

Polynomial time algorithms for dual volume sampling

Authored by:

Suvrit Sra
Stefanie Jegelka
Chengtao Li

Abstract

We study dual volume sampling, a method for selecting k columns from an $n \times m$ short and wide matrix ($n \leq k \leq m$) such that the probability of selection is proportional to the volume spanned by the rows of the induced submatrix. This method was proposed by Avron and Boutsidis (2013), who showed it to be a promising method for column subset selection and its multiple applications. However, its wider adoption has been hampered by the lack of polynomial time sampling algorithms. We remove this hindrance by developing an exact (randomized) polynomial time sampling algorithm as well as its derandomization. Thereafter, we study dual volume sampling via the theory of real stable polynomials and prove that its distribution satisfies the “Strong Rayleigh” property. This result has numerous consequences, including a provably fast-mixing Markov chain sampler that makes dual volume sampling much more attractive to practitioners. This sampler is closely related to classical algorithms for popular experimental design methods that are to date lacking theoretical analysis but are known to empirically work well.

1 Paper Body

A variety of applications share the core task of selecting a subset of columns from a short, wide matrix A with n rows and $m \geq n$ columns. The criteria for selecting these columns typically aim at preserving information about the span of A while generating a well-conditioned submatrix. Classical and recent examples include experimental design, where we select observations or experiments [38]; preconditioning for solving linear systems and constructing low-stretch spanning trees (here A is a version of the node-edge incidence matrix and we select edges in a graph) [6, 4]; matrix approximation [11, 13, 24]; feature selection in k -means clustering [10, 12]; sensor selection [25] and graph signal processing [14, 41]. In this work, we study a randomized approach that holds promise for all of these applications. This approach relies on sampling columns of A according

to a probability distribution defined over its submatrices: the probability of selecting a set S of k columns from A , with $n \leq k \leq m$, is $P(S; A) = \det(AS A_i^{\perp} S) / \det(AS A_i^{\perp} S)$,

(1.1)

where AS is the submatrix consisting of the selected columns. This distribution is reminiscent of volume sampling, where $k \leq n$ columns are selected with probability proportional to the determinant $\det(A_i^{\perp} S AS)$ of a $k \times k$ matrix, i.e., the squared volume of the parallelepiped spanned by the selected columns. (Volume sampling does not apply to $k \leq n$ as the involved determinants vanish.) In contrast, $P(S; A)$ uses the determinant of an $n \times n$ matrix and uses the volume spanned by the rows formed by the selected columns. Hence we refer to $P(S; A)$ -sampling as dual volume sampling (DVS). Contributions. Despite the ostensible similarity between volume sampling and DVS, and despite the many practical implications of DVS outlined below, efficient algorithms for DVS are not known and were raised as open questions in [6]. In this work, we make two key contributions: • We develop polynomial-time randomized sampling algorithms and their derandomization for DVS. Surprisingly, our proofs require only elementary (but involved) matrix manipulations. 31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

• We establish that $P(S; A)$ is a Strongly Rayleigh measure [8], a remarkable property that captures a specific form of negative dependence. Our proof relies on the theory of real stable polynomials, and the ensuing result implies a provably fast-mixing, practical MCMC sampler. Moreover, this result implies concentration properties for dual volume sampling. In parallel with our work, [16] also proposed a polynomial time sampling algorithm that works efficiently in practice. Our work goes on to further uncover the hitherto unknown Strong Rayleigh property of DVS, which has important consequences, including those noted above. 1.1 Connections and implications. The selection of $k \leq n$ columns from a short and wide matrix has many applications. Our algorithms for DVS hence have several implications and connections; we note a few below. Experimental design. The theory of optimal experiment design explores several criteria for selecting the set of columns (experiments) S . Popular choices are $S \in \arg\min_{S \subseteq \{1, \dots, m\}} J(AS)$, with

$$\begin{aligned} J(AS) &= \|A^{\perp} S\|_F^2 = \|AS A_i^{\perp} S\|_F^2 \\ J(AS) &= \|A^{\perp} S\|_2^2 \text{ (E-optimal design)}, J(AS) = \\ &1 / \|A^{\perp} S\|_F^2 \text{ (A-optimal design)}, \\ &1 / \log \det(AS A_i^{\perp} S) \text{ (D-optimal design)}. \end{aligned} \quad (1.2)$$

