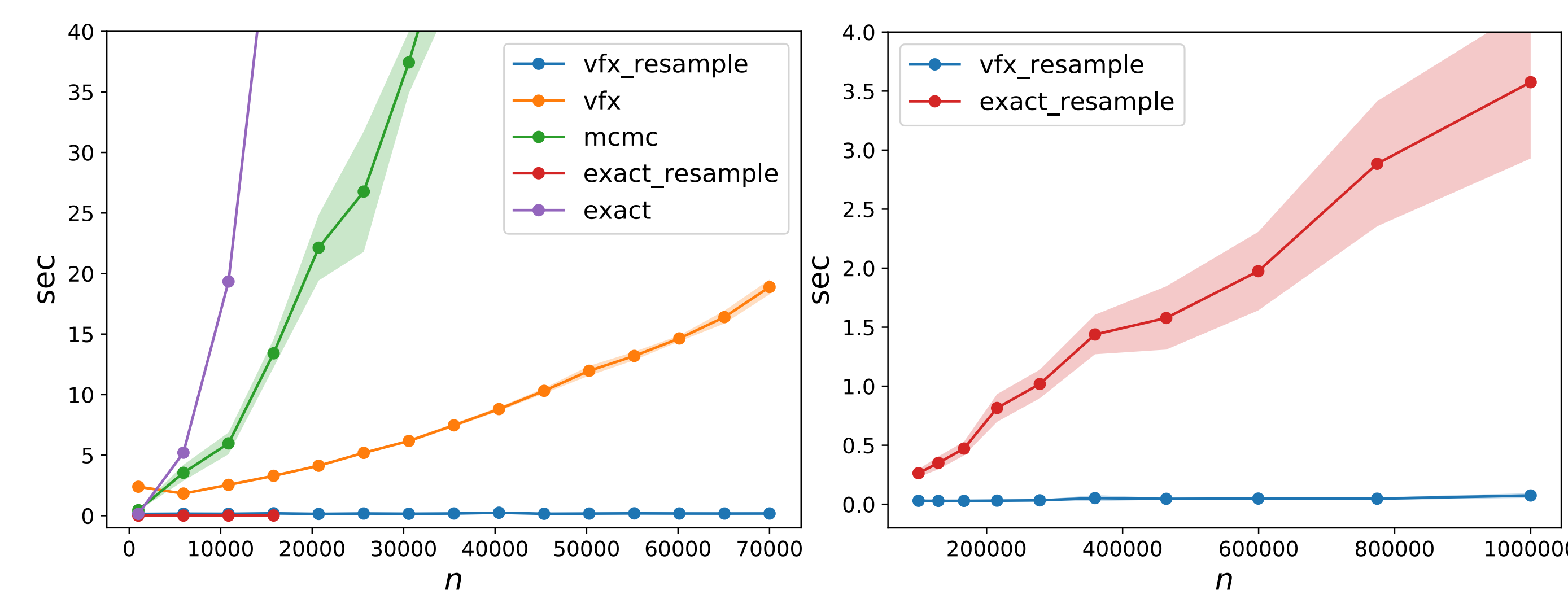
Thm. We can compute α^* in $\tilde{O}(n \cdot \text{poly}(k))$ such that

$$\text{Mode}(|S_{\alpha^*}|) = k, \quad \mathbf{P}(|S_{\alpha^*}| = k) \geq \Omega(k^{-1/2})$$

inducing at most $O(\sqrt{k})$ rejections.

Experiments

Sampling from MNIST8M with $n \in \{1..10^6\}$



	DPP-VFX		eigendecomp. (exact)		MCMC	
	first	succ.	first	succ.	first	succ.
$n = 15000$	4	1	54 (10x)	1 (1x)	13 (4x)	13 (13x)
$n = 70000$	19	1	DNF	DNF	175 (9x)	175 (175x)

Runtime in seconds and corresponding (speedup), DNF = Did Not Finish