



HAL
open science

Un processus ponctuel déterminantal pour la sélection d'attributs

Ayoub Belhadji, Rémi Bardenet, Pierre Chainais

► **To cite this version:**

Ayoub Belhadji, Rémi Bardenet, Pierre Chainais. Un processus ponctuel déterminantal pour la sélection d'attributs. GRETSI 2019, Aug 2019, Lille, France. hal-02439935

HAL Id: hal-02439935

<https://hal.science/hal-02439935v1>

Submitted on 14 Jan 2020

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.

Un processus ponctuel déterminantal pour la sélection d’attributs

Ayoub BELHADJI, Rémi BARDENET, Pierre CHAINAIS
Université de Lille, CNRS, Centrale Lille, Inria, UMR 9189 – CRISTAL,
59651 Villeneuve d’Ascq, France

ayoub.belhadji@centralelille.fr, remi.bardenet@gmail.com
pierre.chainais@centralelille.fr

Résumé – La réduction de dimension est une tâche récurrente de l’analyse des signaux en grande dimension. Les axes principaux de l’analyse en composantes principales sont difficilement interprétables. Lorsqu’on souhaite préserver l’interprétabilité des dimensions, la sélection d’attributs est préférable, mais implique a priori une optimisation combinatoire très coûteuse. Cet article contourne la difficulté en proposant un nouvel algorithme de sélection aléatoire d’attributs, ou de manière équivalente, de colonnes de la matrice des données. On utilise pour cela un processus ponctuel déterminantal (PPD) sur les indices des colonnes, dont le noyau est contrôlé par la structure des données.

Abstract – Dimension reduction is often a preamble to the analysis of high-dimensional signals. Principal component analysis projects data onto directions of large variance, but these directions are difficult to interpret. Feature selection avoids this loss of interpretability, at the price of an expensive combinatorial optimization. To avoid paying this price, we propose a novel algorithm for random feature selection, i.e., random column subset selection in the design matrix. The crux is a determinantal point process on the column indices, with a data-dependent kernel.

1 Introduction

La réduction de dimension est une tâche récurrente de l’analyse de données. Les techniques les plus connues sont l’analyse en composantes principales (ACP) et la sélection d’attributs. L’analyse en composantes principales projette les données sur les directions de plus grande variance, les directions principales. Ces nouveaux attributs, représentés par les axes principaux, sont difficilement interprétables comparés aux attributs initiaux. Lorsqu’on souhaite préserver la sémantique de ces derniers, on pratique plutôt la sélection d’attributs.

Si l’on appelle $\mathbf{X} \in \mathbb{R}^{N \times d}$ la matrice rassemblant les données sous la forme de N observations exprimées sur d attributs, la sélection d’attributs consiste à choisir un sous-ensemble S d’indices de colonnes de la matrice \mathbf{X} . Indépendamment de la tâche d’apprentissage à réaliser ensuite, on peut mesurer la qualité d’une sélection S d’attributs à l’aide de l’erreur d’approximation $\|\mathbf{X} - \Pi_S \mathbf{X}\|$; ici $\|\cdot\|$ est une norme matricielle, et Π_S est la projection sur l’espace engendré par les colonnes de \mathbf{X} d’indice S . Trouver naïvement k colonnes qui optimisent ce critère requiert d’explorer les C_d^k combinaisons possibles, ce qui est habituellement hors d’atteinte quand $d \gg 1$. Une alternative à l’exploration combinatoire est de sous-échantillonner aléatoirement des colonnes. Les qualités attendues d’un tel algorithme sont alors (i) une erreur d’approximation moyenne faible, i.e. proche de celle de l’ACP de cardinal comparable et (ii) un coût d’échantillonnage polynomial.