Here, A^{\perp} denotes the Moore-Penrose pseudoinverse of A , and the minimization ranges over all S such that AS has full row rank n . A-optimal design, for instance, is statistically optimal for linear regression [38]. •

Finding an optimal solution for these design problems is NP-hard; and most discrete algorithms use local search [33]. Avron and Boutsidis [6, Theorem 3.1] show that dual volume sampling yields an approximation guarantee for both A- and E-optimal design: if S is sampled from $P(S; A)$, then $\|A^{\perp} S\|_2^2 \leq (n+1) \|A^{\perp} S\|_F^2$.

$\lambda_1 \text{ATc} \text{ } 2 \text{ } R(n \text{ } 2 \text{ } p \text{ } 1 \text{ } 6 \text{ } C=6 \text{ } 4$
 $\text{ } 2 \text{ } 1 \text{ } (\text{AT} \text{ }) + "$
 $0 \text{ } \dots$
 p
 $r(\text{AT} \text{ }))?(m \text{ } -T \text{ } -)$
 0
 \dots
 $, 3$
 $7 \text{ } \dots 7 \text{ } Q \lambda_1 \text{AT} \text{ } 2 \text{ } Rr(\text{AT} \text{ })?(m \text{ } c \text{ } 5 \text{ } \dots$
 $1 \text{ } 2 \text{ } 2 \text{ } (\text{AT} \text{ }) + "$
 \dots
 $-T \text{ } -)$
 $.$

λ_1 Let $Q_B = \text{diag}(\lambda_1^2(B))Q$ of $B \lambda_1 B$ where $Q_B \text{ } 2 \text{ } R \text{ } -Tc \text{ } -?r(B)$. More B be $?$
 $?$ the eigenvalue decomposition $\lambda_1 \lambda_1 ? \lambda_1 ?$ over, let $W = ITc ; C$ and $= ek -T$
 $- r(B) (W ((Q_B) Q_B)W \lambda_1)$. Then the marginal probability of T in DVS is
 $hQ \text{ } i \text{ } hQ \text{ } i \text{ } r(\text{AT} \text{ }) \text{ } 2 \text{ } r(B) \text{ } 2 \text{ } i \text{ } (\text{AT} \text{ }) ? j \text{ } (B) ? i=1 \text{ } j=1 \text{ } P (T ? S; A) = . \text{ } ZA$

We prove Theorem 2 via a perturbation argument that connects DVS to volume sampling. Specifically, observe that for $?$ $\lambda_1 0$ and $-S - n$ it holds that
 $! ? \lambda_1 ? A \text{ } S \text{ } S \text{ } \lambda_1 \text{ } n \text{ } k \text{ } \lambda_1 \text{ } n \text{ } k \text{ } p \text{ } \det(AS \text{ } AS + "In) = " \det(AS \text{ } AS + "Ik) = " \det$
 $p . (2.1) " (Im)S " (Im)S$ Carefully letting $?$ $! 0$ bridges volumes with $?dual?$
 volumes. The technical remainder of the proof further relates this equality to singular values, and exploits properties of characteristic polynomials. A similar argument yields an alternative proof of Lemma 1. We show the proofs in detail in Appendix A and B respectively. Complexity. The numerator of $P (T ? S; A)$ in Theorem 2 requires $O(mn^2)$ time to compute the first term, $O(mn^2)$ to compute the second and $O(m^3)$ to compute the third. The denominator takes $O(mn^2)$ time, amounting in a total time of $O(m^3)$ to compute the marginal probability.

