



**HAL**  
open science

## A determinantal point process for column subset selection

Ayoub Belhadji, R. Bardenet, Pierre Chainais

► **To cite this version:**

Ayoub Belhadji, R. Bardenet, Pierre Chainais. A determinantal point process for column subset selection. *Journal of Machine Learning Research*, 2020, 21 (197), pp.1-62. hal-01966298

**HAL Id: hal-01966298**

**<https://hal.science/hal-01966298>**

Submitted on 28 Dec 2018

**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 determinantal point process for column subset selection

Ayoub Belhadji<sup>1\*</sup>, Rémi Bardenet<sup>1</sup>, Pierre Chainais<sup>1</sup>

<sup>1</sup> Univ. Lille, CNRS, Centrale Lille, UMR 9189 - CRIStAL, 59651 Villeneuve d'Ascq, France

## Abstract

Dimensionality reduction is a first step of many machine learning pipelines. Two popular approaches are principal component analysis, which projects onto a small number of well chosen but non-interpretable directions, and feature selection, which selects a small number of the original features. Feature selection can be abstracted as a numerical linear algebra problem called the column subset selection problem (CSSP). CSSP corresponds to selecting the best subset of columns of a matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$ , where *best* is often meant in the sense of minimizing the approximation error, i.e., the norm of the residual after projection of  $\mathbf{X}$  onto the space spanned by the selected columns. Such an optimization over subsets of  $\{1, \dots, d\}$  is usually impractical. One workaround that has been vastly explored is to resort to polynomial-cost, random subset selection algorithms that favor small values of this approximation error. We propose such a randomized algorithm, based on sampling from a projection determinantal point process (DPP), a repulsive distribution over a fixed number  $k$  of indices  $\{1, \dots, d\}$  that favors diversity among the selected columns. We give bounds on the ratio of the expected approximation error for this DPP over the optimal error of PCA. These bounds improve over the state-of-the-art bounds of *volume sampling* when some realistic structural assumptions are satisfied for  $\mathbf{X}$ . Numerical experiments suggest that our bounds are tight, and that our algorithms have comparable performance with the *double phase* algorithm, often considered to be the practical state-of-the-art. Column subset selection with DPPs thus inherits the best of both worlds: good empirical performance and tight error bounds.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Notation</b>	<b>4</b>
<b>3</b>	<b>Related Work</b>	<b>5</b>
3.1	Rank revealing QR decompositions . . . . .	5
3.2	Length square importance sampling and additive bounds . . . . .	6
3.3	$k$ -leverage scores sampling and multiplicative bounds . . . . .	6
3.4	The geometric interpretation of the $k$ -leverage scores . . . . .	9
3.5	Negative correlation: volume sampling and the double phase algorithm	9
3.6	Excess risk in sketched linear regression . . . . .	10

---

\*Corresponding author: [ayoub.belhadji@centralelille.fr](mailto:ayoub.belhadji@centralelille.fr)

<b>4</b>	<b>Determinantal Point Processes</b>	<b>12</b>
4.1	Definitions . . . . .	12
4.2	Sampling from a DPP and a $k$ -DPP . . . . .	13
4.3	Motivations for column subset selection using projection DPPs . . . . .	14
<b>5</b>	<b>Main Results</b>	<b>15</b>
5.1	Multiplicative bounds in spectral and Frobenius norm . . . . .	15
5.2	Bounds for the excess risk in sketched linear regression . . . . .	17
<b>6</b>	<b>Numerical experiments</b>	<b>17</b>
6.1	Toy datasets . . . . .	18
6.1.1	Generating toy datasets . . . . .	18
6.1.2	volume sampling vs projection DPP . . . . .	18
6.2	Real datasets . . . . .	22
6.3	Discussion . . . . .	23
<b>7</b>	<b>Conclusion</b>	<b>23</b>
<b>A</b>	<b>Another interpretation of the <math>k</math>-leverage scores</b>	<b>30</b>
<b>B</b>	<b>Majorization and Schur convexity</b>	<b>30</b>
<b>C</b>	<b>Principal angles and the Cosine Sine decomposition</b>	<b>32</b>
C.1	Principal angles . . . . .	32
C.2	The Cosine Sine decomposition . . . . .	32
<b>D</b>	<b>Proofs</b>	<b>34</b>
D.1	Technical lemmas . . . . .	34
D.2	Proof of Proposition 16 . . . . .	35
D.3	Proof of Proposition 17 . . . . .	36
D.3.1	Frobenius norm bound . . . . .	36
D.3.2	Spectral norm bound . . . . .	39
D.4	Proof of Theorem 18 . . . . .	39
D.4.1	Frobenius norm bound . . . . .	40
D.4.2	Spectral norm bound . . . . .	42
D.4.3	Bounding the probability of rejection . . . . .	43
D.5	Proof of Proposition 20 . . . . .	43
<b>E</b>	<b>Generating orthogonal matrices with prescribed leverage scores</b>	<b>43</b>
E.1	Definitions and statement of the problem . . . . .	44
E.2	Related work . . . . .	45
E.3	The restricted Gelfand-Tsetlin polytope . . . . .	46
E.4	Our algorithm . . . . .	49

## 1 Introduction

Datasets come in always larger dimensions, and dimension reduction is thus often one of the first steps in any machine learning pipeline. Two of the most widespread strate-

gies are principal component analysis (PCA) and feature selection. PCA projects the data in directions of large variance, called principal components. While the initial features (the canonical coordinates) generally have a direct interpretation, principal components are linear combinations of these original variables, which makes them hard to interpret. On the contrary, using a selection of original features will preserve interpretability when it is desirable. Once the data are gathered in an  $N \times d$  matrix, of which each row is an observation encoded by  $d$  features, feature selection boils down to selecting columns of  $\mathbf{X}$ . Independently of what comes after feature selection in the machine learning pipeline, a common performance criterion for feature selection is the approximation error in some norm, that is, the norm of the residual after projecting  $\mathbf{X}$  onto the subspace spanned by the selected columns. Optimizing such a criterion over subsets of  $\{1, \dots, d\}$  requires exhaustive enumeration of all possible subsets, which is prohibitive in high dimension. One alternative is to use a polynomial-cost, random subset selection strategy that favors small values of the criterion.

This rationale corresponds to a rich literature on randomized algorithms for column subset selection (Deshpande and Vempala, 2006; Drineas et al., 2008; Boutsidis et al., 2011). A prototypical example corresponds to sampling  $s$  columns of  $\mathbf{X}$  i.i.d. from a multinomial distribution of parameter  $\mathbf{p} \in \mathbb{R}^d$ . This parameter  $\mathbf{p}$  can be the squared norms of each column (Drineas et al., 2004), for instance, or the more subtle  $k$ -leverage scores (Drineas et al., 2008). While the former only takes  $\mathcal{O}(dN^2)$  time to evaluate, it comes with loose guarantees; see Section 3.2. The  $k$ -leverage scores are more expensive to evaluate, since they call for a truncated SVD of order  $k$ , but they come with tight bounds on the ratio of their expected approximation error over that of PCA.

To minimize approximation error, the subspace spanned by the selected columns should be as large as possible. Simultaneously, the number of selected columns should be as small as possible, so that intuitively, diversity among the selected columns is desirable. The column subset selection problem (CSSP) then becomes a question of designing a discrete point process over the column indices  $\{1, \dots, d\}$  that favors diversity in terms of directions covered by the corresponding columns of  $\mathbf{X}$ . Beyond the problem of designing such a point process, guarantees on the resulting approximation error are desirable. Since, given a target dimension  $k \leq d$  after projection, PCA provides the best approximation in Frobenius or spectral norm, it is often used a reference: a good CSS algorithm preserves interpretability of the  $c$  selected features while guaranteeing an approximation error not much worse than that of rank- $k$  PCA, all of this with  $c$  not much larger than  $k$ .

In this paper, we introduce and analyze a new randomized algorithm for selecting  $k$  diverse columns. Diversity is ensured using a determinantal point process (DPP). DPPs can be viewed as the kernel machine of point processes; they were introduced by Macchi (1975) in quantum optics, and their use widely spread after the 2000s in random matrix theory (Johansson, 2005), machine learning (Kulesza et al., 2012), spatial statistics (Lavancier et al., 2015), and Monte Carlo methods (Bardenet and Hardy, 2016), among others. In a sense, the DPP we propose is a nonindependent generalization of the multinomial sampling with  $k$ -leverage scores of (Boutsidis et al., 2009). It further naturally connects to volume sampling, the CSS algorithm that has the best error bounds (Deshpande et al., 2006). We give error bounds for DPP

sampling that exploit sparsity and decay properties of the  $k$ -leverage scores, and outperform volume sampling when these properties hold. Our claim is backed up by experiments on toy and real datasets.

The paper is organized as follows. Section 2 introduces our notation. Section 3 is a survey of column subset selection, up to the state of the art to which we later compare. In Section 4, we discuss determinantal point processes and their connection to volume sampling. Section 5 contains our main results, in the form of both classical bounds on the approximation error and risk bounds when CSS is a prelude to linear regression. In Section 6, we numerically compare CSS algorithms, using in particular a routine that samples random matrices with prescribed  $k$ -leverage scores.

## 2 Notation

We use  $[n]$  to denote the set  $\{1, \dots, n\}$ , and  $[n : m]$  for  $\{n, \dots, m\}$ . We use bold capitals  $\mathbf{A}, \mathbf{X}, \dots$  to denote matrices. For a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and subsets of indices  $I \subset [m]$  and  $J \subset [n]$ , we denote by  $\mathbf{A}_{I,J}$  the submatrix of  $\mathbf{A}$  obtained by keeping only the rows indexed by  $I$  and the columns indexed by  $J$ . When we mean to take all rows or  $\mathbf{A}$ , we write  $\mathbf{A}_{:,J}$ , and similarly for all columns. We write  $\text{rk}(\mathbf{A})$  for the rank of  $\mathbf{A}$ , and  $\sigma_i(\mathbf{A})$ ,  $i = 1, \dots, \text{rk}(\mathbf{A})$  for its singular values, ordered decreasingly. Sometimes, we will need the vectors  $\Sigma(\mathbf{A})$  and  $\Sigma(\mathbf{A})^2$  the vectors of  $\mathbb{R}^d$  with respective entries  $\sigma_i(\mathbf{A})$  and  $\sigma_i^2(\mathbf{A})$ ,  $i = 1, \dots, \text{rk}(\mathbf{A})$ . Similarly, when  $\mathbf{A}$  can be diagonalized,  $\Lambda(\mathbf{A})$  (and  $\Lambda(\mathbf{A})^2$ ) are vectors with the decreasing eigenvalues (squared eigenvalues) of  $\mathbf{A}$  as entries.

The spectral norm of  $\mathbf{A}$  is  $\|\mathbf{A}\|_2 = \sigma_1(\mathbf{A})$ , while the Frobenius norm of  $\mathbf{A}$  is defined by

$$\|\mathbf{A}\|_{\text{Fr}} = \sqrt{\sum_{i=1}^{\text{rk}(\mathbf{A})} \sigma_i(\mathbf{A})^2}.$$

For  $\ell \in \mathbb{N}$ , we need to introduce the  $\ell$ -th elementary symmetric polynomial on  $L \in \mathbb{N}$  variables, that is

$$e_\ell(X_1, \dots, X_L) = \sum_{\substack{T \subset [L] \\ |T| = \ell}} \prod_{j \in T} X_j. \quad (1)$$

Finally, we follow [Ben-Israel \(1992\)](#) and denote spanned volumes by

$$\text{Vol}_q(\mathbf{A}) = \sqrt{e_q(\sigma_1(\mathbf{A})^2, \dots, \sigma_{\text{rk}(\mathbf{A})}(\mathbf{A})^2)}, \quad q = 1, \dots, \text{rk}(\mathbf{A}).$$

Throughout the paper,  $\mathbf{X}$  will always denote an  $N \times d$  matrix that we think of as the original data matrix, of which we want to select  $k \leq d$  columns. Unless otherwise specified,  $r$  is the rank of  $\mathbf{X}$ , and matrices  $\mathbf{U}, \Sigma, \mathbf{V}$  are reserved for the SVD of  $\mathbf{X}$ , that is,

$$\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^T \quad (2)$$

$$= \left[ \mathbf{U}_k \mid \mathbf{U}_{r-k} \right] \left[ \begin{array}{c|c} \Sigma_k & \mathbf{0} \\ \hline \mathbf{0} & \Sigma_{r-k} \end{array} \right] \left[ \begin{array}{c} \mathbf{V}_k^T \\ \hline \mathbf{V}_{r-k}^T \end{array} \right], \quad (3)$$

where  $\mathbf{U} \in \mathbb{R}^{N \times r}$  and  $\mathbf{V} \in \mathbb{R}^{d \times r}$  are orthogonal, and  $\Sigma \in \mathbb{R}^{r \times r}$  is diagonal. The diagonal entries of  $\Sigma$  are denoted by  $\sigma_i = \sigma_i(\mathbf{X})$ ,  $i = 1, \dots, r$ , and we assume they

are in decreasing order. We will also need the blocks given in (3), where we separate blocks of size  $k$  corresponding to the largest  $k$  singular values. To simplify notation, we abusively write  $\mathbf{U}_k$  for  $\mathbf{U}_{:, [k]}$  and  $\mathbf{V}_k$  for  $\mathbf{V}_{:, [k]}$  in (3), among others. Though they will be introduced and discussed at length in Section 3.3, we also recall here that we note  $\ell_i^k = \|\mathbf{V}_{[k], i}\|_2^2$  the so-called  $k$ -leverage score of the  $i$ -th column of  $\mathbf{X}$ .

We need some notation for the selection of columns. Let  $S \subset [d]$  be such that  $|S| = k$ , and let  $\mathbf{S} \in \{0, 1\}^{d \times k}$  be the corresponding sampling matrix:  $\mathbf{S}$  is defined by  $\forall \mathbf{M} \in \mathbb{R}^{N \times d}$ ,  $\mathbf{M}\mathbf{S} = \mathbf{M}_{:, S}$ . In the context of column selection, it is often referred to  $\mathbf{X}\mathbf{S} = \mathbf{X}_{:, S}$  as  $\mathbf{C}$ . We set for convenience  $\mathbf{Y}_{:, S}^\top = (\mathbf{Y}_{:, S})^\top$ .

The result of column subset selection will usually be compared to the result of PCA. We denote by  $\Pi_k \mathbf{X}$  the best rank- $k$  approximation to  $\mathbf{X}$ . The sense of *best* can be understood either in Frobenius or spectral norm, as both give the same result. On the other side, for a given subset  $S \subset [d]$  of size  $|S| = s$  and  $\nu \in \{2, \text{Fr}\}$ , let

$$\Pi_{S, k}^\nu \mathbf{X} = \arg \min_A \|\mathbf{X} - A\|_\nu$$

where the minimum is taken over all matrices  $\mathbf{A} = \mathbf{X}_{:, S} \mathbf{B}$  such that  $\mathbf{B} \in \mathbb{R}^{s \times d}$  and  $\text{rk } \mathbf{B} \leq k$ ; in words, the minimum is taken over matrices of rank at most  $k$  that lie in the column space of  $\mathbf{C} = \mathbf{X}_{:, S}$ . When  $|S| = k$ , we simply write  $\Pi_S^\nu \mathbf{X} = \Pi_{S, k}^\nu \mathbf{X}$ . In practice, the Frobenius projection can be computed as  $\Pi_S^{\text{Fr}} \mathbf{X} = \mathbf{C}\mathbf{C}^+ \mathbf{X}$ , yet there is no simple expression for  $\Pi_S^2 \mathbf{X}$ . However,  $\Pi_S^{\text{Fr}} \mathbf{X}$  can be used as an approximation of  $\Pi_S^2 \mathbf{X}$  since

$$\|\mathbf{X} - \Pi_S^2 \mathbf{X}\|_2 \leq \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_2 \leq \sqrt{2} \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_2, \quad (4)$$

see (Boutsidis et al., 2011, Lemma 2.3).

### 3 Related Work

In this section, we review the main results about column subset selection.

#### 3.1 Rank revealing QR decompositions

The first  $k$ -CSSP algorithm can be traced back to the article of Golub (1965) on pivoted QR factorization. This work introduced the concept of Rank Revealing QR factorization (RRQR). The original motivation was to calculate a well-conditioned QR factorization of a matrix  $\mathbf{X}$  that reveals its numerical rank.

**Definition 1** *Let  $\mathbf{X} \in \mathbb{R}^{N \times d}$  and  $k \in \mathbb{N}$  ( $k \leq d$ ). A RRQR factorization of  $\mathbf{X}$  is a 3-tuple  $(\mathbf{\Pi}, \mathbf{Q}, \mathbf{R})$  with  $\mathbf{\Pi} \in \mathbb{R}^{d \times d}$  a permutation matrix,  $\mathbf{Q} \in \mathbb{R}^{N \times d}$  an orthogonal matrix, and  $\mathbf{R} \in \mathbb{R}^{d \times d}$  a triangular matrix, such that  $\mathbf{X}\mathbf{\Pi} = \mathbf{Q}\mathbf{R}$ ,*

$$\frac{\sigma_k(\mathbf{X})}{p_1(k, d)} \leq \sigma_{\min}(\mathbf{R}_{[k], [k]}) \leq \sigma_k(\mathbf{X}), \quad (5)$$

and

$$\sigma_{k+1}(\mathbf{X}) \leq \sigma_{\max}(\mathbf{R}_{[k+1:d], [k+1:d]}) \leq p_2(k, d) \sigma_{k+1}(\mathbf{X}), \quad (6)$$

where  $p_1(k, d)$  and  $p_2(k, d)$  are controlled.

In practice, a RRQR factorization algorithm interchanges pairs of columns and updates or builds a QR decomposition on the fly. The link between RRQR factorization and  $k$ -CSSP was first discussed by [Boutsidis, Mahoney, and Drineas \(2009\)](#). The structure of a RRQR factorization indeed gives a deterministic selection of a subset of  $k$  columns of  $\mathbf{X}$ . More precisely, if we take  $\mathbf{C}$  to be the first  $k$  columns of  $\mathbf{X}\mathbf{\Pi}$ ,  $\mathbf{C}$  is a subset of columns of  $\mathbf{X}$  and  $\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_2 = \|\mathbf{R}_{[k+1:r], [k+1:r]}\|_2$ . By (6), any RRQR algorithm thus provides provable guarantees in spectral norm for  $k$ -CSSP.

Following ([Golub, 1965](#)), many papers gave algorithms that improved on  $p_1(k, d)$  and  $p_2(k, d)$  in Definition 1. Table 1 sums up the guarantees of the original algorithm of ([Golub, 1965](#)) and the state-of-the-art algorithms of [Gu and Eisenstat \(1996\)](#). Note the dependency of  $p_2(k, d)$  on the dimension  $d$  through the term  $\sqrt{d - k}$ ; this term is common for guarantees in spectral norm for  $k$ -CSSP. We refer to ([Boutsidis et al., 2009](#)) for an exhaustive survey on RRQR factorization.

### 3.2 Length square importance sampling and additive bounds

[Drineas, Frieze, Kannan, Vempala, and Vinay \(2004\)](#) proposed a randomized CSS algorithm based on i.i.d. sampling  $s$  indices  $S = \{i_1, \dots, i_s\}$  from a multinomial distribution of parameter  $\mathbf{p}$ , where

$$p_j = \frac{\|\mathbf{X}_{:,j}\|_2^2}{\|\mathbf{X}\|_{\text{Fr}}^2}, j \in [d]. \quad (7)$$

Let  $\mathbf{C} = \mathbf{X}_{:,S}$  be the corresponding submatrix. First, we note that some columns of  $\mathbf{X}$  may appear more than once in  $\mathbf{C}$ . Second, ([Drineas et al., 2004](#), Theorem 3) states that

$$\mathbb{P} \left( \|\mathbf{X} - \Pi_{S,k}^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \leq \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}}^2 + 2(1 + \sqrt{8 \log(\frac{2}{\delta})}) \sqrt{\frac{k}{s}} \|\mathbf{X}\|_{\text{Fr}}^2 \right) \geq 1 - \delta. \quad (8)$$

Equation (8) is a high-probability additive upper bound for  $\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2$ . The drawback of such bounds is that they can be very loose if the first  $k$  singular values of  $\mathbf{X}$  are large compared to  $\sigma_{k+1}$ . For this reason, multiplicative approximation bounds have been considered.

### 3.3 $k$ -leverage scores sampling and multiplicative bounds

[Drineas, Mahoney, and Muthukrishnan \(2008\)](#) proposed an algorithm with provable multiplicative upper bound using multinomial sampling, but this time according to  $k$ -leverage scores.

Algorithm	$p_2(k, d)$	Complexity	References
Pivoted QR	$2^k \sqrt{d - k}$	$\mathcal{O}(dNk)$	( <a href="#">Golub and Van Loan, 1996</a> )
Strong RRQR (Alg. 3)	$\sqrt{(d - k)k + 1}$	not polynomial	( <a href="#">Gu and Eisenstat, 1996</a> )
Strong RRQR (Alg. 4)	$\sqrt{f^2(d - k)k + 1}$	$\mathcal{O}(dNk \log_f(d))$	( <a href="#">Gu and Eisenstat, 1996</a> )

Table 1: Examples of some RRQR algorithms and their theoretical performances.

**Definition 2 ( $k$ -leverage scores)** Let  $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top \in \mathbb{R}^{N \times d}$  be the SVD of  $\mathbf{X}$ . We note  $\mathbf{V}_k = \mathbf{V}_{:, [k]}$  the first  $k$  columns of  $\mathbf{V}$ . For  $j \in [d]$ , the  $k$ -leverage score of the  $j$ -th column of  $\mathbf{X}$  is defined by

$$\ell_j^k = \sum_{i=1}^k V_{i,j}^2. \quad (9)$$

In particular, it holds

$$\sum_{j \in [d]} \ell_j^k = \sum_{j \in [d]} \|(\mathbf{V}_k^\top)_{:,j}\|_2^2 = \text{Tr}(\mathbf{V}_k \mathbf{V}_k^\top) = k, \quad (10)$$

since  $\mathbf{V}_k$  is an orthogonal matrix. Therefore, one can consider the multinomial distribution on  $[d]$  with parameters

$$p_j = \frac{\ell_j^k}{k}, j \in [d]. \quad (11)$$

This multinomial is called the  $k$ -leverage scores distribution.