Cet article propose un nouvel algorithme de sélection aléatoire d’attributs qui possède ces deux qualités. On utilise pour cela un processus ponctuel déterminantal (PPD) sur les indices des colonnes de \mathbf{X} . Les PPD sont des distributions de

probabilité qui produisent des ensembles de points en interaction répulsive. Un PPD est paramétré par une matrice que l’on appelle noyau. En choisissant un noyau lié à la décomposition en valeurs singulières de \mathbf{X} , nous obtenons alors des garanties fortes sur l’erreur moyenne d’approximation de la sélection d’attributs. Les résultats détaillés sont à retrouver dans l’article long [8].

Le reste du présent article s’organise comme suit. En Partie 2, on présente quelques-uns des algorithmes de sélection existants. En Partie 3, on introduit les PPD. En Partie 4, on introduit un PPD particulier défini via la décomposition en valeurs singulières de \mathbf{X} et on liste quelques-unes de ses propriétés théoriques. Enfin, en Partie 5, on valide ces résultats par des simulations.

2 La sélection aléatoire de colonnes

Dans cette partie, on rappelle quelques algorithmes de l’état de l’art après avoir expliqué la forme des résultats théoriques.

2.1 Comment et à quoi se comparer

On compare les algorithmes de sélection de colonnes à l’ACP qui utiliserait un nombre k de directions principales. L’ACP résume les données \mathbf{X} par leur meilleure approximation $\Pi_k \mathbf{X}$ de rang k . L’adjectif « meilleur » s’entend ici au sens de la norme $\|\cdot\|_2$ ou de la norme de Frobenius $\|\cdot\|_{\text{Fr}}$, deux normes matricielles usuelles [7] donnant le même minimiseur

$$\Pi_k \mathbf{X} = \arg \min_{\text{rg } \mathbf{P} = k} \|\mathbf{X} - \mathbf{P}\|_2 = \arg \min_{\text{rg } \mathbf{P} = k} \|\mathbf{X} - \mathbf{P}\|_{\text{Fr}}. \quad (1)$$

En sélection de colonnes, on souhaite comparer \mathbf{X} à sa projection sur l’espace engendré par les colonnes sélectionnées.

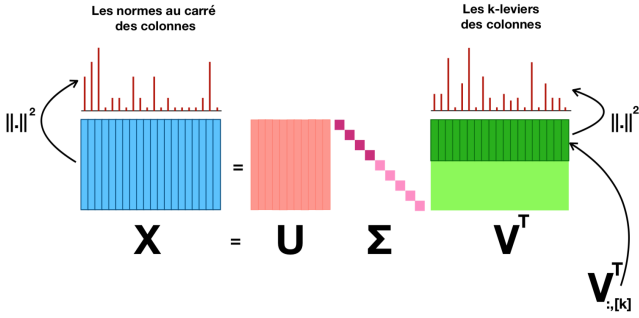


FIGURE 1 – La différence entre les k -leviers et les normes au carré des colonnes.

Concrètement, pour un ensemble $S \subset [d] = \{1, \dots, d\}$ d'indices de colonnes, on restreint le minimum de (1) aux produits $P = X_{:,S}B$, avec B de taille $|S| \times d$. On note alors $\Pi_{S,k}^2 X$ et $\Pi_{S,k}^{\text{Fr}} X$ les minimiseurs obtenus pour la norme 2 et de Frobenius, respectivement. Notons que contrairement à la PCA, le choix de la norme a une influence sur la projection.

L'enjeu de la sélection de colonnes est de trouver S de petit cardinal tel que $\|X - \Pi_{S,k}^\nu X\|_\nu$ est comparable à $\|X - \Pi_k X\|_\nu$, pour $\nu \in \{2, \text{Fr}\}$.

2.2 Échantillonnage multinomial

L'algorithme de sélection de colonnes le plus simple consiste à tirer chaque colonne indépendamment avec une probabilité p_i , qui traduit l'importance de la colonne $i \in [d]$ dans la matrice X . Une possibilité est d'utiliser directement les normes des colonnes [4]. Une autre possibilité, plus efficace, est de tirer parti des k -leviers, ou k -leverage scores en anglais [1].

