



**HAL**  
open science

# A Faster Sampler for Discrete Determinantal Point Processes

Simon Barthelme, Nicolas Tremblay, Pierre-Olivier Amblard

► **To cite this version:**

Simon Barthelme, Nicolas Tremblay, Pierre-Olivier Amblard. A Faster Sampler for Discrete Determinantal Point Processes. AISTATS 2023 - 26th International Conference on Artificial Intelligence and Statistics, Apr 2023, Valencia, Spain. pp.5582-5592. hal-03835312v2

**HAL Id: hal-03835312**

**<https://hal.science/hal-03835312v2>**

Submitted on 22 Feb 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A Faster Sampler for Discrete Determinantal Point Processes

Simon Barthelmé, Nicolas Tremblay, and Pierre-Olivier Amblard

**Abstract.** Discrete Determinantal Point Processes (DPPs) have a wide array of potential applications for subsampling datasets. They are however held back in some cases by the high cost of sampling. In the worst-case scenario, the sampling cost scales as  $\mathcal{O}(n^3)$  where  $n$  is the number of elements of the ground set. A popular workaround to this prohibitive cost is to sample DPPs defined by low-rank kernels. In such cases, the cost of standard sampling algorithms scales as  $\mathcal{O}(np^2 + nm^2)$  where  $m$  is the (average) number of samples of the DPP (usually  $m \ll n$ ) and  $p$  the rank of the kernel used to define the DPP ( $m \leq p \leq n$ ). The first term,  $\mathcal{O}(np^2)$ , comes from a SVD-like step. We focus here on the second term of this cost,  $\mathcal{O}(nm^2)$ , and show that it can be brought down to  $\mathcal{O}(nm + m^3 \log m)$  without loss on the sampling's exactness. In practice, we observe very substantial speedups compared to the classical algorithm as soon as  $n > 1,000$ . The algorithm described here is a close variant of the standard algorithm for sampling continuous DPPs, and uses rejection sampling. In the specific case of projection DPPs, we also show that any additional sample can be drawn in time  $\mathcal{O}(m^3 \log m)$ . Finally, an interesting by-product of the analysis is that a realisation from a DPP is typically contained in a subset of size  $\mathcal{O}(m \log m)$  formed using leverage score i.i.d. sampling.

Discrete Determinantal Point Processes have been advocated as a way of subsampling large datasets, because they produce samples that preserve some of the diversity of the original dataset [1]. One impediment to their broad adoption in practice lies in their computational cost; aside from some special cases (*e.g.*, random spanning forests [2]), exact sampling of a DPP with a large number of elements is rather expensive.

In this manuscript, we show that a simple modification of the standard algorithm yields a substantial improvement, without loss on the algorithm's exactness. The modification we suggest consists in using a form of rejection sampling. The idea is not very original, and indeed appears in works by [3] for continuous DPPs, and more recently by [4] in the context of experimental design. What is striking is how effective this modification can be, when sampling discrete DPPs, especially given how easy it is to implement.

Here and throughout, let  $n$  designate the size of the ground set the DPP draws from, and  $m$  be the (average) size of the subsample produced by the DPP (usually  $m \ll n$ ). In the worst-case, the cost of producing a sample may be as high as  $\mathcal{O}(n^3)$ , as it requires a full diagonalisation of the kernel [5]. However, in the more realistic context of low-rank kernels and using standard exact sampling algorithms, this figure drops to  $\mathcal{O}(np^2 + nm^2)$  where  $p$  is the rank of the kernel ( $m \leq p \leq n$ ) [6]. We improve this to  $\mathcal{O}(np^2 + nm + m^3 \log m)$ . We readily see that, even though this is an improvement for sampling any DPP, the closer is  $p$  to  $m$ , the more substantial the improvement in practice, as  $np^2$  stays the headline complexity. We identify three popular and general use-cases for which rejection sampling brings substantial speed-ups compared to the classical algorithm:

- 1) (*very significant speed-up:  $p = m$  and orthogonalisation is already computed*) sample a DPP with kernel  $\mathbf{K} = \mathbf{Q}\mathbf{Q}^\top$ , where  $\mathbf{Q} \in \mathbb{R}^{n \times m}$  is orthonormal ( $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$ ) and given (for instance, the DPPs used by [7]). As there is no orthogonalisation to compute, the total cost using rejection

---

All three authors are with CNRS, Univ Grenoble-Alpes, Gipsa-lab, France.

sampling is  $\mathcal{O}(nm + m^3 \log m)$ , substantially faster than the usual cost in  $\mathcal{O}(nm^2)$ .

- 2) (*significant speed-up:  $p = m$  and orthogonalisation has yet to be computed*) in some cases, the orthonormal basis  $\mathbf{Q}$  is not known from the start. A popular context is when one wishes to sample a fixed-size L-ensemble of size  $m$  with  $\mathbf{L} = \mathbf{V}\mathbf{V}^\top$  with  $\mathbf{V} \in \mathbb{R}^{n \times m}$  a matrix of features. In this case, one i/ first computes an orthonormal basis  $\mathbf{Q}$  of the span of  $\mathbf{V}$ , before ii/ sampling a DPP with kernel  $\mathbf{K} = \mathbf{Q}\mathbf{Q}^\top$  (as in the previous case). Step i/ involves, e.g., a QR decomposition. Even though the cost of QR scales theoretically as  $\mathcal{O}(nm^2)$ , it is highly efficient (and parallelisable) in modern hardware such that the bottleneck in previous state-of-the-art is step ii/. Our improvement of step ii/ thus also has practical (possibly very large) speed-ups in this context. In addition, there are special cases of feature matrices  $\mathbf{V}$  for which computing an orthogonal basis has cost less than  $\mathcal{O}(nm^2)$ ; increasing further the potential benefits of our approach. This is the case, e.g., for some classes of sparse  $\mathbf{V}$  [8].
- 3) (*moderate speed-up:  $p$  equals a few times  $m$* ) Same context as 2/ but in the case where the feature matrix  $\mathbf{V}$  is of size  $n \times p$  with  $p > m$ . In this case, step i/ is to compute the SVD of  $\mathbf{V}$ . If  $p$  is too large, this will be the dominant step and our improvement over step ii/ will be negligible. However, in popular cases where  $p$  is only a few times  $m$ , the speed-up is appreciable (see Section 3 for details).

Moreover, in contexts where one needs several realisations of the same DPP, step i/ is computed once, and step ii/ as many times as the number of samples needed; such that our improvement over step ii/ becomes that much more useful.

*Organisation of the paper.* Section 1 briefly introduces the main objects and the state-of-the-art. Section 2 describes our algorithm and its runtime, and Section 3 presents empirical results. A corollary of our result states that DPPs are typically contained in a i.i.d. sample of size  $\mathcal{O}(m \log m)$ . Section 4 discusses this fact and offers concluding remarks.

## 1. Background and notation

### 1.1. Discrete DPPs

For more background on discrete DPPs, we refer readers to [1] and [9]. Discrete DPPs are a specific instance of a *discrete point process*. We say  $\mathcal{X}$  is a discrete point process on a ground set  $\Omega$ , if it is a random subset of  $\Omega$ . Without loss of generality, we let  $\Omega = \{1, \dots, n\}$  so that  $\mathcal{X}$  is a random subset of indices. Also, for two matrices  $\mathbf{A}$  and  $\mathbf{B}$  of same size, the notation  $\mathbf{A} \preceq \mathbf{B}$  means that  $\mathbf{B} - \mathbf{A}$  is positive semidefinite.

**Definition 1.1** (DPP).  $\mathcal{X}$  is a DPP on  $\Omega$  with marginal kernel  $\mathbf{K} \in \mathbb{R}^{n \times n}$  such that  $\mathbf{0} \preceq \mathbf{K} \preceq \mathbf{I}$ , noted  $\mathcal{X} \sim \text{DPP}(\mathbf{K})$ , if for all fixed subsets  $S \subseteq \Omega$ , we have

$$(1) \quad p(S \subseteq \mathcal{X}) = \det \mathbf{K}_S$$

Here  $\mathbf{K}_S$  is the principal submatrix of  $\mathbf{K}$  with indices given by  $S$ . We shall use ‘‘Matlab’’ notation, where  $\mathbf{K}_{A,B}$  denotes the submatrix with row indices  $A$  and column indices  $B$ ,  $\mathbf{K}_{A,:}$  means all columns and  $\mathbf{K}_{:,B}$  all rows.