**Theorem 3 (Drineas et al., 2008, Theorem 3)** *If the number  $s$  of sampled columns satisfies*

$$s \geq \frac{4000k^2}{\epsilon^2} \log\left(\frac{1}{\delta}\right), \quad (12)$$

then, under the  $k$ -leverage scores distribution,

$$\mathbb{P}\left(\|\mathbf{X} - \Pi_{S,k}^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \leq (1 + \epsilon)\|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}}^2\right) \geq 1 - \delta. \quad (13)$$

Drineas et al. (2008) also considered replacing multinomial with Bernoulli sampling, still using the  $k$ -leverage scores. The expected number of columns needed for (13) to hold is then lowered to  $\mathcal{O}\left(\frac{k \log k}{\epsilon^2}\right)$ . A natural question is then to understand how low the number of columns can be, while still guaranteeing a multiplicative bound like (13). A partial answer has been given by Deshpande and Vempala (2006).

**Proposition 4 (Deshpande and Vempala, 2006, Proposition 4)** *Given  $\epsilon > 0$ ,  $k, d \in \mathbb{N}$  such that  $d\epsilon \geq 2k$ , there exists a matrix  $\mathbf{X}^\epsilon \in \mathbb{R}^{kd \times k(d+1)}$  such that for any  $S \subset [d]$ ,*

$$\|\mathbf{X}^\epsilon - \Pi_{S,k}^{\text{Fr}} \mathbf{X}^\epsilon\|_{\text{Fr}}^2 \geq (1 + \epsilon)\|\mathbf{X}^\epsilon - \mathbf{X}_k^\epsilon\|_{\text{Fr}}^2. \quad (14)$$

This suggests that the lower bound for the number of columns is  $k/\epsilon$ , at least in the worst case sense of Proposition 4. Interestingly, the  $k$ -leverage scores distribution of the matrix  $\mathbf{X}^\epsilon$  in the proof of Proposition 4 is uniform, so that  $k$ -leverage score sampling boils down to simple uniform sampling.

To match the lower bound of Deshpande and Vempala (2006), Boutsidis, Drineas, and Magdon-Ismail (2011) proposed a greedy algorithm to select columns. This algorithm is inspired by the sparsification of orthogonal matrices proposed in (Batson et al., 2009). The full description of this family of algorithms is beyond the scope of this article. We only recall one of the results of the article.



**Theorem 5 (Boutsidis et al., 2011, Theorem 1.5)** *There exists a randomized algorithm  $\mathcal{A}$  that select at most  $c = \frac{2k}{\epsilon}(1 + o(1))$  columns of  $\mathbf{X}$  such that*

$$\mathbb{E}_{\mathcal{A}} \|\mathbf{X} - \Pi_{S,k}^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \leq (1 + \epsilon) \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}}^2. \quad (15)$$

Finally, a deterministic algorithm based on  $k$ -leverage score sampling was proposed by Papailiopoulos, Kyrillidis, and Boutsidis (2014). The algorithm selects the  $c(\theta)$  columns of  $\mathbf{X}$  with the largest  $k$ -leverage scores, where

$$c(\theta) \in \arg \min_u \left( \sum_{i=1}^u \ell_i^k > \theta \right), \quad (16)$$

and  $\theta$  is a free parameter that controls the approximation error. To guarantee that there exists a matrix of rank  $k$  in the subspace spanned by the selected columns, Papailiopoulos et al. (2014) assume that

$$0 \leq k - \theta < 1. \quad (17)$$

Loosely speaking, this condition is satisfied for a low value of  $c(\theta)$  if the  $k$ -leverage scores (after ordering) are decreasing rapidly enough. The authors give empirical evidence that this condition is satisfied by a large proportion of real datasets.

**Theorem 6 (Papailiopoulos et al., 2014, Theorem 2)** *Let  $\epsilon = k - \theta \in [0, 1)$ , letting  $S$  index the columns with the  $c(\theta)$  largest  $k$ -leverage scores,*

$$\|\mathbf{X} - \Pi_{S,k}^{\nu} \mathbf{X}\|_{\nu} \leq \frac{1}{1 - \epsilon} \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\nu}, \quad \nu \in \{2, \text{Fr}\}. \quad (18)$$

*In particular, if  $\epsilon \in [0, \frac{1}{2}]$ ,*

$$\|\mathbf{X} - \Pi_{S,k}^{\nu} \mathbf{X}\|_{\nu} \leq (1 + 2\epsilon) \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\nu}, \quad \nu \in \{2, \text{Fr}\}. \quad (19)$$

Furthermore, they proved that if the  $k$ -leverage scores decay like a power law, the number of columns needed to obtain a multiplicative bound can actually be smaller than  $\frac{k}{\epsilon}$ .

**Theorem 7 (Papailiopoulos et al., 2014, Theorem 3)** *Assume, for  $\eta > 0$ ,*

$$\ell_i^k = \frac{\ell_1^k}{i^{\eta+1}}. \quad (20)$$

*Let  $\epsilon = k - \theta \in [0, 1)$ , then*

$$c(\theta) = \max \left\{ \left( \frac{4k}{\epsilon} \right)^{\frac{1}{\eta+1}} - 1, \left( \frac{4k}{\eta\epsilon} \right)^{\frac{1}{\eta}}, k \right\}. \quad (21)$$

This complements the fact that the worst case example in Proposition 4 had uniform  $k$ -leverage scores. Loosely speaking, matrices with fast decaying  $k$ -leverage scores can be efficiently subsampled.

### 3.4 The geometric interpretation of the $k$ -leverage scores

The  $k$ -leverage scores can be given a geometric interpretation, the generalization of which serves as a first motivation for our work.

For  $i \in [d]$ , let  $\mathbf{e}_i$  be the  $i$ -th canonical basis vector of  $\mathbb{R}^d$ . Let further  $\theta_i$  be the angle between  $\mathbf{e}_i$  and the subspace  $\mathcal{P}_k = \text{Span}(\mathbf{V}_k)$ , and denote by  $\Pi_{\mathcal{P}_k} \mathbf{e}_i$  the orthogonal projection of  $\mathbf{e}_i$  onto the subspace  $\mathcal{P}_k$ . Then

$$\cos^2(\theta_i) := \frac{(\mathbf{e}_i, \Pi_{\mathcal{P}_k} \mathbf{e}_i)^2}{\|\Pi_{\mathcal{P}_k} \mathbf{e}_i\|^2} = (\mathbf{e}_i, \Pi_{\mathcal{P}_k}(\mathbf{e}_i)) = (\mathbf{e}_i, \sum_{j=1}^k V_{i,j} \mathbf{V}_{:,j}) = \sum_{j=1}^k V_{i,j}^2 = \ell_i^k. \quad (22)$$

A large  $k$ -leverage score  $\ell_i^k$  thus indicates that  $\mathbf{e}_i$  is almost aligned with  $\mathcal{P}_k$ . Selecting columns with large  $k$ -leverage scores as in (Drineas et al., 2008) can thus be interpreted as replacing the principal eigenspace  $\mathcal{P}_k$  by a subspace that must contain  $k$  of the original coordinate axes. Intuitively, a closer subspace to the original  $\mathcal{P}_k$  would be obtained by selecting columns *jointly* rather than independently, considering the angle with  $\mathcal{P}_k$  of the subspace spanned by these columns. More precisely, consider  $S \subset [d]$ ,  $|S| = k$ , and denote  $\mathcal{P}_S = \text{Span}(\mathbf{e}_j, j \in S)$ . A natural definition of the cosine between  $\mathcal{P}_k$  and  $\mathcal{P}_S$  is in terms of the so-called *principal angles* (Golub and Van Loan, 1996, Section 6.4.4); see Appendix C. In particular, Proposition 27 in Appendix C yields

$$\cos^2(\mathcal{P}_k, \mathcal{P}_S) = \text{Det}(\mathbf{V}_{S, [k]})^2. \quad (23)$$

This paper is about sampling  $k$  columns proportionally to (23).

In Appendix A, we contribute a different interpretation of  $k$ -leverage scores and volumes, which relates them to the length-square distribution of Section 3.2.

### 3.5 Negative correlation: volume sampling and the double phase algorithm

In this section, we survey algorithms that randomly sample exactly  $k$  columns from  $\mathbf{X}$ , unlike the multinomial sampling schemes of Sections 3.2 and 3.3, which typically require more than  $k$  columns.

Deshpande, Rademacher, Vempala, and Wang (2006) obtained a multiplicative bound on the expected approximation error, with only  $k$  columns, using so-called *volume sampling*.

**Theorem 8 (Deshpande et al., 2006)** *Let  $S$  be a random subset of  $[d]$ , chosen with probability*

$$\mathbb{P}_{\text{VS}}(S) = Z \text{Det}(\mathbf{X}_{:,S}^T \mathbf{X}_{:,S}) \mathbb{1}_{\{|S|=k\}}, \quad (24)$$

where  $Z = \sum_{|S|=k} \text{Det}(\mathbf{X}_{:,S}^T \mathbf{X}_{:,S})$ . Then

$$\mathbb{E}_{\text{VS}} \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \leq (k+1) \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}}^2 \quad (25)$$

and

$$\mathbb{E}_{\text{VS}} \|\mathbf{X} - \Pi_S^2 \mathbf{X}\|_2^2 \leq (d-k)(k+1) \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}}^2. \quad (26)$$

Note that the bound for the spectral norm was proven in (Deshpande et al., 2006) for the Frobenius projection, that is, they bound  $\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_2$ . The bound (26) easily follows from (4). Later, sampling according to (24) was shown to be doable in polynomial time (Deshpande and Rademacher, 2010). Using a worst case example, Deshpande et al. (2006) proved that the  $k + 1$  factor in (25) cannot be improved.

**Proposition 9 (Deshpande et al., 2006)** *Let  $\epsilon > 0$ . There exists a  $(k + 1) \times (k + 1)$  matrix  $\mathbf{X}^\epsilon$  such that for every subset  $S$  of  $k$  columns of  $\mathbf{X}^\epsilon$ ,*

$$\|\mathbf{X}^\epsilon - \Pi_S^{\text{Fr}} \mathbf{X}^\epsilon\|_{\text{Fr}}^2 > (1 - \epsilon)(k + 1)\|\mathbf{X}^\epsilon - \Pi_k \mathbf{X}^\epsilon\|_{\text{Fr}}^2. \quad (27)$$

We note that there has been recent interest in a similar but different distribution called *dual volume sampling* (Avron and Boutsidis, 2013; Li et al., 2017a; Dereziński and Warmuth, 2018), sometimes also termed volume sampling. The main application of dual VS is row subset selection of a matrix  $\mathbf{X}$  for linear regression on label budget constraints.

(Boutsidis et al., 2009) proposed a  $k$ -CSSP algorithm, called *double phase*, that combines ideas from multinomial sampling and RRQR factorization. The motivating idea is that the theoretical performance of RRQR factorizations depends on the dimension through a factor  $\sqrt{d - k}$ ; see Table 1. To improve on that, the authors propose to first reduce the dimension  $d$  to  $c$  by preselecting a large number of columns  $c > k$  using multinomial sampling from the  $k$ -leverage scores distribution, as in Section 3.3. Then only, they perform a RRQR factorization of the reduced matrix  $\mathbf{V}_k^\top \mathbf{S}_1 \mathbf{D}_1 \in \mathbb{R}^{k \times c}$ , where  $\mathbf{S}_1 \in \mathbb{R}^{d \times c}$  is the sampling matrix of the multinomial phase and  $\mathbf{D}_1 \in \mathbb{R}^{c \times c}$  is a scaling matrix.

**Theorem 10 (Boutsidis et al., 2009)** *Let  $S$  be the output of the double phase algorithm with  $c = \Theta(k \log k)$ . Then*

$$\mathbb{P}_{\text{DPh}} \left( \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}} \leq \Theta(k \log^{\frac{1}{2}} k) \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}} \right) \geq 0.8. \quad (28)$$

$$\mathbb{P}_{\text{DPh}} \left( \|\mathbf{X} - \Pi_S^2 \mathbf{X}\|_2 \leq \Theta(k \log^{\frac{1}{2}} k) \|\mathbf{X} - \Pi_k \mathbf{X}\|_2 + \Theta(k^{\frac{3}{4}} \log^{\frac{1}{4}} k) \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}} \right) \geq 0.8. \quad (29)$$

Note that the spectral norm bound was proven for a slightly different distribution in the randomized phase. Furthermore this bound was proved in (Deshpande et al., 2006) for  $\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_2$  but using (4) the bound (29) follows. The constants  $\Theta(k \log^{\frac{1}{2}} k)$  and  $\Theta(k^{\frac{3}{4}} \log^{\frac{1}{4}} k)$  in the bounds (28) and (29) depends on  $c$  the number of pre-selected columns in the randomized step. In practice, the choice of the parameter  $c$  of the randomized pre-selection phase has an influence on the quality of the approximation. We refer to (Boutsidis et al., 2009) for details.

### 3.6 Excess risk in sketched linear regression

So far, we have focused on approximation bounds in spectral or Frobenius norm for  $\mathbf{X} - \Pi_{S,k}^\nu \mathbf{X}$ . This is a reasonable measure of error as long as it is not known what

the practitioner wants to do with the submatrix  $\mathbf{X}_{:,S}$ . In this section, we assume that the ultimate goal is to perform linear regression of some  $\mathbf{y} \in \mathbb{R}^N$  onto  $\mathbf{X}$ . Other measures of performance then become of interest, such as the excess risk incurred by regressing onto  $\mathbf{X}_{:,S}$  rather than  $\mathbf{X}$ . We use here the framework of [Slawski \(2018\)](#), further assuming well-specification for simplicity.

For every  $i \in [N]$ , assume  $y_i = \mathbf{X}_{i,:} \mathbf{w}^* + \xi_i$ , where the noises  $\xi_i$  are i.i.d. with mean 0 and variance  $v$ . For a given estimator  $\mathbf{w} = \mathbf{w}(\mathbf{X}, \mathbf{y})$ , its excess risk is defined as

$$\mathcal{E}(\mathbf{w}) = \mathbb{E}_{\xi} \left[ \frac{\|\mathbf{X} \mathbf{w}^* - \mathbf{X} \mathbf{w}\|_2^2}{N} \right]. \quad (30)$$

In particular, it is easy to show that the ordinary least squares (OLS) estimator  $\hat{\mathbf{w}} = \mathbf{X}^+ \mathbf{y}$  has excess risk

$$\mathcal{E}(\hat{\mathbf{w}}) = v \times \frac{\text{rk}(\mathbf{X})}{N}. \quad (31)$$

Selecting  $k$  columns indexed by  $S$  in  $\mathbf{X}$  prior to performing linear regression yields  $\mathbf{w}_S = (\mathbf{X} \mathbf{S})^+ \mathbf{y} \in \mathbb{R}^k$ . We are interested in the excess risk of the corresponding sparse vector

$$\hat{\mathbf{w}}_S := \mathbf{S} \mathbf{w}_S = \mathbf{S} (\mathbf{X} \mathbf{S})^+ \mathbf{y} \in \mathbb{R}^d$$

which has all coordinates zero, except those indexed by  $S$ .

**Proposition 11 (Theorem 9, [Mor-Yosef and Avron, 2018](#))** *Let  $S \subset [d]$ , such that  $|S| = k$ . Let  $(\theta_i(S))_{i \in [k]}$  be the principal angles between  $\text{Span } \mathbf{S}$  and  $\text{Span } \mathbf{V}_k$ , see [Appendix C](#). Then*

$$\mathcal{E}(\hat{\mathbf{w}}_S) \leq \frac{1}{N} \left( 1 + \max_{i \in [k]} \tan^2 \theta_i(S) \right) \|\mathbf{w}^*\|^2 \sigma_{k+1}^2 + \frac{vk}{N}. \quad (32)$$

Compared to the excess risk (31) of the OLS estimator, the second term of the right-hand side of (32) replaces  $\text{rk} \mathbf{X}$  by  $k$ . But the price is the first term of the right-hand side of (32), which we loosely term *bias*. To interpret this bias term, we first look at the excess risk of the principal component regressor (PCR)

$$\mathbf{w}_k^* \in \arg \min_{\mathbf{w} \in \text{Span } \mathbf{V}_k} \mathbb{E}_{\xi} [\|\mathbf{y} - \mathbf{X} \mathbf{w}\|^2 / N]. \quad (33)$$

**Proposition 12 (Corollary 11, [Mor-Yosef and Avron, 2018](#))**

$$\mathcal{E}(\mathbf{w}_k^*) \leq \frac{\|\mathbf{w}^*\|^2 \sigma_{k+1}^2}{N} + \frac{vk}{N}. \quad (34)$$

The right-hand side of (34) is almost that of (32), except that the bias term in the CSS risk (32) is larger by a factor that measures how well the subspace spanned by  $S$  is aligned with the principal eigenspace  $\mathbf{V}_k$ . This makes intuitive sense: the performance of CSS will match PCR if selecting columns yields almost the same eigenspace.

The excess risk (32) is yet another motivation to investigate DPPs for column subset selection. We shall see in [Section 5.2](#) that the expectation of (32) under a well-chosen DPP for  $S$  has a particularly simple bias term.

## 4 Determinantal Point Processes

In this section, we introduce discrete determinantal point processes (DPPs) and the related  $k$ -DPPs, of which volume sampling is an example. DPPs were introduced by [Macchi \(1975\)](#) as probabilistic models for beams of fermions in quantum optics. Since then, DPPs have been thoroughly studied in random matrix theory ([Johansson, 2005](#)), and have more recently been adopted in machine learning ([Kulesza et al., 2012](#)), spatial statistics [Lavancier et al. \(2015\)](#), and Monte Carlo methods ([Bardenet and Hardy, 2016](#)).

### 4.1 Definitions

For all the definitions in this section, we refer the reader to ([Kulesza et al., 2012](#)). Recall that  $[d] = \{1, \dots, d\}$ .

**Definition 13 (DPP)** *Let  $\mathbf{K} \in \mathbb{R}^{d \times d}$  be a positive semi-definite matrix. A random subset  $Y \subset [d]$  is drawn from a DPP of marginal kernel  $\mathbf{K}$  if and only if*

$$\forall S \subset [d], \quad \mathbb{P}(S \subset Y) = \text{Det}(\mathbf{K}_S), \quad (35)$$

where  $\mathbf{K}_S = [\mathbf{K}_{i,j}]_{i,j \in S}$ . We take as a convention  $\text{Det}(\mathbf{K}_\emptyset) = 1$ .

For a given matrix  $\mathbf{K}$ , it is not obvious that (35) consistently defines a point process. One sufficient condition is that  $\mathbf{K}$  is symmetric and its spectrum is in  $[0, 1]$ ; see ([Macchi, 1975](#)) and ([Soshnikov, 2000](#)) [Theorem 3]. In particular, when the spectrum of  $\mathbf{K}$  is included in  $\{0, 1\}$ , we call  $\mathbf{K}$  a projection kernel and the corresponding DPP a *projection DPP*<sup>1</sup>. Letting  $r$  be the number of unit eigenvalues of its kernel, samples from a projection DPP have fixed cardinality  $r$  with probability 1 ([Hough et al., 2005](#), Lemma 17).

For symmetric kernels  $\mathbf{K}$ , a DPP can be seen as a *repulsive* distribution, in the sense that for all  $i, j \in [d]$ ,

$$\mathbb{P}(\{i, j\} \subset Y) = \mathbf{K}_{i,i} \mathbf{K}_{j,j} - \mathbf{K}_{i,j}^2 \quad (36)$$

$$= \mathbb{P}(\{i\} \subset Y) \mathbb{P}(\{j\} \subset Y) - \mathbf{K}_{i,j}^2 \quad (37)$$

$$\leq \mathbb{P}(\{i\} \subset Y) \mathbb{P}(\{j\} \subset Y). \quad (38)$$

Besides projection DPPs, there is another natural way of using a kernel matrix to define a random subset of  $[d]$  with prespecified cardinality  $k$ .

**Definition 14 ( $k$ -DPP)** *Let  $\mathbf{L} \in \mathbb{R}^{d \times d}$  be a positive semi-definite matrix. A random subset  $Y \subset [d]$  is drawn from a  $k$ -DPP of kernel  $\mathbf{L}$  if and only if*

$$\forall S \subset [d], \quad \mathbb{P}(Y = S) \propto \mathbb{1}_{\{|S|=k\}} \text{Det}(\mathbf{L}_S) \quad (39)$$

where  $\mathbf{L}_S = [\mathbf{L}_{i,j}]_{i,j \in S}$ .

DPPs and  $k$ -DPPs are closely related but different objects. For starters,  $k$ -DPPs are always well defined, provided  $\mathbf{L}$  has a nonzero minor of size  $k$ .