Définition 1 (k -leviers). Soit $X = U\Sigma V^T \in \mathbb{R}^{N \times d}$ la décomposition en valeurs singulières de X . Pour $j \in [d]$, le k -levier de la j -ème colonne de X est défini par

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

La Définition 1 est illustrée Figure 1. Le k -levier ℓ_j^k est une mesure de l'alignement du j -ème axe canonique avec les k premières directions principales, cf. Figure 2. Nous reviendrons sur cette interprétation en Partie 4.1. En utilisant une distribution multinomiale de paramètre ℓ^k pour sélectionner des colonnes, on favorise donc des colonnes qui sont individuellement proches des axes principaux.

Théorème 1 ([1]). On tire s colonnes i.i.d. de distribution multinomiale de paramètre ℓ^k . Si $s \geq 3200k^2\epsilon^{-2}$, alors

$$\mathbb{P} \left(\|X - \Pi_{S,k}^{\text{Fr}} X\|_{\text{Fr}}^2 \leq (1 + \epsilon) \|X - \Pi_k X\|_{\text{Fr}}^2 \right) \geq 0.7. \quad (3)$$

L'échantillonnage multinomial par k -leviers permet donc d'avoir des erreurs d'approximation relatives à l'ACP de rang k arbitrairement petites, mais au prix de conserver un nombre de colonnes largement supérieur à k .

2.3 Échantillonnage volumique

Un autre algorithme consiste à sélectionner exactement k colonnes de façon non indépendante pour éviter toute redondance. C'est le cas de l'échantillonnage volumique, ou *Volume Sampling* (VS).

Théorème 2 ([3]). Soit S un sous-ensemble de $[d]$ tiré selon

$$\mathbb{P}_{\text{VS}}(S) \propto \text{Det}(X_{:,S}^T X_{:,S}) \mathbb{1}_{\{|S|=k\}}. \quad (4)$$

Alors

$$\mathbb{E}_{\text{VS}} \|X - \Pi_{S,k}^{\text{Fr}} X\|_{\text{Fr}}^2 \leq (k+1) \|X - \Pi_k X\|_{\text{Fr}}^2, \quad (5)$$

et

$$\mathbb{E}_{\text{VS}} \|X - \Pi_{S,k}^2 X\|_2^2 \leq (d-k)(k+1) \|X - \Pi_k X\|_2^2. \quad (6)$$

L'échantillonnage multinomial en Partie 2.2 sélectionnait les colonnes les plus importantes de la matrice X sans tenir compte de leur complémentarité éventuelle. L'échantillonnage volumique favorise quant à lui les sous-ensembles S qui correspondent à des colonnes *diverses* dans \mathbb{R}^N . Pour comprendre cette diversité, supposons que les colonnes de X sont des échantillons i.i.d. de moyenne nulle. L'équation (4) dit alors que l'échantillonnage volumique sélectionne S proportionnellement au déterminant de la matrice de covariance des colonnes $X_{:,S}$ indexées par S .

Nous analysons maintenant l'échantillonnage volumique à travers la lentille des processus ponctuels déterminantaux.

3 Processus ponctuels déterminantaux

On rappelle dans cette section quelques notions de la théorie des processus ponctuels déterminantaux discrets (PPD; [5]).

Définition 2 (PPD). Soit $K \in \mathbb{R}^{d \times d}$ une matrice semi-définie positive. Un sous-ensemble aléatoire noté $Y \subset [d]$ est tiré d'un PPD de noyau marginal K si et seulement si

$$\forall S \subset [d], \quad \mathbb{P}(S \subset Y) = \text{Det}(K_S), \quad (7)$$

avec $K_S = [K_{i,j}]_{i,j \in S}$. Par convention $\text{Det}(K_\emptyset) = 1$.