**Definition 1.2** (Projection DPP). A projection DPP is a DPP whose kernel is a projection matrix (ie.  $\mathbf{K}^2 = \mathbf{K}$ ).

If  $\mathbf{K}$  is a projection matrix, it can be written as  $\mathbf{Q}\mathbf{Q}^\top$  where  $\mathbf{Q} \in \mathbb{R}^{n \times m}$  is an orthonormal basis for span  $\mathbf{K}$  ( $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$  and span  $\mathbf{Q} = \text{span } \mathbf{K}$ ). Note that any orthonormal basis for  $\mathbf{K}$  is enough,  $\mathbf{Q}$  need not be a basis of eigenvectors.

Projection DPPs are important because of the following mixture decomposition, due to [5].

**Theorem 1.3** (mixture representation). *Let  $\mathcal{X} \sim \text{DPP}(\mathbf{K})$ , and  $\mathbf{K} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$  the eigendecomposition of  $\mathbf{K}$ , with  $\mathbf{\Lambda}$  the diagonal matrix of eigenvalues  $\{\lambda_j\}_{j=1,\dots,n}$  and  $\mathbf{U} = (\mathbf{u}_1|\mathbf{u}_2|\dots|\mathbf{u}_n)$  the matrix of eigenvectors. Then the following process produces a sample from  $\mathcal{X}$ :*

- 1) *Sample a subset  $\mathcal{Y}$  of eigenvectors by including each eigenvector  $\mathbf{u}_j$  with probability  $\lambda_j$*
- 2) *Form the projection kernel  $\mathbf{K}_{\mathcal{Y}} = \mathbf{U}_{:, \mathcal{Y}}(\mathbf{U}_{:, \mathcal{Y}})^\top$*
- 3) *Sample  $\mathcal{X} \sim \text{DPP}(\mathbf{K}_{\mathcal{Y}})$*

The cost of sampling a DPP when following this recipe equals the cost of computing the eigendecomposition of  $\mathbf{K}$  ( $\mathcal{O}(np^2)$  with  $p$  the rank of  $\mathbf{K}$ ) followed by the cost of sampling a projection DPP (step 3). It is the latter step that we focus on here.

## 1.2. Fixed-size DPPs

The cardinal of a DPP is in general random. Such varying-sized samples are not practical in many applications, which led [10] to define fixed-size DPPs<sup>1</sup>

**Definition 1.4** (Fixed-size DPP). *A fixed size DPP of size  $m$  is a DPP  $\mathcal{X}$  conditioned on  $|\mathcal{X}| = m$ .*

To sample a fixed-size DPP with kernel  $\mathbf{K}$ , one follows the same recipe as in Theorem 1.3 except for the first step that is replaced by:

- 1) *Sample a subset  $\mathcal{Y}$  of eigenvectors by including each eigenvector  $\mathbf{u}_j$  with probability  $\lambda_j$ ; conditioned on  $|\mathcal{Y}| = m$ .*

Performing such a conditioned sampling can be done by Algorithm 8 of [1], which works by explicitly computing elementary polynomials. If  $n$  and/or  $m$  are too large, numerical instabilities usually arise with this method, and [11] propose a way to stabilize this conditioned sampling.

## 1.3. L-ensembles and fixed-size L-ensembles

L-ensembles are a subclass of DPPs popular in machine learning applications because of their intuitive definition:

**Definition 1.5** (L-ensemble). *Let  $\mathbf{L}$  be a positive semi-definite matrix.  $\mathcal{X}$  is a L-ensemble on  $\Omega$  if for all  $X \subseteq \Omega$*

$$(2) \quad p(\mathcal{X} = X) = \frac{1}{\det(\mathbf{I} + \mathbf{L})} \det(\mathbf{L}_X)$$

L-ensembles are specified via their likelihood function, Eq. (2), which states that those subsets of  $\Omega$  where the submatrix  $\mathbf{L}_X$  is well-conditioned, are preferred. Intuitively, if  $L_{i,j}$  represents a similarity between items  $i$  and  $j$  of  $\Omega$ , then the L-ensemble favours subsets of  $\Omega$  that hold dissimilar items.

As with DPPs, one defines fixed-size L-ensembles as:

**Definition 1.6** (Fixed-size L-ensemble). *A fixed size L-ensemble of size  $m$  is a L-ensemble  $\mathcal{X}$  conditioned on  $|\mathcal{X}| = m$ .*

<sup>1</sup>They are often called k-DPPs in the literature, but we prefer “fixed-size DPPs” in order not to overload the symbol  $k$  too much.

(Fixed-size) L-ensembles are (fixed-size) DPPs with kernel  $\mathbf{K} = (\mathbf{I} + \mathbf{L})^{-1} \mathbf{L}$  [1, 9], and the mixture representation thus applies.

To conclude this section, we have seen that all (fixed-size) L-ensembles and more generally all (fixed-size) DPPs have a mixture representation that divides the sampling algorithm in two steps: i/ a diagonalisation step that costs  $\mathcal{O}(np^2)$ , ii/ a step consisting of sampling a projection DPP. Step ii/ is known<sup>2</sup> to cost  $\mathcal{O}(nm^2)$  in the literature. The purpose of this paper is to show that the cost of this second step can always (and easily) be reduced to  $\mathcal{O}(nm + m^3 \log m)$ .

#### 1.4. State-of-the-art

Various directions have been explored when designing fast samplers for discrete DPPs. Some have focused on bypassing the eigendecomposition of  $\mathbf{K}$  or  $\mathbf{L}$  [12–14]. If  $\mathbf{K}$  is a sparse matrix, then the algorithms in [12] can be quite advantageous compared to standard  $\mathcal{O}(n^3)$  algorithms. These algorithms are difficult to adapt to L-ensembles in an efficient manner (for instance, they cannot take advantage of sparsity in  $\mathbf{L}$ ). Random spanning forests [2] are an example of a discrete DPP where a fast sampler (Wilson’s algorithm, [15]) is available, and sparsity in  $\mathbf{L}$  seems to play a role. However, Wilson’s algorithm does not extend readily to L-ensembles with arbitrary structures.

Another direction for generic DPP samplers is to give up on exactness. Approximate samplers are available, based on Markov Chain Monte Carlo methods. The most recent results in that direction are in [16], where the authors show that given some preprocessing there are MCMC samplers that run in  $\tilde{\mathcal{O}}(m^\omega)$ , where the  $\tilde{\mathcal{O}}$  is shorthand for  $\mathcal{O}(\cdot)$  “up to logarithmic factors”, and  $\omega$  is the exponent of matrix-multiplication time, which for practical values of  $m$  is effectively 3. The preprocessing consists essentially in estimating the inclusion probabilities, also known as the leverage scores, and its runtime is given by [16] as  $\tilde{\mathcal{O}}(nm^{\omega-1})$ . Our results are essentially the same (preprocessing in  $\mathcal{O}(nm^2)$ , sampling in  $\tilde{\mathcal{O}}(m^3)$ ), but our sampler is exact.

Also, two papers [17, 18] extend the tree-based approach of [19] to obtain both approximate and exact samplers with complexity slightly larger than our proposal. They are also more complicated to implement. However, they can handle non-symmetric DPPs, which the algorithm given here cannot do.

Finally, the use of rejection sampling is not new in the context of DPP sampling, since algorithms for sampling continuous DPPs use this strategy out of necessity [3]. More recently, [4] describe a similar algorithm to ours, in the context of volume sampling for experimental design. Our contribution compared to [4] is to i/ lay out a more refined analysis: better bound on the total number of proposals in Theorem 2.1, novel investigation in the theoretical implications of this theorem in Section 4, ii/ keep an eye on practical implementations: see Sections 2.3 and 3. What we would like to stress is that rejection sampling leads to an algorithm that is much faster in practice, but no more complicated to implement, than the traditional discrete sampler.