2.2 Sampling The marginal probabilities derived above directly yield a polynomial-time exact DVS algorithm. $!$ Instead of k -sets, we sample ordered k -tuples $S = (s_1, \dots, s_k) \text{ } 2 \text{ } [m]^k$. We denote the k -tuple $!$ variant of the DVS distribution by $P (?; A): Yk ! ! 1 P ((sj = ij)_{kj=1} ; A) = P (\{i_1, \dots, i_k\}; A) = P (sj = ij -s_1 = i_1, \dots, s_{j-1} = i_{j-1}; A)$. $j=1 \text{ } k! !$ Sampling S is now straightforward. At the j th step we sample s_j via $P (sj = ij -s_1 = i_1, \dots, s_{j-1} = i_{j-1}; A)$; these probabilities are easily obtained from the marginals in Theorem 2. Corollary 3. Let $T = \{i_1, \dots, i_{t+1}\}$, and $P (T ? S; A)$ as in Theorem 2. Then, $! P (T [\{i\} ? S; A) P (st = i; A -s_1 = i_1, \dots, s_{t+1} = i_{t+1}) = . (k \text{ } t + 1) P (T ? S; A)$

As a result, it is possible to draw an exact dual volume sample in time $O(km^4)$.

The full proof may be found in the appendix. The running time claim follows since the sampling algorithm invokes $O(mk)$ computations of marginal probabilities, each costing $O(m^3)$ time. 3

Remark A potentially more efficient approximate algorithm could be derived by noting the relations between volume sampling and DVS. Specifically, we add

a small perturbation to DVS as in Equation 2.1 to transform it into a volume sampling problem, and apply random projection for more efficient volume sampling as in [17]. Please refer to Appendix C for more details. 2.3

Derandomization

Next, we derandomize the above sampling algorithm to deterministically select a subset that satisfies the bound (1.3) for the Frobenius norm, thereby answering another question in [6]. The key insight for derandomization is that conditional expectations can be computed in polynomial time, given the marginals in Theorem 2: ! Corollary 4. Let $(i_1, \dots, i_t) \in [m]^t$ be such that the marginal distribution satisfies $P(s_1 = i_1, \dots, s_t = i_t; A) \geq 0$. The conditional expectation can be expressed as $\mathbb{E}[P_0(\{i_1, \dots, i_t\} \mid S) \mid S] = P(S; A[n]_{\{j\}}) \mid_{j=1}^t \mathbb{E}[kAS \mid kF \mid s_1 = i_1, \dots, s_t = i_t = , P_0(\{i_1, \dots, i_t\} \mid S) \mid S] = P(S; A)$ where P_0 are the unnormalized marginal distributions, and it can be computed in $O(nm^3)$ time. We show the full derivation in Appendix D. ! Corollary 4 enables a greedy derandomization procedure. Starting with the empty tuple $S_0 = ;$, in ! the i th iteration, we greedily select $j \in [n] \setminus \{i_1, \dots, i_t\}$ such that $\mathbb{E}[kAS \mid kF \mid (s_1, \dots, s_i) = S_{i-1} \cup j]$ and append !!! it to our selection: $S_i = S_{i-1} \cup j$. The final set is the non-ordered version S_k of S . Theorem 5 shows that this greedy procedure succeeds, and implies a deterministic version of the bound (1.3). Theorem 5. The greedy derandomization selects a column set S satisfying $kAS \mid kF \leq$

$$\begin{aligned} & m \cdot k \\ & n+1 \leq kAS \mid kF \leq n+1 \\ & kAS \mid kF \leq \\ & n(m \cdot n + 1) \leq kAS \mid kF \leq k \cdot n+1 \end{aligned}$$

In the proof, we construct a greedy algorithm. In each iteration, the algorithm computes, for each column that has not yet been selected, the expectation conditioned on this column being included in the current set. Then it chooses the element with the lowest conditional expectation to actually be added to the current set. This greedy inclusion of elements will only decrease the conditional expectation, thus retaining the bound in Theorem 5. The detailed proof is deferred to Appendix E. Complexity. Each iteration of the greedy selection requires $O(nm^3)$ to compute $O(m)$ conditional expectations. Thus, the total running time for k iterations is $O(knm^4)$. The approximation bound for the spectral norm is slightly worse than that in (1.3), but is of the same order if $k = O(n)$.