<sup>1</sup>All projection DPPs in this paper have symmetric kernels

## 4.2 Sampling from a DPP and a $k$ -DPP

Let  $\mathbf{K} \in \mathbb{R}^{d \times d}$  be a symmetric, positive semi-definite matrix, with eigenvalues in  $[0, 1]$ , so that  $\mathbf{K}$  is the marginal kernel of a DPP on  $[d]$ . Let us diagonalize it as  $\mathbf{K} = \mathbf{V} \text{Diag}(\lambda_i) \mathbf{V}^\top$ . [Hough et al. \(2005\)](#) established that sampling from the DPP with kernel  $\mathbf{K}$  can be done by (i) sampling independent Bernoullis  $B_i, i = 1, \dots, d$ , with respective parameters  $\lambda_i$ , (ii) forming the submatrix  $\mathbf{V}_{:,B}$  of  $\mathbf{V}$  corresponding to columns  $i$  such that that  $B_i = 1$ , and (iii) sampling from the projection DPP with kernel

$$\mathbf{K}_{\text{proj}} = \mathbf{V}_{:,B} \mathbf{V}_{:,B}^\top.$$

The only nontrivial step is sampling from a projection DPP, for which we give pseudocode in Figure 1; see ([Hough et al., 2005](#), Theorem 7) or ([Kulesza et al., 2012](#), Theorem 2.3) for a proof. For a survey of variants of the algorithm, we also refer to ([Tremblay et al., 2018](#)) and the documentation of the DPPy toolbox<sup>2</sup> ([Gautier et al., 2018](#)). For our purposes, it is enough to remark that general DPPs are mixtures of projection DPPs of different ranks, and that the cardinality of a general DPP is a sum of independent Bernoulli random variables.

```

PROJECTIONDPP( $\mathbf{K}_{\text{proj}} = \mathbf{V} \mathbf{V}^\top$ )
1    $Y \leftarrow \emptyset$ 
2    $\mathbf{W} \leftarrow \mathbf{V}$ 
3   while  $\text{rk}(\mathbf{W}) > 0$ 
4       Sample  $i$  from  $\Omega$  with probability  $\propto \|\mathbf{W}_{i,:}\|_2^2$  ▷ Chain rule
5        $Y \leftarrow Y \cup \{i\}$ 
6        $\mathbf{V} \leftarrow \mathbf{V}_\perp$  an orthonormal basis of  $\text{Span}(\mathbf{V} \cap e_i^\perp)$ 
7   return  $Y$ 

```

Figure 1: Pseudocode for sampling from a DPP of marginal kernel  $\mathbf{K}$ .

The next proposition establishes that  $k$ -DPPs also are mixtures of projection DPPs.

**Proposition 15** ([Kulesza et al. \(2012, Section 5.2.2\)](#)) *Let  $Y$  be a random subset of  $[d]$  sampled from a  $k$ -DPP with kernel  $\mathbf{L}$ . We further assume that  $\mathbf{L}$  is symmetric, we denote its rank by  $r$  and its diagonalization by  $\mathbf{L} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$ . Finally, let  $k \leq r$ . It holds*

$$\mathbb{P}(Y = S) = \sum_{\substack{T \subset [r] \\ |T|=k}} \mu_T \left[ \frac{1}{k!} \text{Det} \left( \mathbf{V}_{T,S} \mathbf{V}_{T,S}^\top \right) \right] \quad (40)$$

where

$$\mu_T = \frac{\prod_{i \in T} \lambda_i}{\sum_{\substack{U \subset [r] \\ |U|=k}} \prod_{i \in U} \lambda_i}. \quad (41)$$

<sup>2</sup><http://github.com/guilgautier/DPPy>

Each mixture component in square brackets in (40) is a projection DPP with cardinality  $k$ . Sampling a  $k$ -DPP can thus be done by (i) sampling a multinomial distribution with parameters (41), and (ii) sampling from the corresponding projection DPP using the algorithm in Figure 1. The main difference between  $k$ -DPPs and DPPs is that all mixture components in (40) have the same cardinality  $k$ . In particular, projection DPPs are the only DPPs that are also  $k$ -DPPs.

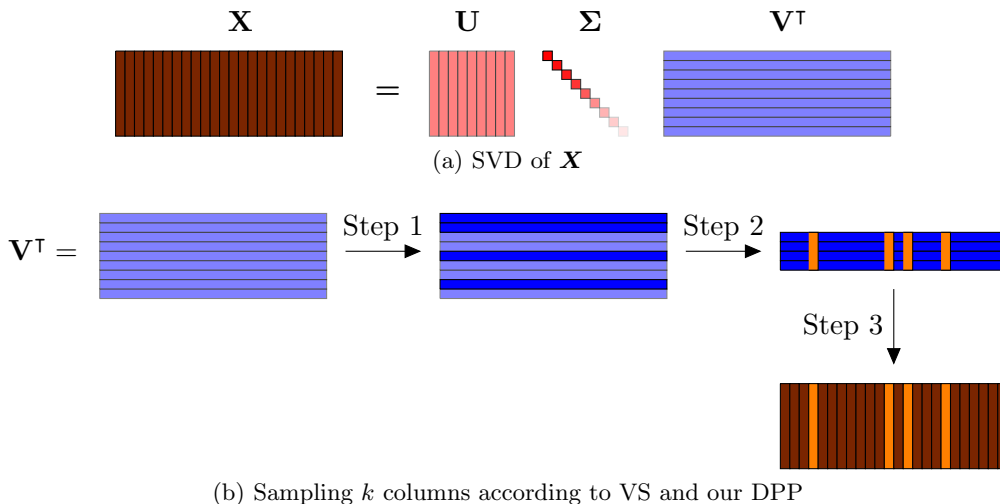


Figure 2: A graphical depiction of the sampling algorithms for volume sampling (VS) and the DPP with marginal kernel  $\mathbf{V}_k \mathbf{V}_k^\top$ . (a) Both algorithms start with an SVD. (b) In Step 1, VS randomly selects  $k$  rows of  $\mathbf{V}^\top$ , while the DPP always picks the first  $k$  rows. Step 2 is the same for both algorithms: jointly sample  $k$  columns of the subsampled  $\mathbf{V}^\top$ , proportionally to their squared volume. Step 3 is simply the extraction of the corresponding columns of  $\mathbf{X}$ .

A fundamental example of  $k$ -DPP is volume sampling, as defined in Section 3.5. Its kernel is the Gram matrix of the data  $\mathbf{L} = \mathbf{X}^\top \mathbf{X}$ . In general,  $\mathbf{L}$  is not an orthogonal projection, so that volume sampling is not a DPP.

### 4.3 Motivations for column subset selection using projection DPPs

volume sampling has been successfully used for column subset selection, see Section 3.5. Our motivation to investigate projection DPPs instead of volume sampling is twofold.

Following (40), volume sampling can be seen as a mixture of projection DPPs indexed by  $T \subset [d], |T| = k$ , with marginal kernels  $\mathbf{K}_T = \mathbf{V}_{:,T} \mathbf{V}_{:,T}^\top$  and mixture weights  $\mu_T \propto \prod_{i \in T} \sigma_i^2$ . The component with the highest weight thus corresponds to the  $k$  largest singular values, that is, the projection DPP with marginal kernel  $\mathbf{K} := \mathbf{V}_k \mathbf{V}_k^\top$ . This paper is about column subset selection using precisely this DPP. Alternately, we could motivate the study of this DPP by remarking that its marginals  $\mathbb{P}(i \subset Y)$  are the  $k$ -leverage scores introduced in Section 3.3. Since  $\mathbf{K}$  is symmetric, this DPP can be seen as a repulsive generalization of leverage score sampling.

Finally, we recap the difference between volume sampling and the DPP with kernel  $\mathbf{K}$  with a graphical depiction in Figure 2 of the two procedures to sample

from them that we introduced in Section 4.2. Figure 2 is another illustration of the decomposition of volume sampling as a mixture of projection DPPs.

## 5 Main Results

In this section, we prove bounds for  $\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^\nu \mathbf{X}\|_\nu$  under the projection DPP of marginal kernel  $\mathbf{K} = \mathbf{V}_k \mathbf{V}_k^\top$  presented in Section 4. Throughout, we compare our bounds to the state-of-the-art bounds of volume sampling obtained by [Deshpande et al. \(2006\)](#); see Theorem 8 and Section 3.5. For clarity, we defer the proofs of our results from this section to Appendix D.

### 5.1 Multiplicative bounds in spectral and Frobenius norm

Let  $S$  be a random subset of  $k$  columns of  $\mathbf{X}$  chosen with probability:

$$\mathbb{P}_{\text{DPP}}(S) = \text{Det}(\mathbf{V}_{S,[k]})^2. \quad (42)$$

First, without any further assumption, we have the following result.

**Proposition 16** *Under the projection DPP of marginal kernel  $\mathbf{V}_k \mathbf{V}_k^\top$ , it holds*

$$\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^\nu \mathbf{X}\|_\nu^2 \leq k(d+1-k) \|\mathbf{X} - \Pi_k \mathbf{X}\|_\nu^2, \quad \nu \in \{2, \text{Fr}\}. \quad (43)$$

For the spectral norm, the bound is practically the same as that of volume sampling (26). However, our bound for the Frobenius norm is worse than (25) by a factor  $(d-k)$ . In the rest of this section, we sharpen our bounds by taking into account the sparsity level of the  $k$ -leverage scores and the decay of singular values.

In terms of sparsity, we first replace the dimension  $d$  in (43) by the number  $p \in [d]$  of nonzero  $k$ -leverage scores

$$p = |\{i \in [d], \mathbf{V}_{i,[k]} \neq \mathbf{0}\}|. \quad (44)$$

To quantify the decay of the singular values, we define the flatness parameter

$$\beta = \sigma_{k+1}^2 \left( \frac{1}{d-k} \sum_{j \geq k+1} \sigma_j^2 \right)^{-1}. \quad (45)$$

In words,  $\beta \in [1, d-k]$  measures the flatness of the spectrum of  $\mathbf{X}$  above the cut-off at  $k$ . Indeed, (45) is the ratio of the largest term in a sum to that sum. The closer  $\beta$  is to 1, the more similar the terms in the sum in the denominator of (45). At the extreme,  $\beta = d-k$  when  $\sigma_{k+1}^2 > 0$  while  $\sigma_j^2 = 0, \forall j \geq k+2$ .

**Proposition 17** *Under the projection DPP of marginal kernel  $\mathbf{V}_k \mathbf{V}_k^\top$ , it holds*

$$\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^2 \mathbf{X}\|_2^2 \leq k(p-k) \|\mathbf{X} - \Pi_k \mathbf{X}\|_2^2 \quad (46)$$

and

$$\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \leq \left( 1 + \beta \frac{p-k}{d-k} \right) \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}}^2. \quad (47)$$



The bound in (46) compares favorably with volume sampling (26) since the dimension  $d$  has been replaced by the sparsity level  $p$ . For  $\beta$  close to 1, the bound in (47) is better than the bound (25) of volume sampling since  $(p - k)/(d - k) \leq 1$ . Again, the sparser the  $k$ -leverage scores, the smaller the bounds.

Now, one could argue that, in practice, sparsity is never exact: it can well be that  $p = d$  while there still are a lot of small  $k$ -leverage scores. We will demonstrate in Section 6 that the DPP still performs better than volume sampling in this setting, which Proposition 17 doesn't reflect. We introduce two ideas to further tighten the bounds of Proposition 17. First, we define an effective sparsity level in the vein of Papailiopoulos et al. (2014), see Section 3.3. Second, we condition the DPP on a favourable event with controlled probability.

**Theorem 18** *Let  $\pi$  be a permutation of  $[d]$  such that leverage scores are reordered*

$$\ell_{\pi_1}^k \geq \ell_{\pi_2}^k \geq \dots \geq \ell_{\pi_d}^k. \quad (48)$$

For  $\delta \in [d]$ , let  $T_\delta = [\pi_\delta, \dots, \pi_d]$ . Let  $\theta > 1$  and

$$p_{\text{eff}}(\theta) = \min \left\{ q \in [d] \mid \sum_{i \leq q} \ell_{\pi_i}^k \geq k - 1 + \frac{1}{\theta} \right\}. \quad (49)$$

Finally, let  $\mathcal{A}_\theta$  be the event  $\{S \cap T_{p_{\text{eff}}(\theta)} = \emptyset\}$ . Then, on the one hand,

$$\mathbb{P}_{\text{DPP}}(\mathcal{A}_\theta) \geq \frac{1}{\theta}, \quad (50)$$

and, on the other hand,

$$\mathbb{E}_{\text{DPP}} [\|\mathbf{X} - \Pi_S^2 \mathbf{X}\|_2^2 \mid \mathcal{A}_\theta] \leq (p_{\text{eff}}(\theta) - k + 1)(k - 1 + \theta) \|\mathbf{X} - \Pi_k \mathbf{X}\|_2^2 \quad (51)$$

and

$$\mathbb{E}_{\text{DPP}} [\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \mid \mathcal{A}_\theta] \leq \left( 1 + \beta \frac{(p_{\text{eff}}(\theta) + 1 - k)}{d - k} (k - 1 + \theta) \right) \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}}^2. \quad (52)$$

In Theorem 18, the effective sparsity level  $p_{\text{eff}}(\theta)$  replaces the sparsity level  $p$  of Proposition 17. The key is to condition on  $S$  not containing any index of column with too small a  $k$ -leverage score, that is, the event  $\mathcal{A}_\theta$ . In practice, this is achieved by rejection sampling: we repeatedly and independently sample  $S \sim \text{DPP}(\mathbf{K})$  until  $S \cap T_{p_{\text{eff}}(\theta)} = \emptyset$ .

The caveat of any rejection sampling procedure is a potentially large number of samples required before acceptance. But in the present case, Equation (50) guarantees that the expectation of that number of samples is less than  $\theta$ . The free parameter  $\theta$  thus interestingly controls both the “energy” threshold in (49), and the complexity of the rejection sampling. The approximation bounds suggest picking  $\theta$  close to 1, which implies a compromise with the value of  $p_{\text{eff}}(\theta)$  that should not be too large either. We have empirically observed that the performance of the DPP is relatively insensitive to the choice of  $\theta$ .

## 5.2 Bounds for the excess risk in sketched linear regression

In Section 3.6, we surveyed bounds on the excess risk of ordinary least squares estimators that relied on a subsample of the columns of  $\mathbf{X}$ . Importantly, the generic bound (32) of Mor-Yosef and Avron (2018) has a bias term that depends on the maximum squared tangent of the principal angles between  $\text{Span}(\mathbf{S})$  and  $\text{Span}(\mathbf{V}_k)$ . When  $|S| = k$ , this quantity is hard to control without making strong assumptions on the matrix  $\mathbf{V}_k$ . But it turns out that, in expectation under the same DPP as in Section 5.1, this bias term drastically simplifies.

**Proposition 19** *We use the notation of Section 3.6. Under the projection DPP with marginal kernel  $\mathbf{V}_k \mathbf{V}_k^\top$ , it holds*

$$\mathbb{E}_{\text{DPP}} [\mathcal{E}(\mathbf{w}_S)] \leq (1 + k(p - k)) \frac{\|\mathbf{w}^*\|^2 \sigma_{k+1}^2}{N} + \frac{vk}{N}. \quad (53)$$

The sparsity level  $p$  appears again in the bound (53): The sparser the  $k$ -leverage scores distribution, the smaller the bias term. The bound (53) only features an additional  $(1 + k(p - k))$  factor in the bias term, compared to the bound obtained by Mor-Yosef and Avron (2018) for PCR, see Proposition 12. Loosely speaking, this factor is to be seen as the price we accept to pay in order to get more interpretable features than principal components in the linear regression problem. Finally, a natural question is to investigate the choice of  $k$  to minimize the bound in (53), but this is out of the scope of this paper.

As in Theorem 18, for practical purposes, it can be desirable to bypass the need for the exact sparsity level  $p$  in Proposition 19. We give a bound that replaces  $p$  with the effective sparsity level  $p_{\text{eff}}(\theta)$  introduced in (49).

**Theorem 20** *Using the notation of Section 3.6 for linear regression, and of Theorem 18 for leverage scores and their indices, it holds*

$$\mathbb{E}_{\text{DPP}} [\mathcal{E}(\hat{\mathbf{w}}_S) \mid \mathcal{A}_\theta] \leq [1 + (k - 1 + \theta)(p_{\text{eff}}(\theta) - k + 1)] \frac{\|\mathbf{w}^*\|^2 \sigma_{k+1}^2}{N} + \frac{vk}{N}. \quad (54)$$

In practice, the same rejection sampling routine as in Theorem 18 can be used to sample conditionally on  $\mathcal{A}_\theta$ . Finally, to the best of our knowledge, bounding the excess risk in linear regression has not been investigated under volume sampling.

In summary, we have obtained two sets of results. We have proven a set of multiplicative bounds in spectral and Frobenius norm for  $\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^\nu \mathbf{X}\|_\nu$ ,  $\nu \in \{2, \text{Fr}\}$ , under the projection DPP of marginal kernel  $\mathbf{K} = \mathbf{V}_k \mathbf{V}_k^\top$ , see Propositions 16 & 17 and Theorem 18. As far as the linear regression problem is concerned, we have proven bounds for the excess risk in sketched linear regression, see Proposition 19 and Theorem 20.

## 6 Numerical experiments

In this section, we empirically compare our algorithm to the state of the art in column subset selection. In Section 6.1, the projection DPP with kernel  $\mathbf{K} = \mathbf{V}_k \mathbf{V}_k^\top$  and

volume sampling are compared on toy datasets. In Section 6.2, several column subset selection algorithms are compared to the projection DPP on four real datasets from genomics and text processing. In particular, the numerical simulations demonstrate the favorable influence of the sparsity of the  $k$ -leverage scores on the performance of our algorithm both on toy datasets and real datasets. Finally, we packaged all CSS algorithms in this section in a Python toolbox<sup>3</sup>.

## 6.1 Toy datasets

This section is devoted to comparing the expected approximation error  $\mathbb{E}\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2$  for the projection DPP and volume sampling. We focus on the Frobenius norm to avoid effects due to different choices of the projection  $\Pi_S^{\text{Fr}}$ , see (4).

In order to be able to evaluate the expected errors *exactly*, we generate matrices of low dimension ( $d = 20$ ) so that the subsets of  $[d]$  can be exhaustively enumerated. Furthermore, to investigate the role of leverage scores and singular values on the performance of CSS algorithms, we need to generate datasets  $\mathbf{X}$  with prescribed spectra and  $k$ -leverage scores.

### 6.1.1 Generating toy datasets

Recall that the SVD of  $\mathbf{X} \in \mathbb{R}^{N \times d}$  reads  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ , where  $\mathbf{\Sigma}$  is a diagonal matrix and  $\mathbf{U}$  and  $\mathbf{V}$  are orthogonal matrices. To sample a matrix  $\mathbf{X}$ ,  $\mathbf{U}$  is first drawn from the Haar measure of  $\mathcal{O}_N(\mathbb{R})$ , then  $\mathbf{\Sigma}$  is chosen among a few deterministic diagonal matrices that illustrate various spectral properties. Sampling the matrix  $\mathbf{V}$  is trickier. The first  $k$  columns of  $\mathbf{V}$  are structured as follows: the number of non vanishing rows of  $\mathbf{V}_k$  is equal to  $p$  and the norms of the nonvanishing rows are prescribed by a vector  $\boldsymbol{\ell}$ . By considering the matrix  $\mathbf{K} = \mathbf{V}_k \mathbf{V}_k^\top$ , generating  $\mathbf{V}$  boils down to the simulation of an Hermitian matrix with prescribed diagonal and spectrum (in this particular case the spectrum is included in  $\{0, 1\}$ ). For this reason, we propose an algorithm that takes as input a leverage scores profile  $\boldsymbol{\ell}$  and a spectrum  $\boldsymbol{\sigma}^2$ , and outputs a corresponding random orthogonal matrix  $\mathbf{X}$ ; see Appendix E. This algorithm is a randomization<sup>4</sup> of the algorithm proposed by Fickus, Mixon, Poteet, and Strawn (2011b). Finally, the matrix  $\mathbf{V}_k \in \mathbb{R}^{d \times k}$  is completed by applying the Gram-Schmidt procedure to  $d - k$  additional i.i.d. unit Gaussian vectors, resulting in a matrix  $\mathbf{V} \in \mathbb{R}^{d \times d}$ . Figure 3 summarizes the algorithm proposed to generate matrices  $\mathbf{X}$  with a  $k$ -leverage scores profile  $\boldsymbol{\ell}$  and a sparsity level  $p$ .

### 6.1.2 volume sampling vs projection DPP

This section sums up the results of numerical simulations on toy datasets. The number of observations is fixed to  $N = 100$ , the dimension to  $d = 20$ , and the number of selected columns to  $k \in \{3, 5\}$ . Singular values of are chosen from the following profiles: a spectrum with a cutoff called the projection spectrum,

$$\boldsymbol{\Sigma}_{k=3, \text{proj}} = 100 \sum_{i=1}^3 \mathbf{e}_i \mathbf{e}_i^\top + 0.1 \sum_{i=4}^{20} \mathbf{e}_i \mathbf{e}_i^\top,$$

<sup>3</sup><http://github.com/AyoubBelhadji/CSSPy>

<sup>4</sup><http://github.com/AyoubBelhadji/FrameBuilder>

MATRIXGENERATOR( $\ell, \Sigma$ )	
1	Sample $\mathbf{U}$ from the Haar measure $\mathbb{O}_N(\mathbb{R})$ .
2	Pick $\Sigma$ a diagonal matrix.
3	Pick $p \in [k + 1 : d]$ .
4	Pick a $k$ -leverage-scores profile $\ell \in \mathbb{R}_+^d$ with a sparsity level $p$ .
5	Generate a matrix $\mathbf{V}_k$ with the $k$ -leverage-scores profile $\ell$ .
6	Extend the matrix $\mathbf{V}_k$ to an orthogonal matrix $\mathbf{V}$ .
7	<b>return</b> $\mathbf{X} \leftarrow \mathbf{U}\Sigma\mathbf{V}^\top$

Figure 3: The pseudocode of the algorithm generating a matrix  $\mathbf{X}$  with prescribed profile of  $k$ -leverage scores.

$$\Sigma_{k=5,\text{proj}} = 100 \sum_{i=1}^5 \mathbf{e}_i \mathbf{e}_i^\top + 0.1 \sum_{i=6}^{20} \mathbf{e}_i \mathbf{e}_i^\top.$$

and a smooth spectrum

$$\Sigma_{k=3,\text{smooth}} = 100\mathbf{e}_1 \mathbf{e}_1^\top + 10\mathbf{e}_2 \mathbf{e}_2^\top + \mathbf{e}_3 \mathbf{e}_3^\top + 0.1 \sum_{i=4}^{20} \mathbf{e}_i \mathbf{e}_i^\top,$$

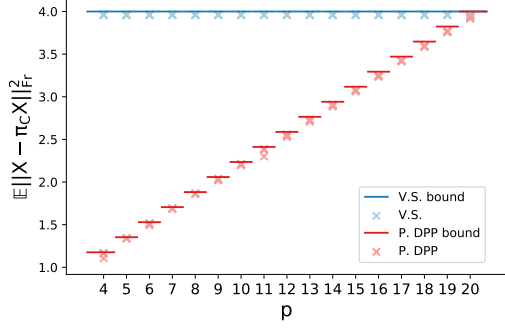
$$\Sigma_{k=5,\text{smooth}} = 10000\mathbf{e}_1 \mathbf{e}_1^\top + 1000\mathbf{e}_2 \mathbf{e}_2^\top + 100\mathbf{e}_3 \mathbf{e}_3^\top + 10\mathbf{e}_4 \mathbf{e}_4^\top + \mathbf{e}_5 \mathbf{e}_5^\top + 0.1 \sum_{i=6}^{20} \mathbf{e}_i \mathbf{e}_i^\top.$$

Note that all profiles satisfy  $\beta = 1$ ; see (45). In each experiment, for each spectrum, we sample 200 independent leverage scores profiles that satisfy the sparsity constraints from a Dirichlet distribution with concentration parameter 1 and equal means. For each leverage scores profile, we sample a matrix  $\mathbf{X}$  from the algorithm in Appendix E.

