

# Determinantal Point Processes for Coresets

(Journées MAS, Rouen, August 2022)

Nicolas Tremblay, Simon Barthelmé, Pierre-Olivier Amblard

CNRS, GIPSA-lab, Univ. Grenoble-Alpes, France



## Illustration and context

### Coresets: definition and iid theorem

- Coresets

- Sensitivities

- A classical iid coreset result

### DPPs for Coresets

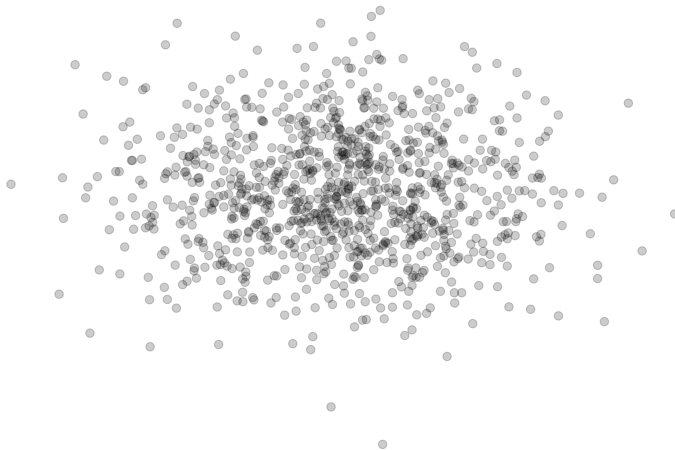
- Determinantal Point Processes

- A theoretical point-of-view

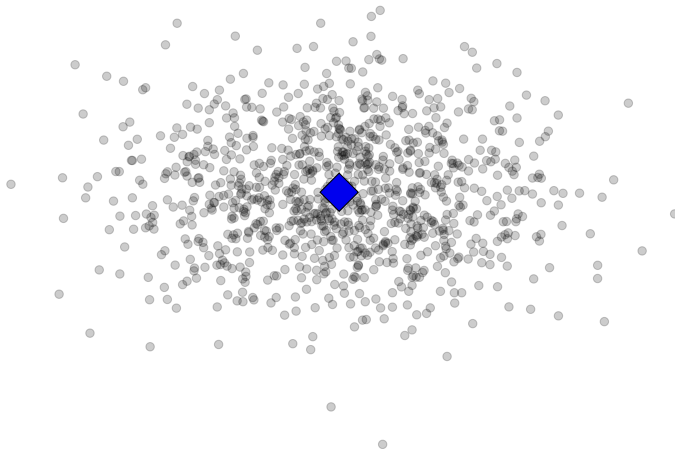
- A practical point-of-view

## Conclusion

## Task: find the mean



## Task: find the mean



## Task: find the mean

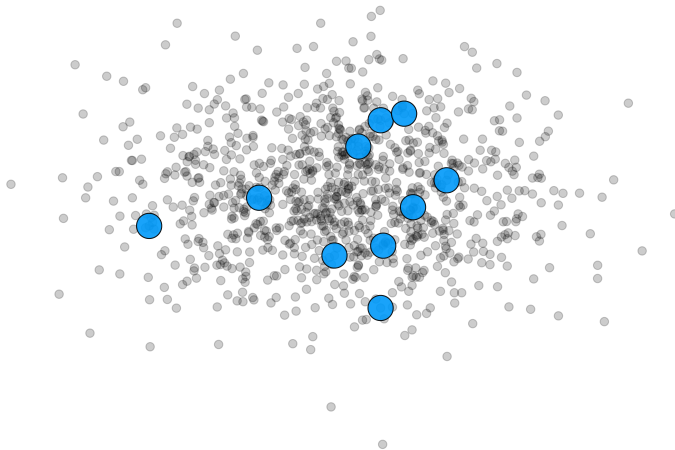


Figure: Example of iid uniform sampling

## Task: find the mean

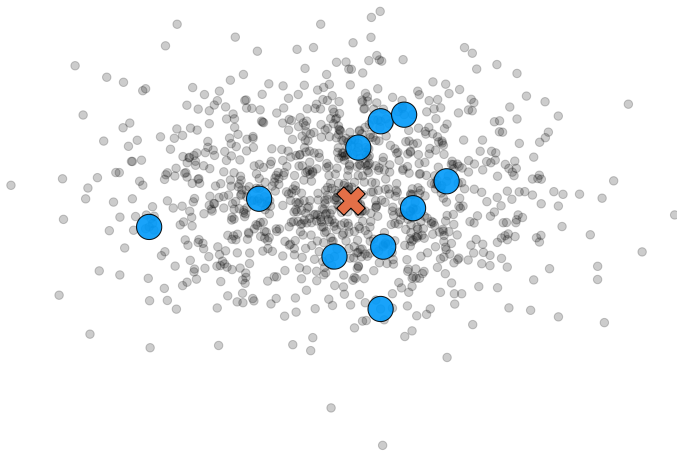


Figure: Example of iid uniform sampling

## Task: find the mean

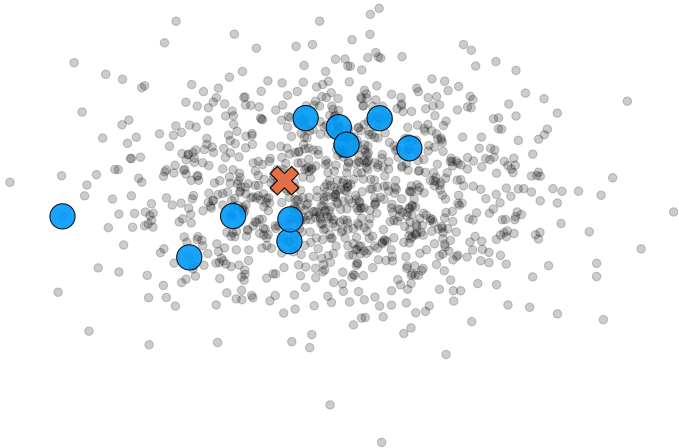


Figure: Example of iid uniform sampling

## Task: find the mean

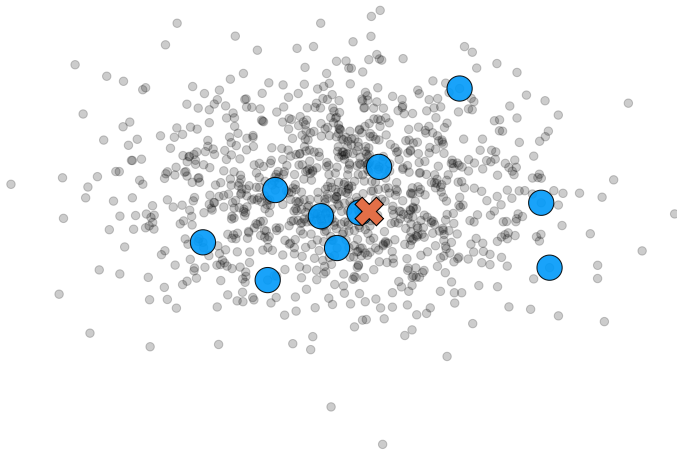


Figure: Example of iid uniform sampling



## Task: find the mean

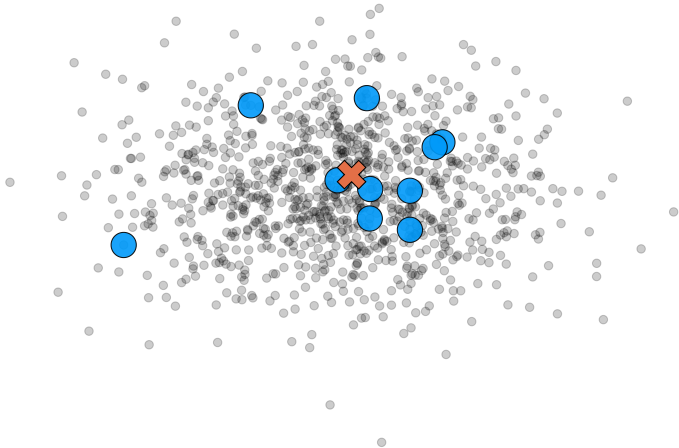


Figure: Example of iid uniform sampling

## Task: find the mean

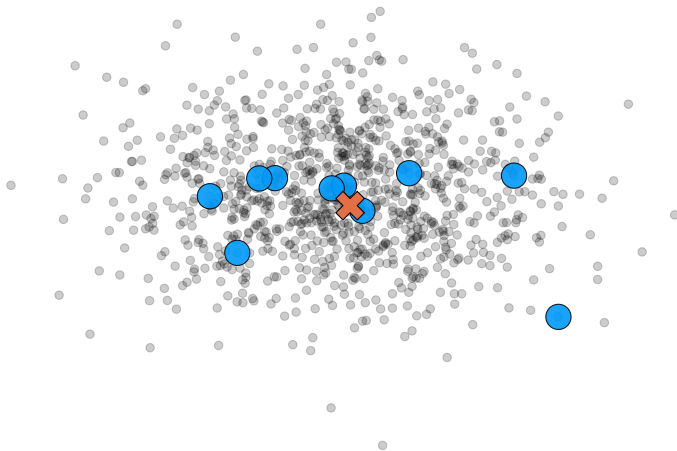


Figure: Example of iid uniform sampling

## Task: find the mean

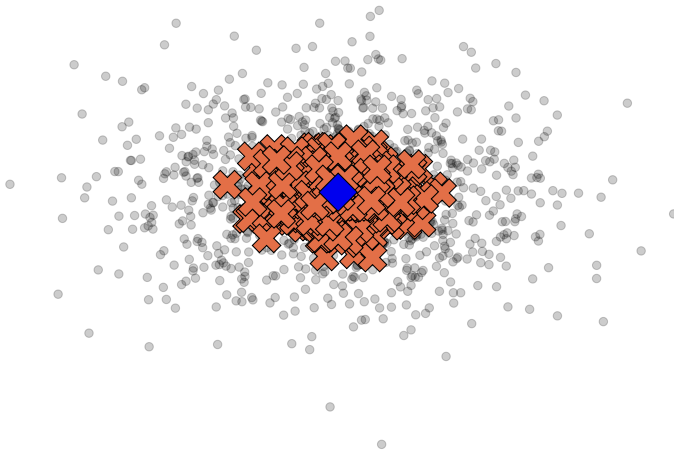


Figure: Uniform iid estimations of the mean

## Task: find the mean

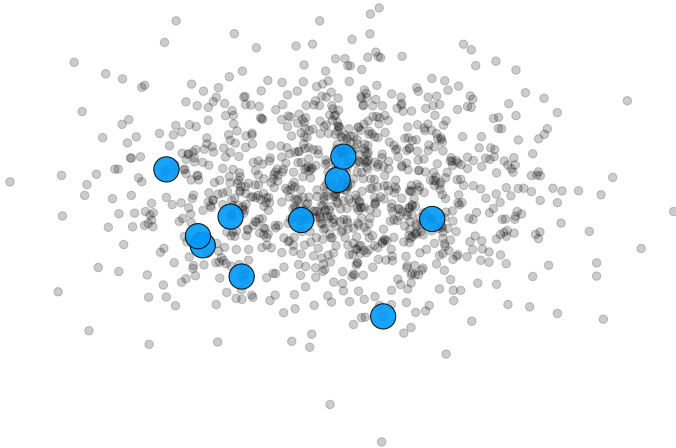


Figure: Example of a smarter iid sampling (sensitivity-based importance sampling)