3

Strong Rayleigh Property and Fast Markov Chain Sampling

Next, we investigate DVS more deeply and discover that it possesses a remarkable structural property, namely, the Strongly Rayleigh (SR) [8] property. This property has proved remarkably fruitful in a variety of recent contexts, including recent progress in approximation algorithms [23], fast sampling [2, 27], graph sparsification [22, 39], extensions to the Kadison-Singer problem [1], and certain concentration of measure results [37], among others. For DVS, the SR property has two major consequences: it leads to a fast mixing practical MCMC sampler, and it implies results on concentration of measure. Strongly Rayleigh

measures. SR measures were introduced in the landmark paper of Borcea et al. [8], who develop a rich theory of negatively associated measures. In particular, we say that a probability measure $\mu : 2^{[n]} \rightarrow \mathbb{R}^+$ is negatively associated if $\mu(F \cap G) \leq \mu(F)\mu(G)$ for $F, G \subseteq [n]$ increasing functions on 2 with disjoint support. This property reflects a ‘repelling’ nature of μ , a property that occurs more broadly across probability, combinatorics, physics, and other fields; see [36, 8, 42] and references therein. The negative association property turns out to be quite subtle in general; the class of SR measures captures a strong notion of negative association and provides a framework for analyzing such measures. 4

Specifically, SR measures are defined via their connection to real stable polynomials [36, 8, 42]. A multivariate polynomial $f \in \mathbb{C}[z]$ where $z \in \mathbb{C}^m$ is called real stable if all its coefficients are real and $f(z) \neq 0$ whenever $\text{Im}(z_i) \leq 0$ for $1 \leq i \leq m$. A measure is called an SR measure if its multivariate generating polynomial $f_\mu(z) := \sum_{S \subseteq [n]} \mu(S) \prod_{i \in S} z_i$ is real stable. Notable examples of SR measures are Determinantal Point Processes [31, 29, 9, 26], balanced matroids [19, 37], Bernoulli conditioned on their sum, among others. It is known (see [8, pg. 523]) that the class of SR measures is exponentially larger than the class of determinantal measures. 3.1

Strong Rayleigh Property of DVS

Theorem 6 establishes the SR property for DVS and is the main result of this section. Here and in the following, we use the notation $z \in \mathbb{C}^m$.

Theorem 6. Let $A \in \mathbb{R}^{n \times m}$ and $n \leq k \leq m$. Then the multiaffine polynomial $X \mapsto \det(A_S + \sum_{i \in S} z_i A_{i, :})$ is real stable.

$$(3.1) \quad \det(A_S + \sum_{i \in S} z_i A_{i, :})$$

is real stable. Consequently, $P(S; A)$ is an SR measure. The proof of Theorem 6 relies on key properties of real stable polynomials and SR measures established in [8]. Essentially, the proof demonstrates that the generating polynomial of $P(S; A)$ can be obtained by applying a few carefully chosen stability preserving operations to a polynomial that we know to be real stable. Stability, although easily destroyed, is closed under several operations noted in the important proposition below. Proposition 7 (Prop. 2.1 [8]). Let $f : \mathbb{C}^m \rightarrow \mathbb{C}$ be a stable polynomial. The following properties preserve stability: (i) Substitution: $f(z_1, \dots, z_m)$ for $z_i \in \mathbb{C}$; (ii) Differentiation: $\frac{\partial}{\partial z_i} f(z_1, \dots, z_m)$ for any $i \in [m]$; (iii) Diagonalization: $f(z, z, z, \dots, z)$ is stable, and hence $f(z, z, \dots, z)$; and (iv) Inversion: $z \mapsto 1/z$ for $z \neq 0$. In addition, we need the following two propositions for proving Theorem 6. Proposition 8 (Prop. 2.4 [7]). Let B be Hermitian, $z \in \mathbb{C}^m$ and A_i ($1 \leq i \leq m$) be Hermitian semidefinite matrices. Then, the following polynomial is stable: $X \mapsto \det(B + \sum_{i=1}^m z_i A_i)$. (3.2)

Proposition 9. For $n \leq k \leq m$ and $L := A A^T$, we have $\det(A_S + \sum_{i \in S} z_i A_{i, :}) = \det(L_S + \sum_{i \in S} z_i L_{i, :})$.