## 2. Sampling via accept-reject

In this section the goal is to formulate and analyse an algorithm for sampling a projection DPP  $\mathcal{X} \sim DPP(\mathbf{Q}\mathbf{Q}^\top)$  with  $\mathbf{Q} \in \mathbb{R}^{n \times m}$ , verifying  $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$ . The first exact such algorithm was described by [5], and adapted in [1] to the discrete

---

<sup>2</sup>This is an average (resp. deterministic) cost for DPPs (resp. fixed-size DPPs) for which  $m$  refers to the average (resp. desired) size of the sample.

case. The first efficient version appeared in [6]. It is effectively a variant of the Gram-Schmidt algorithm.

For completeness we give an easy derivation of this classical algorithm in the next section (Section 2.1), and the notation will serve to describe our own variant, in Section 2.2.

## 2.1. State-of-the-art algorithm

A projection DPP has size  $m$  almost surely [1]. We shall sample the  $m$  elements successively. Let  $\boldsymbol{\mathcal{X}} = (x_1, \dots, x_m)$  be an ordered version of  $\mathcal{X}$ ; we can go from  $\boldsymbol{\mathcal{X}}$  to  $\mathcal{X}$  by forgetting the order and from  $\mathcal{X}$  to  $\boldsymbol{\mathcal{X}}$  by ordering randomly. The latter can be achieved by picking a first item uniformly from  $\mathcal{X}$ , then a second, then a third etc. The sampling algorithm proceeds via the following decomposition:

$$(3) \quad p(\boldsymbol{\mathcal{X}}) = p(x_1)p(x_2|x_1)p(x_3|x_1, x_2) \dots p(x_m|x_1 \dots x_{m-1})$$

The algorithm samples  $x_1$  first, then  $x_2$  given  $x_1$  has been selected, etc. The law of  $x_1$  is the law of the first element of  $\boldsymbol{\mathcal{X}}$ , a randomly ordered version of  $\mathcal{X}$ . That is equivalent to  $x_1$  being sampled uniformly at random from  $\mathcal{X}$ , and so:

$$p(x_1 = i) = \frac{1}{m}p(i \in \mathcal{X}) = \frac{K_{i,i}}{m}$$

We can similarly obtain the law of  $x_2$  given  $x_1$ , as two elements drawn randomly (without replacement) from  $\mathcal{X}$ :

$$\begin{aligned} p(x_2 = i|x_1 = j) &= \frac{p(x_1 = j, x_2 = i)}{p(x_1 = j)} \\ &= \frac{p(i \in \mathcal{X}, j \in \mathcal{X})}{m(m-1)} \times \frac{m}{p(j \in \mathcal{X})} \\ &= \frac{1}{(m-1)} \frac{\det \mathbf{K}_{\{i,j\}}}{K_{j,j}} \end{aligned}$$

The formula for determinants of block matrices yields:

$$\frac{\det \mathbf{K}_{\{i,j\}}}{K_{j,j}} = K_{i,i} - \frac{K_{i,j}^2}{K_{j,j}}$$

and so:

$$p(x_2 = i|x_1 = j) = \frac{1}{(m-1)} \left( K_{i,i} - \frac{K_{i,j}^2}{K_{j,j}} \right)$$

For the general term in the chain rule decomposition (Eq. (3)), the same reasoning applies. We obtain:

$$(4) \quad \begin{aligned} p(x_t = i | \boldsymbol{\mathcal{X}}_{1:(t-1)} = X_t) \\ = \frac{1}{(m-t)} \left( K_{i,i} - \mathbf{K}_{i, X_t} (\mathbf{K}_{X_t})^{-1} \mathbf{K}_{X_t, i} \right) \end{aligned}$$

Eq. (4) is enough to give us a sampling algorithm, since at each step we have an explicit (discrete) probability distribution to sample from. However, implementing Eq. (4) naïvely, we would be computing a matrix inverse at each step, which would turn out to be quite expensive for large  $m$ . To get an efficient algorithm a bit more work is needed.

Let us reexpress Eq. (4) in terms of  $\mathbf{Q}$ . Recalling  $\mathbf{K} = \mathbf{Q}\mathbf{Q}^\top$ , we obtain:

$$(5) \quad \begin{aligned} p(x_t = i | \boldsymbol{\mathcal{X}}_{1:(t-1)} = X_t) \\ = \frac{1}{(m-t)} \left( K_{ii} - \mathbf{Q}_{i,:} \mathbf{M}_t(\mathbf{Q}_{i,:})^\top \right) \end{aligned}$$

where  $\mathbf{M}_t = (\mathbf{Q}_{X_t,:})^\top (\mathbf{Q}_{X_t,:}(\mathbf{Q}_{X_t,:})^\top)^{-1} \mathbf{Q}_{X_t,:}$  is a projection matrix ( $\mathbf{M}_t^2 = \mathbf{M}_t$ ) of size  $m$  and rank  $|X_t| = t - 1$ , and so can be rewritten

$$\mathbf{M}_t = \sum_{i=1}^{t-1} \mathbf{s}_i \mathbf{s}_i^\top$$

where  $\mathbf{s}_1 \dots \mathbf{s}_t$  form an orthonormal basis for  $\text{span } \mathbf{M}_t = \text{span } \mathbf{Q}_{X_t,:}^\top$ , the linear subspace spanned by the rows of  $\mathbf{Q}$  selected so far. A first source of computational savings comes from realising that  $\mathbf{M}_t$  can be computed iteratively via the Gram-Schmidt process. Notice that  $\mathbf{M}_t = \mathbf{M}_{t-1} + \mathbf{s}_{t-1} \mathbf{s}_{t-1}^\top$ , and  $\mathbf{M}_{t-1}$  spans  $\text{span } \mathbf{Q}_{X_{t-1},:}^\top$ . We obtain  $\mathbf{s}_{t-1}$  via Gram-Schmidt: first we compute the residual

$$\mathbf{z}_{t-1} = (\mathbf{I} - \mathbf{M}_{t-1}) \mathbf{Q}_{x_{t-1},:}^\top$$

and then we normalise:

$$\mathbf{s}_{t-1} = \frac{\mathbf{z}_{t-1}}{\|\mathbf{z}_{t-1}\|}$$

At each step  $t$  this costs  $\mathcal{O}(m(t-1))$  operations, and we will need to do this  $m-1$  times at a total cost of  $\mathcal{O}(m^3)$ .

Next, we can show that the probability distribution we sample from at step  $t$  can be easily obtained from the one we had at step  $t-1$ . It is more convenient to write this using unnormalised versions of the densities. Let  $\pi^{(1)}(i) = mp(x_1 = i) = K_{i,i}$ . Next, we define:

$$\begin{aligned} \pi^{(2)}(i) &= (m-1)p(x_2 = i | x_1 = j) \\ &= K_{i,i} - \frac{K_{i,j}^2}{K_{i,i}} = \pi^{(1)}(i) - \frac{K_{i,j}^2}{K_{i,i}} \end{aligned}$$

Note that we have suppressed the dependency on the past in the notation  $\pi^{(2)}(i)$ : it is to be understood as the (unnormalised) density we draw from at the second step of the algorithm. In the general case, we define:

$$(6) \quad \pi^{(t)}(i) = (m-t+1) p(x_t = i | \mathcal{X}_{1:(t-1)} = X_t)$$

Injecting Eq. (4) and Eq. (5), we find

$$\begin{aligned} \pi^{(t)}(i) &= \pi^{(1)}(i) - \mathbf{Q}_{i,:} \mathbf{M}_t (\mathbf{Q}_{i,:})^\top \\ &= \pi^{(1)}(i) - \mathbf{Q}_{i,:} (\mathbf{M}_{t-1} + \mathbf{s}_{t-1} \mathbf{s}_{t-1}^\top) (\mathbf{Q}_{i,:})^\top \\ (7) \quad &= \pi^{(t-1)}(i) - (\mathbf{Q}_{i,:} \mathbf{s}_{t-1})^2 \end{aligned}$$

All we need to do at each step of the algorithm is to

- 1) pick an item  $i$  according to  $\pi^{(t)}$
- 2) perform a step of the Gram-Schmidt algorithm to update  $\mathbf{M}_t$  based on the new vector  $\mathbf{Q}_{i,:}$
- 3) Update the probability distribution to  $\pi^{(t+1)}$  according to Eq. (7)

Sampling from a discrete distribution (step 1 above) can be done at cost  $\mathcal{O}(n)$ , and is needed  $m$  times, for a total cost of  $\mathcal{O}(nm)$ . We have already established that the cost of the Gram-Schmidt algorithm is  $\mathcal{O}(m^3)$ . It is the update to the probability distribution that is the most costly, with each step costing  $\mathcal{O}(nm)$  ( $n$  dot products in  $\mathbb{R}^m$ ) for a total of  $\mathcal{O}(nm^2)$ . Since  $n > m$  the cost of the algorithm scales as  $\mathcal{O}(nm^2)$ . We show pseudo-code for this standard algorithm as Alg. 1.