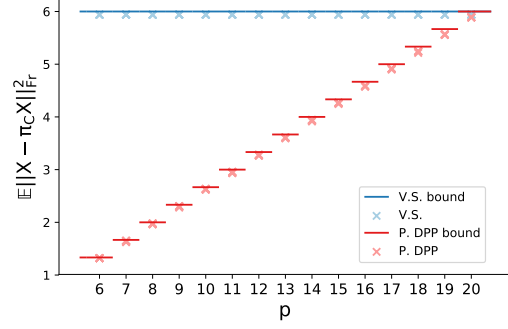
Figure 4 compares, on the one hand, the theoretical bounds in Theorem 8 for volume sampling and Proposition 17 for the projection DPP, to the numerical evaluation of the expected error for sampled toy datasets on the other hand. The x-axis indicates various sparsity levels  $p$ . The unit on the y-axis is the error of PCA. There are 400 crosses on each subplot: each of the 200 matrices appears once for both algorithms. The 200 matrices are spread evenly across the values of  $p$ .

The VS bounds in  $(k + 1)$  are independent of  $p$ . They appear to be tight for projection spectra, and looser for smooth spectra. For the projection DPP, the bound  $(k + 1) \frac{p-k}{d-k}$  is linear in  $p$ , and can be much lower than the bound of VS. The numerical evaluations of the error also suggest that this DPP bound is tight for a projection spectrum, and looser in the smooth case. In both cases, the bound is representative of the actual behavior of the algorithm.

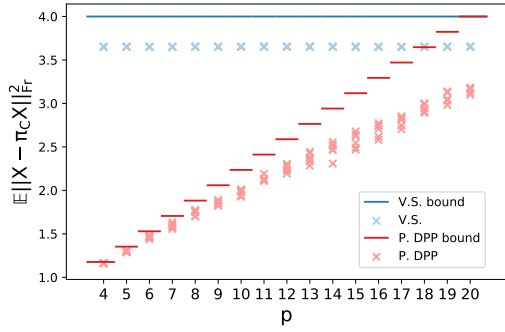
The bottom row of Figure 4 displays the same results for identity spectra, again for  $k = 3$  and  $k = 5$ . This setting is extremely nonsparse and represents an arbitrarily bad scenario where even PCA would not make much practical sense. Both VS and DPP sampling perform constantly badly, and all crosses superimpose at  $y = 1$ , which indicates the PCA error. In this particular case, our linear bound in



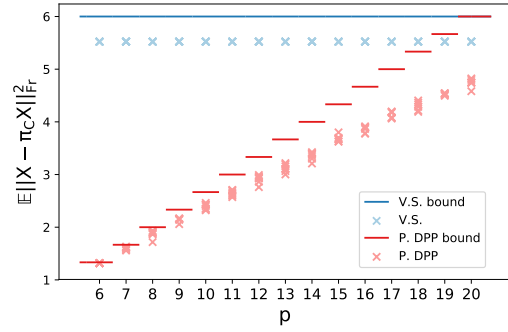
(a)  $\Sigma_{3,\text{proj}}, k = 3$



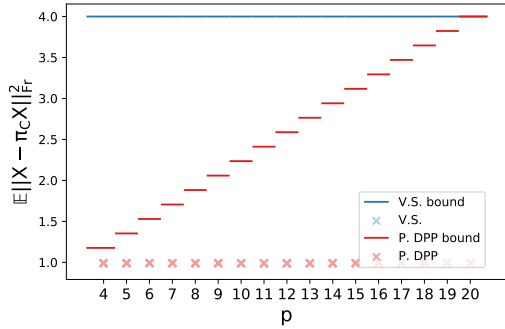
(d)  $\Sigma_{5,\text{proj}}, k = 5$



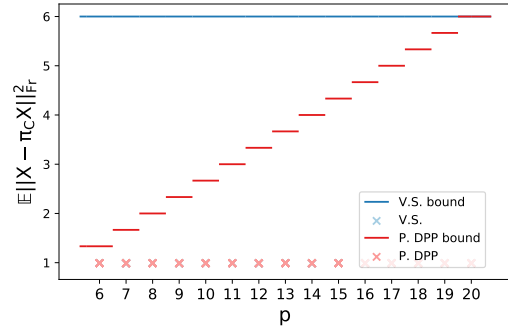
(b)  $\Sigma_{3,\text{smooth}}, k = 3$



(e)  $\Sigma_{5,\text{smooth}}, k = 5$



(d)  $I_{20}, k = 3$



(f)  $I_{20}, k = 5$

Figure 4: Realizations and bounds for  $\mathbb{E}\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2$  as a function of the sparsity level  $p$ .

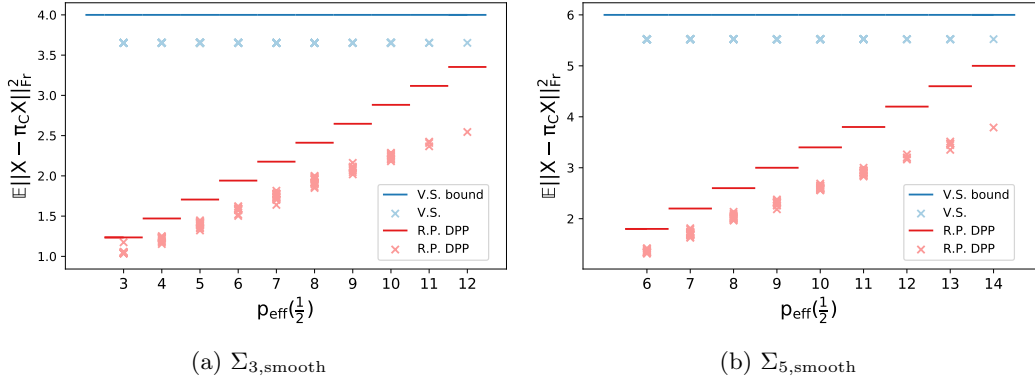


Figure 5: Realizations and bounds for  $\mathbb{E}\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2$  as a function of the effective sparsity level  $p_{\text{eff}}(\frac{1}{2})$ .

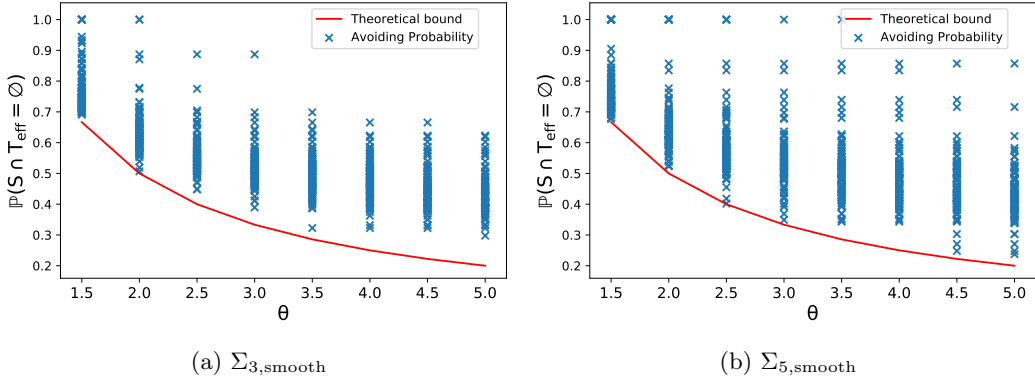


Figure 6: Realizations and bounds for the avoiding probability  $\mathbb{P}(S \cap T_{p_{\text{eff}}(\theta)} = \emptyset)$  in Theorem 18 as a function of  $\theta$ .

$p$  is not representative of the actual behavior of the error. This observation can be explained for volume sampling using Theorem 25, which states that the expected squared error under VS is Schur-concave, and is thus minimized for flat spectra. We have no similar result for the projection DPP.

Figure 5 provides a similar comparison for the two smooth spectra  $\Sigma_{3,\text{smooth}}$  and  $\Sigma_{5,\text{smooth}}$ , but this time using the effective sparsity level  $p_{\text{eff}}(\theta)$  introduced in Theorem 18. We use  $\theta = 1/2$ ; qualitatively, we have observed the results to be robust to the choice of  $\theta$ . The 200 sampled matrices are now unevenly spread across the  $x$ -axis, since we do not control  $p_{\text{eff}}(\theta)$ . Note finally that the DPP here is conditioned on the event  $\{S \cap T_{p_{\text{eff}}(\theta)} = \emptyset\}$ , and sampled using an additional rejection sampling routine as detailed below Theorem 18.

For the DPP, the bound is again linear on the effective sparsity level  $p_{\text{eff}}(\frac{1}{2})$ , and can again be much lower than the VS bound. The behavior of both VS and the projection DPP are similar to the exact sparsity setting of Figure 4: the DPP has uniformly better bounds and actual errors, and the bound reflects the actual behavior, if relatively loosely when  $p_{\text{eff}}(1/2)$  is large.

Figure 6 compares the theoretical bound in Theorem 18 for the avoiding probability  $\mathbb{P}(S \cap T_{p_{\text{eff}}(\theta)} = \emptyset)$  with 200 realizations, as a function of  $\theta$ . More precisely, we drew 200 matrices  $\mathbf{X}$ , and then for each  $\mathbf{X}$ , we computed exactly – by enumeration – the value  $\mathbb{P}(S \cap T_{p_{\text{eff}}(\theta)} = \emptyset)$  for all values of  $\theta$ . The only randomness is thus in the sampling of  $\mathbf{X}$ , not the evaluation of the probability. The results suggest again that the bound is relatively tight.

## 6.2 Real datasets

This section compares the empirical performances of several column subset selection algorithms on the datasets in Table 2.

Dataset	Application domain	$N \times d$	References
Colon	genomics	$62 \times 2000$	(Alon et al., 1999)
Leukemia	genomics	$72 \times 7129$	(Golub et al., 1999)
Basehock	text processing	$1993 \times 4862$	(Li et al., 2017b)
Relathe	text processing	$1427 \times 4322$	(Li et al., 2017b)

Table 2: Datasets used in the experimental section.

These datasets are illustrative of two extreme situations regarding the sparsity of the  $k$ -leverage scores. For instance, the dataset Basehock has a very sparse profile of  $k$ -leverage scores, while the dataset Colon has a quasi-uniform distribution of  $k$ -leverage scores, see Figures 7 (a) & (b).

We consider the following algorithms presented in Section 3: 1) the projection DPP with marginal kernel  $\mathbf{K} = \mathbf{V}_k \mathbf{V}_k^\top$ , 2) volume sampling, using the implementation proposed by Kulesza and Taskar (2011), 3) deterministically picking the largest  $k$ -leverage scores, 4) pivoted QR as in (Golub, 1965), although the only known bounds for this algorithm are for the spectral norm, and 5) double phase, with  $c$  manually tuned to optimize the performance, usually around  $c \approx 10k$ .

The rest of Figure 7 sums up the empirical results of the previously described algorithms on the Colon and Basehock datasets. Figures 7 (c) & (d) illustrate the results of the five algorithms in the following setting. An ensemble of 50 subsets are sampled from each algorithm. We give the corresponding boxplots for the Frobenius errors, on Colon and Basehock respectively. We observe that the increase in performance using projection DPP compared to volume sampling is more important for the Basehock dataset than for the Colon dataset: this improvement can be explained by the sparsity of the  $k$ -leverage scores as predicted by our approximation bounds. Deterministic methods (largest leverage scores and pivoted QR) perform well compared with other algorithms on the Basehock dataset; in contrast, they display very bad performances on the Colon dataset. The double phase algorithm has the best results on both datasets. However its theoretical guarantees cannot predict such an improvement, as noted in Section 3. The performance of the projection DPP is comparable to those Double Phase and makes it a close second, with a slightly larger gap on the Colon dataset. We emphasize that our approximation bounds are sharp compared to numerical observations.

Figures 7 (e) & (f) show results obtained using a classical boosting technique for randomized algorithms. We repeat 20 times: sample 50 subsets and take the

best subset selection. Displayed boxplots are for these 20 best results. The same comments apply as without boosting.

Figure 8 calls for similar comments, comparing this time the datasets Relathe (with concentrated profile of  $k$ -leverage scores) and Leukemia (with almost uniform profile of  $k$ -leverage scores).

### 6.3 Discussion

The performance of our algorithm has been compared to state-of-the-art column subset selection algorithms. We emphasize that the theoretical performances of the proposed approach take into account the sparsity of the  $k$ -leverage scores as in Proposition 17 or their fast decrease as in Proposition 18, and that the bounds are in good agreement with the actual behavior of the algorithm. In contrast, state-of-the-art algorithms like volume sampling have looser bounds and worse performances, or like double phase display great performance but have overly pessimistic theoretical bounds.

## 7 Conclusion

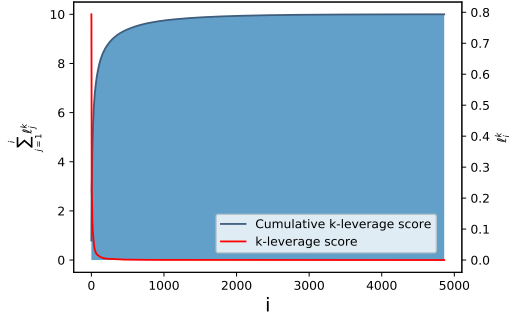
We have proposed, analyzed, and empirically investigated a new randomized column subset selection (CSS) algorithm. The crux of our algorithm is a discrete determinantal point process (DPP) that selects a diverse set of  $k$  columns of a matrix  $\mathbf{X}$ . This DPP is tailored to CSS through its parametrization by the marginal kernel  $\mathbf{K} = \mathbf{V}_k \mathbf{V}_k^\top$ , where  $\mathbf{V}_k$  are the first  $k$  right singular vectors of the matrix  $\mathbf{X}$ . This specific kernel is related to volume sampling, the state-of-the-art for CSS guarantees in Frobenius and spectral norm.

We have identified generic conditions on the matrix  $\mathbf{X}$  under which our algorithm has bounds that improve on volume sampling. In particular, our bounds highlight the importance of the sparsity and the decay of the  $k$ -leverage scores on the approximation performance of our algorithm. This resonates with the compressed sensing literature. We have further numerically illustrated this relation to the sparsity and decay of the  $k$ -leverage scores using toy and real datasets. In these experiments, our algorithm performs comparably to the so-called double phase algorithm, which is the empirical state-of-the-art for CSS despite more conservative theoretical guarantees than volume sampling. Thus, our DPP sampling inherits both favorable bounds and increased empirical performance under sparsity or fast decay of the  $k$ -leverage scores.

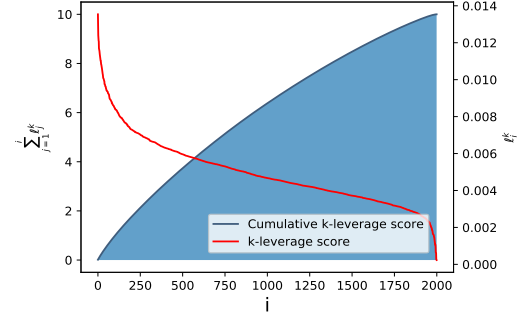
In terms of computational cost, our algorithms scale with the cost of finding the  $k$  first right singular vectors, which is currently the main bottleneck. In the line of (Drineas et al., 2012) and (Boutsidis et al., 2011), where the authors estimates the  $k$ -leverage scores using random projections, we plan to investigate the impact of random projections to estimate the full matrix  $\mathbf{K}$  on the approximation guarantees of our algorithms.

Although often studied as an independent task, in practice CSS is often a prelude to a learning algorithm. We have considered linear regression and we have given a bound on the excess risk of a regression performed on the selected columns only. In particular, sparsity and decay of the  $k$ -leverage scores are again involved: the more

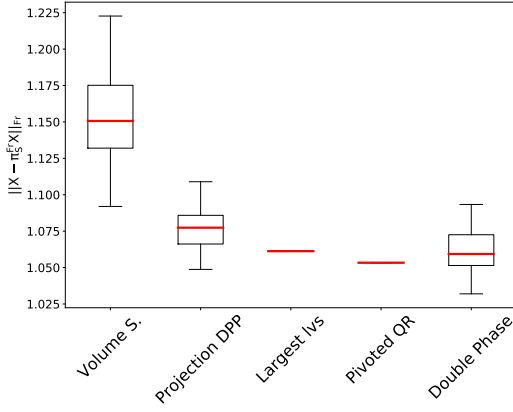




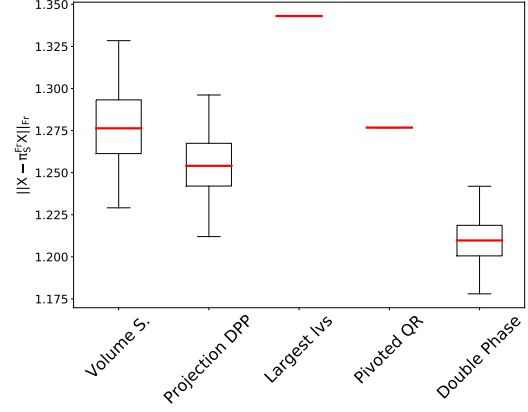
(a)  $k$ -leverage scores profile for the dataset Basehock ( $k=10$ ).



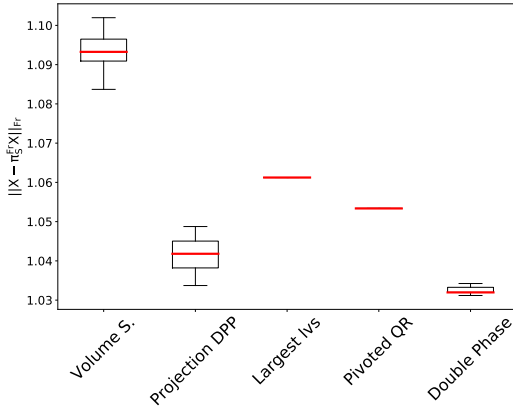
(b)  $k$ -leverage scores profile for the dataset Colon ( $k=10$ ).



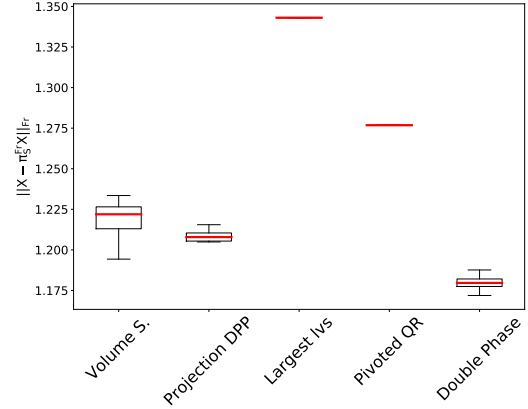
(c) Boxplots of  $\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}$  on a batch of 50 samples for the five algorithms on the dataset Basehock ( $k=10$ ).



(d) Boxplots of  $\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}$  on a batch of 50 samples for the five algorithms on the dataset Colon ( $k=10$ ).

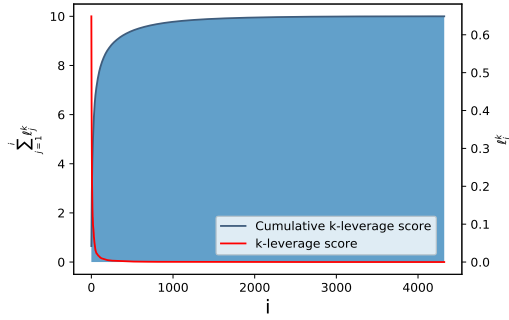


(e) Boxplots of  $\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}$  on a batch of 50 samples for the boosting of randomized algorithms on the dataset Basehock ( $k=10$ ).

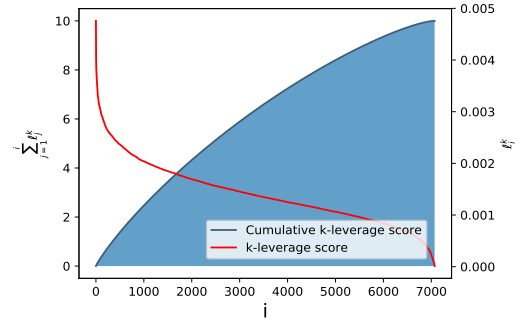


(f) Boxplots of  $\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}$  on a batch of 50 samples for the boosting of randomized algorithms on the dataset Colon ( $k=10$ ).

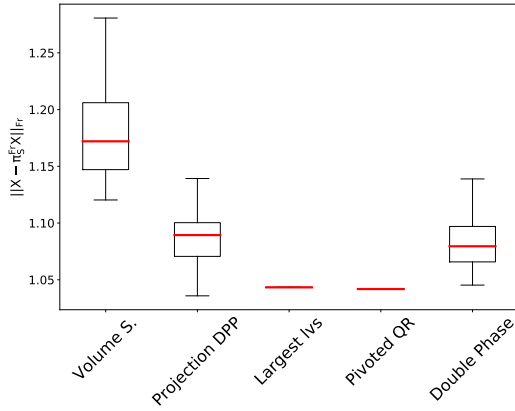
Figure 7: Comparison of several column subset selection algorithms for two datasets: Basehock and Colon.