## Task: find the mean

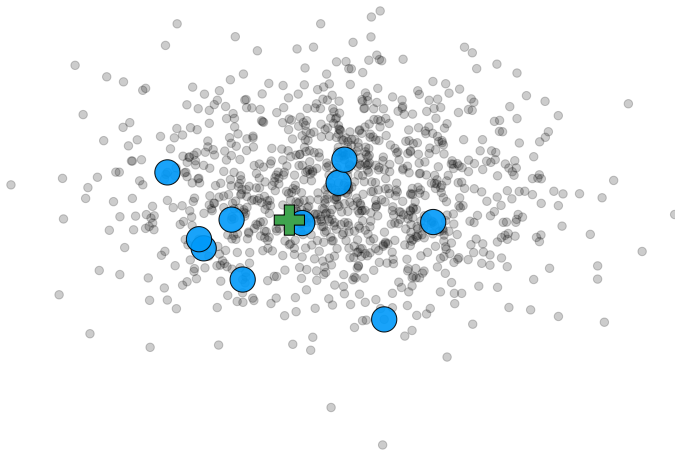


Figure: Example of a smarter iid sampling (sensitivity-based importance sampling)

## Task: find the mean

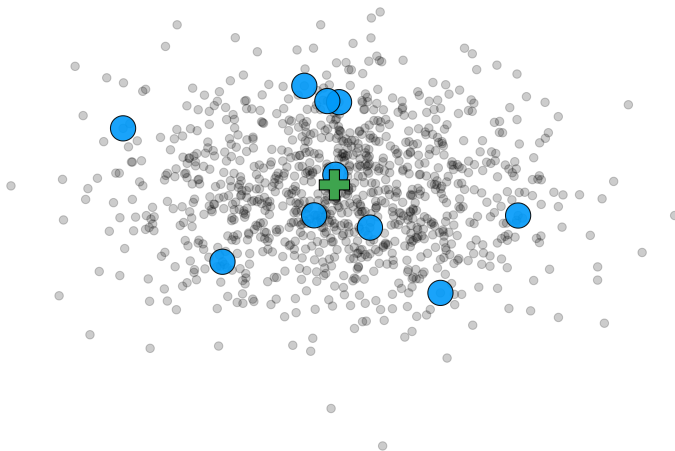


Figure: Example of a smarter iid sampling (sensitivity-based importance sampling)

## Task: find the mean

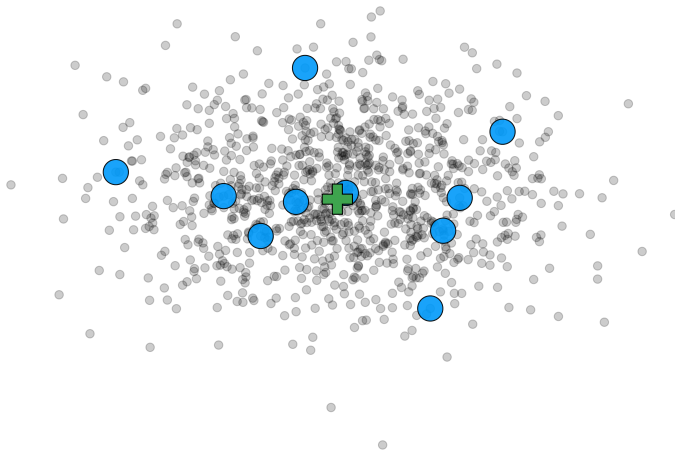


Figure: Example of a smarter iid sampling (sensitivity-based importance sampling)

## Task: find the mean

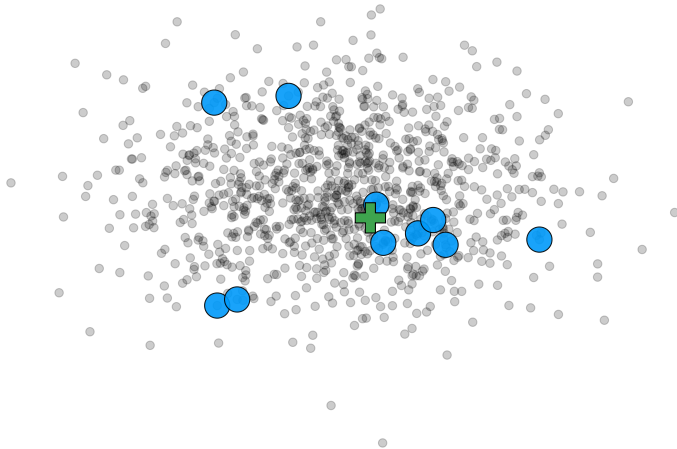


Figure: Example of a smarter iid sampling (sensitivity-based importance sampling)



## Task: find the mean

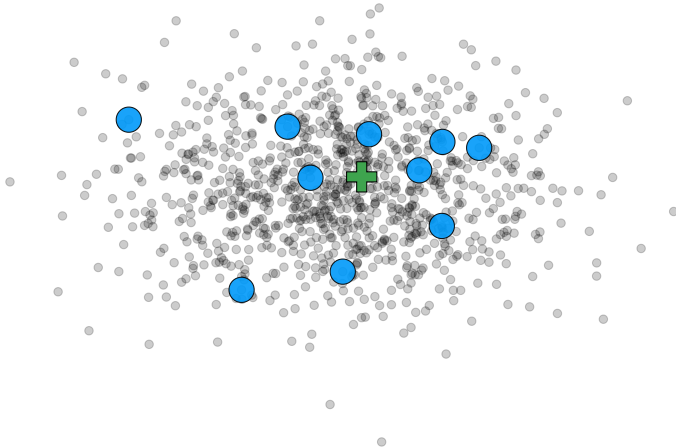


Figure: Example of a smarter iid sampling (sensitivity-based importance sampling)

## Task: find the mean

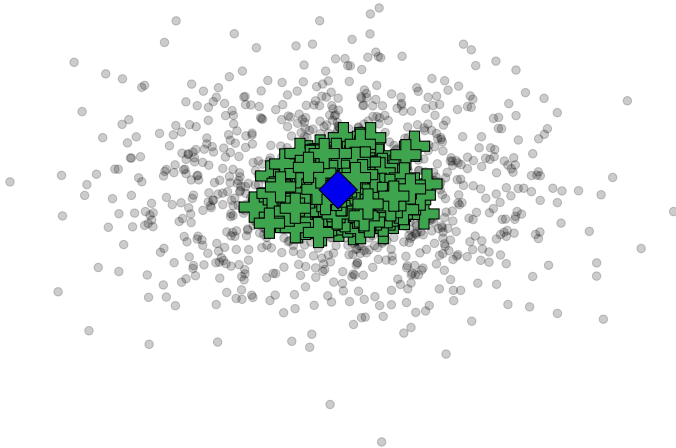


Figure: Smarter iid estimations of the mean (sensitivity-based importance sampling)

## Task: find the mean

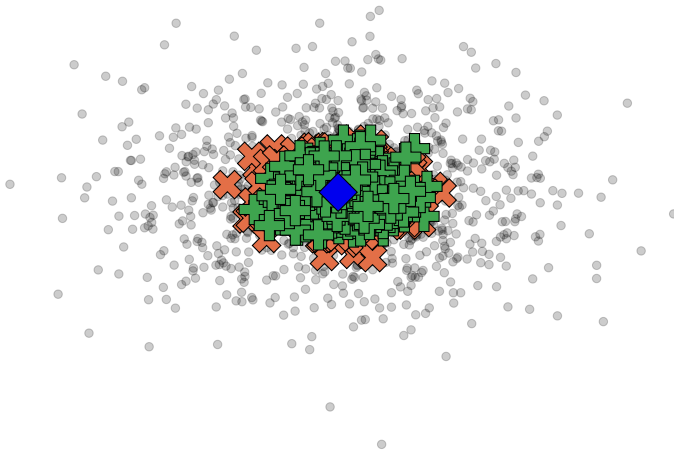


Figure: Comparison of both estimators: **variance reduction**

## Task: find the mean

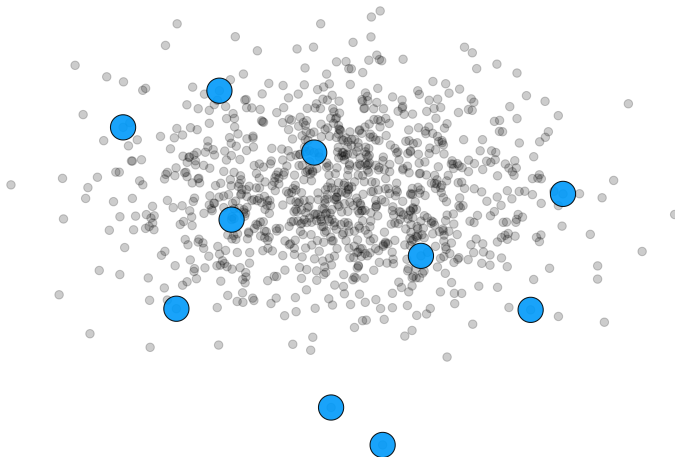


Figure: Example of DPP sampling (using Alg. 1 of the paper)