In the next section, we move on to the core of our contribution, showing that this  $\mathcal{O}(nm^2)$  cost can be reduced to  $\mathcal{O}(nm + m^3 \log m)$  via rejection sampling.

---

**Algorithm 1:** Sampling from a projection DPP  $\mathcal{X} \sim DPP(\mathbf{K} = \mathbf{Q}\mathbf{Q}^\top)$ , standard algorithm

---

Initialise  $\pi(x) \leftarrow \sum_{j=1}^m Q_{x,j}^2$ ,  $\mathcal{X} = \emptyset$ , Gram-Schmidt basis  $\mathbf{S} = []$  ;  
**foreach**  $t \in 1 \dots m$  **do**  
    Sample  $x$  from  $\frac{\pi(x)}{m-(t-1)}$ , add to  $\mathcal{X}$  ;  
    Compute residual  $\mathbf{z}_t = (\mathbf{I} - \mathbf{S}\mathbf{S}^\top)(\mathbf{Q}_{x,:})^\top$  ;  
    Add column  $\mathbf{s}_t = \frac{\mathbf{z}_t}{\|\mathbf{z}_t\|}$  to  $\mathbf{S}$  ;  
    Compute  $\mathbf{v} = \mathbf{Q}\mathbf{s}_t \in \mathbb{R}^n$  ;  
    Update probabilities  $\pi(x) \leftarrow \pi(x) - (v_x)^2$  ;  
**end**

---

## 2.2. Using rejection sampling

As discussed above, the most expensive part of alg. 1, lies in updating the probability distribution to sample from (last step of the *for* loop). It turns out that the accept-reject method lets us bypass this step.

To briefly recall the rejection sampling idea, suppose we have an unnormalised density  $\pi(x)$  (the target) we wish to draw from, and a proposal  $q(x)$ , also unnormalised, but that we know how to draw from, and is not too far from  $\pi$ . Further,  $q(x)$  has support at least as wide as  $\pi(x)$ , and upper bounds it ( $\pi(x) \leq q(x)$  over the support). We may then draw  $x$  from  $q(x)$  and accept the sample with probability  $\frac{\pi(x)}{q(x)}$ . The accepted samples then have density  $\pi(x)$ . If  $q(x)$  is a good bound for  $\pi(x)$ , then the rejection sampler will be quite efficient. In the limit where  $\pi = q$ , the acceptance probability goes to 1. If on the other hand  $q(x)$  is quite loose, the acceptance probability may be bad.

In the DPP sampling algorithm, we need to sample from  $\pi^{(1)}$ , then  $\pi^{(2)}$ , etc. up to  $\pi^{(m)}$ . Our proposal is to compute  $\pi^{(1)}$  exactly for all  $n$  entries, then use  $\pi^{(1)}$  as proposal distribution for the rest of the sequence.

Recall that  $\pi^{(t)}$  is the unnormalised density defined in Eq. (6). From the recursion in Eq. (7), we have that

$$\pi^{(t)}(x) \geq \pi^{(t+1)}(x)$$

for all  $x$  and  $1 \leq t \leq m-1$ . This is true in particular for  $\pi^{(1)}$ , which can therefore be used as a proposal distribution for all subsequent  $\pi^{(t)}$ . We can directly compute the probability of accepting a proposed sample. At step  $t$ , we sample  $x$  from the normalised density  $\frac{1}{m}\pi^{(1)}(x)$  and accept it with probability  $\frac{\pi^{(t)}(x)}{\pi^{(1)}(x)}$ . The acceptance probability equals:

$$(8) \quad \rho_t = \sum_{i=1}^n \frac{\pi^{(t)}(i)}{\pi^{(1)}(i)} \frac{\pi^{(1)}(i)}{m} = \frac{1}{m} \sum_{i=1}^n \pi^{(t)}(i) = \frac{(m-t+1)}{m}$$

This probability decreases at each step of the algorithm, but at the final step it is still positive and equals  $\frac{1}{m}$ .

Let us outline the proposed algorithm. First, one computes every entry of  $\pi^{(1)}$ . For this we use the following formula

$$(9) \quad \pi^{(1)}(i) = K_{i,i} = \sum_{j=1}^n Q_{ij}^2$$

The computation is equivalent to computing the norm of each row of  $\mathbf{Q}$ , at cost  $\mathcal{O}(nm)$ . Because we need to sample from  $\pi^{(1)}$  repeatedly, it pays to use Walker's



alias method [20] (see also Chapter III.4 of [21]). Given a preprocessing cost of  $\mathcal{O}(n)$ , the alias method gives us all subsequent samples at cost  $\mathcal{O}(1)$  instead of  $\mathcal{O}(n)$ .

At the first step we sample our first item from  $\pi^{(1)}$ . At step 2, and all subsequent steps, we use rejection sampling, which involves computing the ratio

$$(10) \quad \frac{\pi^{(t)}(x)}{\pi^{(1)}(x)} = 1 - \frac{1}{\pi^{(1)}(x)} \sum_{j=1}^{t-1} (\mathbf{Q}_{x,:} \cdot \mathbf{s}_j)^2$$

by Eq. (7). Computing the acceptance ratio has cost  $\mathcal{O}(m(t-1))$  at step  $t$ , the cost of  $t-1$  dot products in  $\mathbb{R}^m$ . We do this repeatedly until a proposal is accepted, at which point we need to perform a Gram-Schmidt step to update  $\mathbf{M}_t$  to  $\mathbf{M}_{t+1}$ . We then move on to the next iteration, or stop if  $t = m$ .

We summarise the whole process as Alg. 2. To recapitulate the different computational costs:

- Preprocessing cost: computing  $\pi^{(1)}$  for all entries comes at cost  $\mathcal{O}(nm)$  and setting up Walker's alias method at cost  $\mathcal{O}(n)$
- The Gram-Schmidt process (computing  $\mathbf{z}_t$  then  $\mathbf{s}_t$ ) costs  $\mathcal{O}(tm)$  at step  $t$ . Summing this figure for  $t = 1$  to  $m$  gives a cost of  $\mathcal{O}(m^3)$
- We now need to compute the average cost of the *while* loop. At step  $t$ , the rejection sampler has probability  $\rho_t$  of succeeding, given by Eq. (8). The number of proposals  $R_t$  that are required until acceptance is thus a random variable that follows a geometric distribution with success probability  $\rho_t$ . One thus has:

$$(11) \quad \mathbb{E}(R_t) = \frac{1}{\rho_t} = \frac{m}{m-t+1}.$$

Since computing the acceptance ratio costs  $\mathcal{O}(mt)$  for each trial, the expected cost of the *while* loop at step  $t$  scales as  $\mathcal{O}\left(\frac{m^2 t}{m-t+1}\right)$ . Summing this figure over  $t$ :

$$\sum_{t=1}^m \frac{m^2 t}{m-t+1} \leq m^3 \sum_{t=1}^m \frac{1}{m-t+1}$$

yields a total expected cost scaling as <sup>3</sup>  $\mathcal{O}(m^3 \log m)$ .