Étant donnée une matrice K , (7) ne permet pas nécessairement de définir correctement une distribution de probabilité. Une condition suffisante est que K soit symétrique et de spectre inclus dans $[0, 1]$. Si la matrice K est symétrique, le PPD correspondant est une distribution répulsive dans le sens suivant : pour tout $i, j \in [d]$,

$$\mathbb{P}(\{i, j\} \subset Y) = K_{i,i}K_{j,j} - K_{i,j}^2 \quad (8)$$

$$= \mathbb{P}(\{i\} \subset Y) \mathbb{P}(\{j\} \subset Y) - K_{i,j}^2 \quad (9)$$

$$\leq \mathbb{P}(\{i\} \subset Y) \mathbb{P}(\{j\} \subset Y). \quad (10)$$

Dans le cas extrême où le spectre de K est inclus dans $\{0, 1\}$, le noyau et le PPD correspondants sont dits *de projection*. Les PPD de projection ont la particularité de tirer des échantillons Y qui ont toujours le même cardinal avec probabilité 1, égal au nombre des valeurs propres de K égales à 1. Ces PPD

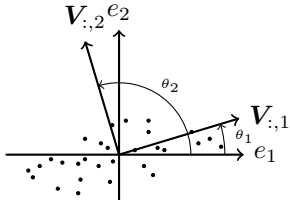


FIGURE 2 – Interprétation géométrique des k -leviers dans le cas $k = 1$: $\ell_1^k = \cos^2 \theta_1$, $\ell_2^k = \cos^2 \theta_2$.

de projection apparaissent en filigrane dans l'échantillonnage volumique de la Partie 2.3, qui peut-être vu comme un mélange de PPD de projection. En effet, si on note $\lambda_1 \geq \dots \geq \lambda_d$ les valeurs propres de $\mathbf{X}^\top \mathbf{X}$ et r le rang de \mathbf{X} , on a [5] :

$$\mathbb{P}_{\text{VS}}(Y = S) = \sum_{\substack{T \subset [r] \\ |T|=k}} \mu_T \left[\frac{1}{k!} \text{Det} \left(\mathbf{V}_{T,S} \mathbf{V}_{T,S}^\top \right) \right] \quad (11)$$

avec

$$\mu_T \propto \prod_{i \in T} \lambda_i, \quad \sum_{\substack{T \subset [r] \\ |T|=k}} \mu_T = 1. \quad (12)$$

Lorsque d est grand devant N , le coût du calcul de la matrice \mathbf{V} est en Nd^2 . Une fois \mathbf{V} calculée, le coût d'échantillonnage d'un PPD du mélange (11) est $\mathcal{O}(dk^2)$. Une analyse plus détaillée du coût d'échantillonnage est fournie dans [5] ou [2].

4 Un PPD bien choisi

On considère l'algorithme de sélection de colonnes suivant : on tire les colonnes S avec une probabilité

$$\mathbb{P}(S) \propto \text{Det} \mathbf{V}_{S,[k]} (\mathbf{V}_{S,[k]})^\top \quad (13)$$

Cet échantillonnage correspond au PPD de projection de noyau $\mathbf{K} = \mathbf{V}_{:, [k]} (\mathbf{V}_{:, [k]})^\top$, qui est aussi le PPD de plus grand poids dans le mélange (11).

4.1 Interprétation géométrique

Pour $j \in [d]$, soit \mathbf{e}_j le j -ème vecteur de la base canonique de \mathbb{R}^d et soit θ_j l'angle entre \mathbf{e}_j et le sous-espace $\mathcal{P}_k = \text{Vect}(\mathbf{V}_{:, [k]})$. Soit $\Pi_{\mathcal{P}_k} \mathbf{e}_j$ la projection orthogonale de \mathbf{e}_j sur \mathcal{P}_k . Alors

$$\cos^2(\theta_j) := \frac{(\mathbf{e}_j, \Pi_{\mathcal{P}_k} \mathbf{e}_j)^2}{\|\Pi_{\mathcal{P}_k} \mathbf{e}_j\|^2} = (\mathbf{e}_j, \Pi_{\mathcal{P}_k}(\mathbf{e}_j)) \quad (14)$$

$$= (\mathbf{e}_j, \sum_{i=1}^k V_{j,i} \mathbf{V}_{:,i}) = \sum_{i=1}^k V_{j,i}^2 = \ell_j^k. \quad (15)$$