## Task: find the mean

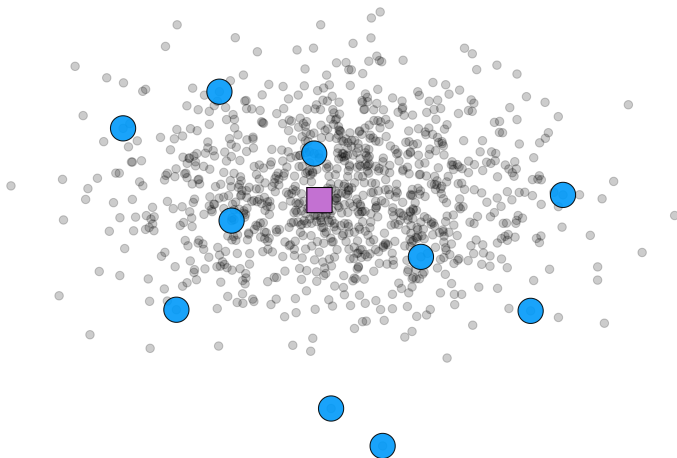


Figure: Example of DPP sampling (using Alg. 1 of the paper)

## Task: find the mean

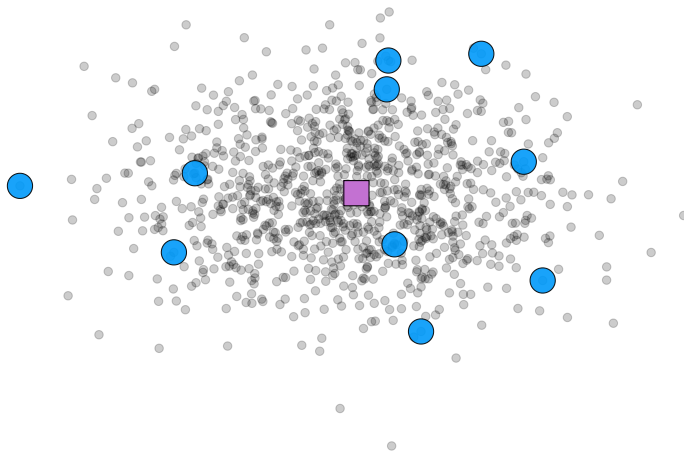


Figure: Example of DPP sampling (using Alg. 1 of the paper)

## Task: find the mean

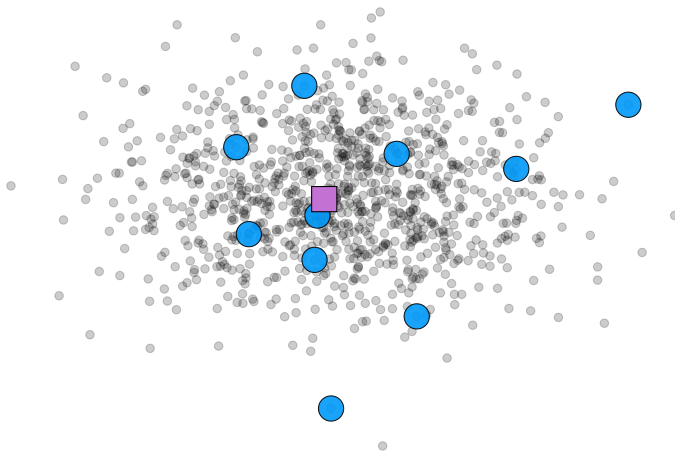


Figure: Example of DPP sampling (using Alg. 1 of the paper)

## Task: find the mean

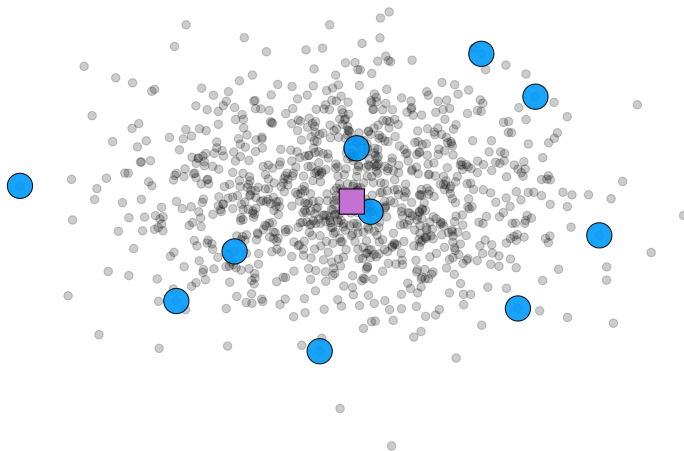


Figure: Example of DPP sampling (using Alg. 1 of the paper)



## Task: find the mean

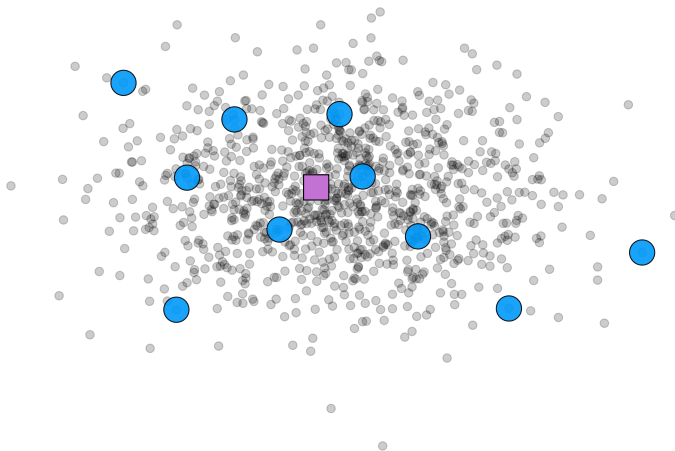


Figure: Example of DPP sampling (using Alg. 1 of the paper)

## Task: find the mean

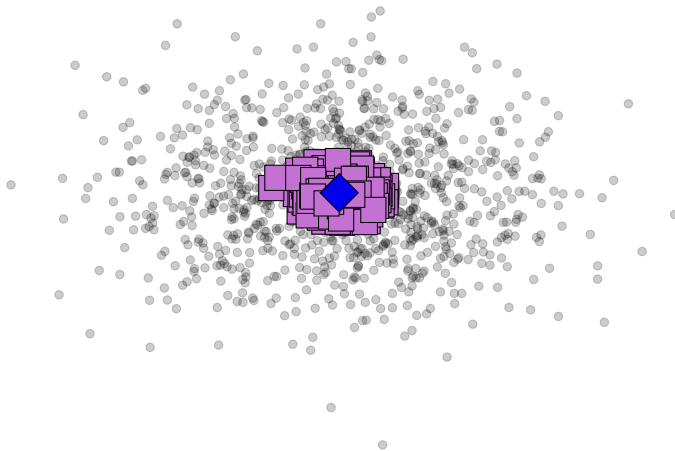


Figure: DPP estimations of the mean (using Alg. 1 of the paper)

## Task: find the mean

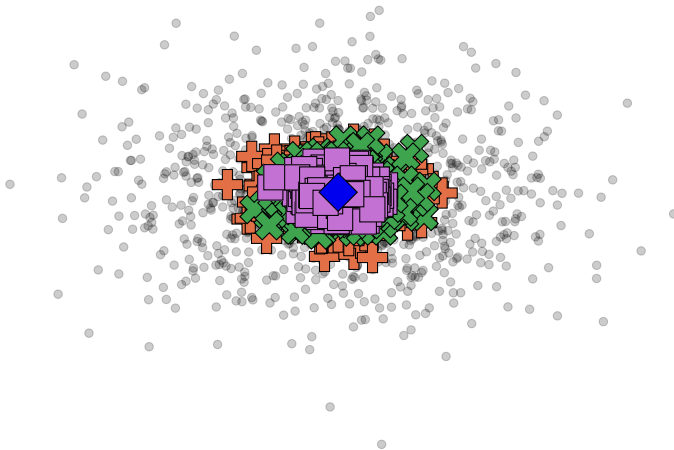
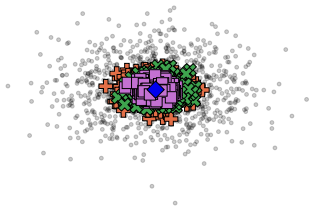


Figure: Comparison of all estimators: **more variance reduction**

## Context and goal of coresets



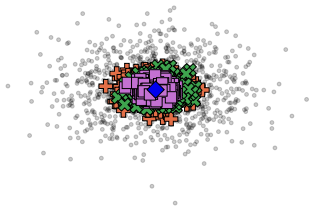
### Context

- (very) large  $n$  (size of dataset)
- a precise task at hand

---

<sup>1</sup>Munteanu and Schwiegelshohn, *Coresets-Methods and History: A Theoreticians Design Pattern...*, KI, 2017

## Context and goal of coresets



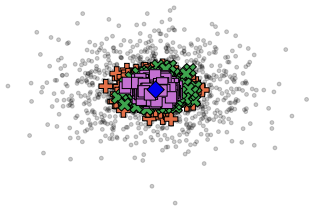
### Context

- (very) large  $n$  (size of dataset)
- a precise task at hand

### Goal

- a coreset is a **tiny** (size indep./polylog of  $n$ ) sample of the data **for the task** at hand
- a coreset has **provable guarantees** on the error made
- a coreset should be **sampled faster** than solving the task on the original data (!)

## Context and goal of coresets



### Context

- (very) large  $n$  (size of dataset)
- a precise task at hand

### Goal

- a coreset is a **tiny** (size indep./polylog of  $n$ ) sample of the data **for the task** at hand
- a coreset has **provable guarantees** on the error made
- a coreset should be **sampled faster** than solving the task on the original data (!)

### State-of-the-art<sup>1</sup>

- Verifying all 3 points is very challenging. The state-of-the-art usually comes with
  - a provable algorithm but very expensive
  - some work-arounds more affordable, still provable, but (much) less efficient
  - some heuristics inspired by these provable algorithms
- coresets under research: deterministic, iid random, multi-task, streaming, ...

<sup>1</sup>Munteanu and Schwiegelshohn, *Coresets-Methods and History: A Theoreticians Design Pattern...*, KI, 2017

## Illustration and context

### Coresets: definition and iid theorem

- Coresets

- Sensitivities

- A classical iid coreset result

### DPPs for Coresets

- Determinantal Point Processes

- A theoretical point-of-view

- A practical point-of-view

## Conclusion

## A generic class of problems

- Consider a dataset  $\mathcal{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ , say:  $n$  points in dimension  $d$ .
- Let  $\Theta$  be a parameter space and consider cost functions of the form:

$$L(\mathcal{X}, \theta) = \sum_{i=1}^n f(\mathbf{x}_i, \theta)$$

where  $f : \mathcal{X}, \Theta \rightarrow \mathbb{R}^+$ .