Tallying everything we obtain the following theorem.

**Theorem 2.1.** *Alg. 2 samples a projection DPP, with an expected runtime scaling as  $\mathcal{O}(nm + m^3 \log m)$ . Also, any additional sample from the same DPP can be obtained in an extra  $\mathcal{O}(m^3 \log m)$  expected runtime.*

*Moreover, these expected runtimes are representative. Indeed,  $\forall \delta \in (0, \frac{1}{2})$ , the total number of proposals  $R = \sum_{t=1}^m R_t$  satisfies, with probability greater than  $1 - \delta$ :*

$$R \leq 2m \log m + 3m \log \frac{1}{\delta}$$

*Proof.* The fact that Alg. 2 samples a projection DPP in  $\mathcal{O}(nm + m^3 \log m)$  expected runtime is proven above the Theorem's statement. The fact that any additional sample from the same DPP only costs an extra  $\mathcal{O}(m^3 \log m)$  expected runtime comes from the observation that all initialisation steps (the computation of  $\pi^{(1)}$  and the setting-up cost of Walkers' algorithm) have already been computed for the first sample. To obtain any extra sample, one only needs to run the *for* loop once more, costing  $\mathcal{O}(m^3 \log m)$ .

The final statement relates to concentration properties of  $R$ . For  $m = 1$ , it is trivial

---

<sup>3</sup> $\sum_{t=1}^m \frac{1}{m-t+1} = \sum_{t=1}^m \frac{1}{t}$  scales as  $\mathcal{O}(\log m)$ : see, e.g., Chapter 6 of [22]

---

**Algorithm 2:** Sampling from a projection DPP  $\mathcal{X} \sim DPP(\mathbf{K} = \mathbf{Q}\mathbf{Q}^\top)$  with rejection sampling

---

Initialise  $\pi^{(1)}(x) \leftarrow \sum_{j=1}^m Q_{x,j}^2$ ,  $\mathcal{X} \leftarrow \emptyset$ , Gram-Schmidt basis  $\mathbf{S} \leftarrow []$  ;  
 Initialise alias table for Walker's alias algorithm to sample from  $\pi^{(1)}$ .  
**foreach**  $t \in 1 \dots m$  **do**  
   accept  $\leftarrow$  false ;  
   **while** not accept **do**  
     Draw  $x$  from  $\pi^{(1)}$  using the alias method;  
     Compute acceptance ratio  $r = 1 - \frac{1}{\pi^{(1)}(x)} \sum_{j=1}^{t-1} (\mathbf{Q}_{x,:} \cdot \mathbf{s}_j)^2$  ;  
     **if**  $\text{rand}() < r$  **then**  
       | accept  $\leftarrow$  true  
     **end**  
   **end**  
   Add  $x$  to  $\mathcal{X}$  ;  
   Compute residual  $\mathbf{z}_t = (\mathbf{I} - \mathbf{S}\mathbf{S}^\top)(\mathbf{Q}_{x,:})^\top$  ;  
   Add column  $\mathbf{s}_t = \frac{\mathbf{z}_t}{\|\mathbf{z}_t\|}$  to  $\mathbf{S}$  ;  
**end**

---

(log refers to the natural logarithm in the result) as  $R = 1$  with probability 1 in this case. We now show the result for  $m \geq 2$ . At step  $t$  of Alg. 2 the number of proposals  $R_t$  is a random geometric variable with parameter  $(m - t + 1)/m$ . Note that all the  $R_t$ 's are independent. We study here the behavior of  $R = \sum_{t=1}^m R_t$ , the total number of rejection sampling steps in the whole course of Alg. 2. The following one-tailed upper bound is drawn from Thm 2.3 in [23]:

$$\forall \lambda \geq 1, \quad p(R \geq \lambda E[R]) \leq \frac{1}{\lambda} \left(1 - \frac{1}{m}\right)^{(\lambda-1-\log \lambda)E[R]}$$

Let  $\delta \in (0, \frac{1}{2})$ . We look for  $\lambda$  large enough such that  $p(R \geq \lambda E[R]) \leq \delta$ , that is:

$$-\log \lambda + (\lambda - 1 - \log \lambda)E[R] \log \left(1 - \frac{1}{m}\right) \leq \log \delta$$

As  $\lambda \geq 1$ ,  $-\log \lambda \leq 0$ . It thus suffices to seek  $\lambda$  verifying:

$$(12) \quad (\lambda - 1 - \log \lambda)E[R] \log \left(1 - \frac{1}{m}\right) \leq \log \delta$$

Note that  $E[R] \log \left(1 - \frac{1}{m}\right)$  is negative so we seek a lower bound of  $\lambda - 1 - \log \lambda$ . One has<sup>4</sup>:

$$\forall \lambda \geq 1, \quad \frac{1}{3}\lambda - \log \frac{3}{2} \leq \lambda - 1 - \log \lambda$$

such that Eq. (12) is verified provided that:

$$\lambda \geq 3 \left( \log \frac{3}{2} - \frac{\log \delta^{-1}}{E[R] \log \left(1 - \frac{1}{m}\right)} \right).$$

---

<sup>4</sup>In fact, the function  $\lambda - 1 - \log \lambda$  is convex for all  $\lambda \geq 1$  and thus lower-bounded by all its tangents. The one we use is the tangent in  $3/2$ . Other choices lead to other constants in the result.

Stated differently, setting  $\lambda$  to this lower bound implies  $p(R \leq \lambda E[R]) \geq 1 - \delta$ . All is left to show is that  $\forall m \geq 2$ :

$$(13) \quad 3 \left( \log \frac{3}{2} - \frac{\log \delta^{-1}}{E[R] \log \left(1 - \frac{1}{m}\right)} \right) E[R] \leq 2m \log m + 3m \log \delta^{-1}.$$

For this, we use two upper bounds. The first one is

$$\forall m > 1, \quad -\frac{1}{\log \left(1 - \frac{1}{m}\right)} \leq m - 1/2.$$

The second one is the following bound on  $E[R]$ . As  $E[R_t] = m/(m-t+1)$ , one has (see, *e.g.*, Chapter 6 of [22]):

$$\frac{E[R]}{m} = \sum_{t=1}^m \frac{1}{m-t+1} = \sum_{t=1}^m \frac{1}{t} = \gamma + \psi(m) + \frac{1}{m}$$

where  $\gamma \approx 0.577$  is Euler's constant and  $\psi(m)$  the digamma function. Now, a known bound on  $\psi(m)$  is  $\psi(m) \leq \log m - \frac{1}{2m}$ , which gives  $E[R] \leq m(\log m + \gamma) + \frac{1}{2}$ . Applying these two upper bounds yields Eq. (13).  $\square$

### 2.3. Some refinements

Alg. 2 works well enough as is but there are some refinements that can further reduce the computational cost.

#### 2.3.1. Preprocessing for general DPPs

In some applications we require several samples from the same DPP, and algorithms have been described that trade higher set-up cost for a lower cost per sample (see, *e.g.*, [19]). If the target DPP is a projection DPP, then setting up Alg. 2 for repeated sampling could not be easier, as stated in Thm 2.1: computing  $\pi^{(1)}$  and setting up the alias table is part of the preprocessing, so the first sample from the DPP costs  $\mathcal{O}(nm + m^3 \log m)$  but after that the cost is just  $\mathcal{O}(m^3 \log m)$  per sample.

If the target DPP is not a projection DPP, then one has to use the mixture representation (Thm 1.3): draw a random set of eigenvectors  $\mathcal{Y}$  and run Alg. 2 with  $\mathbf{Q} = \mathbf{U}_{\mathcal{Y}, \cdot}$ . Since the kernel changes every time, so does  $\pi^{(1)}$ , and it cannot be computed as part of pre-processing. However, the kernels encountered in practice tend to have rapidly decreasing eigenvalues, so that the variance in  $\mathcal{Y}$  is quite small and the DPP is close to a projection DPP. Without getting into too much detail, it is possible to pre-compute the partial sum

$$\widehat{\pi^{(1)}}(i) = \sum_{j \in \hat{\mathcal{Y}}} U_{i,j}^2$$

for some highly likely subset of  $\mathcal{Y}$ , denoted here  $\hat{\mathcal{Y}}$ .  $\pi^{(1)}$  for the actual sampled  $\mathcal{Y}$  can be obtained efficiently by removing the extra entries and adding the missing ones. The alias table can be computed from scratch. This type of preprocessing brings down the cost to  $\mathcal{O}(ns + m^3 \log m)$  per sample, where  $m = E(|\mathcal{Y}|)$  and  $s$  is the expected size of the symmetric difference between  $\hat{\mathcal{Y}}$  and  $\mathcal{Y}$ .