Proof. Let $Y = \text{Diag}([y_i]_{i=1}^m)$ be a diagonal matrix. Using the Cauchy-Binet identity we have $X \mapsto \det(A_S + \sum_{i \in S} z_i A_{i, :}) = \det((A_S + \sum_{i \in S} z_i A_{i, :})^T Y) = \det(A_S^T Y + \sum_{i \in S} z_i A_{i, :}^T Y)$.

$$-T = -n, T \text{ ?}[m]$$

Thus, when $Y = IS$, the (diagonal) indicator matrix for S , we obtain AY
 $A_i = AS A_i S$. Consequently, in the summation above only terms with $T \text{ ? } S$
survive, yielding $X \det(AS A_i \det(A_i \det(LT, T)) = \text{en}(LS, S) \cdot S = T \det(AT)$
 $= -T = -n, T \text{ ? } S$

$$-T = -n, T \text{ ? } S$$

We are now ready to sketch the proof of Theorem 6. Proof. (Theorem 6).
Notationally, it is more convenient to prove that the "complement" polynomial
 $P_{Sc}^{pc}(z) := -S = -k, S \text{ ?}[m] \det(AS A_i$ is stable; subsequently, an application
of Prop. 7-(iv) yields $S \text{ ? } z$ stability of (3.1). Using matrix notation $W = \text{Diag}(w_1$
 $, \dots, w_m)$, $Z = \text{Diag}(z_1, \dots, z_m)$, our starting stable polynomial (this
stability follows from Prop. 8) is $h(z, w) := \det(L + W + Z)$, which can be
expanded as $X \det(h(z, w)) = \det(WS + LS) \text{ ? } z \text{ ? } S = S \text{ ?}[m]$

$$5$$

$$w \text{ ? } C_m, z \text{ ? } C_m,$$

$$S \text{ ?}[m]$$

$$\text{?}X$$

$$T \text{ ? } S$$

$$\text{? } wST \det(LT, T) \text{ ? } z \text{ ? } S.$$

Thus, $h(z, w)$ is real stable in $2m$ variables, indexed below by S and R where
 $R := ST$. Instead of the form above, We can sum over $S, R \text{ ? } [m]$ but then
have to constrain the support to the case when $Sc T = ;$ and $Sc R = ;$. In other
words, we may write (using Iverson-brackets $J \text{ ? } K$) $X h(z, w) = JSc R = ; Sc T$
 $= ;K \det(LT, T) \text{ ? } z \text{ ? } Sc wR \cdot (3.3) S, R \text{ ?}[m]$

Next, we truncate polynomial (3.3) at degree $(m k) + (k n) = m n$ by restrict-
ing $-Sc [R = m n$. By [8, Corollary 4.18] this truncation preserves stability,
whence $X H(z, w) := JSc R = ;K \det(LSR, SR) \text{ ? } z \text{ ? } Sc wR, S, R \text{ ?}[m] -Sc [R = m$
 n

is also stable. Using Prop. 7-(iii), setting $w_1 = \dots = w_m = y$ retains
stability; thus $X g(z, y) := H(z, (y, y, \dots, y)) = JSc R = ;K \det(LSR, SR)$
 $\text{? } z \text{ ? } Sc y -R = -\{z\} S, R \text{ ?}[m] -Sc [R = m n$

$$m \text{ times}$$

$$=$$

$$X \text{ ?}X$$

$$S \text{ ?}[m]$$

$$-T = -n, T \text{ ? } S$$

$$\text{? } \det(LT, T) \text{ ? } y -S -$$

$$-T - Sc$$

$$z$$

$$=$$

$$X$$

$$\text{en}(LS, S) \text{ ? } y -S -$$

$$n \text{ ? } Sc$$

$$z,$$

$$S \text{ ?}[m]$$

is also stable. Next, differentiating $g(z, y)$, k times with respect to y and evaluating at 0 preserves stability (Prop. 7-(ii) and (i)). In doing so, only terms corresponding to $|S| = k$ survive, resulting in