- A classical ML objective: find

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} L(\mathcal{X}, \theta).$$



## Examples falling in this class of problems

Find  $\theta^* = \operatorname{argmin}_{\theta \in \Theta} L(\mathcal{X}, \theta)$  where  $L(\mathcal{X}, \theta) = \sum_{i=1}^n f(x_i, \theta)$ , and  $f : \mathcal{X}, \Theta \rightarrow \mathbb{R}^+$ .

## Examples falling in this class of problems

Find  $\theta^* = \operatorname{argmin}_{\theta \in \Theta} L(\mathcal{X}, \theta)$  where  $L(\mathcal{X}, \theta) = \sum_{i=1}^n f(\mathbf{x}_i, \theta)$ , and  $f : \mathcal{X}, \Theta \rightarrow \mathbb{R}^+$ .

- ***k*-means** The *k*-means objective is to find *k* centroids  $\{\mathbf{c}_\ell\}_{\ell=1, \dots, k}$  in  $\mathbb{R}^d$  such that  $L(\mathcal{X}, \{\mathbf{c}_\ell\}) = \sum_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}, \{\mathbf{c}_\ell\})$  is minimal, with

$$f(\mathbf{x}, \{\mathbf{c}_\ell\}) = \min_{\mathbf{c}_\ell} \|\mathbf{x} - \mathbf{c}_\ell\|^2$$

## Examples falling in this class of problems

Find  $\theta^* = \operatorname{argmin}_{\theta \in \Theta} L(\mathcal{X}, \theta)$  where  $L(\mathcal{X}, \theta) = \sum_{i=1}^n f(\mathbf{x}_i, \theta)$ , and  $f : \mathcal{X}, \Theta \rightarrow \mathbb{R}^+$ .

- ***k*-means** The *k*-means objective is to find *k* centroids  $\{\mathbf{c}_\ell\}_{\ell=1, \dots, k}$  in  $\mathbb{R}^d$  such that  $L(\mathcal{X}, \{\mathbf{c}_\ell\}) = \sum_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}, \{\mathbf{c}_\ell\})$  is minimal, with

$$f(\mathbf{x}, \{\mathbf{c}_\ell\}) = \min_{\mathbf{c}_\ell} \|\mathbf{x} - \mathbf{c}_\ell\|^2 \quad (\geq 0).$$

## Examples falling in this class of problems

Find  $\theta^* = \operatorname{argmin}_{\theta \in \Theta} L(\mathcal{X}, \theta)$  where  $L(\mathcal{X}, \theta) = \sum_{i=1}^n f(\mathbf{x}_i, \theta)$ , and  $f : \mathcal{X}, \Theta \rightarrow \mathbb{R}^+$ .

- ***k*-means** The *k*-means objective is to find *k* centroids  $\{\mathbf{c}_\ell\}_{\ell=1, \dots, k}$  in  $\mathbb{R}^d$  such that  $L(\mathcal{X}, \{\mathbf{c}_\ell\}) = \sum_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}, \{\mathbf{c}_\ell\})$  is minimal, with

$$f(\mathbf{x}, \{\mathbf{c}_\ell\}) = \min_{\mathbf{c}_\ell} \|\mathbf{x} - \mathbf{c}_\ell\|^2 \quad (\geq 0).$$

- **Other examples:** linear regression, logistic regression, *k*-median, low-rank approx. of matrices, etc.

## Coresets: definition

- Consider  $\{\mathcal{S}, \{\omega_s > 0\}\}$  a weighted sample of  $\mathcal{X}$  and its associated estimated cost

$$\hat{L}(\mathcal{S}, \theta) = \sum_{\mathbf{s} \in \mathcal{S}} \omega_s f(\mathbf{s}, \theta)$$

## Coresets: definition

- Consider  $\{S, \{\omega_s > 0\}\}$  a weighted sample of  $\mathcal{X}$  and its associated estimated cost

$$\hat{L}(S, \theta) = \sum_{s \in S} \omega_s f(s, \theta)$$

- **Definition** ( $\epsilon$ -coreset) *A weighted sample  $S$  is an  $\epsilon$ -coreset of  $\mathcal{X}$  wrt  $L$  if:*

$$\forall \theta \in \Theta \quad (1 - \epsilon)L(\mathcal{X}, \theta) \leq \hat{L}(S, \theta) \leq (1 + \epsilon)L(\mathcal{X}, \theta)$$

## Coresets: definition

- Consider  $\{S, \{\omega_s > 0\}\}$  a weighted sample of  $\mathcal{X}$  and its associated estimated cost

$$\hat{L}(S, \theta) = \sum_{s \in S} \omega_s f(s, \theta)$$

- **Definition** ( $\epsilon$ -coreset) *A weighted sample  $S$  is an  $\epsilon$ -coreset of  $\mathcal{X}$  wrt  $L$  if:*

$$\forall \theta \in \Theta \quad (1 - \epsilon)L(\mathcal{X}, \theta) \leq \hat{L}(S, \theta) \leq (1 + \epsilon)L(\mathcal{X}, \theta)$$

## Coresets: definition

- Consider  $\{\mathcal{S}, \{\omega_s > 0\}\}$  a weighted sample of  $\mathcal{X}$  and its associated estimated cost

$$\hat{L}(\mathcal{S}, \theta) = \sum_{s \in \mathcal{S}} \omega_s f(s, \theta)$$

- Definition** ( $\epsilon$ -coreset) *A weighted sample  $\mathcal{S}$  is an  $\epsilon$ -coreset of  $\mathcal{X}$  wrt  $L$  if:*

$$\forall \theta \in \Theta \quad (1 - \epsilon)L(\mathcal{X}, \theta) \leq \hat{L}(\mathcal{S}, \theta) \leq (1 + \epsilon)L(\mathcal{X}, \theta)$$

- Denote by  $\hat{\theta}^*$  the argmin of  $\hat{L}$ :  $\hat{\theta}^* = \underset{\theta \in \Theta}{\operatorname{argmin}} \hat{L}(\mathcal{S}, \theta)$ .



## Coresets: definition

- Consider  $\{\mathcal{S}, \{\omega_s > 0\}\}$  a weighted sample of  $\mathcal{X}$  and its associated estimated cost

$$\hat{L}(\mathcal{S}, \theta) = \sum_{s \in \mathcal{S}} \omega_s f(s, \theta)$$

- **Definition** ( $\epsilon$ -coreset) *A weighted sample  $\mathcal{S}$  is an  $\epsilon$ -coreset of  $\mathcal{X}$  wrt  $L$  if:*

$$\forall \theta \in \Theta \quad (1 - \epsilon)L(\mathcal{X}, \theta) \leq \hat{L}(\mathcal{S}, \theta) \leq (1 + \epsilon)L(\mathcal{X}, \theta)$$

- Denote by  $\hat{\theta}^*$  the argmin of  $\hat{L}$ :  $\hat{\theta}^* = \operatorname{argmin}_{\theta \in \Theta} \hat{L}(\mathcal{S}, \theta)$ .
- Why are coresets interesting?  
If  $L$  has a clear global minimum (up to an  $\epsilon$  factor) in  $\theta^*$ , then  $\hat{\theta}^* \simeq \theta^*$ :

## Coresets: definition

- Consider  $\{\mathcal{S}, \{\omega_s > 0\}\}$  a weighted sample of  $\mathcal{X}$  and its associated estimated cost

$$\hat{L}(\mathcal{S}, \theta) = \sum_{s \in \mathcal{S}} \omega_s f(s, \theta)$$

- **Definition** ( $\epsilon$ -coreset) *A weighted sample  $\mathcal{S}$  is an  $\epsilon$ -coreset of  $\mathcal{X}$  wrt  $L$  if:*

$$\forall \theta \in \Theta \quad (1 - \epsilon)L(\mathcal{X}, \theta) \leq \hat{L}(\mathcal{S}, \theta) \leq (1 + \epsilon)L(\mathcal{X}, \theta)$$

- Denote by  $\hat{\theta}^*$  the argmin of  $\hat{L}$ :  $\hat{\theta}^* = \operatorname{argmin}_{\theta \in \Theta} \hat{L}(\mathcal{S}, \theta)$ .

- Why are coresets interesting?

If  $L$  has a clear global minimum (up to an  $\epsilon$  factor) in  $\theta^*$ , then  $\hat{\theta}^* \simeq \theta^*$ :

$$\hat{L}(\mathcal{S}, \hat{\theta}^*)$$

## Coresets: definition

- Consider  $\{\mathcal{S}, \{\omega_s > 0\}\}$  a weighted sample of  $\mathcal{X}$  and its associated estimated cost

$$\hat{L}(\mathcal{S}, \theta) = \sum_{s \in \mathcal{S}} \omega_s f(s, \theta)$$

- Definition** ( $\epsilon$ -coreset) A weighted sample  $\mathcal{S}$  is an  $\epsilon$ -coreset of  $\mathcal{X}$  wrt  $L$  if:

$$\forall \theta \in \Theta \quad (1 - \epsilon)L(\mathcal{X}, \theta) \leq \hat{L}(\mathcal{S}, \theta) \leq (1 + \epsilon)L(\mathcal{X}, \theta)$$

- Denote by  $\hat{\theta}^*$  the argmin of  $\hat{L}$ :  $\hat{\theta}^* = \operatorname{argmin}_{\theta \in \Theta} \hat{L}(\mathcal{S}, \theta)$ .

- Why are coresets interesting?

If  $L$  has a clear global minimum (up to an  $\epsilon$  factor) in  $\theta^*$ , then  $\hat{\theta}^* \simeq \theta^*$ :

$$\hat{L}(\mathcal{S}, \hat{\theta}^*) \leq \hat{L}(\mathcal{S}, \theta^*)$$

## Coresets: definition

- Consider  $\{S, \{\omega_s > 0\}\}$  a weighted sample of  $\mathcal{X}$  and its associated estimated cost

$$\hat{L}(S, \theta) = \sum_{s \in S} \omega_s f(s, \theta)$$

- Definition** ( $\epsilon$ -coreset) *A weighted sample  $S$  is an  $\epsilon$ -coreset of  $\mathcal{X}$  wrt  $L$  if:*

$$\forall \theta \in \Theta \quad (1 - \epsilon)L(\mathcal{X}, \theta) \leq \hat{L}(S, \theta) \leq (1 + \epsilon)L(\mathcal{X}, \theta)$$

- Denote by  $\hat{\theta}^*$  the argmin of  $\hat{L}$ :  $\hat{\theta}^* = \operatorname{argmin}_{\theta \in \Theta} \hat{L}(S, \theta)$ .