#### 2.3.2. Caching computations and updating the proposal distribution

Clearly, Alg. 2 has some wasted computation, since all the computations done when a proposal  $x$  is rejected are performed again should  $x$  come up a second time. When  $n$  is small, or when  $\pi^{(1)}$  has low entropy, this may indeed happen several times. Caching is one way of reducing the amount of redundant computations that are

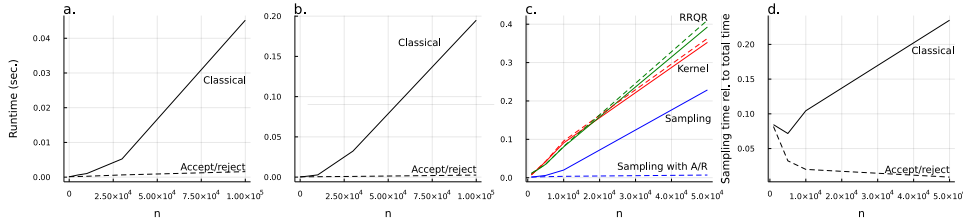


Figure 1. Left panels: median time needed to sample a projection DPP, using the standard approach (Alg. 1) vs. A/R (Alg. 2). For two values of  $m$ : a)  $m = 30$ , b)  $m = 60$ . Note that here the time taken to compute an orthogonal basis is not taken into account (see text). Right panels: simulation of a full workload that includes feature generation and orthogonalisation (see text). c) time taken by each step in the computation as a function of  $n$ . Solid lines are runtimes when sampling using the classical method, dashed, when using A/R. Note that the first two steps (labeled “kernel” and “RRQR”, in red and green respectively) are identical. A/R sampling becomes beneficial at around  $n = 1,000$  and is orders of magnitude faster at  $n = 100,000$ . d) Sampling time as a function of total computation time.

performed. Going back to the recursive formula for  $\pi^{(t)}$  (Eq. (7)), we see that it is cheaper to compute  $\pi^{(t)}$  from  $\pi^{(t-1)}$  than it is to compute it from scratch. A reasonable caching strategy is then to keep track for every point  $x$  of the last density evaluation performed for that point. If  $x$  comes up again, the evaluation of the acceptance ratio is simplified.

Another natural idea is to update the proposal distribution over the course of the algorithm (instead of sticking with  $\pi^{(1)}$  throughout). The most basic version is that any  $x$  that has already been selected has an acceptance probability of 0, so we may as well not suggest them. Another is that points similar to a selected point are quite unlikely to come up further down, and so it may be worth computing  $\pi^{(t)}$  for these neighbours to tighten the bound. Finally, we may combine this idea with the caching idea, which provides a better bound for every point that has ever been suggested. Unfortunately this runs against the difficulty of updating the alias table in Walker’s algorithm, which one would have to compute from scratch at every update (at cost  $\mathcal{O}(n)$ ). A better way would be to use a binary tree representation [21], which can be updated at cost  $\mathcal{O}(\log n)$  and provides samples also at cost  $\mathcal{O}(\log n)$ . The implementation complexity increases a lot however, and we have not pursued this further. As we shall see below, Alg. 2 is quite fast in practice, and implementation time may be better invested in feature computation and orthogonalisation.

### 3. Empirical results

We compare the Accept/Reject algorithm (Alg. 2) to its classical counterpart (Alg. 1) for different values of  $n$  and  $m$ . Both algorithms are implemented in the Julia language and are publicly available<sup>5</sup>. For each value of  $n$ , we sample a random projection matrix of size  $n \times m$  (via QR decomposition of a matrix with Gaussian entries). We then pre-compute the leverage scores, and run each algorithm 100 times. Fig. 1 a) and b) show the measured median runtimes. All tests are run on a 2017 Linux laptop with i7-8550U Intel CPU and 8 Go of RAM.

<sup>5</sup>we’ve added a folder in our DPP.jl repository containing the code necessary to reproduce the figures, available here: [https://github.com/dahtah/DPP.jl/tree/main/misc/sampling\\_paper](https://github.com/dahtah/DPP.jl/tree/main/misc/sampling_paper). In addition, the A/R sampler is available as part of the DPP.jl software package, <https://github.com/dahtah/DPP.jl/>.

For very small values of  $n$ , the classical algorithm is faster, which can be explained by the efficiency of BLAS calls. At each step the whole conditional distribution is computed, the main cost being a matrix multiplication (i.e., a BLAS call), which benefits from efficient multithreaded code. However, the different asymptotic scalings ( $\mathcal{O}(nm^2)$  vs.  $\mathcal{O}(nm)$ ) soon makes the classical algorithm uncompetitive. The cross-over point in our simulations is at around  $n = 1,000$ , and by  $n = 10,000$  the difference is stark. This illustrates the “very significant speed-up” scenario described in the introduction:  $p = m$  and orthogonalisation is already computed. We now move on to illustrate a less favorable case related to the “moderate speed-up” scenario of the introduction:  $p$  equals a few times  $m$ .

In many cases, sampling a projection DPP is only one of the steps in a process that involves feature computation and orthogonalisation. For instance, in [24], a DPP based on the Gaussian kernel is used to produce a subset of the data suitable for running clustering algorithms. Starting from  $n$  points,  $\mathbf{x}_1, \dots, \mathbf{x}_n$  in  $\mathbb{R}^d$  they use a DPP with L-ensemble given by  $L_{ij} = \gamma\kappa(\mathbf{x}, \mathbf{y}) = \gamma \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{y}\|^2\right)$ , where  $\gamma$  is a tuning parameter that determines the expected size. Producing a sample from this exact DPP requires the eigendecomposition of  $\mathbf{L}$  which is impractical; however,  $\mathbf{L}$  is numerically low-rank for relevant values of  $\sigma$  and this can be exploited. In [24] a low-rank approximation of  $\mathbf{L}$  is used, based on random Fourier features [25] followed by a SVD, bringing down the total cost to  $\mathcal{O}(nm^2)$ . We now sketch (without any formal justification) another procedure which gives comparable results at lower cost.

Gaussian kernel matrices have rapidly decaying spectra (see, e.g., [26]), which implies in particular that DPPs sampled from a Gaussian L-ensemble are well approximated by projection DPPs with kernels  $\mathbf{P}_m = \mathbf{U}_m \mathbf{U}_m^\top$  where  $\mathbf{P}_m$  projects onto the dominant eigenspace of  $\mathbf{L}$  of order  $m$ . Thus, all we need is a good basis for the dominant eigenspace. Methods from randomised linear algebra offer good practical tools (“range finders”) to obtain a basis for such a space [27]. For these simulations, we used the following approximation:

- 1) Select (and compute)  $5m$  columns uniformly from  $\mathbf{L}$ . Call this matrix  $\mathbf{A}$ . We call this the “kernel step”.
- 2) Use Rank-Revealing QR (RRQR, [28]) and random projections, as implemented in the Julia package LowRankApprox.jl, to produce  $\mathbf{Q}$ , an orthonormal matrix of size  $n \times m$  that approximates the image of  $\mathbf{A}$ . We call this the “RRQR step”.
- 3) Sample a DPP with projection kernel  $\mathbf{Q}\mathbf{Q}^\top$  using either the classical or the A/R algorithm.