$$\begin{aligned} & \frac{\partial^k g(z, y)}{\partial y^k} \bigg|_{y=0} \\ &= \frac{1}{k!} \sum_{|S|=k} \det(A_S) z^S, \end{aligned}$$

which is just $p(z)$ (up to a constant); here, the last equality follows from Prop. 9. This establishes stability of $p(z)$ and hence of $p(z)$. Since $p(z)$ is in addition multiaffine, it is the generating polynomial of an SR measure, completing the proof. 3.2

Implications: MCMC

The SR property of $P(S; A)$ established in Theorem 6 implies a fast mixing Markov chain for sampling S . The states for the Markov chain are all sets of cardinality k . The chain starts with a randomly-initialized active set S , and in each iteration we swap an element $s_{in} \in S$ with an element $s_{out} \notin S$ with a specific probability determined by the probability of the current and proposed set. The stationary distribution of this chain is the one induced by DVS, by a simple detailed-balance argument. The chain is shown in Algorithm 1.

Algorithm 1 Markov Chain for Dual Volume Sampling
Input: $A \in \mathbb{R}^{n \times m}$ the matrix of interest, k the target cardinality, T the number of steps
Output: $S \sim P(S; A)$
Initialize $S \subseteq [m]$ such that $|S| = k$ and $\det(A_S) \neq 0$ for $i = 1$ to T do
draw $b \in \{0, 1\}$ uniformly
if $b = 1$ then
Pick $s_{in} \in S$ and $s_{out} \notin S$ uniformly randomly
 $\alpha = \min(1, \frac{\det(A_{S \setminus \{s_{in}\} \cup \{s_{out}\}})}{\det(A_S)})$
 $S \leftarrow S \setminus \{s_{in}\} \cup \{s_{out}\}$ with probability α
end if
end for

The convergence of the markov chain is measured via its mixing time: The mixing time of the chain indicates the number of iterations t that we must perform (starting from S_0) before we can consider S_t as an approximately valid sample from $P(S; A)$. Formally, if $S_0(t)$ is the total variation distance between the distribution of S_t and $P(S; A)$ after t steps, then $\tau_{S_0}(\epsilon) := \min\{t :$

$$S_0(t) \leq \epsilon\}$$

is the mixing time to sample from a distribution ϵ -close to $P(S; A)$ in terms of total variation distance. We say that the chain mixes fast if τ_{S_0} is polynomial

$P(S; A) / \exp\{\log \det(AS A_i \text{ to zero essen} S) / \}$ with temperature parameter β . Driving β tially recovers Fedorov's method, while our results imply fast mixing for $\beta = 1$, together with approximation guarantees. Through this lens, simulated annealing may be viewed as initializing Fedorov's method with the fast-mixing sampler. In practice, we observe that letting $\beta \rightarrow 1$ improves the approximation results, which opens interesting questions for future work.

4

Experiments

We report selection performance of DVS on real regression data (CompAct, CompAct(s), Abalone and Bank32NH1) for experimental design. We use 4,000 samples from each dataset for estimation. We compare against various baselines, including uniform sampling (Unif), leverage score sampling (Lev) [30], predictive length sampling (PL) [45], the sampling (Smpl)/greedy (Greedy) selection methods in [43] and Fedorov's exchange algorithm [20]. We initialize the MCMC sampler with Kmeans++ [5] for DVS and run for 10,000 iterations, which empirically yields selections that are 1

<http://www.dcc.fc.up.pt/~ltorgo/Regression/DataSets.html>

7

sufficiently good. We measure performances via (1) the prediction error $\|X - \hat{X}\|_2$, (2) running times. Figure 1 shows the results for these three measures with sample sizes k varying from 60 to 200. Further experiments (including for the interpolation $\beta \rightarrow 1$), may be found in the appendix. Unif Lev PL Smpl Greedy DVS Fedorov