- Why are coresets interesting?

If  $L$  has a clear global minimum (up to an  $\epsilon$  factor) in  $\theta^*$ , then  $\hat{\theta}^* \simeq \theta^*$ :

$$\hat{L}(S, \hat{\theta}^*) \leq \hat{L}(S, \theta^*) \leq (1 + \epsilon)L(\mathcal{X}, \theta^*)$$

## Coresets: definition

- Consider  $\{\mathcal{S}, \{\omega_s > 0\}\}$  a weighted sample of  $\mathcal{X}$  and its associated estimated cost

$$\hat{L}(\mathcal{S}, \theta) = \sum_{s \in \mathcal{S}} \omega_s f(s, \theta)$$

- Definition** ( $\epsilon$ -coreset) *A weighted sample  $\mathcal{S}$  is an  $\epsilon$ -coreset of  $\mathcal{X}$  wrt  $L$  if:*

$$\forall \theta \in \Theta \quad (1 - \epsilon)L(\mathcal{X}, \theta) \leq \hat{L}(\mathcal{S}, \theta) \leq (1 + \epsilon)L(\mathcal{X}, \theta)$$

- Denote by  $\hat{\theta}^*$  the argmin of  $\hat{L}$ :  $\hat{\theta}^* = \underset{\theta \in \Theta}{\operatorname{argmin}} \hat{L}(\mathcal{S}, \theta)$ .

- Why are coresets interesting?

If  $L$  has a clear global minimum (up to an  $\epsilon$  factor) in  $\theta^*$ , then  $\hat{\theta}^* \simeq \theta^*$ :

$$(1 - \epsilon)L(\mathcal{X}, \hat{\theta}^*) \leq \hat{L}(\mathcal{S}, \hat{\theta}^*) \leq \hat{L}(\mathcal{S}, \theta^*) \leq (1 + \epsilon)L(\mathcal{X}, \theta^*)$$

## Coresets: definition

- Consider  $\{\mathcal{S}, \{\omega_s > 0\}\}$  a weighted sample of  $\mathcal{X}$  and its associated estimated cost

$$\hat{L}(\mathcal{S}, \theta) = \sum_{s \in \mathcal{S}} \omega_s f(s, \theta)$$

- Definition** ( $\epsilon$ -coreset) *A weighted sample  $\mathcal{S}$  is an  $\epsilon$ -coreset of  $\mathcal{X}$  wrt  $L$  if:*

$$\forall \theta \in \Theta \quad (1 - \epsilon)L(\mathcal{X}, \theta) \leq \hat{L}(\mathcal{S}, \theta) \leq (1 + \epsilon)L(\mathcal{X}, \theta)$$

- Denote by  $\hat{\theta}^*$  the argmin of  $\hat{L}$ :  $\hat{\theta}^* = \operatorname{argmin}_{\theta \in \Theta} \hat{L}(\mathcal{S}, \theta)$ .

- Why are coresets interesting?

If  $L$  has a clear global minimum (up to an  $\epsilon$  factor) in  $\theta^*$ , then  $\hat{\theta}^* \simeq \theta^*$ :

$$(1 - \epsilon)L(\mathcal{X}, \theta^*) \leq (1 - \epsilon)L(\mathcal{X}, \hat{\theta}^*) \leq \hat{L}(\mathcal{S}, \hat{\theta}^*) \leq \hat{L}(\mathcal{S}, \theta^*) \leq (1 + \epsilon)L(\mathcal{X}, \theta^*)$$

## Coresets: definition

- Consider  $\{\mathcal{S}, \{\omega_s > 0\}\}$  a weighted sample of  $\mathcal{X}$  and its associated estimated cost

$$\hat{L}(\mathcal{S}, \theta) = \sum_{s \in \mathcal{S}} \omega_s f(s, \theta)$$

- Definition** ( $\epsilon$ -coreset) *A weighted sample  $\mathcal{S}$  is an  $\epsilon$ -coreset of  $\mathcal{X}$  wrt  $L$  if:*

$$\forall \theta \in \Theta \quad (1 - \epsilon)L(\mathcal{X}, \theta) \leq \hat{L}(\mathcal{S}, \theta) \leq (1 + \epsilon)L(\mathcal{X}, \theta)$$

- Denote by  $\hat{\theta}^*$  the argmin of  $\hat{L}$ :  $\hat{\theta}^* = \operatorname{argmin}_{\theta \in \Theta} \hat{L}(\mathcal{S}, \theta)$ .

- Why are coresets interesting?

If  $L$  has a clear global minimum (up to an  $\epsilon$  factor) in  $\theta^*$ , then  $\hat{\theta}^* \simeq \theta^*$ :

$$\begin{aligned} (1 - \epsilon)L(\mathcal{X}, \theta^*) &\leq (1 - \epsilon)L(\mathcal{X}, \hat{\theta}^*) \leq \hat{L}(\mathcal{S}, \hat{\theta}^*) \leq \hat{L}(\mathcal{S}, \theta^*) \leq (1 + \epsilon)L(\mathcal{X}, \theta^*) \\ &\Downarrow \\ (1 - \epsilon)L(\mathcal{X}, \theta^*) &\leq \hat{L}(\mathcal{S}, \hat{\theta}^*) \leq (1 + \epsilon)L(\mathcal{X}, \theta^*) \end{aligned}$$

# Coresets: illustration on the 1-means problem

- Data  $\mathcal{X}$





## Coresets: illustration on the 1-means problem

- Data  $\mathcal{X}$
- Cost function:

$$L(\mathcal{X}, \theta) = \sum_{i=1}^n \|\mathbf{x}_i - \theta\|^2$$



## Coresets: illustration on the 1-means problem

- Data  $\mathcal{X}$
- Cost function:

$$L(\mathcal{X}, \theta) = \sum_{i=1}^n \|\mathbf{x}_i - \theta\|^2$$

- Optimal  $\theta$ :

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} L(\mathcal{X}, \theta)$$



## Coresets: illustration on the 1-means problem

- Data  $\mathcal{X}$

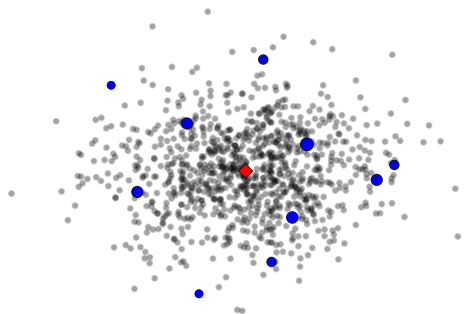
- Cost function:

$$L(\mathcal{X}, \theta) = \sum_{i=1}^n \|\mathbf{x}_i - \theta\|^2$$

- Optimal  $\theta$ :

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} L(\mathcal{X}, \theta)$$

- A weighted subset  $\mathcal{S}$



## Coresets: illustration on the 1-means problem

- Data  $\mathcal{X}$
- Cost function:

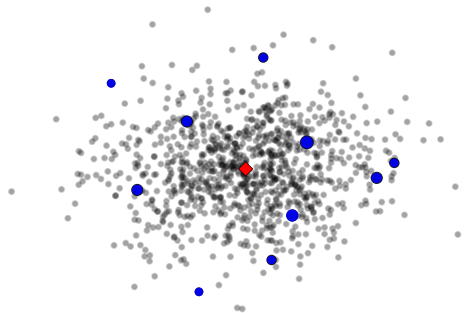
$$L(\mathcal{X}, \theta) = \sum_{i=1}^n \|\mathbf{x}_i - \theta\|^2$$

- Optimal  $\theta$ :

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} L(\mathcal{X}, \theta)$$

- A weighted subset  $\mathcal{S}$
- Estimated cost function:

$$\hat{L}(\mathcal{S}, \theta) = \sum_{\mathbf{s} \in \mathcal{S}} \omega_{\mathbf{s}} \|\mathbf{s} - \theta\|^2$$



## Coresets: illustration on the 1-means problem

- Data  $\mathcal{X}$
- Cost function:

$$L(\mathcal{X}, \theta) = \sum_{i=1}^n \|\mathbf{x}_i - \theta\|^2$$

- Optimal  $\theta$ :

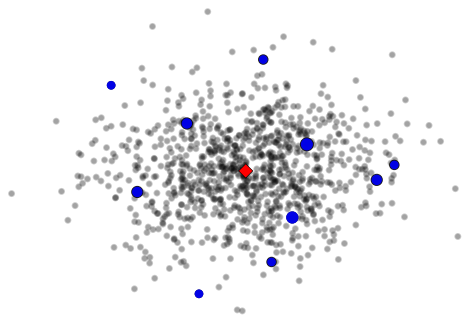
$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} L(\mathcal{X}, \theta)$$

- A weighted subset  $\mathcal{S}$
- Estimated cost function:

$$\hat{L}(\mathcal{S}, \theta) = \sum_{\mathbf{s} \in \mathcal{S}} \omega_{\mathbf{s}} \|\mathbf{s} - \theta\|^2$$

- $\mathcal{S}$  is a  $\epsilon$ -coreset if:

$$\forall \theta \quad \left| \frac{\hat{L}}{L} - 1 \right| \leq \epsilon$$



## Coresets: illustration on the 1-means problem

- Data  $\mathcal{X}$
- Cost function:

$$L(\mathcal{X}, \theta) = \sum_{i=1}^n \|\mathbf{x}_i - \theta\|^2$$

- Optimal  $\theta$ :

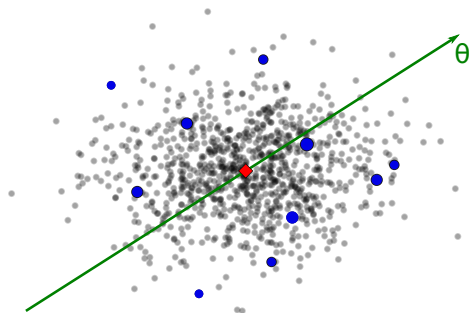
$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} L(\mathcal{X}, \theta)$$

- A weighted subset  $\mathcal{S}$
- Estimated cost function:

$$\hat{L}(\mathcal{S}, \theta) = \sum_{\mathbf{s} \in \mathcal{S}} \omega_{\mathbf{s}} \|\mathbf{s} - \theta\|^2$$

- $\mathcal{S}$  is a  $\epsilon$ -coreset if:

$$\forall \theta \quad \left| \frac{\hat{L}}{L} - 1 \right| \leq \epsilon$$



## Coresets: illustration on the 1-means problem

- Data  $\mathcal{X}$
- Cost function:

$$L(\mathcal{X}, \theta) = \sum_{i=1}^n \|x_i - \theta\|^2$$

- Optimal  $\theta$ :