Un k -levier élevé indique que \mathbf{e}_j est presque aligné avec \mathcal{P}_k . La sélection des colonnes qui ont un k -levier élevé comme dans le Théorème 1 peut donc être vue comme le remplacement du sous-espace principal \mathcal{P}_k par un sous-espace formé par des colonnes qui sont individuellement alignées avec \mathcal{P}_k . Pour éviter la redondance entre colonnes, on peut mesurer la ressemblance entre \mathcal{P}_k et $\mathcal{P}_S = \text{Vect}(\mathbf{e}_j, j \in S)$ d'une façon naturelle à l'aide des *angles principaux* [7, Section 6.4.4]. Formellement,

$$\cos^2(\mathcal{P}_k, \mathcal{P}_S) := \text{Det}(\mathbf{V}_{S,[k]})^2. \quad (16)$$

Ainsi, le PPD de noyau (4) est une généralisation naturelle du tirage multinomial avec les k -leviers : selon (7), chaque colonne est toujours marginalement tirée proportionnellement à son k -levier, mais on ajoute la contrainte que le sous-ensemble \mathcal{P}_S soit une bonne approximation de \mathcal{P}_k .

4.2 Bornes sous hypothèse de parcimonie

Avant de présenter les résultats, on définit la parcimonie des k -leviers comme le nombre de k -leviers non-nuls :

$$p = |\{j \in [d], \ell_j^k \neq 0\}|. \quad (17)$$

La décroissance des valeurs singulières $\sigma_1 \geq \dots \geq \sigma_r$ de \mathbf{X} est caractérisée par

$$\beta = \sigma_{k+1}^2 \left(\frac{1}{d-k} \sum_{j \geq k+1} \sigma_j^2 \right)^{-1} \in [1, d-k]. \quad (18)$$

Ces quantités contrôlent la qualité de l'approximation sous la loi du PPD de noyau marginal (4). La parcimonie p des k -leviers reflète le nombre de coordonnées initiales qui influent le calcul des axes principaux. La constante β reflète l'importance de σ_{k+1}^2 relativement à la moyenne des valeurs singulières au carré $(\sigma_i^2)_{i \geq k+1}$. Une grande valeur de β indique une coupure entre les directions principales d'ordre $> k+1$ relativement à celles d'ordre $\leq k+1$. Notons que cette constante β n'intervient ci-après que dans les garanties en norme de Frobenius. Nous énonçons maintenant nos résultats.

Proposition 1. *Sous la loi du PPD de noyau marginal \mathbf{K} ,*

$$\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_{S,k}^2 \mathbf{X}\|_2^2 \leq (1 + (p-k)k) \|\mathbf{X} - \Pi_k \mathbf{X}\|_2^2 \quad (19)$$

et

$$\mathbb{E}_{\text{DPP}} \|\mathbf{X} - \Pi_{S,k}^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \leq \left(1 + \beta \frac{p-k}{d-k} k \right) \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}}^2. \quad (20)$$

Notre borne spectrale (19) est meilleure que la borne spectrale (6) de l'échantillonnage volumique d'un facteur $(p-k)/(d-k)$. Plus le profil des k -leviers est parcimonieux, meilleure est la qualité de l'approximation. Quand il n'y a pas de coupure de spectre ($\beta \approx 1$), la borne (20) est meilleure que la borne de Frobenius (5) de l'échantillonnage volumique.