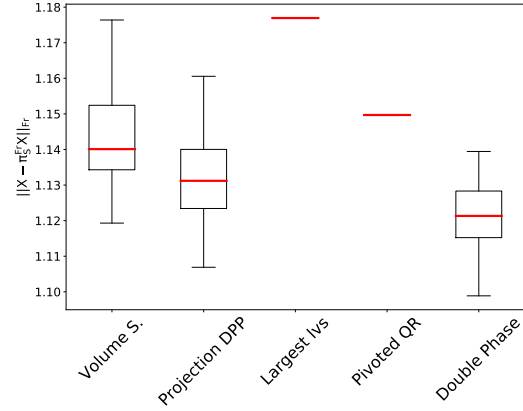
(a)  $k$ -leverage scores profile for the dataset Relathe ( $k=10$ ).



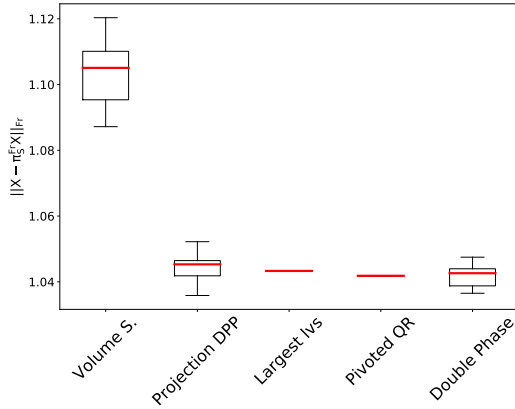
(b)  $k$ -leverage scores profile for the dataset Leukemia ( $k=10$ ).



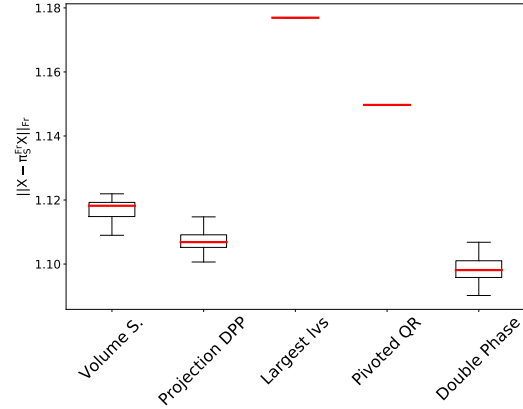
(c) Boxplots of  $\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}$  on a batch of 50 samples for the five algorithms on the dataset Relathe ( $k=10$ ).



(d) Boxplots of  $\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}$  on a batch of 50 samples for the five algorithms on the dataset Leukemia ( $k=10$ ).



(e) Boxplots of  $\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}$  on a batch of 50 samples for the boosting of randomized algorithms on the dataset Relathe ( $k=10$ ).



(f) Boxplots of  $\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}$  on a batch of 50 samples for the boosting of randomized algorithms on the dataset Leukemia ( $k=10$ ).

Figure 8: Comparison of several column subset selection algorithms for two datasets: Relathe and Leukemia.

localized the  $k$ -leverage scores, the smaller the excess risk bounds. Such an analysis of the excess risk in regression highlights the interest of the proposed approach since it would be difficult to conduct for both volume sampling or the double phase algorithms. Future work in this direction includes investigating the importance of the sparsity of the  $k$ -leverage scores on the performance of other learning algorithms such as spectral clustering or support vector machines.

Finally, in our experimental section, we used an adhoc randomized algorithm inspired by (Fickus et al., 2011b) to sample toy datasets with a prescribed profile of  $k$ -leverage scores. An interesting question would be to characterize the distribution of the output of our algorithm. In particular, sampling from the uniform measure on the set of symmetric matrices with prescribed spectrum and leverage scores is still an open problem (Dhillon et al., 2005).

## Acknowledgments

AB and RB acknowledge support from ANR grant BoB (ANR-16-CE23-0003), and all authors acknowledge support from ANR grant BNPSI (ANR-13-BS03-0006).

## References

- U. Alon, N. Barkai, D. A. Notterman, K. Gish, S. Ybarra, D. Mack, and A. J. Levine. Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proceedings of the National Academy of Sciences*, 96(12):6745–6750, 1999.
- H. Avron and C. Boutsidis. Faster subset selection for matrices and applications. *SIAM Journal on Matrix Analysis and Applications*, 34(4):1464–1499, 2013.
- R. Bardenet and A. Hardy. Monte Carlo with Determinantal Point Processes. *ArXiv e-prints*, May 2016.
- Y. Baryshnikov. GUEs and queues. *Probability Theory and Related Fields*, 119(2): 256–274, 2001.
- J. D. Batson, D. A. Spielman, and N. Srivastava. Twice-ramanujan sparsifiers. In *Proceedings of the Forty-first Annual ACM Symposium on Theory of Computing*, STOC '09, pages 255–262, New York, NY, USA, 2009. ACM. ISBN 978-1-60558-506-2. doi: 10.1145/1536414.1536451. URL <http://doi.acm.org/10.1145/1536414.1536451>.
- Adi Ben-Israel. A volume associated with  $m \times n$  matrices. *Linear Algebra and its Applications*, 167:87 – 111, 1992. ISSN 0024-3795. doi: [http://dx.doi.org/10.1016/0024-3795\(92\)90340-G](http://dx.doi.org/10.1016/0024-3795(92)90340-G). URL <http://www.sciencedirect.com/science/article/pii/002437959290340G>.
- Å. Björck and G. H. Golub. Numerical methods for computing angles between linear subspaces. *Mathematics of computation*, 27(123):579–594, 1973.
- C. Boutsidis, M. W. Mahoney, and P. Drineas. An improved approximation algorithm for the column subset selection problem. In *Proceedings of the Twentieth*

- Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '09, pages 968–977, Philadelphia, PA, USA, 2009. Society for Industrial and Applied Mathematics. URL <http://dl.acm.org/citation.cfm?id=1496770.1496875>.
- C. Boutsidis, P. Drineas, and M. Magdon-Ismail. Near optimal column-based matrix reconstruction. In *Proceedings of the 2011 IEEE 52Nd Annual Symposium on Foundations of Computer Science*, FOCS '11, pages 305–314, Washington, DC, USA, 2011. IEEE Computer Society. ISBN 978-0-7695-4571-4. doi: 10.1109/FOCS.2011.21. URL <http://dx.doi.org/10.1109/FOCS.2011.21>.
- M. Dereziński and M. K. Warmuth. Reverse iterative volume sampling for linear regression. *arXiv preprint arXiv:1806.01969*, 2018.
- A. Deshpande and L. Rademacher. Efficient volume sampling for row/column subset selection. In *Proceedings of the 2010 IEEE 51st Annual Symposium on Foundations of Computer Science*, FOCS '10, pages 329–338, Washington, DC, USA, 2010. IEEE Computer Society. ISBN 978-0-7695-4244-7. doi: 10.1109/FOCS.2010.38. URL <http://dx.doi.org/10.1109/FOCS.2010.38>.
- A. Deshpande and S. Vempala. Adaptive sampling and fast low-rank matrix approximation. In *Proceedings of the 9th International Conference on Approximation Algorithms for Combinatorial Optimization Problems, and 10th International Conference on Randomization and Computation*, APPROX'06/RANDOM'06, pages 292–303, Berlin, Heidelberg, 2006. Springer-Verlag. ISBN 3-540-38044-2, 978-3-540-38044-3. doi: 10.1007/11830924\_28. URL [http://dx.doi.org/10.1007/11830924\\_28](http://dx.doi.org/10.1007/11830924_28).
- A. Deshpande, L. Rademacher, S. Vempala, and G. Wang. Matrix approximation and projective clustering via volume sampling. In *Proceedings of the Seventeenth Annual ACM-SIAM Symposium on Discrete Algorithm*, SODA '06, pages 1117–1126, Philadelphia, PA, USA, 2006. Society for Industrial and Applied Mathematics. ISBN 0-89871-605-5. URL <http://dl.acm.org/citation.cfm?id=1109557.1109681>.
- I. Dhillon, R. Heath, M. Sustik, and J. Tropp. Generalized finite algorithms for constructing hermitian matrices with prescribed diagonal and spectrum. *SIAM Journal on Matrix Analysis and Applications*, 27(1):61–71, 2005. doi: 10.1137/S0895479803438183. URL <https://doi.org/10.1137/S0895479803438183>.
- P. Drineas, A. Frieze, R. Kannan, S. Vempala, and V. Vinay. Clustering large graphs via the singular value decomposition. *Mach. Learn.*, 56(1-3):9–33, June 2004. ISSN 0885-6125. doi: 10.1023/B:MACH.0000033113.59016.96. URL <https://doi.org/10.1023/B:MACH.0000033113.59016.96>.
- P. Drineas, M. W. Mahoney, and S. Muthukrishnan. Relative-error CUR matrix decompositions. *SIAM Journal on Matrix Analysis and Applications*, 30(2):844–881, 2008.
- P. Drineas, M. Magdon-Ismail, M. W. Mahoney, and D. P. Woodruff. Fast approximation of matrix coherence and statistical leverage. *Journal of Machine Learning Research*, 13(Dec):3475–3506, 2012.

- M. Fickus, D. G. Mixon, and M. J. Poteet. Frame completions for optimally robust reconstruction. 2011a. doi: 10.1117/12.891813.
- M. Fickus, D. G. Mixon, M. J. Poteet, and N. Strawn. Constructing all self-adjoint matrices with prescribed spectrum and diagonal, 2011b.
- G. Gautier, R. Bardenet, and M. Valko. DPPy: Sampling determinantal point processes with Python. *arXiv preprint arXiv:1809.07258*, 2018.
- G. H. Golub. Numerical methods for solving linear least squares problems. *Numer. Math.*, 7(3):206–216, June 1965. ISSN 0029-599X. doi: 10.1007/BF01436075. URL <http://dx.doi.org/10.1007/BF01436075>.
- G. H. Golub and C. F. Van Loan. *Matrix Computations (3rd Ed.)*. Johns Hopkins University Press, Baltimore, MD, USA, 1996. ISBN 0-8018-5414-8.
- T. R. Golub, D. K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J. P. Mesirov, H. Coller, M. L. Loh, J. R. Downing, M. A. Caligiuri, et al. Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *science*, 286(5439):531–537, 1999.
- M. Gu and S. C. Eisenstat. Efficient algorithms for computing a strong rank-revealing QR factorization. *SIAM J. Sci. Comput.*, 17(4):848–869, July 1996. ISSN 1064-8275. doi: 10.1137/0917055. URL <http://dx.doi.org/10.1137/0917055>.
- V. Guruswami and A. K. Sinop. Optimal column-based low-rank matrix reconstruction. In *Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete Algorithms*, pages 1207–1214. SIAM, 2012.
- A. Horn. Doubly stochastic matrices and the diagonal of a rotation matrix. *American Journal of Mathematics*, 76(3):620–630, 1954.
- J. B. Hough, M. Krishnapur, Y. Peres, and B. Virág. Determinantal processes and independence. 2005. doi: 10.1214/154957806000000078.
- K. Johansson. Random matrices and determinantal processes. *ArXiv Mathematical Physics e-prints*, October 2005.
- A. Kulesza and B. Taskar. k-DPPs: Fixed-size determinantal point processes. In Lise Getoor and Tobias Scheffer, editors, *ICML*, pages 1193–1200. Omnipress, 2011. URL <http://dblp.uni-trier.de/db/conf/icml/icml2011.html#KuleszaT11>.
- A. Kulesza, B. Taskar, et al. Determinantal point processes for machine learning. *Foundations and Trends® in Machine Learning*, 5(2–3):123–286, 2012.
- F. Lavancier, J. Møller, and E. Rubak. Determinantal point process models and statistical inference. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 77(4):853–877, 2015.
- C. Li, S. Jegelka, and S. Sra. Polynomial time algorithms for dual volume sampling. In *Advances in Neural Information Processing Systems*, pages 5038–5047, 2017a.

- J. Li, K. Cheng, S. Wang, F. Morstatter, R. P. Trevino, J. Tang, and H. Liu. Feature selection: A data perspective. *ACM Computing Surveys (CSUR)*, 50(6):94, 2017b.
- P. Ma, M. W. Mahoney, and B. Yu. A statistical perspective on algorithmic leveraging. *The Journal of Machine Learning Research*, 16(1):861–911, 2015.
- O Macchi. The coincidence approach to stochastic point processes. 7:83–122, 03 1975.
- A. W. Marshall, I. Olkin, and B. C. Arnold. *Inequalities: Theory of Majorization and its Applications*, volume 143. Springer, second edition, 2011. doi: 10.1007/978-0-387-68276-1.
- L. Mor-Yosef and H. Avron. Sketching for principal component regression, 2018.
- D. Papailiopoulos, A. Kyrillidis, and C. Boutsidis. Provable deterministic leverage score sampling. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14*, pages 997–1006, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-2956-9. doi: 10.1145/2623330.2623698. URL <http://doi.acm.org/10.1145/2623330.2623698>.
- G. Raskutti and M. W. Mahoney. A statistical perspective on randomized sketching for ordinary least-squares. *The Journal of Machine Learning Research*, 17(1):7508–7538, 2016.
- M. Slawski. On principal components regression, random projections, and column subsampling. *Electronic Journal of Statistics*, 12(2):3673–3712, 2018.
- A. Soshnikov. Determinantal random point fields. *Russian Mathematical Surveys*, 55:923–975, October 2000. doi: 10.1070/RM2000v055n05ABEH000321.
- N. Tremblay, S. Barthelmé, and P.-O. Amblard. Optimized algorithms to sample determinantal point processes. *arXiv preprint arXiv:1802.08471*, 2018.

## A Another interpretation of the $k$ -leverage scores

For  $i \in [d]$ , the SVD of  $\mathbf{X}$  yields

$$\mathbf{X}_{:,i} = \sum_{\ell=1}^r V_{i,\ell} \mathbf{f}_\ell, \quad (55)$$

where  $\mathbf{f}_\ell = \sigma_\ell \mathbf{U}_{:, \ell}$ ,  $\ell \in [r]$ , are orthogonal. Thus

$$\mathbf{X}_{:,i}^\top \mathbf{f}_j = V_{i,j} \|\mathbf{f}_j\|^2 = V_{i,j} \sigma_j^2. \quad (56)$$

Then

$$\frac{V_{i,j}}{\|\mathbf{X}_{:,i}\|} = \frac{\mathbf{X}_{:,i}^\top \mathbf{f}_j}{\sigma_j \|\mathbf{X}_{:,i}\| \|\mathbf{f}_j\|} =: \frac{\cos \eta_{i,j}}{\sigma_j}, \quad (57)$$

where  $\eta_{i,j} \in [0, \pi/2]$  is the angle formed by  $\mathbf{X}_{:,i}$  and  $\mathbf{f}_j$ . Finally, (56) also yields

$$\ell_i^k = \|\mathbf{X}_{:,i}\|^2 \sum_{j=1}^k \frac{\cos^2 \eta_{i,j}}{\sigma_j^2}. \quad (58)$$

Compared to the length-square distribution in Section 3.2,  $k$ -leverage scores thus favour columns that are aligned with the principal features. The weight  $1/\sigma_j^2$  corrects the fact that features associated with large singular values are typically aligned with more columns. One could also imagine more arbitrary weights  $w_j/\sigma_j^2$  in lieu of  $1/\sigma_j^2$ , or, equivalently, modified  $k$ -leverage scores

$$\ell_i^k(\mathbf{w}) = \sum_{j=1}^k w_j V_{i,j}^2.$$

However, the projection DPP with marginal kernel  $\mathbf{K} = \mathbf{V}_k \mathbf{V}_k^\top$  that we study in this paper is invariant to such reweightings. Indeed, for any  $S \subset [d]$  of cardinality  $k$ ,

$$\text{Det} \left[ \mathbf{V}_{S,[k]} \text{Diag}(\mathbf{w}_{[k]}) \mathbf{V}_{[k],S}^\top \right] = \text{Det}(\mathbf{V}_{S,[k]})^2 \prod_{j \in [k]} w_j^2 \propto \text{Det}(\mathbf{V}_{S,[k]})^2. \quad (59)$$

Such a scaling is thus not a free parameter in  $\mathbf{K}$ .

## B Majorization and Schur convexity

This section recalls some definitions and results from the theory of majorization and the notions of Schur-convexity and Schur-concavity. We refer to (Marshall et al., 2011) for further details. In this section, a subset  $\mathcal{D} \subset \mathbb{R}^d$  is a symmetric domain if  $\mathcal{D}$  is stable under coordinate permutations. Furthermore, a function  $f$  defined on a symmetric domain  $\mathcal{D}$  is called symmetric if it is stable under coordinate permutations.

**Definition 21** Let  $\mathbf{p}, \mathbf{q} \in \mathbb{R}_+^d$ .  $\mathbf{p}$  is said to majorize  $\mathbf{q}$  according to Schur order and we note  $\mathbf{q} \prec_S \mathbf{p}$  if

$$\begin{cases} q_{i_1} \leq p_{j_1} \\ q_{i_1} + q_{i_2} \leq p_{j_1} + p_{j_2} \\ \dots \\ \sum_{k=1}^{d-1} q_{i_k} \leq \sum_{k=1}^{d-1} p_{j_k} \\ \sum_{k=1}^d q_{i_k} = \sum_{k=1}^d p_{j_k} \end{cases} \quad (60)$$

where  $\mathbf{p}, \mathbf{q}$  are reordered so that  $p_{i_d} \leq \dots \leq p_{i_1}$  and  $q_{j_d} \leq \dots \leq q_{j_1}$ .

The majorization order has an algebraic characterization using doubly stochastic matrices first proven by Hardy, Littlewood, and Polya in 1929.

**Proposition 22 (Theorem B.2, Marshall et al., 2011)** The vector  $\mathbf{p}$  majorizes the vector  $\mathbf{q}$  if and only if there exists a  $d \times d$  doubly stochastic matrix  $\Pi$  such that  $\mathbf{q} = \mathbf{p}\Pi$ .

**Example 1** Let  $\mathbf{p} = (3, 0, 0)$  and  $\mathbf{q} = (1, 1, 1)$ . We check easily that  $\mathbf{p}$  majorizes  $\mathbf{q}$ . Note that we can 'redistribute'  $\mathbf{p}$  over  $\mathbf{q}$  as follows:  $\mathbf{q} = \frac{1}{3}\mathbf{J}\mathbf{p}$ , where  $\mathbf{J}$  is a  $3 \times 3$  matrix of ones. The matrix  $\Pi = \frac{1}{3}\mathbf{J}$  is a doubly stochastic matrix.

Schur order compares two vectors using multiple inequalities. To avoid such cumbersome calculations, a scalar metric of inequality in a vector is desired. This is possible using the notion of Schur-convex/concave function.

**Definition 23** Let  $f$  be a function on a symmetric domain  $\mathcal{D} \subset \mathbb{R}_+^d$ .  $f$  is said to be Schur convex if

$$\forall \mathbf{p}, \mathbf{q} \in \mathbb{R}_+^d, \mathbf{q} \prec_S \mathbf{p} \implies f(\mathbf{q}) \leq f(\mathbf{p}). \quad (61)$$

$f$  is said to be Schur concave if

$$\forall \mathbf{p}, \mathbf{q} \in \mathbb{R}_+^d, \mathbf{q} \prec_S \mathbf{p} \implies f(\mathbf{q}) \geq f(\mathbf{p}). \quad (62)$$

**Proposition 24 (Theorem A.3, Marshall et al., 2011)** Let  $f$  be a symmetric function defined on  $\mathbb{R}_+^d$ , let  $\mathcal{D}$  be a permutation-symmetric domain in  $\mathbb{R}_+^d$  and suppose that

$$\forall x_i, x_j \in \mathbb{R}_+, (x_i - x_j) \left( \frac{\partial f}{\partial x_i} - \frac{\partial f}{\partial x_j} \right) > 0 \quad (63)$$

then

$$\forall \mathbf{p}, \mathbf{q} \in \mathcal{D}, \mathbf{q} \prec_S \mathbf{p} \implies f(\mathbf{q}) \leq f(\mathbf{p}), \quad (64)$$

and  $f$  is Schur convex.

We get a similar result for Schur concavity by switching the orders in the previous proposition.

**Theorem 25 (Theorem 3.1, Guruswami and Sinop, 2012)** Let  $\mathbf{X} \in \mathbb{R}^{N \times d}$ , and let  $\boldsymbol{\sigma} \in \mathbb{R}^d$  the vector containing the squares of the singular values of  $\mathbf{X}$ . The function

$$\boldsymbol{\sigma} \mapsto \mathbb{E}_{\text{VS}} \|\mathbf{X} - \Pi_S \mathbf{X}\|_{\text{Fr}}^2 = (k+1) \frac{e_k(\boldsymbol{\sigma})}{e_{k-1}(\boldsymbol{\sigma})} \quad (65)$$

is Schur-concave.



## C Principal angles and the Cosine Sine decomposition

### C.1 Principal angles

This section surveys the notion of principal angles between subspaces, see (Golub and Van Loan, 1996, Section 6.4.3) for details.

**Definition 26** Let  $\mathcal{P}, \mathcal{Q}$  be two subspaces in  $\mathbb{R}^d$ . Let  $p = \dim \mathcal{P}$  and  $q = \dim \mathcal{Q}$  and assume that  $q \leq p$ . To define the vector of principal angles  $\boldsymbol{\theta} \in [0, \pi/2]^q$  between  $\mathcal{P}$  and  $\mathcal{Q}$ , let

$$\cos(\theta_1) = \max \left\{ \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}; \quad \mathbf{x} \in \mathcal{P}, \mathbf{y} \in \mathcal{Q} \right\} \quad (66)$$

be the cosine of the smallest angle between a vector of  $\mathcal{P}$  and a vector of  $\mathcal{Q}$ , and let  $(\mathbf{x}_1, \mathbf{y}_1) \in \mathcal{P} \times \mathcal{Q}$  be a pair of vectors realizing the maximum. For  $i \in [2, q]$ , define successively

$$\cos(\theta_i) = \max \left\{ \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}; \quad \mathbf{x} \in \mathcal{P}, \mathbf{y} \in \mathcal{Q}; \mathbf{x} \perp \mathbf{x}_j, \mathbf{y} \perp \mathbf{y}_j, \forall j \in [1 : i - 1] \right\} \quad (67)$$

and denote  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{P} \times \mathcal{Q}$  such that  $\cos(\theta_i) = \mathbf{x}_i^T \mathbf{y}_i$ .