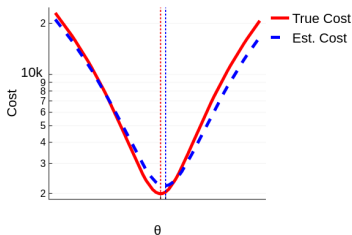
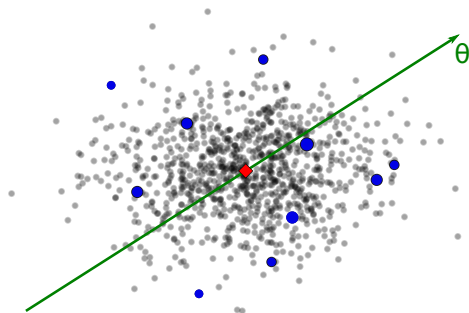
$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} L(\mathcal{X}, \theta)$$

- A weighted subset  $\mathcal{S}$
- Estimated cost function:

$$\hat{L}(\mathcal{S}, \theta) = \sum_{s \in \mathcal{S}} \omega_s \|s - \theta\|^2$$

- $\mathcal{S}$  is a  $\epsilon$ -coreset if:

$$\forall \theta \quad \left| \frac{\hat{L}}{L} - 1 \right| \leq \epsilon$$



## Coresets: illustration on the 1-means problem

- Data  $\mathcal{X}$
- Cost function:

$$L(\mathcal{X}, \theta) = \sum_{i=1}^n \|x_i - \theta\|^2$$

- Optimal  $\theta$ :

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} L(\mathcal{X}, \theta)$$

- A weighted subset  $\mathcal{S}$
- Estimated cost function:

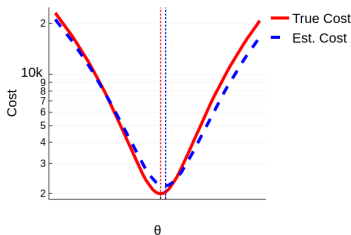
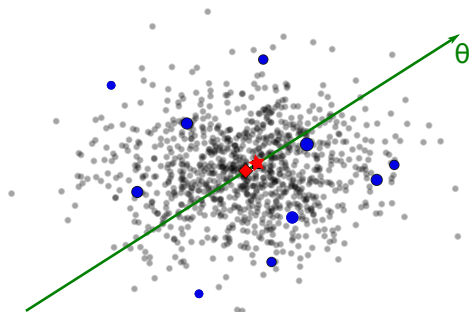
$$\hat{L}(\mathcal{S}, \theta) = \sum_{s \in \mathcal{S}} \omega_s \|s - \theta\|^2$$

- $\mathcal{S}$  is a  $\epsilon$ -coreset if:

$$\forall \theta \quad \left| \frac{\hat{L}}{L} - 1 \right| \leq \epsilon$$

- Estimated optimal  $\theta$ :

$$\hat{\theta}^* = \operatorname{argmin}_{\theta \in \Theta} \hat{L}(\mathcal{S}, \theta)$$





## Sensitivity: a central object

- The **sensitivity** of a datapoint  $\mathbf{x}_i \in \mathcal{X}$  with respect to  $f : \mathcal{X}, \Theta \rightarrow \mathbb{R}^+$  is:

$$\sigma_i = \max_{\theta \in \Theta} \frac{f(\mathbf{x}_i, \theta)}{L(\mathcal{X}, \theta)} \in [0, 1].$$

The total sensitivity is denoted  $\mathfrak{S} = \sum_{i=1}^n \sigma_i$ .

## Sensitivity: a central object

- The **sensitivity** of a datapoint  $\mathbf{x}_i \in \mathcal{X}$  with respect to  $f : \mathcal{X}, \Theta \rightarrow \mathbb{R}^+$  is:

$$\sigma_i = \max_{\theta \in \Theta} \frac{f(\mathbf{x}_i, \theta)}{L(\mathcal{X}, \theta)} \in [0, 1].$$

The total sensitivity is denoted  $\mathfrak{S} = \sum_{i=1}^n \sigma_i$ .

- In general, the sensitivity is unknown analytically. In the paper, we managed to compute the analytic sensitivities for two simple cases: 1-means and linear regression

## Random coresets

- Random context: suppose  $\mathcal{S}$  is a *random* subset  $\mathcal{S} \subset \mathcal{X}$  (possibly with repetitions)
- Importance sampling notations:
  - Define  $\epsilon_i$  the *random variable* counting the number of times  $\mathbf{x}_i$  is in  $\mathcal{S}$
  - To each element  $\mathbf{x}_i$  associate the weight  $\omega_i = \frac{1}{\mathbb{E}(\epsilon_i)}$

- One has:

$$\hat{L}(\mathcal{S}, \theta) = \sum_{i=1}^n f(\mathbf{x}_i, \theta) \frac{\epsilon_i}{\mathbb{E}(\epsilon_i)}.$$

- By construction,  $\hat{L}$  is an unbiased estimator of  $L$ :

$$\mathbb{E} \left( \hat{L}(\mathcal{S}, \theta) \right) = \sum_{i=1}^n f(\mathbf{x}_i, \theta) = L(\mathcal{X}, \theta).$$

## A classical iid coreset theorem<sup>1</sup>

- Let  $\mathbf{p} \in [0, 1]^n$  be a probability distribution over all datapoints  $\mathcal{X}$  with  $p_i$  the probability of sampling  $\mathbf{x}_i$  and  $\sum_i p_i = 1$ .

---

<sup>1</sup>Langberg and Schulman, *Universal  $\epsilon$ -approximators for integrals*, SIAM, 2010

## A classical iid coreset theorem<sup>1</sup>

- Let  $\mathbf{p} \in [0, 1]^n$  be a probability distribution over all datapoints  $\mathcal{X}$  with  $p_i$  the probability of sampling  $\mathbf{x}_i$  and  $\sum_i p_i = 1$ .
- Draw  $\mathcal{S}$ :  $m$  iid samples with replacement according to  $\mathbf{p}$ .

---

<sup>1</sup>Langberg and Schulman, *Universal  $\epsilon$ -approximators for integrals*, SIAM, 2010

## A classical iid coreset theorem<sup>1</sup>

- Let  $\mathbf{p} \in [0, 1]^n$  be a probability distribution over all datapoints  $\mathcal{X}$  with  $p_i$  the probability of sampling  $\mathbf{x}_i$  and  $\sum_i p_i = 1$ .
- Draw  $\mathcal{S}$ :  $m$  iid samples with replacement according to  $\mathbf{p}$ .
- Associate importance sampling weights to each element of  $\mathcal{S}$ .

---

<sup>1</sup>Langberg and Schulman, *Universal  $\epsilon$ -approximators for integrals*, SIAM, 2010

## A classical iid coreset theorem<sup>1</sup>

- Let  $\mathbf{p} \in [0, 1]^n$  be a probability distribution over all datapoints  $\mathcal{X}$  with  $p_i$  the probability of sampling  $\mathbf{x}_i$  and  $\sum_i p_i = 1$ .
- Draw  $\mathcal{S}$ :  $m$  iid samples with replacement according to  $\mathbf{p}$ .
- Associate importance sampling weights to each element of  $\mathcal{S}$ .
- **Theorem** The weighted sample  $\mathcal{S}$  is a  $\epsilon$ -coreset with high probability if:

$$m \geq \mathcal{O} \left( \frac{d'}{\epsilon^2} \left( \max_i \frac{\sigma_i}{p_i} \right)^2 \right),$$

where  $d'$  is the pseudo-dimension of  $\Theta$  (a generalization of the Vapnik-Chervonenkis dimension).

---

<sup>1</sup>Langberg and Schulman, *Universal  $\epsilon$ -approximators for integrals*, SIAM, 2010

## A classical iid coreset theorem<sup>1</sup>

- Let  $\mathbf{p} \in [0, 1]^n$  be a probability distribution over all datapoints  $\mathcal{X}$  with  $p_i$  the probability of sampling  $\mathbf{x}_i$  and  $\sum_i p_i = 1$ .
- Draw  $\mathcal{S}$ :  $m$  iid samples with replacement according to  $\mathbf{p}$ .
- Associate importance sampling weights to each element of  $\mathcal{S}$ .
- **Theorem** The weighted sample  $\mathcal{S}$  is a  $\epsilon$ -coreset with high probability if:

$$m \geq \mathcal{O} \left( \frac{d'}{\epsilon^2} \left( \max_i \frac{\sigma_i}{p_i} \right)^2 \right),$$

where  $d'$  is the pseudo-dimension of  $\Theta$  (a generalization of the Vapnik-Chervonenkis dimension).

- The optimal probability distribution minimizing the rhs is  $p_i = \sigma_i / \mathfrak{S}$ .

---

<sup>1</sup>Langberg and Schulman, *Universal  $\epsilon$ -approximators for integrals*, SIAM, 2010



## A classical iid coreset theorem<sup>1</sup>

- Let  $\mathbf{p} \in [0, 1]^n$  be a probability distribution over all datapoints  $\mathcal{X}$  with  $p_i$  the probability of sampling  $\mathbf{x}_i$  and  $\sum_i p_i = 1$ .
- Draw  $\mathcal{S}$ :  $m$  iid samples with replacement according to  $\mathbf{p}$ .
- Associate importance sampling weights to each element of  $\mathcal{S}$ .
- **Theorem** The weighted sample  $\mathcal{S}$  is a  $\epsilon$ -coreset with high probability if:

$$m \geq \mathcal{O} \left( \frac{d'}{\epsilon^2} \left( \max_i \frac{\sigma_i}{p_i} \right)^2 \right),$$

where  $d'$  is the pseudo-dimension of  $\Theta$  (a generalization of the Vapnik-Chervonenkis dimension).