We set  $m = 100$  and time each step. This results in a total runtime of around 1.2 sec. at  $n = 10^5$  with the A/R sampler, which challenges the notion that DPPs are very slow to sample from. With this procedure, the time spent sampling the actual DPP goes up to 20% of total time for the classical algorithm at  $n = 10^5$ , but using the A/R sampler sampling time becomes negligible. See Fig. 1 c) and d) to see how these times vary with  $n$ . This indicates that for some computations the implementation effort may be better allocated to speeding up the linear algebra and feature computation part rather than the sampling part. In this particular instance, step (1) at least could be sped up by exploiting parallelism, or the GPU, which we did not attempt.

#### 4. Discussion and perspectives

On top of the improvement on the sampling time of DPPs, our results imply the following intriguing by-product. Let  $\mathcal{X}$  be a projection DPP of size  $m$ . A set of  $\mathcal{O}(m \log m)$  points sampled i.i.d. from the inclusion probability distribution

$p(i \in \mathcal{X}) = \pi^{(1)}(i)/m$  (also known as *leverage scores*) contains with high probability a realisation from the DPP. The consequences of this fact are worth discussing. First, let us put the result a bit more formally.

**Definition 4.1.** *Let  $\mathcal{X}$  be a DPP on  $\Omega$  and  $\mathcal{Y} \subseteq \Omega$ . We call  $\Phi(\mathcal{Y})$  a thinning algorithm if it returns a subset of  $\mathcal{Y}$ . Moreover, we say  $\Phi(\mathcal{Y})$  is successful when it returns a realisation from  $\mathcal{X}$ .*

**Corollary 4.2.** *Let  $\mathcal{X}$  be a projection DPP of size  $m$ , and  $\mathcal{Y}$  be a set of i.i.d. points sampled with replacement with probability proportional to the leverage scores:  $p(i \in \mathcal{X}) = \pi^{(1)}(i)/m$ . Let  $\delta \in (0, 1/2)$ . A simple modification of Alg. 2 gives a thinning algorithm  $\Phi$  that verifies:  $\Phi(\mathcal{Y})$  is successful with probability greater than  $1 - \delta$  provided that  $|\mathcal{Y}| \geq 2m \log m + 3m \log(\frac{1}{\delta})$ .*

*Proof.* Let  $\mathcal{Y}$  be drawn i.i.d. with replacement from the leverage score distribution  $\pi^{(1)}$ .  $\Phi$  is the following simple modification of Alg. 2. Instead of drawing a new proposal  $x$  using the alias method at the beginning of the *while* loop as in Alg. 2, draw uniformly and without replacement from  $\mathcal{Y}$ . If  $\Phi$  finishes before emptying  $\mathcal{Y}$ , then it is successful. If  $\mathcal{Y}$  is empty and  $\Phi$  is not terminated, then it fails. The probability that  $\Phi$  succeeds is thus equal to the probability that  $|\mathcal{Y}|$  is larger than the number of proposals  $R$  of Alg. (2). By Thm 2.1, setting  $|\mathcal{Y}| = 2m \log m + 3m \log(\frac{1}{\delta})$  yields  $p(|\mathcal{Y}| \geq R) \leq 1 - \delta$  and ends the proof.  $\square$

Note that this is a substantial improvement over the work of [14], which gives this result only for  $|\mathcal{Y}| \geq \mathcal{O}(m^2)$  i.i.d. points.

A natural question is to ask if this result is optimal: can we find a thinning algorithm that succeeds with high probability for even smaller i.i.d sets? The answer is no in general:

**Proposition 4.3.** *Corollary 4.2 is optimal in the following sense. Let  $\mathcal{X}$  and  $\mathcal{Y}$  be as previously. There does not exist a generic thinning algorithm able to succeed with fixed non-null probability if  $|\mathcal{Y}| = o(m \log m)$ .*

*Proof.* We show the proposition by exhibiting a type of DPP for which there does not exist a thinning algorithm that succeeds with a non-null probability if  $|\mathcal{Y}|$  is asymptotically smaller than  $m \log m$ .

It is well-known in the folklore that a form of *stratified sampling* is a special case of projection DPPs. In stratified sampling, we partition the ground set  $\Omega$  into  $m$  classes, and sample an item uniformly from each segment of the partition. To simplify the argument, assume  $\Omega$  can be cut into  $m$  subsets of equal size, and define vector  $e_j$  as the (normalised) indicator of segment  $j$ , i.e.  $e_j(i) = \sqrt{\frac{1}{m}}$  if item  $i$  is in segment  $j$  and 0 otherwise. Let  $\mathbf{E} = [e_1 \dots e_m]$ . Then it is easy to show that stratified sampling is equivalent to a DPP  $\mathcal{X}$  with marginal kernel  $\mathbf{K} = \mathbf{E}\mathbf{E}^\top$ . Since  $\mathbf{E}^\top \mathbf{E} = \mathbf{I}$ , the DPP in question is a projection DPP.

Because of the nature of stratified sampling, we know that  $\mathcal{X}$  contains a point from each one of the segments, and that  $p(i \in \mathcal{X}) = \frac{m}{n}$  for all  $i$ . Now in order for any thinning algorithm to produce a stratified sample, the i.i.d. sample  $\mathcal{Y}$  needs to contain at least one point from each segment. Let  $l \stackrel{\text{def}}{=} |\mathcal{Y}|$ : how large does  $l$  need to be so that  $\mathcal{Y}$  contains at least one point from each segment with probability at least  $\alpha$  ( $\alpha > 0$  fixed)? This is an instance of the coupon collector's problem. Assume that at each time  $t$  we add a ball to one of  $m$  urns with equal probability, and call  $T$  the smallest  $t$  such all urns have at least one ball. We show that for any  $l(m) = o(m \log m)$ ,  $\lim_{m \rightarrow \infty} p(T < l(m)) = 0$ .

To do this, we need an upper bound for  $p(T < l(m))$ . One could work with results from [29] for instance. However, we prefer an elegant line of proof inspired

by a contribution of a stackexchange user called “cardinal”<sup>6</sup>.  $T$  can be viewed as a sum of  $m$  geometric variables:  $T = \sum_{i=1}^m T_i$ , where  $T_i$  is the time at which  $i$  urns have at least one ball. All the  $T_i$ ’s are independent geometric random variables with success probability  $p_i = 1 - \frac{i-1}{m}$ . Indeed, the same representation is obtained by considering Alg. 2 in the special case of stratified sampling (each iteration fills one urn). Now, Markov’s inequality gives:

$$\forall s > 0, \quad p(T < l) = p(\exp(-sT) > \exp(-sl)) \\ \leq \exp(sl)E(\exp(-sT))$$

Since  $T$  is a sum of independent geometric variables,  $E(\exp(-sT))$  is easy to compute<sup>7</sup>:

$$E(\exp(-sT)) = \prod_{i=1}^m \frac{i}{m(e^s - 1) + i}$$

Picking  $s = \frac{1}{m}$ , we obtain:

$$p(T < l) \leq \exp^{\frac{l}{m}} \prod_{i=1}^m \frac{i}{m(e^{1/m} - 1) + i}$$

Since  $e^{\frac{1}{m}} \geq 1 + \frac{1}{m}$ , we upper bound the right-hand side to:

$$p(T < l) \leq \exp^{\frac{l}{m}} \prod_{i=1}^m \frac{i}{1 + i} = \frac{\exp(\frac{l}{m})}{m + 1}$$

Thus, any choice of sample size  $l(m)$  such that  $\frac{\exp(\frac{l(m)})}{m+1}$  goes to 0 in the limit is asymptotically too small (the probability of success goes to 0). Noting that  $\frac{\exp(\frac{l}{m})}{m+1} = o(1)$  is equivalent to  $l(m) = o(m \log m)$  yields the claim.  $\square$

This transition occurring at  $\mathcal{O}(m \log m)$  calls for discussion, and paves the way to future interesting lines of research. First of all, Corollary 4.2 shows, from an original angle, that the repulsiveness of DPPs is weak. Indeed, other repulsive processes such as hard-core processes cannot verify such property in all generality. For instance, in the high density limit of a hard-sphere model, the probability that the position of  $m$  non-overlapping spheres can be found within a set of only  $\mathcal{O}(m \log m)$  iid points drawn uniformly, tends to 0. In addition, these results ask the following question: in what cases should one pay the extra cost of sampling  $m$  elements from a DPP, rather than simply sampling  $\mathcal{O}(m \log m)$  elements i.i.d. from the leverage score distribution? Of course, when the objective is to sample a diverse set, such as in search engines, it is always worthwhile to sample the DPP. However, in the case of integration [30, 31]; or in the case of coresets [24], the answer is not so clear and requires further investigation.