Note that although the so-called principal vectors  $(\mathbf{x}_i, \mathbf{y}_i)_{i \in [q]}$  are not uniquely defined by (66) and (67), the principal angles  $\boldsymbol{\theta}$  are uniquely defined, see (Björck and Golub, 1973). The following result confirms this, while also providing a way to compute  $\boldsymbol{\theta}$ .

**Proposition 27 (Björck and Golub, 1973, Ben-Israel, 1992)** Let  $\mathcal{P}$  and  $\mathcal{Q}$  and  $\boldsymbol{\theta}$  be as in Definition 26. Let  $\mathbf{P} \in \mathbb{R}^{d \times p}$ ,  $\mathbf{Q} \in \mathbb{R}^{d \times q}$  be two orthogonal matrices, whose columns are orthonormal bases of  $\mathcal{P}$  and  $\mathcal{Q}$ , respectively. Then

$$\forall i \in [q], \quad \cos(\theta_i) = \sigma_i(\mathbf{Q}^T \mathbf{P}). \quad (68)$$

In particular

$$\text{Vol}_q^2(\mathbf{Q}^T \mathbf{P}) = \prod_{i \in [q]} \cos^2(\theta_i). \quad (69)$$

An important case for our work arises when  $q = k$ ,  $\mathbf{Q} = \mathbf{V} \in \mathbb{R}^{d \times k}$ , and  $\mathbf{P} = \mathbf{S} \in \mathbb{R}^{d \times k}$  is a sampling matrix. The left-hand side of (69) then equals  $\text{Det}(\mathbf{V}_{:,S})^2$ .

### C.2 The Cosine Sine decomposition

The Cosine Sine (CS) decomposition is useful for the study of the relative position of two subspaces. It generalizes the notion of cosine, sine and tangent to subspaces.

**Proposition 28 (Golub and Van Loan, 1996)** Let  $q \leq d/2$  and  $\mathbf{Q} = \begin{bmatrix} \mathbf{Q}_1 \\ \mathbf{Q}_2 \end{bmatrix}$  be a  $d \times q$  orthogonal matrix, where  $\mathbf{Q}_1 \in \mathbb{R}^{q \times q}$  and  $\mathbf{Q}_2 \in \mathbb{R}^{(d-q) \times q}$ . Assume that  $\mathbf{Q}_1$  is non singular, then there exist orthogonal matrices  $\mathbf{Y} \in \mathbb{R}^{d \times q}$  and

$$\mathbf{W} = \left[ \begin{array}{c|c} \mathbf{W}_1 & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{W}_2 \end{array} \right] \in \mathbb{R}^{d \times d}, \quad (70)$$

and a matrix

$$\boldsymbol{\Sigma} = \begin{bmatrix} \mathcal{C} \\ \mathcal{S} \\ \mathbf{0} \end{bmatrix} \in \mathbb{R}^{d \times q}, \quad (71)$$

such that

$$\mathbf{Q} = \mathbf{W} \boldsymbol{\Sigma} \mathbf{Y}^T, \quad (72)$$

where  $\mathbf{W}_1 \in \mathbb{R}^{q \times q}$  and  $\mathbf{W}_2 \in \mathbb{R}^{d-q \times d-q}$ , and  $\mathcal{C}, \mathcal{S} \in \mathbb{R}^{q \times q}$  are diagonal matrices satisfying the identity  $\mathcal{C}^2 + \mathcal{S}^2 = \mathbb{1}_q$ . In particular, each block  $\mathbf{Q}_i$  factorizes as

$$\begin{aligned} \mathbf{Q}_1 &= \mathbf{W}_1 \mathcal{C} \mathbf{Y}^T \\ \mathbf{Q}_2 &= \mathbf{W}_2 \begin{bmatrix} \mathcal{S} \\ \mathbf{0} \end{bmatrix} \mathbf{Y}^T. \end{aligned} \quad (73)$$

The CS decomposition is defined for every orthogonal matrix. An important case is when  $\mathbf{Q}$  is the product of an orthogonal matrix  $\mathbf{V} \in \mathbb{R}^{d \times d}$  and a sampling matrix  $\mathbf{S} \in \mathbb{R}^{d \times k}$ , that is  $\mathbf{Q} = \mathbf{V}^\top \mathbf{S}$ .

**Corollary 29** *Let  $\mathbf{V} \in \mathbb{R}^{d \times d}$  be an orthogonal matrix and  $\mathbf{S} \in \mathbb{R}^{d \times k}$  be a sampling matrix. Let*

$$\mathbf{Q} = \mathbf{V}^\top \mathbf{S} = \begin{bmatrix} \mathbf{V}_k^\top \mathbf{S} \\ \mathbf{V}_{d-k}^\top \mathbf{S} \end{bmatrix} \quad (74)$$

be a  $d \times k$  orthogonal matrix, with  $\text{Det}(\mathbf{V}_k^\top \mathbf{S})^2 > 0$ . Let further  $\mathbf{Z}_S = \mathbf{V}_{d-k}^\top \mathbf{S} (\mathbf{V}_k^\top \mathbf{S})^{-1}$ . Then

$$\text{Tr}(\mathbf{Z}_S \mathbf{Z}_S^\top) \leq \sum_{i \in [k]} \tan^2(\theta_i(S)). \quad (75)$$

**Proof** In the case  $k \leq d/2$ , Proposition 28 applied to the matrix  $\mathbf{Q} = \mathbf{V}^\top \mathbf{S}$  with  $\mathbf{Q}_1 = \mathbf{V}_k^\top \mathbf{S}$  and  $\mathbf{Q}_2 = \mathbf{V}_{d-k}^\top \mathbf{S}$  yields

$$\mathbf{Q}_1 = \mathbf{W}_1 \mathcal{C} \mathbf{Y}^T \quad (76)$$

$$\mathbf{Q}_2 = \mathbf{W}_2 \begin{bmatrix} \mathcal{S} \\ \mathbf{0} \end{bmatrix} \mathbf{Y}^T. \quad (77)$$

Thus, the diagonal matrix  $\mathcal{C}$  contains the singular values of the matrix  $\mathbf{V}_k^\top \mathbf{S}$  that are cosines of the principal angles  $(\theta_i(S))_{i \in [k]}$  between  $\text{Span}(\mathbf{V}_k)$  and  $\text{Span}(\mathbf{S})$  thanks to Proposition 27.

The identity  $\mathcal{C}^2 + \mathcal{S}^2 = \mathbb{1}_k$  and the fact that  $\theta_i(S) \in [0, \frac{\pi}{2}]$  imply that the (diagonal) elements of  $\mathcal{S}$  are equal to the sines of the principal angles between  $\text{Span}(\mathbf{V}_k)$  and  $\text{Span}(\mathbf{S})$ . Let  $\mathcal{T} = \mathcal{S} \mathcal{C}^{-1}$ .  $\mathcal{T} \in \mathbb{R}^{k \times k}$  is a diagonal matrix containing the tangents of the principal angles  $(\theta_i(S))_{i \in [k]}$ . Using (76), we get

$$\mathbf{Z}_S = \mathbf{V}_{d-k}^\top \mathbf{S} (\mathbf{V}_k^\top \mathbf{S})^{-1} = \mathbf{W}_2 \begin{bmatrix} \mathcal{S} \\ \mathbf{0} \end{bmatrix} \mathbf{Y}^\top \mathbf{Y} \mathcal{C}^{-1} \mathbf{W}_1^\top = \mathbf{W}_2 \begin{bmatrix} \mathcal{S} \\ \mathbf{0} \end{bmatrix} \mathcal{C}^{-1} \mathbf{W}_1^\top = \mathbf{W}_2 \begin{bmatrix} \mathcal{S} \mathcal{C}^{-1} \\ \mathbf{0} \end{bmatrix} \mathbf{W}_1^\top. \quad (78)$$

Then,

$$\text{Tr}(\mathbf{Z}_S \mathbf{Z}_S^\top) = \text{Tr}(\mathbf{W}_2 \begin{bmatrix} \mathcal{T}^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{W}_2^\top) = \sum_{i \in [k]} \tan^2(\theta_i(S)). \quad (79)$$

■

## D Proofs

### D.1 Technical lemmas

We start with two useful lemmas borrowed from the literature.

**Lemma 30 (Lemma 3.1, Boutsidis et al., 2011)** *Let  $S \subset [d]$ , then*

$$\|\mathbf{X} - \Pi_{S,k}^\nu \mathbf{X}\|_\nu^2 \leq \|\mathbf{E}(\mathbf{I} - \mathbf{P}_S)\|_\nu^2, \quad \nu \in \{2, \text{Fr}\}, \quad (80)$$

where  $\mathbf{E} = \mathbf{X} - \Pi_k \mathbf{X}$  and  $\mathbf{P}_S = \mathbf{S}(\mathbf{V}_k^\top \mathbf{S})^{-1} \mathbf{V}_k^\top$ . Furthermore,

$$\|\mathbf{X} - \Pi_{S,k}^\nu \mathbf{X}\|_\nu^2 \leq \frac{1}{\sigma_k^2(\mathbf{V}_{S,[k]})} \|\mathbf{X} - \Pi_k \mathbf{X}\|_\nu^2, \quad \nu \in \{2, \text{Fr}\}. \quad (81)$$

The following lemma was first proven by [Deshpande et al., 2006](#), and later rephrased.

**Lemma 31 (Lemma 11, Deshpande and Rademacher, 2010)** *Let  $\mathbf{V} \in \mathbb{R}^{k \times d}$ ,  $r = \text{rk}(\mathbf{V})$  and  $\ell \in [1 : r]$ . Then*

$$\sum_{S \subset [d], |S|=\ell} e_\ell(\Sigma(\mathbf{V}_{:,S})^2) = e_\ell(\Sigma(\mathbf{V})^2) \quad (82)$$

where  $e_\ell$  is the  $\ell$ -th elementary symmetric polynomial on  $r$  variables, see [Section 2](#).

Elementary symmetric polynomials play an important role in the proof of [Proposition 18](#), in particular their interplay with the Schur order; see [Appendix B](#) for definitions.

**Lemma 32** *Let  $\phi, \psi : \mathbb{R}_+^d \rightarrow \mathbb{R}_+$  be defined by*

$$\phi : \boldsymbol{\sigma} \mapsto \frac{e_{k-1}(\boldsymbol{\sigma})}{e_k(\boldsymbol{\sigma})} \quad (83)$$

and

$$\psi : \boldsymbol{\sigma} \mapsto e_k(\boldsymbol{\sigma}). \quad (84)$$

Then both functions are symmetric,  $\phi$  is Schur-convex, and  $\psi$  is Schur-concave.

**Proof** [of [Lemma 32](#)] Let  $i, j \in [r], i \neq j$ . Let  $\sigma_i, \sigma_j \in \mathbb{R}_+$ , it holds

$$\begin{aligned} (\sigma_i - \sigma_j)(\partial_i \phi(\boldsymbol{\sigma}) - \partial_j \phi(\boldsymbol{\sigma})) &= (\sigma_i - \sigma_j) \left( -\frac{1}{\sigma_i^2} + \frac{1}{\sigma_j^2} \right) \\ &= \frac{(\sigma_i - \sigma_j)^2 (\sigma_i + \sigma_j)}{\sigma_i^2 \sigma_j^2} \geq 0, \end{aligned}$$

so that  $\phi$  is Schur-convex by [Proposition 24](#). Similarly,

$$\begin{aligned} (\sigma_i - \sigma_j)(\partial_i \psi(\boldsymbol{\sigma}) - \partial_j \psi(\boldsymbol{\sigma})) &= (\sigma_i - \sigma_j) \left( \prod_{\ell \neq i} \sigma_\ell - \prod_{\ell \neq j} \sigma_\ell \right) \\ &= -(\sigma_i - \sigma_j)^2 \prod_{\ell \neq i, j} \sigma_\ell \geq 0, \end{aligned}$$

so that  $\psi$  is Schur-concave by [Proposition 24](#). ■

Elementary symmetric polynomials also interact nicely with “marginalizing” sums.

**Lemma 33** Let  $\mathbf{V}$  be a real  $k \times d$  matrix and let  $r = \text{rk}(\mathbf{V})$ . Denote by  $p$  the number of non zero columns of  $\mathbf{V}$ . Then for all  $k \leq r + 1$ ,

$$\sum_{\substack{S \subset [d], |S|=k \\ \text{Vol}_k(\mathbf{V}_{:,S})^2 > 0}} \sum_{\substack{T \subset [S] \\ |T|=k-1}} e_{k-1}(\Sigma(\mathbf{V}_{:,T})^2) \leq (p - k + 1)e_{k-1}(\Sigma(\mathbf{V})^2). \quad (85)$$

*A fortiori*,

$$\sum_{\substack{S \subset [d], |S|=k \\ \text{Vol}_k(\mathbf{V}_{:,S})^2 > 0}} \sum_{\substack{T \subset [S] \\ |T|=k-1}} e_{k-1}(\Sigma(\mathbf{V}_{:,T})^2) \leq (d - k + 1)e_{k-1}(\Sigma(\mathbf{V})^2). \quad (86)$$

**Proof** [of Lemma 33] For  $T \subset [d]$ ,  $|T| = k - 1$ ,

$$\begin{aligned} \Omega_1(T) &= \{S \subset [d] : |S| = k, T \subset S, \forall i \in S, \mathbf{V}_{:,i} \neq \mathbf{0}\} \\ \Omega_2(T) &= \{S \subset [d] : |S| = k, T \subset S, \text{Vol}_k(\mathbf{V}_{:,S})^2 > 0\}. \end{aligned}$$

Note that  $\Omega_2(T) \subset \Omega_1(T)$  so that

$$\begin{aligned} \sum_{\substack{S \subset [d], |S|=k \\ \text{Vol}_k(\mathbf{V}_{:,S})^2 > 0}} \sum_{\substack{T \subset S \\ |T|=k-1}} e_{k-1}(\Sigma(\mathbf{V}_{:,T})^2) &= \sum_{\substack{T \subset [d] \\ |T|=k-1}} \sum_{S \in \Omega_2(T)} e_{k-1}(\Sigma(\mathbf{V}_{:,T})^2) \\ &\leq \sum_{\substack{T \subset [d] \\ |T|=k-1}} \sum_{S \in \Omega_1(T)} e_{k-1}(\Sigma(\mathbf{V}_{:,T})^2). \end{aligned}$$

The set  $\Omega_1(T)$  has at most  $(p - k + 1)$  elements so that

$$\sum_{\substack{T \subset [d] \\ |T|=k-1}} \sum_{S \in \Omega_1(T)} e_{k-1}(\Sigma(\mathbf{V}_{:,T})^2) \leq (p - k + 1) \sum_{\substack{T \subset [d] \\ |T|=k-1}} e_{k-1}(\Sigma(\mathbf{V}_{:,T})^2). \quad (87)$$

Lemma 31 for  $\ell = k - 1$  further yields

$$(p - k + 1) \sum_{\substack{T \subset [d] \\ |T|=k-1}} e_{k-1}(\Sigma(\mathbf{V}_{:,T})^2) \leq (p - k + 1) e_{k-1}(\Sigma(\mathbf{V})^2). \quad (88)$$

■

## D.2 Proof of Proposition 16

First, Lemma 30 yields

$$\begin{aligned} \sum_{S \subset [d], |S|=k} \text{Det}(\mathbf{V}_{S,[k]})^2 \|\mathbf{X} - \Pi_S^\nu \mathbf{X}\|_\nu^2 &\leq \sum_{S \subset [d], |S|=k} \frac{1}{\sigma_k^2(\mathbf{V}_{S,[k]})} \text{Det}(\mathbf{V}_{S,[k]})^2 \|\mathbf{X} - \Pi_k \mathbf{X}\|_\nu^2 \\ &= \|\mathbf{X} - \Pi_k \mathbf{X}\|_\nu^2 \sum_{S \subset [d], |S|=k} \prod_{\ell=1}^{k-1} \sigma_\ell^2(\mathbf{V}_{S,[k]}), \quad (89) \end{aligned}$$

where the last equality follows from

$$\text{Det}(\mathbf{V}_{S,[k]})^2 = \prod_{\ell=1}^k \sigma_{\ell}^2(\mathbf{V}_{S,[k]}). \quad (90)$$

By definition of the polynomial  $e_{k-1}$ , it further holds

$$\prod_{\ell=1}^{k-1} \sigma_{\ell}^2(\mathbf{V}_{S,[k]}) \leq e_{k-1}(\Sigma(\mathbf{V}_{S,[k]})^2), \quad (91)$$

so that (89) leads to

$$\sum_{S \subset [d], |S|=k} \text{Det}(\mathbf{V}_{S,[k]})^2 \|\mathbf{X} - \Pi_S^{\nu} \mathbf{X}\|_{\nu}^2 \leq \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\nu}^2 \sum_{S \subset [d], |S|=k} e_{k-1}(\Sigma(\mathbf{V}_{S,[k]})^2). \quad (92)$$

Now, Lemma 31 applied to the matrix  $\mathbf{V}_{S,[k]}^{\top}$  gives

$$e_{k-1}(\Sigma(\mathbf{V}_{S,[k]})^2) = \sum_{T \subset S, |T|=k-1} e_{k-1}(\Sigma(\mathbf{V}_{T,[k]})^2), \quad (93)$$

Therefore, Lemma 33 yields

$$\sum_{S \subset [d], |S|=k} e_{k-1}(\Sigma(\mathbf{V}_{S,[k]})^2) \leq (d-k+1) \sum_{T \subset [d], |T|=k-1} e_{k-1}(\Sigma(\mathbf{V}_{T,[k]})^2). \quad (94)$$

Using Lemma 31 and the fact that  $\mathbf{V}_k$  is orthogonal, we finally write

$$\sum_{T \subset [d], |T|=k-1} e_{k-1}(\Sigma(\mathbf{V}_{T,[k]})^2) = e_{k-1}(\Sigma(\mathbf{V}_k)^2) = k. \quad (95)$$

Plugging (95) into (94), and then into (92) concludes the proof of Proposition 16.

### D.3 Proof of Proposition 17

We first prove the Frobenius norm bound, which requires more work. The spectral bound is easier and uses a subset of the arguments for the Frobenius norm.

#### D.3.1 Frobenius norm bound

Recall that  $\mathbf{E} = \mathbf{X} - \Pi_k \mathbf{X}$ . We start with Lemma 30:

$$\begin{aligned} \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 &\leq \|\mathbf{E}(\mathbf{I} - \mathbf{P}_S)\|_{\text{Fr}}^2 \\ &\leq \|\mathbf{E}\|_{\text{Fr}}^2 + \text{Tr}(\mathbf{E}^{\top} \mathbf{E} \mathbf{P}_S \mathbf{P}_S^{\top}) - 2 \text{Tr}(\mathbf{P}_S^{\top} \mathbf{E}^{\top} \mathbf{E}). \end{aligned} \quad (96)$$

Since  $\mathbf{E}^{\top} \mathbf{E} = \mathbf{V}_{r-k} \Sigma_{r-k}^2 \mathbf{V}_{r-k}^{\top}$  and  $\mathbf{P}_S = \mathbf{S}(\mathbf{V}_k^{\top} \mathbf{S})^{-1} \mathbf{V}_k^{\top}$ ,

$$\begin{aligned} \text{Tr}(\mathbf{P}_S^{\top} \mathbf{E}^{\top} \mathbf{E}) &= \text{Tr} \left( \mathbf{V}_k ((\mathbf{V}_k^{\top} \mathbf{S})^{\top})^{-1} \mathbf{S}^{\top} \mathbf{V}_{r-k} \Sigma_{r-k} \mathbf{V}_{r-k}^{\top} \right) \\ &= \text{Tr} \left( \mathbf{V}_{r-k}^{\top} \mathbf{V}_k ((\mathbf{V}_k^{\top} \mathbf{S})^{\top})^{-1} \mathbf{S}^{\top} \mathbf{V}_{r-k} \Sigma_{r-k} \right) \\ &= 0, \end{aligned} \quad (97)$$

where the last equality follows from  $\mathbf{V}_{r-k}^\top \mathbf{V}_k = \mathbf{0}$ . Therefore, (96) becomes

$$\|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \leq \|\mathbf{E}\|_{\text{Fr}}^2 + \text{Tr}(\mathbf{E}^\top \mathbf{E} \mathbf{P}_S \mathbf{P}_S^\top). \quad (98)$$

Taking expectations,

$$\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \leq \|\mathbf{E}\|_{\text{Fr}}^2 + \sum_{S \subset [d], |S|=k} \text{Det}(\mathbf{V}_{S,[k]})^2 \text{Tr}(\mathbf{E}^\top \mathbf{E} \mathbf{P}_S \mathbf{P}_S^\top). \quad (99)$$

Proposition 27 expresses  $\text{Det}(\mathbf{V}_{S,[k]})^2$  as a function of the principal angles  $(\theta_i(S))$  between  $\text{Span}(\mathbf{V}_k)$  and  $\text{Span}(\mathbf{S})$ , namely

$$\text{Det}(\mathbf{V}_{S,[k]})^2 = \prod_{i \in [k]} \cos^2(\theta_i(S)). \quad (100)$$

The remainder of the proof is in two steps. First, we bound the second factor in the sum in the right-hand side of (99) with a similar geometric expression. This allows trigonometric manipulations. Second, we work our way back to elementary symmetric polynomials of spectra, and we conclude after some simple algebra.