- The optimal probability distribution minimizing the rhs is  $p_i = \sigma_i / \mathfrak{S}$ .
- In this case,  $\mathcal{S}$  is a  $\epsilon$ -coreset with high probability if:

$$m \geq \mathcal{O} \left( \frac{d' \mathfrak{S}^2}{\epsilon^2} \right).$$

---

<sup>1</sup>Langberg and Schulman, *Universal  $\epsilon$ -approximators for integrals*, SIAM, 2010

## A classical iid coreset theorem<sup>1</sup>

- Let  $\mathbf{p} \in [0, 1]^n$  be a probability distribution over all datapoints  $\mathcal{X}$  with  $p_i$  the probability of sampling  $\mathbf{x}_i$  and  $\sum_i p_i = 1$ .
- Draw  $\mathcal{S}$ :  $m$  iid samples with replacement according to  $\mathbf{p}$ .
- Associate importance sampling weights to each element of  $\mathcal{S}$ .
- **Theorem** The weighted sample  $\mathcal{S}$  is a  $\epsilon$ -coreset with high probability if:

$$m \geq \mathcal{O} \left( \frac{d'}{\epsilon^2} \left( \max_i \frac{\sigma_i}{p_i} \right)^2 \right),$$

where  $d'$  is the pseudo-dimension of  $\Theta$  (a generalization of the Vapnik-Chervonenkis dimension).

- The optimal probability distribution minimizing the rhs is  $p_i = \sigma_i / \mathfrak{S}$ .
- In this case,  $\mathcal{S}$  is a  $\epsilon$ -coreset with high probability if:

$$m \geq \mathcal{O} \left( \frac{d' \mathfrak{S}^2}{\epsilon^2} \right).$$

- In the  $k$ -means setting,  $\mathfrak{S} = \mathcal{O}(k)$ ,  $d' = kd \log k$ , yielding  $m \geq \mathcal{O} \left( \frac{dk^3 \log k}{\epsilon^2} \right)$ .

---

<sup>1</sup>Langberg and Schulman, *Universal  $\epsilon$ -approximators for integrals*, SIAM, 2010

## In practice?

- In practice, computing the sensitivity is either i/ impossible, or ii/ costs more than solving the initial problem on the whole data set.

---

<sup>1</sup>see Feldman and Langberg, *A unified framework for . . .*, ACM symp. on Theory of computing, 2011  
or Bachem et al., *Practical Coreset Constructions for Machine Learning*, Arxiv, 2017

<sup>2</sup>Braverman et al., *New frameworks for offline and streaming coreset constructions*, Arxiv, 2016

## In practice?

- In practice, computing the sensitivity is either i/ impossible, or ii/ costs more than solving the initial problem on the whole data set.
- To circumvent this problem, upper bounds (easier to estimate) are used<sup>1</sup>.
- Even then, finding algorithms that estimate useful upper bounds faster than the time needed to solve the problem on the whole dataset, remains a challenge.
- *N.B.* Those are not the current best sensitivity-based iid theorems<sup>2</sup>

---

<sup>1</sup>see Feldman and Langberg, *A unified framework for . . .*, ACM symp. on Theory of computing, 2011  
or Bachem et al., *Practical Coreset Constructions for Machine Learning*, Arxiv, 2017

<sup>2</sup>Braverman et al., *New frameworks for offline and streaming coreset constructions*, Arxiv, 2016

## Illustration and context

### Coresets: definition and iid theorem

- Coresets

- Sensitivities

- A classical iid coreset result

### DPPs for Coresets

- Determinantal Point Processes

- A theoretical point-of-view

- A practical point-of-view

## Conclusion

## DPPs in a nutshell

Determinantal point processes (or DPP) are:

- *random* processes that *induce diversity*.
- *tractable*.
- used for three main purposes:
  - i/ *produce diverse samples* of a large database
  - ii/ *use as a tool* in a variety of SP/ML contexts
  - iii/ *characterize* various observed phenomena.

i/ This sample diversity can be directly useful<sup>12</sup>:

summary generation:



<sup>1</sup>left figure: from Kulesza and Taskar, *DPPs for machine learning*, Found. and Trends in ML, 2013

<sup>2</sup>right figure: from G. Gautier's slides [guilgautier.github.io/pdfs/GaBaVa17\\_slides.pdf](https://github.com/guilgautier/pdfs/GaBaVa17_slides.pdf)

i/ This sample diversity can be directly useful<sup>12</sup>:

summary generation:



search engines / recommendation:



<sup>1</sup>left figure: from Kulesza and Taskar, *DPPs for machine learning*, Found. and Trends in ML, 2013

<sup>2</sup>right figure: from G. Gautier's slides [guilgautier.github.io/pdfs/GaBaVa17\\_slides.pdf](http://guilgautier.github.io/pdfs/GaBaVa17_slides.pdf)



i/ This sample diversity can be directly useful<sup>12</sup>:

summary generation:



search engines / recommendation:



ii/ DPP samples can also be used as a tool in several SP/ML contexts:

- Monte Carlo integration
- Feature selection problems
- Coresets!
- etc.

<sup>1</sup>left figure: from Kulesza and Taskar, *DPPs for machine learning*, Found. and Trends in ML, 2013

<sup>2</sup>right figure: from G. Gautier's slides [guilgautier.github.io/pdfs/GaBaVa17\\_slides.pdf](http://guilgautier.github.io/pdfs/GaBaVa17_slides.pdf)

## Determinantal Point Processes: formal definition and notations

- Let  $\mathbf{L} \in \mathbb{R}^{n \times n}$  be a positive semi-definite matrix, where  $L_{ij}$  encodes some kind of interaction between  $x_i$  and  $x_j$ ; e.g., the Gaussian kernel  $L_{ij} = \exp^{-\frac{\|x_i - x_j\|^2}{2\tau^2}}$ .

## Determinantal Point Processes: formal definition and notations

- Let  $\mathbf{L} \in \mathbb{R}^{n \times n}$  be a positive semi-definite matrix, where  $L_{ij}$  encodes some kind of interaction between  $x_i$  and  $x_j$ ; e.g., the Gaussian kernel  $L_{ij} = \exp\left(-\frac{\|x_i - x_j\|^2}{2\tau^2}\right)$ .
- Let  $m$  be a fixed integer and  $\mathcal{S}$  a random subset of  $\mathcal{X}$  of size  $m$ .
- We say that  $\mathcal{S}$  is distributed according to a  $m$ -DPP, and write  $\mathcal{S} \sim mDPP(\mathbf{L})$ , if:

$$\mathbb{P}(\mathcal{S} = S) = \begin{cases} 0 & \text{if } |S| \neq m \\ \frac{1}{Z} \det(L_S) & \text{if } |S| = m \end{cases}$$

where  $Z$  is a normalization constant.

## Determinantal Point Processes: formal definition and notations

- Let  $\mathbf{L} \in \mathbb{R}^{n \times n}$  be a positive semi-definite matrix, where  $L_{ij}$  encodes some kind of interaction between  $x_i$  and  $x_j$ ; e.g., the Gaussian kernel  $L_{ij} = \exp\left(-\frac{\|x_i - x_j\|^2}{2\tau^2}\right)$ .
- Let  $m$  be a fixed integer and  $\mathcal{S}$  a random subset of  $\mathcal{X}$  of size  $m$ .
- We say that  $\mathcal{S}$  is distributed according to a  $m$ -DPP, and write  $\mathcal{S} \sim mDPP(\mathbf{L})$ , if:

$$\mathbb{P}(\mathcal{S} = S) = \begin{cases} 0 & \text{if } |S| \neq m \\ \frac{1}{Z} \det(L_S) & \text{if } |S| = m \end{cases}$$

where  $Z$  is a normalization constant.

- Denote by  $\pi_i$  the inclusion probability of  $x_i$ :

$$\pi_i = \mathbb{P}(x_i \in \mathcal{S}) = \frac{1}{Z} \sum_{S \text{ s.t. } i \in S, |S|=m} \det(L_S).$$

- By construction,  $\sum_i \pi_i = m$ .

## Questions

Consider the same class of minimization problems as previously. Say  $\mathcal{S} \sim mDPP(\mathbf{L})$ .

On the theoretical side (forgetting numerical efficiency for now):

- Under what conditions on  $\mathbf{L}$  is  $\mathcal{S}$  an  $\epsilon$ -coreset with high probability?
- Can we *do better* than the iid case? (as strong coresets, but with smaller  $m$ ?)
- What is the *optimal*  $\mathbf{L}$ ?

## Questions

Consider the same class of minimization problems as previously. Say  $\mathcal{S} \sim mDPP(\mathbf{L})$ .

On the theoretical side (forgetting numerical efficiency for now):

- Under what conditions on  $\mathbf{L}$  is  $\mathcal{S}$  an  $\epsilon$ -coreset with high probability?
- Can we *do better* than the iid case? (as strong coresets, but with smaller  $m$ ?)
- What is the *optimal*  $\mathbf{L}$ ?

On the practical side (back IRL where we look for a practical implementation):

- In the iid world, the (sub-optimal but more realistic) strategies based on upper-bounding the sensitivity have a cost linear in  $n$ . For instance in the  $k$ -means context, they have a cost in  $\mathcal{O}(nkd)$ .
- Can we design a coreset algorithm based on DPPs that
  - outperforms in practice its iid counterpart
  - is not *ridiculously* longer than its iid counterpart?

## In theory: two theorems

**A first theorem** (#9 in the paper) states that if  $\pi_i = m \frac{\sigma_i}{\sigma}$  then we recover the iid performance. Frustratingly, this thm

- i/ only proves that DPPs do not fare worse than iid
- ii/ only provides conditions on  $\{\pi_i\}$ , nothing on higher order marginals (concentration tools well adapted to this case are hard to come by)

## In theory: two theorems

**A first theorem** (#9 in the paper) states that if  $\pi_i = m \frac{\sigma_i}{\sigma}$  then we recover the iid performance. Frustratingly, this thm

- i/ only proves that DPPs do not fare worse than iid
- ii/ only provides conditions on  $\{\pi_i\}$ , nothing on higher order marginals (concentration tools well adapted to this case are hard to come by)

**A second theorem** (#14 in the paper)

- Recall the ideal iid probability distribution:  $p_i = \sigma_i / \mathfrak{S}$ .
- Consider a PSD matrix  $L$  verifying:
  1.  $L$  is projective of rank  $m$ :  $L = UU^t$  with  $U \in \mathbb{R}^{n \times m}$  and  $U^tU = I_m$ .
  2.  $\forall i, L_{ii} = mp_i$ .