#### Acknowledgements

We thank the five anonymous reviewers for their helpful comments that led to an improved version of this manuscript. This work was partially supported by the ANR project GRANOLA (ANR-21-CE48-0009), as well as the LabEx PERSYVAL-Lab (ANR-11-LABX-0025-01) and MIAI@Grenoble Alpes (ANR-19-P3IA-0003).

<sup>6</sup>see <https://stats.stackexchange.com/q/7774>

<sup>7</sup>Using  $E(\exp(-sT)) = \prod_{j=1}^m E(\exp(-sT_j)) = \prod_{j=1}^m \frac{p_j \exp^{-s}}{1 - (1-p_j) \exp^{-s}}$  and changing variable  $i \leftarrow m - j + 1$

## . References

- [1] A. Kulesza, B. Taskar *et al.*, “Determinantal point processes for machine learning,” *Foundations and Trends® in Machine Learning*, vol. 5, no. 2–3, pp. 123–286, 2012.
- [2] L. Avena and A. Gaudillièrè, “Two Applications of Random Spanning Forests,” *Journal of Theoretical Probability*, Jul. 2017. [Online]. Available: <http://link.springer.com/10.1007/s10959-017-0771-3>
- [3] F. Lavancier, J. Møller, and E. Rubak, “Determinantal point process models and statistical inference,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 77, no. 4, pp. 853–877, 2015.
- [4] M. Dereziński, K. L. Clarkson, M. W. Mahoney, and M. K. Warmuth, “Minimax experimental design: Bridging the gap between statistical and worst-case approaches to least squares regression,” in *Conference on Learning Theory*. PMLR, 2019, pp. 1050–1069.
- [5] J. B. Hough, M. Krishnapur, Y. Peres, and B. Virág, “Determinantal Processes and Independence,” *Probability Surveys*, vol. 3, pp. 206–229, 2006. [Online]. Available: <http://dx.doi.org/10.1214/154957806000000078>
- [6] J. Gillenwater, “Approximate inference for determinantal point processes,” Ph.D. dissertation, University of Pennsylvania, 2014.
- [7] C. Launay, A. Desolneux, and B. Galerne, “Determinantal point processes for image processing,” *SIAM Journal on Imaging Sciences*, vol. 14, no. 1, pp. 304–348, 2021.
- [8] T. A. Davis, “Algorithm 8xx: Suitesparseqr, a multifrontal multithreaded sparse qr factorization package,” *ACM Trans. Math. Software*, 2008.
- [9] N. Tremblay, S. Barthelmé, K. Usevich, and P.-O. Amblard, “Extended L-ensembles: A new representation for determinantal point processes,” *The Annals of Applied Probability*, vol. 33, no. 1, pp. 613 – 640, 2023. [Online]. Available: <https://doi.org/10.1214/22-AAP1824>
- [10] A. Kulesza and B. Taskar, “k-dpps: fixed-size determinantal point processes,” in *Proceedings of the 28th International Conference on International Conference on Machine Learning*, 2011, pp. 1193–1200.
- [11] S. Barthelmé, P.-O. Amblard, and N. Tremblay, “Asymptotic equivalence of fixed-size and varying-size determinantal point processes,” *Bernoulli*, vol. 25, no. 4B, pp. 3555–3589, 2019. [Online]. Available: <https://doi.org/10.3150/18-BEJ1102>
- [12] J. Poulson, “High-performance sampling of generic determinantal point processes,” *Philosophical Transactions of the Royal Society A*, vol. 378, no. 2166, p. 20190059, 2020.
- [13] C. Launay, B. Galerne, and A. Desolneux, “Exact sampling of determinantal point processes without eigendecomposition,” *Journal of Applied Probability*, vol. 57, no. 4, pp. 1198–1221, 2020.
- [14] M. Dereziński, D. Calandriello, and M. Valko, “Exact sampling of determinantal point processes with sublinear time preprocessing,” *Advances in neural information processing systems*, vol. 32, 2019.
- [15] D. B. Wilson, “Generating random spanning trees more quickly than the cover time,” in *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*, 1996, pp. 296–303. [Online]. Available: <http://dl.acm.org/citation.cfm?id=237880>
- [16] N. Anari, Y. P. Liu, and T.-D. Vuong, “Optimal sublinear sampling of spanning trees and determinantal point processes via average-case entropic independence,” in *2022 IEEE 63rd Annual Symposium on Foundations of Computer Science (FOCS)*, 2022, pp. 123–134.
- [17] I. Han, M. Gartrell, E. Dohmatob, and A. Karbasi, “Scalable mcmc sampling for nonsymmetric determinantal point processes,” in *International Conference on Machine Learning*, 2022, pp. 8213–8229.
- [18] I. Han, M. Gartrell, J. Gillenwater, E. Dohmatob, and A. Karbasi, “Scalable sampling for nonsymmetric determinantal point processes,” in *International Conference on Learning Representations*, 2022.
- [19] J. Gillenwater, A. Kulesza, Z. Mariet, and S. Vassilvtiskii, “A tree-based method for fast repeated sampling of determinantal point processes,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 2260–2268.
- [20] A. J. Walker, “An efficient method for generating discrete random variables with general distributions,” *ACM Transactions on Mathematical Software (TOMS)*, vol. 3, no. 3, pp. 253–256, 1977.
- [21] L. Devroye, *Non-Uniform Random Variate Generation*. Springer-Verlag, 1986.
- [22] M. Abramowitz and I. A. Stegun, *Handbook of mathematical functions with formulas, graphs, and mathematical tables*. US Government printing office, 1964, vol. 55.
- [23] S. Janson, “Tail bounds for sums of geometric and exponential variables,” *Statistics & Probability Letters*, vol. 135, pp. 1–6, 2018.



- [24] N. Tremblay, S. Barthelmé, and P.-O. Amblard, “Determinantal Point Processes for Core-sets.” *Journal of Machine Learning Research*, vol. 20, no. 168, pp. 1–70, 2019.
- [25] A. Rahimi and B. Recht, “Random features for large-scale kernel machines,” in *Advances in Neural Information Processing Systems*, vol. 20, 2007. [Online]. Available: <https://proceedings.neurips.cc/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf>
- [26] A. J. Wathen and S. Zhu, “On spectral distribution of kernel matrices related to radial basis functions,” *Numerical Algorithms*, vol. 70, no. 4, pp. 709–726, Dec 2015.
- [27] P.-G. Martinsson and J. A. Tropp, “Randomized numerical linear algebra: Foundations and algorithms,” *Acta Numerica*, vol. 29, pp. 403–572, 2020.
- [28] T. F. Chan, “Rank revealing qr factorizations,” *Linear algebra and its applications*, vol. 88, pp. 67–82, 1987.
- [29] C. Witt, “Fitness levels with tail bounds for the analysis of randomized search heuristics,” *Information Processing Letters*, vol. 114, no. 1-2, pp. 38–41, 2014.
- [30] R. Bardenet and A. Hardy, “Monte Carlo with determinantal point processes,” *The Annals of Applied Probability*, vol. 30, no. 1, pp. 368 – 417, 2020. [Online]. Available: <https://doi.org/10.1214/19-AAP1504>
- [31] J.-F. Coeurjolly, A. Mazoyer, and P.-O. Amblard, “Monte carlo integration of non-differentiable functions on  $[0, 1]^d$ ,  $iota = 1, \dots, d$ , using a single determinantal point pattern defined on  $[0, 1]^d$ ,” *Electronic Journal of Statistics*, vol. 15, no. 2, pp. 6228 – 6280, 2021. [Online]. Available: <https://doi.org/10.1214/21-EJS1929>