4.3 Cas de la parcimonie approchée

En général, aucun des k -leviers n'est strictement nul et la constante p est égale à d . Toutefois, il est courant qu'un grand nombre de k -leviers soient très petits, voir [8] pour des exemples concrets. De la même manière que l'ACP tire parti d'une dimension effective petite d'un problème, nous souhaitons un résultat qui profite de cette quasi-parcimonie des k -leviers et qui rende compte de l'amélioration observée expérimentalement par rapport à l'échantillonnage volumique.

On définit la parcimonie effective de la façon suivante. Soit π une permutation de $[d]$ telle que les k -leviers soient ordonnés :

$$\ell_{\pi_1}^k \geq \ell_{\pi_2}^k \geq \dots \geq \ell_{\pi_d}^k. \quad (21)$$

Soit $\delta \in [d]$, et $T_\delta = [\pi_\delta, \dots, \pi_d]$. Soit $\theta \geq 1$ et

$$p_{\text{eff}}(\theta) = \min \left\{ q \in [d] \mid \sum_{i \leq q} \ell_{\pi_i}^k \geq k - 1 + \frac{1}{\theta} \right\}. \quad (22)$$

Sachant que par construction, le vecteur ℓ^k somme à k , le paramètre θ contrôle dans (22) le défaut de somme cumulée des $p_{\text{eff}}(\theta)$ plus grands k -leviers.

Proposition 2. Soit \mathcal{A}_θ l'événement $\{S \cap T_{p_{\text{eff}}(\theta)} = \emptyset\}$. Alors,

$$\mathbb{P}_{\text{DPP}}(\mathcal{A}_\theta) \geq \frac{1}{\theta}, \quad (23)$$

$$\mathbb{E}_{\text{DPP}} [\|\mathbf{X} - \Pi_{S,k}^2 \mathbf{X}\|_2^2 \mid \mathcal{A}_\theta] \leq (1 + (p_{\text{eff}}(\theta) - k + 1)(k - 1 + \theta)) \|\mathbf{X} - \Pi_k \mathbf{X}\|_2^2, \quad (24)$$

et

$$\mathbb{E}_{\text{DPP}} [\|\mathbf{X} - \Pi_{S,k}^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2 \mid \mathcal{A}_\theta] \leq \left(1 + \beta \frac{(p_{\text{eff}}(\theta) + 1 - k)}{d - k} (k - 1 + \theta) \right) \|\mathbf{X} - \Pi_k \mathbf{X}\|_{\text{Fr}}^2. \quad (25)$$

Dans la Proposition 2, la parcimonie effective $p_{\text{eff}}(\theta)$ remplace la parcimonie p de la Proposition 1. L'idée est de conditionner le PPD à ne pas contenir les colonnes qui ont des k -leviers faibles, représentées par l'ensemble de « mauvais » indices $T_{p_{\text{eff}}(\theta)}$. En pratique, le conditionnement peut se réaliser à l'aide d'une étape de réjection : on tire des sous ensemble $S \sim \text{DPP}(\mathbf{K})$ jusqu'à réaliser la condition $S \cap T_{p_{\text{eff}}(\theta)} = \emptyset$. Le temps d'attente moyen est alors contrôlé par (23).

Enfin, on a présenté ici des garanties en termes de normes matricielles de l'erreur d'approximation. Lorsqu'on connaît les étapes d'apprentissage qui vont être réalisées après la réduction de dimension, l'objectif peut changer. Nous proposons dans [8] des garanties sur l'erreur de prédiction d'une régression linéaire après sélection de colonnes.