First, for  $S \subset [d], |S| = k$ , let

$$\mathbf{Z}_S = \mathbf{V}_{d-k}^\top \mathbf{S} (\mathbf{V}_k^\top \mathbf{S})^{-1} = \mathbf{V}_{d-k}^\top \mathbf{P}_S \mathbf{V}_k.$$

It allows us to write

$$\text{Tr}(\mathbf{E}^\top \mathbf{E} \mathbf{P}_S \mathbf{P}_S^\top) = \text{Tr}(\mathbf{V}_{d-k} \boldsymbol{\Sigma}_{d-k}^2 \mathbf{V}_{d-k}^\top \mathbf{P}_S \mathbf{P}_S^\top) = \text{Tr}(\boldsymbol{\Sigma}_{r-k}^2 \mathbf{Z}_S \mathbf{Z}_S^\top). \quad (101)$$

However, for real symmetric matrices  $\mathbf{A}$  and  $\mathbf{B}$  with the same size, a simple diagonalization argument yields

$$\text{Tr}(\mathbf{A}\mathbf{B}) \leq \|\mathbf{A}\|_2 \text{Tr}(\mathbf{B}), \quad (102)$$

so that

$$\text{Tr}(\mathbf{E}^\top \mathbf{E} \mathbf{P}_S \mathbf{P}_S^\top) = \text{Tr}(\boldsymbol{\Sigma}_{r-k}^2 \mathbf{Z}_S \mathbf{Z}_S^\top) \leq \sigma_{k+1}^2 \text{Tr}(\mathbf{Z}_S \mathbf{Z}_S^\top). \quad (103)$$

In Appendix C, we characterize  $\text{Tr}(\mathbf{Z}_S \mathbf{Z}_S^\top)$  using principal angles, see (75). This reads

$$\text{Tr}(\mathbf{Z}_S \mathbf{Z}_S^\top) = \sum_{j \in [k]} \tan^2(\theta_j(S)). \quad (104)$$

Combining (99), (103), (100), and (104), we obtain the following intermediate bound

$$\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \leq \|\mathbf{E}\|_{\text{Fr}}^2 + \sigma_{k+1}^2 \sum_{S \subset [d], |S|=k} \left[ \prod_{i \in [k]} \cos^2(\theta_i(S)) \right] \left[ \sum_{j \in [k]} \tan^2(\theta_j(S)) \right]. \quad (105)$$

Distributing the sum and using trigonometric identities, the general term of the sum in (105) becomes

$$\begin{aligned} \left[ \prod_{i \in [k]} \cos^2(\theta_i(S)) \right] \left[ \sum_{j \in [k]} \tan^2(\theta_j(S)) \right] &= \sum_{i \in [k]} (1 - \cos^2(\theta_i(S))) \prod_{j \in [k], j \neq i} \cos^2(\theta_j(S)) \\ &= \sum_{i \in [k]} \prod_{j \in [k], j \neq i} \cos^2(\theta_j(S)) - \sum_{i \in [k]} \prod_{j \in [k]} \cos^2(\theta_j(S)). \end{aligned} \quad (106)$$

The  $(\cos(\theta_j(S)))_{j \in [k]}$  are the singular values of the matrix  $\mathbf{V}_{S,[k]}$  so that

$$\sum_{i \in [k]} \prod_{j \in [k], j \neq i} \cos^2(\theta_j(S)) = e_{k-1}(\Sigma(\mathbf{V}_{S,[k]}))^2, \quad (107)$$

and

$$\prod_{j \in [k]} \cos^2(\theta_j(S)) = e_k(\Sigma(\mathbf{V}_{S,[k]}))^2. \quad (108)$$

Back to (106), one gets

$$\begin{aligned} \left[ \prod_{i \in [k]} \cos^2(\theta_i(S)) \right] \left[ \sum_{j \in [k]} \tan^2(\theta_j(S)) \right] &= e_{k-1}(\Sigma(\mathbf{V}_{S,[k]}))^2 - \sum_{i \in [k]} e_k(\Sigma(\mathbf{V}_{S,[k]}))^2 \\ &= e_{k-1}(\Sigma(\mathbf{V}_{S,[k]}))^2 - k e_k(\Sigma(\mathbf{V}_{S,[k]}))^2. \end{aligned} \quad (109)$$

Thus, plugging (109) back into the intermediate bound (105), it comes

$$\begin{aligned} &\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \\ &\leq \|\mathbf{E}\|_{\text{Fr}}^2 + \sigma_{k+1}^2 \left[ \sum_{\substack{S \subset [d] \\ |S|=k}} e_{k-1}(\Sigma(\mathbf{V}_{S,[k]}))^2 - k \sum_{\substack{S \subset [d] \\ |S|=k}} e_k(\Sigma(\mathbf{V}_{S,[k]}))^2 \right]. \end{aligned} \quad (110)$$

Using Lemma 31 twice, it comes

$$\begin{aligned} &\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \\ &\leq \|\mathbf{E}\|_{\text{Fr}}^2 + \sigma_{k+1}^2 \left[ \sum_{\substack{S \subset [d] \\ |S|=k}} \sum_{\substack{T \subset S \\ |T|=k-1}} e_{k-1}(\Sigma(\mathbf{V}_{T,[k]}))^2 - k e_k(\Sigma(\mathbf{V}_{:, [k]}))^2 \right]. \end{aligned} \quad (111)$$

Lemmas 33 and the identities  $e_{k-1}(\Sigma(\mathbf{V}_{:, [k]}))^2 = k$  and  $e_k(\Sigma(\mathbf{V}_{:, [k]}))^2 = 1$  allow us to conclude

$$\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \leq \|\mathbf{E}\|_{\text{Fr}}^2 + \sigma_{k+1}^2 [(p-k+1)e_{k-1}(\Sigma(\mathbf{V}_{:, [k]}))^2 - k] \quad (112)$$

$$= \|\mathbf{E}\|_{\text{Fr}}^2 + \sigma_{k+1}^2 (p-k)k. \quad (113)$$

By definition of  $\beta$  (45), we have proven (47), i.e.,

$$\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \leq \|\mathbf{E}\|_{\text{Fr}}^2 \left( 1 + \beta \frac{p-k}{d-k} k \right).$$

### D.3.2 Spectral norm bound

The bound in spectral norm is easier to derive. We start from Lemma 30:

$$\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^2 \mathbf{X}\|_2^2 = \sum_{S \subset [d], |S|=k} \text{Det}(\mathbf{V}_{S,[k]})^2 \|\mathbf{X} - \Pi_S \mathbf{X}\|_2^2 \quad (114)$$

$$\leq \|\mathbf{E}\|_2^2 \sum_{\substack{S \subset [d], |S|=k \\ \text{Det}(\mathbf{V}_{S,[k]})^2 > 0}} \prod_{\ell=1}^{k-1} \sigma_\ell^2(\mathbf{V}_{S,[k]}) \quad (115)$$

By definition of  $e_{k-1}$ , it comes

$$\begin{aligned} \mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_S^2 \mathbf{X}\|_2^2 &\leq \|\mathbf{E}\|_2^2 \sum_{\substack{S \subset [d], |S|=k \\ \text{Det}(\mathbf{V}_{S,[k]})^2 > 0}} e_{k-1}(\Sigma(\mathbf{V}_{S,[k]})^2) \\ &\leq (p-k+1) e_{k-1}(\Sigma(\mathbf{V}_{:, [k]})^2) \|\mathbf{E}\|_2^2 \\ &= (p-k+1) k \|\mathbf{E}\|_2^2, \end{aligned}$$

where we again used the double sum trick of (111) and Lemma 33.

### D.4 Proof of Theorem 18

We start with a lemma on evaluations of elementary symmetric polynomials on specific sequences.

**Lemma 34** *Let  $\boldsymbol{\lambda} \in [0, 1]^k$  such that*

$$\begin{cases} \lambda_1 \geq \dots \geq \lambda_k, \\ \Lambda = \sum_{i=1}^k \lambda_i \geq k-1 + \frac{1}{\theta}. \end{cases} \quad (116)$$

*Then, with the functions  $\phi, \psi$  introduced in Lemma 32,*

$$\begin{cases} \psi(\boldsymbol{\lambda}) \geq \frac{1}{\theta}, \\ \phi(\boldsymbol{\lambda}) \leq k-1 + \theta. \end{cases} \quad (117)$$

**Proof** Let  $\hat{\boldsymbol{\lambda}} = (1, \dots, 1, \Lambda - k + 1) \in \mathbb{R}^k$ . Then

$$\begin{cases} \lambda_1 \leq \hat{\lambda}_1 \\ \lambda_1 + \lambda_2 \leq \hat{\lambda}_1 + \hat{\lambda}_2 \\ \dots \\ \sum_{i=1}^{k-1} \lambda_i \leq \sum_{i=1}^{k-1} \hat{\lambda}_i \\ \sum_{i=1}^k \lambda_i = \sum_{i=1}^k \hat{\lambda}_i \end{cases} \quad (118)$$

so that, according to Definition 21,

$$\boldsymbol{\lambda} \prec_S \hat{\boldsymbol{\lambda}}. \quad (119)$$



Lemma 32 ensures the Schur-convexity of  $\phi$  and the Schur-concavity of  $\psi$ , so that

$$\phi(\boldsymbol{\lambda}) \leq \phi(\hat{\boldsymbol{\lambda}}) = k - 1 + \frac{1}{\Lambda - k + 1} \leq k - 1 + \theta,$$

and

$$\psi(\boldsymbol{\lambda}) \geq \psi(\hat{\boldsymbol{\lambda}}) = \Lambda - k + 1 \geq \frac{1}{\theta}.$$

■

#### D.4.1 Frobenius norm bound

Let  $\mathbf{K} = \mathbf{V}_k \mathbf{V}_k^\top$ , and  $\pi$  be a permutation of  $[d]$  that reorders the leverage scores decreasingly,

$$\ell_{\pi_1}^k \geq \ell_{\pi_2}^k \geq \dots \geq \ell_{\pi_d}^k. \quad (120)$$

By construction,  $T_{p_{\text{eff}}} = [\pi_{p_{\text{eff}}}, \dots, \pi_d]$  thus collects the indices of the smallest leverage scores. Finally, denoting by  $\boldsymbol{\Pi} = (\delta_{i,\pi_j})_{(i,j) \in [d] \times [d]}$  the matricial representation of permutation  $\pi$ , we let

$$\mathbf{K}^\pi = \boldsymbol{\Pi} \mathbf{K} \boldsymbol{\Pi}^\top = ((\mathbf{K}_{\pi_i, \pi_j}))_{1 \leq i, j \leq d}.$$

The goal of the proof is to bound

$$\mathbb{E}_{\text{DPP}} \left[ \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \mid S \cap T_{p_{\text{eff}}} = \emptyset \right] = \frac{\sum \text{Det}(\mathbf{V}_{S,[k]})^2 \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2}{\sum \text{Det}(\mathbf{V}_{S,[k]})^2}, \quad (121)$$

where both sums run over subsets  $S \subset [d]$  such that  $|S| = k$  and  $S \cap T_{p_{\text{eff}}} = \emptyset$ . For simplicity, let us write

$$Z_{k,p_{\text{eff}}}(\theta) = \sum_{\substack{S \subset [d], |S|=k \\ S \cap T_{p_{\text{eff}}}(\theta) = \emptyset}} \text{Det}(\mathbf{V}_{S,[k]})^2, \quad (122)$$

$$Y_{k,p_{\text{eff}}}(\theta) = \sum_{\substack{S \subset [d], |S|=k \\ S \cap T_{p_{\text{eff}}}(\theta) = \emptyset}} \text{Det}(\mathbf{V}_{S,[k]})^2 \text{Tr}(\mathbf{Z}_S \mathbf{Z}_S^\top). \quad (123)$$

Following steps (99) to (103) of the previous proof, one obtains

$$\mathbb{E}_{\text{DPP}} \left[ \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \mid S \cap T_{p_{\text{eff}}} = \emptyset \right] \leq \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}}^2 + \sigma_{k+1}^2 \frac{Y_{k,p_{\text{eff}}}(\theta)}{Z_{k,p_{\text{eff}}}(\theta)}. \quad (124)$$

By definition (45) of the flatness parameter  $\beta$ ,

$$\sigma_{k+1}^2 = \beta \frac{1}{d-k} \sum_{j \geq k+1} \sigma_j^2 = \beta \frac{1}{d-k} \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}}^2. \quad (125)$$

Then, it remains to upper bound the ratio  $Y_{k,p_{\text{eff}}}(\theta)/Z_{k,p_{\text{eff}}}(\theta)$  in (124), which is the important part of the proof. We first evaluate  $Z_{k,p_{\text{eff}}}(\theta)$  and then bound  $Y_{k,p_{\text{eff}}}(\theta)$ .

The matrix  $\mathbf{P}\mathbf{V}_k \in \mathbb{R}^{d \times k}$  has its rows ordered by decreasing leverage scores. Let  $\tilde{\mathbf{V}}_{p_{\text{eff}}(\theta)}^\pi \in \mathbb{R}^{p_{\text{eff}}(\theta) \times k}$  be the submatrix corresponding to the first  $p_{\text{eff}}(\theta)$  rows of  $\mathbf{P}\mathbf{V}_k$ . Let also

$$\hat{\mathbf{V}}_{p_{\text{eff}}(\theta)}^\pi = \begin{pmatrix} \tilde{\mathbf{V}}_{\pi, p_{\text{eff}}(\theta)} \\ \mathbf{0}_{d-p_{\text{eff}}(\theta), k} \end{pmatrix}$$

be padded with zeros. Then

$$\mathbf{K}_{p_{\text{eff}}(\theta)}^\pi = \left[ \begin{array}{c|c} \tilde{\mathbf{V}}_{\pi, p_{\text{eff}}(\theta)} \tilde{\mathbf{V}}_{\pi, p_{\text{eff}}(\theta)}^\top & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right] = \hat{\mathbf{V}}_{p_{\text{eff}}(\theta)}^\pi (\hat{\mathbf{V}}_{p_{\text{eff}}(\theta)}^\pi)^\top \in \mathbb{R}^{d \times d}. \quad (126)$$

The nonzero block of  $\mathbf{K}_{p_{\text{eff}}(\theta)}^\pi$  is a submatrix of  $\mathbf{K}^\pi$ , and  $\text{rk } \mathbf{K}^\pi = \text{rk } \mathbf{K} = k$ . Hence  $\mathbf{K}_{p_{\text{eff}}(\theta)}^\pi$  has at most  $k$  nonzero eigenvalues

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k \geq 0 = \lambda_{k+1} = \dots = \lambda_d. \quad (127)$$

Therefore,

$$e_k(\Lambda(\mathbf{K}_{p_{\text{eff}}(\theta)}^\pi)) = \sum_{\substack{T \subset [d] \\ |T|=k}} \prod_{j \in T} \lambda_j = \prod_{i \in [k]} \lambda_i. \quad (128)$$

Note moreover that

$$\forall \ell \in [k], \quad e_\ell(\Sigma(\hat{\mathbf{V}}_{\pi, p_{\text{eff}}(\theta)}^2)) = e_\ell(\Lambda(\mathbf{K}_{p_{\text{eff}}(\theta)}^\pi)). \quad (129)$$

By construction,

$$Z_{k, p_{\text{eff}}(\theta)} = \sum_{\substack{S \subset [d], |S|=k \\ S \cap T_{p_{\text{eff}}(\theta)} = \emptyset}} \text{Det}(\mathbf{V}_{S, [k]})^2 = \sum_{S \subset [d], |S|=k} \text{Det} \left[ \left( \hat{\mathbf{V}}_{p_{\text{eff}}(\theta)}^\pi \right)_{S, :} \right]^2 \quad (130)$$

Then, Lemma 31 yields

$$Z_{k, p_{\text{eff}}(\theta)} = e_k(\Sigma(\hat{\mathbf{V}}_{\pi, p_{\text{eff}}(\theta)}^2)) = e_k(\Lambda(\mathbf{K}_{p_{\text{eff}}(\theta)}^\pi)) = \prod_{i \in [k]} \lambda_i. \quad (131)$$

Now we bound  $Y_{k, p_{\text{eff}}(\theta)}$ . We use again principal angles and trigonometric identities. Using (104) and (109) above, it holds

$$\begin{aligned} Y_{k, p_{\text{eff}}(\theta)} &= \sum_{\substack{S \subset [d], |S|=k \\ S \cap T_{p_{\text{eff}}(\theta)} = \emptyset}} \text{Det}(\mathbf{V}_{S, [k]})^2 \text{Tr}(\mathbf{Z}_S \mathbf{Z}_S^\top) \\ &= \sum_{\substack{S \subset [d], |S|=k \\ S \cap T_{p_{\text{eff}}(\theta)} = \emptyset}} \prod_{i \in [k]} \cos^2(\theta_i(S)) \sum_{j \in [k]} \tan^2(\theta_j(S)) \\ &= \sum_{\substack{S \subset [d], |S|=k \\ S \cap T_{p_{\text{eff}}(\theta)} = \emptyset}} e_{k-1}(\Sigma(\mathbf{V}_{S, [k]})^2) - k e_k(\Sigma(\mathbf{V}_{S, [k]})^2) \end{aligned} \quad (132)$$

$$= \sum_{S \subset [d], |S|=k} e_{k-1} \left( \Sigma \left( \left[ \hat{\mathbf{V}}_{p_{\text{eff}}(\theta)}^\pi \right]_{S, :} \right)^2 \right) - k e_k \left( \Sigma \left( \left[ \hat{\mathbf{V}}_{p_{\text{eff}}(\theta)}^\pi \right]_{S, :} \right)^2 \right) \quad (133)$$

By Lemma 33 applied to the matrix  $\hat{\mathbf{V}}_{\pi, p_{\text{eff}}(\theta)}$  combined to (130), we get

$$\begin{aligned} Y_{k, p_{\text{eff}}(\theta)} &\leq (p_{\text{eff}}(\theta) - k + 1)e_{k-1}(\Sigma(\hat{\mathbf{V}}_{p_{\text{eff}}(\theta)}^\pi)^2) - k e_k(\Sigma(\hat{\mathbf{V}}_{p_{\text{eff}}(\theta)}^\pi)^2) \\ &\leq (p_{\text{eff}}(\theta) - k + 1)e_{k-1}(\Lambda(\mathbf{K}_{p_{\text{eff}}(\theta)}^\pi)) - k e_k(\Lambda(\mathbf{K}_{p_{\text{eff}}(\theta)}^\pi)) \\ &\leq \left( (p_{\text{eff}}(\theta) - k + 1)\phi(\tilde{\boldsymbol{\lambda}}) - k \right) Z_{k, p_{\text{eff}}(\theta)}. \end{aligned} \quad (134)$$

where  $\tilde{\boldsymbol{\lambda}} = (1, \dots, 1, \text{Tr}(\mathbf{K}_{p_{\text{eff}}(\theta)}^\pi) - k + 1) \in \mathbb{R}^k$ , see Lemma 34. Now, as in the proof of Lemma 34,

$$\phi(\tilde{\boldsymbol{\lambda}}) = k - 1 + \frac{1}{\text{Tr}(\mathbf{K}_{p_{\text{eff}}(\theta)}^\pi) - k + 1} \leq k - 1 + \theta$$

by (49). Thus (134) yields

$$\frac{Y_{k, p_{\text{eff}}(\theta)}}{Z_{k, p_{\text{eff}}(\theta)}} \leq (p_{\text{eff}}(\theta) - k + 1)(k - 1 + \theta) - k \leq (p_{\text{eff}}(\theta) - k + 1)(k - 1 + \theta). \quad (135)$$

Finally, plugging (135) and (125) in (124) concludes the proof of (52).

#### D.4.2 Spectral norm bound

We proceed as for the Frobenius norm, using the notation of Section D.3.1. Lemma 30, Equations (132) and (135) yield

$$\begin{aligned} &\mathbb{E}_{\text{DPP}} \left[ \|\mathbf{X} - \Pi_S^2 \mathbf{X}\|_2^2 \mid S \cap T_{p_{\text{eff}}} = \emptyset \right] \\ &= Z_{k, p_{\text{eff}}(\theta)}^{-1} \sum_{\substack{S \subset [d], |S|=k \\ S \cap T_{p_{\text{eff}}(\theta)} = \emptyset}} \text{Det}(\mathbf{V}_{S, [k]})^2 \|\mathbf{X} - \Pi_S^2 \mathbf{X}\|_2^2, \\ &\leq Z_{k, p_{\text{eff}}(\theta)}^{-1} \|\mathbf{X} - \Pi_k \mathbf{X}\|_2^2 \sum_{\substack{S \subset [d], |S|=k \\ S \cap T_{p_{\text{eff}}(\theta)} = \emptyset, \\ \text{Det}(\mathbf{V}_{S, [k]})^2 > 0}} \prod_{\ell=1}^{k-1} \sigma_\ell^2(\mathbf{V}_{S, [k]}) \\ &\leq Z_{k, p_{\text{eff}}(\theta)}^{-1} \|\mathbf{X} - \Pi_k \mathbf{X}\|_2^2 \sum_{\substack{S \subset [d], |S|=k \\ S \cap T_{p_{\text{eff}}(\theta)} = \emptyset \\ \text{Det}(\mathbf{V}_{S, [k]})^2 > 0}} e_{k-1}(\Sigma(\mathbf{V}_{S, [k]})^2) \\ &\leq \frac{Y_{k, p_{\text{eff}}(\theta)}}{Z_{k, p_{\text{eff}}(\theta)}} \|\mathbf{X} - \Pi_k \mathbf{X}\|_2^2 \\ &\leq (p_{\text{eff}}(\theta) - k + 1)(k - 1 + \theta) \|\mathbf{X} - \Pi_k \mathbf{X}\|_2^2, \end{aligned}$$

which is the claimed spectral bound.