## In theory: two theorems

**A first theorem** (#9 in the paper) states that if  $\pi_i = m \frac{\sigma_i}{\sigma}$  then we recover the iid performance. Frustratingly, this thm

- i/ only proves that DPPs do not fare worse than iid
- ii/ only provides conditions on  $\{\pi_i\}$ , nothing on higher order marginals (concentration tools well adapted to this case are hard to come by)

**A second theorem** (#14 in the paper)

- Recall the ideal iid probability distribution:  $p_i = \sigma_i / \mathfrak{S}$ .
- Consider a PSD matrix  $L$  verifying:
  1.  $L$  is projective of rank  $m$ :  $L = UU^t$  with  $U \in \mathbb{R}^{n \times m}$  and  $U^tU = I_m$ .
  2.  $\forall i, L_{ii} = mp_i$ .
- **Lemma [via Schur-Horn]** *Such a matrix necessarily exists. In general, there are many degrees of freedom left to define  $U$ .*

## In theory: two theorems

**A first theorem** (#9 in the paper) states that if  $\pi_i = m \frac{\sigma_i}{\sigma}$  then we recover the iid performance. Frustratingly, this thm

- i/ only proves that DPPs do not fare worse than iid
- ii/ only provides conditions on  $\{\pi_i\}$ , nothing on higher order marginals (concentration tools well adapted to this case are hard to come by)

**A second theorem** (#14 in the paper)

- Recall the ideal iid probability distribution:  $p_i = \sigma_i / \mathfrak{S}$ .
- Consider a PSD matrix  $L$  verifying:
  1.  $L$  is projective of rank  $m$ :  $L = UU^t$  with  $U \in \mathbb{R}^{n \times m}$  and  $U^tU = I_m$ .
  2.  $\forall i, L_{ii} = mp_i$ .
- **Lemma [via Schur-Horn]** *Such a matrix necessarily exists. In general, there are many degrees of freedom left to define  $U$ .*
- **Theorem [Variance reduction theorem]** *Sample  $S_{iid}$  by drawing  $m$  samples iid from  $p$ . Sample  $S_{dpp} \sim mDPP(L)$ . One has:*

$$\forall \theta \in \Theta \quad \text{Var} \left[ \hat{L}(S_{dpp}, \theta) \right] = \text{Var} \left[ \hat{L}(S_{iid}, \theta) \right] - \frac{m-1}{m} Y$$

where  $Y \geq 0$  depends on intricate frame properties of the lines of  $U$ .

- ⇒ Such a DPP necessarily provides a better coreset than its iid counterpart.
- ⇒ Finding the *best* projective DPP is however out of (our) reach theoretically.

## In practice: heuristics

Sampling from a DPP requires a worst-case  $\mathcal{O}(n^3)$  number of operations. Low-rank DPPs have a more reasonable  $\mathcal{O}(nm^2)$  complexity. We propose two DPP heuristics based on low-rank kernels:

### Alg. 1: Approximate Gaussian kernel (with parameter $\tau > 0$ )

- Compute  $r \geq \mathcal{O}(m)$  Random Fourier Features<sup>1</sup> and obtain  $\Psi \in \mathbb{R}^{n \times r}$  s.t.  $\Psi\Psi^t \in \mathbb{R}^{n \times n}$  approximates the Gaussian kernel
- Sample an  $m$ -DPP from  $\mathbf{L} = \Psi\Psi^t$
- ✓ This runs in  $\mathcal{O}(nm^2 + nmd)$
- ✗  $\tau$  is a (annoying) hyper-parameter.

### Alg. 2: Vandermonde kernel (here for $d = 1$ , can be extended to $d \geq 2$ )

- Compute  $V \in \mathbb{R}^{n \times m}$  the partial Vandermonde matrix  $V_{ij} = x_i^{j-1}$
- Sample an  $m$ -DPP from  $\mathbf{L} = VV^t$  (it is a projective DPP)
- ✓ This runs in  $\mathcal{O}(nm^2)$
- ✓ No hyper-parameter tuning
- ✗ For  $d \geq 2$ , not all values of  $m$  are admissible.

---

<sup>1</sup>Rahimi et al., *Random features for large-scale kernel machines*, NIPS, 2008

<sup>3</sup>Barthelmé, Tremblay, Usevich, Amblard. *Determinantal point processes in the flat limit*, Bernoulli, 2022

## In practice: heuristics

Sampling from a DPP requires a worst-case  $\mathcal{O}(n^3)$  number of operations. Low-rank DPPs have a more reasonable  $\mathcal{O}(nm^2)$  complexity. We propose two DPP heuristics based on low-rank kernels:

### Alg. 1: Approximate Gaussian kernel (with parameter $\tau > 0$ )

- Compute  $r \geq \mathcal{O}(m)$  Random Fourier Features<sup>1</sup> and obtain  $\Psi \in \mathbb{R}^{n \times r}$  s.t.  $\Psi\Psi^t \in \mathbb{R}^{n \times n}$  approximates the Gaussian kernel
- Sample an  $m$ -DPP from  $\mathbf{L} = \Psi\Psi^t$
- ✓ This runs in  $\mathcal{O}(nm^2 + nmd)$
- ×  $\tau$  is a (annoying) hyper-parameter.

### Alg. 2: Vandermonde kernel (here for $d = 1$ , can be extended to $d \geq 2$ )

- Compute  $V \in \mathbb{R}^{n \times m}$  the partial Vandermonde matrix  $V_{ij} = x_i^{j-1}$
  - Sample an  $m$ -DPP from  $\mathbf{L} = VV^t$  (it is a projective DPP)
  - ✓ This runs in  $\mathcal{O}(nm^2)$
  - ✓ No hyper-parameter tuning
  - × For  $d \geq 2$ , not all values of  $m$  are admissible.
- Advertisement: both kernels become equivalent in the flat limit ( $\tau \rightarrow \infty$ )<sup>3</sup>.

---

<sup>1</sup>Rahimi et al., *Random features for large-scale kernel machines*, NIPS, 2008

<sup>3</sup>Barthelmé, Tremblay, Usevich, Amblard. *Determinantal point processes in the flat limit*, Bernoulli, 2022

## In practice: illustration on 1-means

- Data  $\mathcal{X}$ , parameter  $\theta$
- Cost func.

$$L(\mathcal{X}, \theta) = \sum_{i=1}^n \|\mathbf{x}_i - \theta\|^2$$



## In practice: illustration on 1-means

- Data  $\mathcal{X}$ , parameter  $\theta$
- Cost func.

$$L(\mathcal{X}, \theta) = \sum_{i=1}^n \|\mathbf{x}_i - \theta\|^2$$



Compare:

- uniform iid sampling
- sensitivity iid: ideal iid sampling based on exact sensitivities
- $m$ -DPP (heuristic) based on RFFs of the Gaussian  $L$ -ensemble
- Proj Poly DPP (heuristic) based on the partial Vandermonde matrix

## In practice: illustration on 1-means

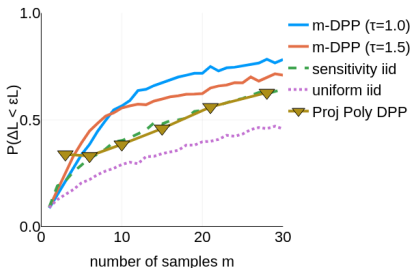
- Data  $\mathcal{X}$ , parameter  $\theta$
- Cost func.

$$L(\mathcal{X}, \theta) = \sum_{i=1}^n \|\mathbf{x}_i - \theta\|^2$$



Compare:

- uniform iid sampling
- sensitivity iid: ideal iid sampling based on exact sensitivities
- $m$ -DPP (heuristic) based on RFFs of the Gaussian  $L$ -ensemble
- Proj Poly DPP (heuristic) based on the partial Vandermonde matrix



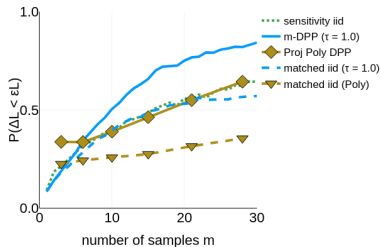
## In practice: illustration on 1-means

- sensitivity iid: as before
- $m$ -DPP ( $\tau = 1$ ): as before
- Proj Poly DPP: as before
- matched iid ( $\tau = 1$ ): iid version of  $m$ -DPP ( $\tau = 1$ )
- matched iid (Poly): iid version of Proj Poly DPP



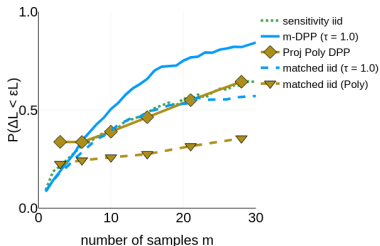
## In practice: illustration on 1-means

- sensitivity iid: as before
- $m$ -DPP ( $\tau = 1$ ): as before
- Proj Poly DPP: as before
- matched iid ( $\tau = 1$ ): iid version of  $m$ -DPP ( $\tau = 1$ )
- matched iid (Poly): iid version of Proj Poly DPP



## In practice: illustration on 1-means

- sensitivity iid: as before
- $m$ -DPP ( $\tau = 1$ ): as before
- Proj Poly DPP: as before
- matched iid ( $\tau = 1$ ): iid version of  $m$ -DPP ( $\tau = 1$ )
- matched iid (Poly): iid version of Proj Poly DPP



As the dimension  $d$  increases:

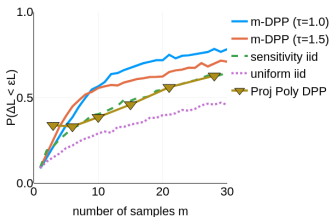


Figure:  $d = 2$

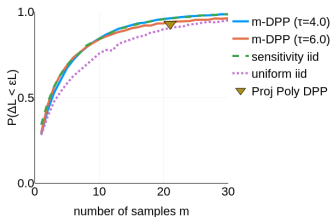


Figure:  $d = 20$

## Conclusion: take home messages

- This is exploratory work on the simple question: can DPPs help create better coresets? If so, how?
- We have a few (mainly frustrating) theorems stating that DPPs do not fare worse than its iid counterpart. The strongest result is the variance reduction theorem.
- We propose 2 DPP-based heuristics (= no provable guarantees), running in  $\mathcal{O}(nm^2)$
- In the  $k$ -means and linear regression problems, these heuristics outperform (= better coresets for a similar computation time) the iid scheme especially:
  - for small  $m$  : to keep the  $\mathcal{O}(nm^2)$  DPP sampling cost under control
  - and small  $d$  : to keep the DPP's repulsiveness significant.
- For (many) more theoretical and experimental details, the paper is available at: <http://jmlr.org/papers/volume20/18-167/18-167.pdf>
- The DPP4Coresets Julia toolbox is available at: <https://gricad-gitlab.univ-grenoble-alpes.fr/tremblan/dpp4coresets.jl> or on my website.