5 Simulations numériques

Cette partie illustre les résultats théoriques à l'aide de simulations numériques sur des jeux de données synthétiques. On compare l'erreur d'approximation $\mathbb{E} \|\mathbf{X} - \Pi_S^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2$ divisée par celle de l'ACP pour l'échantillonnage volumique et l'échantillonnage avec PPD de noyau (4). Pour cela, on génère des matrices \mathbf{X} de petite taille $N \times d = 100 \times 20$, si bien qu'on peut énumérer toutes les combinaisons de $k = 5$ colonnes. On génère 200 telles matrices aléatoires \mathbf{X} avec des profils de k -leviers parcimonieux, à l'aide d'un algorithme dont la description détaillée est fournie dans l'article [8]. Sur la Figure 3, pour chaque matrice on représente deux croix : en bleu l'erreur d'approximation de l'échantillonnage volumique, en rouge celle du PPD. En trait plein bleu, la borne théorique (5) de l'échantillonnage volumique est indépendante du niveau de parcimonie p , tandis que la borne théorique du PPD dépend linéairement de p . On remarque que les bornes

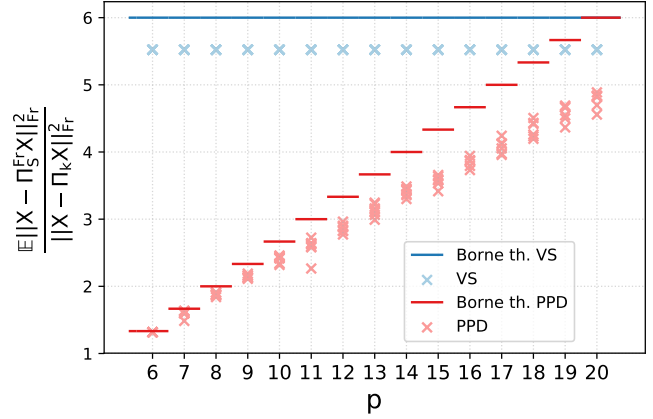


FIGURE 3 – L'erreur d'approximation $\mathbb{E} \|\mathbf{X} - \Pi_{S,k}^{\text{Fr}} \mathbf{X}\|_{\text{Fr}}^2$ et la borne théorique pour l'échantillonnage volumique (bleu) et le PPD (rouge), en fonction de la parcimonie p . La valeur 1 en ordonnée correspond à l'erreur de l'ACP.

théoriques reflètent bien les comportements empiriques des deux approches. La présence de parcimonie est favorable au PPD, comme attendu. D'autres simulations et exemples réels sont disponibles dans [8].

6 Conclusion

On a proposé un nouvel algorithme de sélection aléatoire de colonnes utilisant un processus ponctuel déterminantal inspiré de l'ACP. On a présenté des garanties théoriques de cet algorithme, qui dépassent l'état de l'art dans les cas où les k -leviers sont parcimonieux. Les expériences numériques confirment la précision des bornes obtenues.

Références

- [1] Drineas, P. and Mahoney, M. W. and Muthukrishnan, S. *Relative-error CUR matrix decompositions*. SIAM J MATRIX ANAL A, 2008.
- [2] Tremblay, N. and Barthelme, S. and Amblard, P.O. *Optimized algorithms to sample determinantal point processes*. ArXiv preprint, arXiv :1802.08471, 2018.
- [3] Deshpande, A. and Rademacher, L. and Vempala, S. and Wang, G. *Matrix Approximation and Projective Clustering via Volume Sampling*. Proc. ACM-SIAM SODA, 2006.
- [4] Drineas, P. and Frieze, A. and Kannan, R. and Vempala, S. and Vinay, V. *Clustering Large Graphs via the Singular Value Decomposition*. Mach. Learn., 2004.
- [5] Kulesza, A. and Taskar, B. *Determinantal point processes for machine learning*. FTML, 2012.
- [6] Macchi, O. *The coincidence approach to stochastic point processes*. Adv. Appl. Probab., 1975.
- [7] Golub, G. H. and Van Loan, C. F. *Matrix Computations (3rd Ed.)* Johns Hopkins University Press, 1996
- [8] Belhadji, A. and Bardenet, R. and Chainais, P. *A determinantal point process for column subset selection*. ArXiv preprint, arXiv :1812.09771, 2018.