### D.4.3 Bounding the probability of rejection

Still with the notation of Section D.3.1, (130) yields

$$\begin{aligned} \mathbb{P}(S \cap T_{p_{\text{eff}}}(\theta) = \emptyset) &= \sum_{\substack{S \subset [d], |S|=k \\ S \cap T_{p_{\text{eff}}}(\theta) = \emptyset}} \text{Det}(\mathbf{V}_{S,[k]})^2 \\ &= e_k(\mathbf{K}_{p_{\text{eff}}}^\pi) \end{aligned} \quad (136)$$

$$\begin{aligned} &= \prod_{i \in [k]} \lambda_i \\ &= \psi(\hat{\boldsymbol{\lambda}}). \end{aligned} \quad (137)$$

Lemma 34 concludes the proof since

$$\psi(\hat{\boldsymbol{\lambda}}) \geq \frac{1}{\theta}. \quad (138)$$

## D.5 Proof of Proposition 20

First, Proposition 11 gives

$$\mathcal{E}(\mathbf{w}_S) \leq \frac{(1 + \max_{i \in [k]} \tan^2 \theta_i(S)) \|\mathbf{w}^*\|^2 \sigma_{k+1}^2}{N} + \frac{k}{N} \nu. \quad (139)$$

Now (75) further gives

$$\max_{i \in [k]} \tan^2 \theta_i(S) \leq \sum_{i \in [k]} \tan^2 \theta_i(S) = \text{Tr}(\mathbf{Z}_S \mathbf{Z}_S^\top). \quad (140)$$

The proof now follows the same lines as for the approximation bounds. First, following the lines of Section D.3, , we straightforwardly bound

$$\mathbb{E}_{\text{DPP}} \sum_{i \in [k]} \tan^2(\theta_i(S)) = \sum_{S \subset [d], |S|=k} \prod_{i \in [k]} \cos^2(\theta_i(S)) \sum_{j \in [k]} \tan^2(\theta_j(S)) \quad (141)$$

and obtain (53). In a similar vein, the same lines as in Section D.4 allow bounding

$$\mathbb{E}_{\text{DPP}} \left[ \sum_{i \in [k]} \tan^2(\theta_i(S)) | S \cap T_{p_{\text{eff}}} = \emptyset \right] = \sum_{\substack{S \subset [d], |S|=k \\ S \cap T_{p_{\text{eff}}}(\theta) = \emptyset}} \prod_{i \in [k]} \cos^2(\theta_i(S)) \sum_{j \in [k]} \tan^2(\theta_j(S)). \quad (142)$$

and yield (54).

## E Generating orthogonal matrices with prescribed leverage scores

In this section, we describe an algorithm that samples a random orthonormal matrix with a prescribed profile of  $k$ -leverage scores. This algorithm was used to generate the matrices  $\mathbf{F} = \mathbf{V}_k^\top \in \mathbb{R}^{k \times d}$  for the toy datasets of Section 6. The orthogonality

constraint can be expressed as a condition on the spectrum of the matrix  $\mathbf{K} = \mathbf{V}_k \mathbf{V}_k^\top$ , namely  $\text{Sp}(\mathbf{K}) \subset \{0, 1\}$ . On the other hand, the constraint on the  $k$ -leverage scores can be expressed as a condition on the diagonal of  $\mathbf{K}$ . Thus, the problem of generating an orthogonal matrix with a given profile of  $k$ -leverage scores boils down to enforcing conditions on the spectrum and the diagonal of a symmetric matrix  $\mathbf{K}$ .

### E.1 Definitions and statement of the problem

We denote by  $(\mathbf{f}_i)_{i \in [d]}$  the columns of the matrix  $\mathbf{F}$ . For  $n \in \mathbb{N}$ , we write  $\mathbb{1}_n$  the vector containing ones living in  $\mathbb{R}^n$ , and  $\mathbb{0}_n$  the vector containing zeros living in  $\mathbb{R}^n$ . We say that the vector  $\mathbf{u} \in \mathbb{R}^n$  interlaces on  $\mathbf{v} \in \mathbb{R}^n$  and we denote

$$\mathbf{u} \sqsubseteq \mathbf{v}$$

if  $u_n \leq v_n$  and  $\forall i \in [1 : n - 1]$ ,  $v_{i+1} \leq u_i \leq v_i$ .

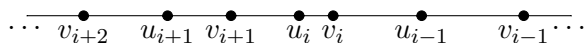


Figure 9: Illustration of the interlacing of  $\mathbf{u}$  on  $\mathbf{v}$ .

**Definition 35** Let  $k, d \in \mathbb{N}$ , with  $k \leq d$ . Let  $\mathbf{F} \in \mathbb{R}^{k \times d}$  be a full rank matrix<sup>5</sup>. Within this section, we denote  $\boldsymbol{\sigma}^2 = (\sigma_1^2, \sigma_2^2, \dots, \sigma_k^2)$  the squares of the nonvanishing singular values of the matrix  $\mathbf{F}$ , and  $\boldsymbol{\ell} = (\ell_1 = \|\mathbf{f}_1\|^2, \ell_2 = \|\mathbf{f}_2\|^2, \dots, \ell_d = \|\mathbf{f}_d\|^2)$  are the squared norms of the columns of  $\mathbf{F}$ , which we assume to be ordered decreasingly:

$$\ell_1 \geq \ell_2 \geq \dots \geq \ell_d.$$

When  $\mathbf{F}$  is orthogonal, we can think of  $\boldsymbol{\ell}$  as a vector of leverage scores.

We are interested in the problem of constructing an orthogonal matrix given its leverage scores.

**Problem 1** Let  $k, d \in \mathbb{N}$ , with  $k \leq d$ , and let  $\boldsymbol{\ell} \in \mathbb{R}_+^d$  such that  $\sum_{i=1}^d \ell_i = k$ . Build a matrix  $\mathbf{F} \in \mathbb{R}^{k \times d}$  such that

$$\text{Sp}(\mathbf{F}^\top \mathbf{F}) = [\mathbb{1}_k, \mathbb{0}_{d-k}], \quad (143)$$

and

$$\text{Diag}(\mathbf{F}^\top \mathbf{F}) = \boldsymbol{\ell}. \quad (144)$$

We actually consider here the generalization of Problem 2 to an arbitrary spectrum.

**Problem 2** Let  $k, d \in \mathbb{N}$ , with  $k \leq d$ , and let  $\boldsymbol{\ell} \in \mathbb{R}_+^d$  such that  $\sum_{i=1}^d \ell_i = \sum_{i=1}^k \sigma_i^2$ . Build a matrix  $\mathbf{F} \in \mathbb{R}^{k \times d}$  such that

$$\text{Sp}(\mathbf{F}^\top \mathbf{F}) = [\boldsymbol{\sigma}^2, \mathbb{0}_{d-k}] =: \hat{\boldsymbol{\sigma}}^2 \quad (145)$$

and

$$\text{Diag}(\mathbf{F}^\top \mathbf{F}) = \boldsymbol{\ell}. \quad (146)$$

<sup>5</sup>A frame, using the definitions of (Fickus et al., 2011a) and (Fickus et al., 2011b).



```

GIVENSALGORITHM( $\ell, \sigma$ )
1    $\mathbf{F} \leftarrow [ \text{Diag}(\sigma) \mid \mathbf{0} ] \in \mathbb{R}^{k \times d}$ 
2   while  $\exists i, j, k \in [d], i < k < j : \|\mathbf{f}_i\|^2 < \ell_i, \|\mathbf{f}_k\|^2 = \ell_k, \|\mathbf{f}_j\|^2 > \ell_j$ 
3     if  $\ell_i - \|\mathbf{f}_i\|^2 \leq \|\mathbf{f}_j\|^2 - \ell_j$ 
4        $\mathbf{F} \leftarrow \mathbf{G}_{i,j}(\theta)\mathbf{F}$ , where  $\|(\mathbf{G}_{i,j}(\theta)\mathbf{F})_i\|^2 = \ell_i$ .
5     else
6        $\mathbf{F} \leftarrow \mathbf{G}_{i,j}(\theta)\mathbf{F}$ , where  $\|(\mathbf{G}_{i,j}(\theta)\mathbf{F})_j\|^2 = \ell_j$ ,
7   return  $\mathbf{F} \in \mathbb{R}^{k \times d}$ .

```

Figure 10: The pseudocode of the algorithm proposed by [Dhillon et al. \(2005\)](#) for generating a matrix given its leverage scores and spectrum by successively applying Givens rotations.

collection of spectra of all minors of  $\mathbf{F} \in \mathcal{M}_{(\ell, \sigma)}$ . This parametrization was introduced by [Fickus et al. \(2011b\)](#), and we recall it in Section E.3. For now, let us simply look at Figure 12, which displays a few outputs of our algorithm for the same input as in Figure 11a. We now obtain different matrices for the same input  $(\ell, \sigma)$ , and these matrices are less structured than the output of Algorithm 10, as required.

### E.3 The restricted Gelfand-Tsetlin polytope

**Definition 38** Recall that  $(\mathbf{f}_i)_{i \in [d]}$  are the columns of the matrix  $\mathbf{F} \in \mathbb{R}^{k \times d}$ . For  $r \in [d]$ , we further define

$$\mathbf{F}_r = \mathbf{F}_{:, [r]} \in \mathbb{R}^{k \times r}, \quad (150)$$

$$\mathbf{C}_r = \sum_{i \in [r]} \mathbf{f}_i \mathbf{f}_i^\top \in \mathbb{R}^{k \times k}, \quad (151)$$

$$\mathbf{G}_r = \mathbf{F}_r^\top \mathbf{F}_r \in \mathbb{R}^{r \times r}. \quad (152)$$

Furthermore, we note for  $r \in [d]$ ,

$$(\lambda_{r,i})_{i \in [k]} = \Lambda(\mathbf{C}_r), \quad (153)$$

$$(\tilde{\lambda}_{r,i})_{i \in [r]} = \Lambda(\mathbf{G}_r). \quad (154)$$

The  $(\lambda_{r,i})_{i \in [k]}$ ,  $r \in [d]$ , are called the outer eigensteps of  $\mathbf{F}$ , and we group them in the matrix

$$\Lambda^{out}(\mathbf{F}) = (\lambda_{r,i})_{i \in [k], r \in [d]} \in \mathbb{R}^{k \times d}.$$

Similarly, the  $(\tilde{\lambda}_{r,i})_{i \in [r]}$  are called inner eigensteps of  $\mathbf{F}$ .

**Example 2** For  $k = 2$ ,  $d = 4$ , consider the full-rank matrix

$$\mathbf{F} = \begin{bmatrix} 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix}, \quad (155)$$

Then

$$\Lambda^{out}(\mathbf{F}) = \begin{bmatrix} 1 & 1 & 2 & 2 \\ 0 & 1 & 1 & 2 \end{bmatrix}. \quad (156)$$

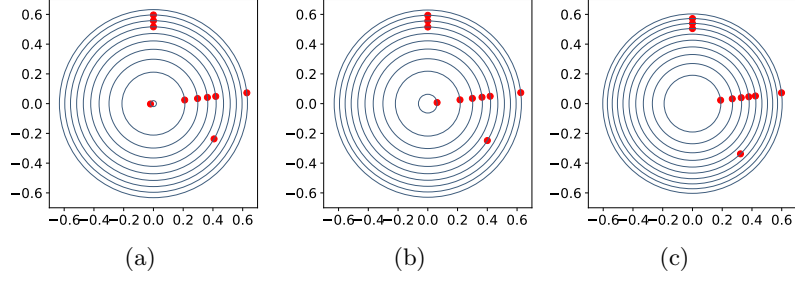


Figure 11: The output of the algorithm in Figure 10 for  $k = 2$ ,  $d = 10$ ,  $\sigma = (1, 1)$ , and three different values of  $\ell$  that each add to  $k$ . Each red dot has coordinates a column of  $\mathbf{F}$ . The blue circles have for radii the prescribed  $(\sqrt{\ell_i})$ .

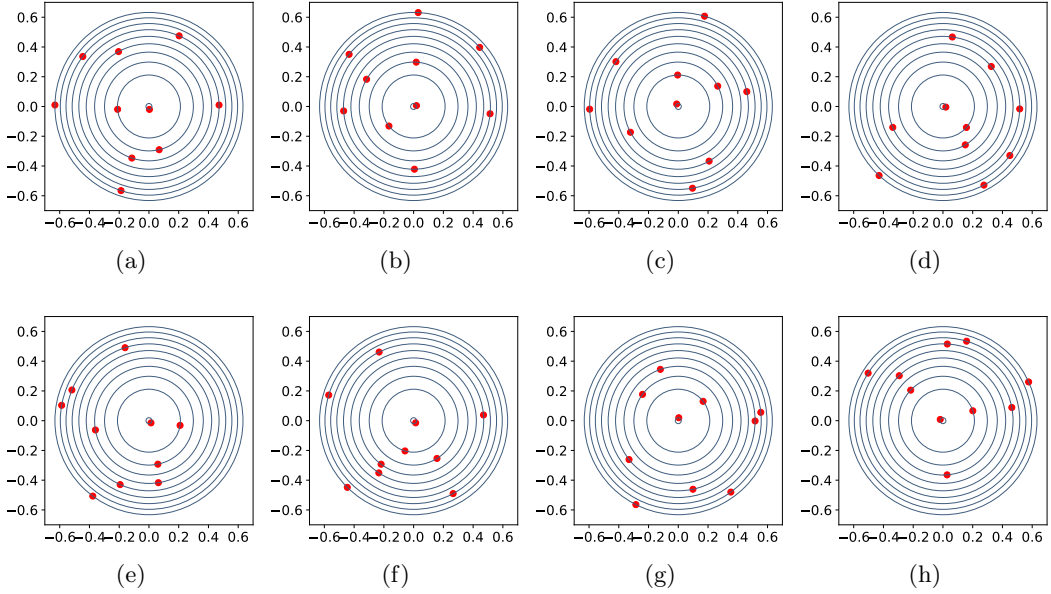


Figure 12: The output of our algorithm for  $k = 2$ ,  $d = 10$ , an input  $\sigma = (1, 1)$ , and  $\ell$  as in Figure 11a. Each red dot has coordinates a column of  $\mathbf{F}$ . The blue circles have for radii the prescribed  $(\sqrt{\ell_i})$ .

**Proposition 39** *The outer eigensteps satisfy the following constraints:*

$$\begin{cases} \forall i \in [k], \lambda_{0,i} = 0 \\ \forall i \in [k], \lambda_{d,i} = \sigma_i^2 \\ \forall r \in [d], (\lambda_{r,:}) \sqsubseteq (\lambda_{r+1,:}) \\ \forall r \in [d], \sum_{i \in [d]} \lambda_{r,i} = \sum_{i \in [r]} \ell_i \end{cases} \quad (157)$$

In other words, the outer eigensteps are constrained to live in a polytope. We define the restricted Gelfand-Tsetlin polytope  $\mathbf{GT}_{(k,d)}(\sigma, \ell)$  to be the subset of  $\mathbb{R}^{k \times d}$  defined by the equations (157). A more graphical summary of the interlacing and sum



$$\begin{array}{ccccccccc}
\ell_1 = \lambda_{1,1} & \blacktriangleleft & \lambda_{2,1} & \blacktriangleleft & \lambda_{3,1} & \cdots & \lambda_{d-1,1} & \blacktriangleleft & \lambda_{d,1} = \sigma_1 \\
+ & \blacktriangleright & + & \blacktriangleright & + & \cdots & + & \blacktriangleright & + \\
0 = \lambda_{1,2} & \blacktriangleleft & \lambda_{2,2} & \blacktriangleleft & \lambda_{3,2} & \cdots & \lambda_{d-1,2} & \blacktriangleleft & \lambda_{d,2} = \sigma_2 \\
+ & \blacktriangleright & + & \blacktriangleright & + & \cdots & + & \blacktriangleright & + \\
0 = \lambda_{1,3} & \blacktriangleleft & \lambda_{2,3} & \blacktriangleleft & \lambda_{3,3} & \cdots & \lambda_{d-1,3} & \blacktriangleleft & \lambda_{d,3} = \sigma_3 \\
\vdots & & \vdots & & \vdots & & \vdots & & \vdots \\
0 = \lambda_{1,k} & \blacktriangleleft & \lambda_{2,k} & \blacktriangleleft & \lambda_{3,k} & \cdots & \lambda_{d-1,k} & \blacktriangleleft & \lambda_{d,k} = \sigma_k \\
\hline
& \ell_1 & \sum_{i \leq 2} \ell_i & \sum_{i \leq 3} \ell_i & \sum_{i \leq d-1} \ell_i & \sum_{i \leq d} \ell_i & & & 
\end{array}$$

Figure 13: The interlacing relationships (157) satisfied by the outer eigensteps of a frame. Thick triangles are used in place of  $\leq$  for improved readability.

constraints is given in Figure 13. The restricted GT polytope<sup>6</sup> allows a parametrization of  $\mathcal{M}_{(\ell, \sigma)}$  by the following reconstruction result.

**Theorem 40 (Theorem 3, Fickus et al., 2011a)** *Every matrix  $\mathbf{F} \in \mathcal{M}_{(\ell, \sigma)}$  can be constructed as follows:*

- pick a valid sequence of outer eigensteps noted  $\Lambda^{\text{out}} \in \mathbf{GT}_{(k,d)}(\sigma, \ell)$ ,
- pick  $\mathbf{f}_1 \in \mathbb{R}^k$  such that

$$\|\mathbf{f}_1\|^2 = \ell_1, \quad (158)$$

- for  $r \in [d]$ , consider the polynomial  $p_r(x) = \prod_{i \in [d]} (x - \lambda_{r,i})$ , and for each  $r \in [d-1]$ , choose  $\mathbf{f}_{r+1} \in \mathbb{R}^k$  such that

$$\forall \lambda \in \{\lambda_{r,i}\}_{i \in [d]}, \|\mathbf{P}_{r,\lambda} \mathbf{f}_{r+1}\|^2 = - \lim_{x \rightarrow \lambda} (x - \lambda) \frac{p_{r+1}(\lambda)}{p_r(\lambda)}, \quad (159)$$

where  $\mathbf{P}_{r,\lambda}$  denotes the orthogonal projection onto the eigenspace  $\text{Ker}(\lambda \mathbb{I}_k - \mathbf{F}_r \mathbf{F}_r^T)$ .

Conversely, any matrix  $\mathbf{F}$  constructed by this process is in  $\mathcal{M}_{(\ell, \sigma)}$ .

Fickus et al. (2011a) propose an algorithm to construct a vector  $\mathbf{f}_r$  satisfying Equation (159). Finally, an algorithm for the construction of a valid sequence of eigensteps  $\Lambda^{\text{out}} \in \mathbf{GT}_{(k,d)}(\sigma, \ell)$  was proposed in (Fickus et al., 2011b). This yields the following constructive result.

**Theorem 41 (Theorem 4.1, Fickus et al., 2011b)** *Every matrix  $\mathbf{F} \in \mathcal{M}(\sigma, \ell)$  can be constructed as follows:*

- Set  $\forall i \in [k], \tilde{\lambda}_{d,i} = \sigma_i^2$ ,
- For  $r \in \{d-1, \dots, 1\}$ , construct  $\{\tilde{\lambda}_{r,\cdot}\}$  as follows. For each  $i \in \{k, \dots, 1\}$ , pick

$$\tilde{\lambda}_{r-1,i} \in [B_{i,r}(\ell, \sigma), A_{i,r}(\ell, \sigma)],$$

<sup>6</sup>Note the difference with the Gelfand-Tsetlin polytope in the random matrix literature (Baryshnikov, 2001), where only the spectrum is constrained, not the diagonal.

```

RANDOM EIGENSTEPS( $\ell, \sigma$ )
1    $\Lambda^{\text{out}} \leftarrow \mathbb{O} \in \mathbb{R}^{k \times d}$ 
2    $\forall i \in [k], \tilde{\lambda}_{d,i} \leftarrow \sigma_i$ 
3   for  $r \in \{d-1, \dots, 1\}$ 
4       for  $i \in \{k, \dots, 1\}$ 
5           Pick  $\tilde{\lambda}_{r-1,i} \sim \mathcal{U}([B_{i,r}(\ell, \sigma), A_{i,r}(\ell, \sigma)])$ 
return  $\Lambda^{\text{out}}$ 

```

Figure 14: The pseudocode of the generator of random valid eigensteps taking as input  $(\ell, \sigma)$ .

where

$$\begin{aligned}
 A_{i,r}(\ell, \sigma) &= \max \left\{ \tilde{\lambda}_{r+1,i+1}, \sum_{t=i}^k \tilde{\lambda}_{r+1,t} - \sum_{t=i+1}^k \tilde{\lambda}_{r,t} - \ell_{r+1} \right\} \\
 B_{i,r}(\ell, \sigma) &= \min \left\{ \tilde{\lambda}_{r+1,i}, \min_{z=1, \dots, i} \left\{ \sum_{t=z}^r \ell_t - \sum_{t=z+1}^i \tilde{\lambda}_{r+1,t} - \sum_{t=i+1}^k \tilde{\lambda}_{r,t} \right\} \right\}.
 \end{aligned} \tag{160}$$

Furthermore, any sequence constructed by this algorithm is a valid sequence of inner eigensteps.

Based on these results we propose an algorithm for the generation of orthogonal random matrices with a given profile of leverage scores.

#### E.4 Our algorithm

We consider a randomization of the algorithm given in Theorem 41. First, we generate a random sequence of valid inner eigensteps  $\Lambda^{\text{in}}$  using Algorithm 14. Then we proceed to the reconstruction a frame that admits  $\Lambda^{\text{in}}$  as a sequence of eigensteps using the Algorithm proposed in (Fickus et al., 2011a).

Note that Equations (158) and (159) admit several solutions. For example, for  $r \in [d]$ , and if  $\mathbf{f}_{r+1}$  satisfies (159),  $-\mathbf{f}_{r+1}$  satisfies this equation too. Fickus et al. (2011a) actually prove that the set of solutions of these equations is invariant under a specific action of the orthogonal group  $\mathbb{O}(\rho(r, k))$  where  $\rho(r, k) \in \mathbb{N}$  nontrivially depends on the eigensteps. In the reconstruction step of our algorithm, we apply a random Haar-distributed orthogonal matrix as soon as such an invariance is provable. Namely, we a random orthogonal matrix sampled from the Haar measure on  $\mathbb{O}(d)$  to the vector  $\mathbf{f}_1$  and, then, we apply an independent random orthogonal matrix sampled from the Haar measure on  $\mathbb{O}(\rho(r, k))$  to each reconstructed vector  $\mathbf{f}_{r+1}$ .

Figure 12 displays a few samples from our algorithm, which display diversity and no apparent structure, as required for a generator of toy datasets. The question of fully characterizing the distribution of the output of our algorithm is an open question.