Error

0.35 0.3 0.25 0.2

100

Running Time

80

0.26

60

0.24

40 20

0.15 150

200

0.22 0.2

0 100

Time-Error Trade-off

0.28

Error

Prediction Error

Seconds

0.4

0.18 100

150

k

k

200
0
10
20 30 Seconds
40
50

Figure 1: Results on the CompAct(s) dataset. Results are the median of 10 runs, except Greedy and Fedorov. Note that Unif, Lev, PL and DVS use less than 1 second to finish experiments. In terms of prediction error, DVS performs well and is comparable with Lev. Its strength compared to the greedy and relaxation methods (Smpl, Greedy, Fedorov) is running time, leading to good time-error tradeoffs. These tradeoffs are illustrated in Figure 1 for $k = 120$. In other experiments (shown in Appendix G) we observed that in some cases, the optimization and greedy methods (Smpl, Greedy, Fedorov) yield better results than sampling, however with much higher running times. Hence, given time-error tradeoffs, DVS may be an interesting alternative in situations where time is a very limited resource and results are needed quickly.

5 Conclusion

In this paper, we study the problem of DVS and develop an exact (randomized) polynomial time sampling algorithm as well as its derandomization. We further study dual volume sampling via the theory of real-stable polynomials and prove that its distribution satisfies the “Strong Rayleigh” property. This result has remarkable consequences, especially because it implies a provably fastmixing Markov chain sampler that makes dual volume sampling much more attractive to practitioners. Finally, we observe connections to classical, computationally more expensive experimental design methods (Fedorov’s method and SA); together with our results here, these could be a first step towards a better theoretical understanding of those methods. Acknowledgement This research was supported by NSF CAREER award 1553284, NSF grant IIS-1409802, DARPA grant N66001-17-1-4039, DARPA FunLoL grant (W911NF-16-1-0551) and a Siebel Scholar Fellowship. The views, opinions, and/or findings contained in this article are those of the author and should not be interpreted as representing the official views or policies, either expressed or implied, of the Defense Advanced Research Projects Agency or the Department of Defense.

2 References

- [1] N. Anari and S. O. Gharan. The Kadison-Singer problem for strongly Rayleigh measures and applications to asymmetric TSP. arXiv:1412.1143, 2014.
- [2] N. Anari and S. O. Gharan. Effective-resistance-reducing flows and asymmetric TSP. In IEEE Symposium on Foundations of Computer Science (FOCS), 2015.
- [3] N. Anari, S. O. Gharan, and A. Rezaei. Monte Carlo Markov chain algorithms for sampling strongly Rayleigh distributions and determinantal point processes. In COLT, pages 23–26, 2016.
- [4] M. Arioli and I. S. Duff. Precondi-

tioning of linear least-squares problems by identifying basic variables. *SIAM J. Sci. Comput.*, 2015. [5] D. Arthur and S. Vassilvitskii. k-means++: The advantages of careful seeding. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, 2007. 8

[6] H. Avron and C. Boutsidis. Faster subset selection for matrices and applications. *SIAM Journal on Matrix Analysis and Applications*, 34(4):1464?1499, 2013. [7] J. Borcea and P. Br?nd?n. Applications of stable polynomials to mixed determinants: Johnson?s conjectures, unimodality, and symmetrized Fischer products. *Duke Mathematical Journal*, pages 205?223, 2008. [8] J. Borcea, P. Br?nd?n, and T. Liggett. Negative dependence and the geometry of polynomials. *Journal of the American Mathematical Society*, 22:521?567, 2009. [9] A. Borodin. Determinantal point processes. *arXiv:0911.1153*, 2009. [10] C. Boutsidis and M. Magdon-Ismail. Deterministic feature selection for k-means clustering. *IEEE Transactions on Information Theory*, pages 6099?6110, 2013. [11] C. Boutsidis, M. W. Mahoney, and P. Drineas. An improved approximation algorithm for the column subset selection problem. In *SODA*, pages 968?977, 2009. [12] C. Boutsidis, A. Zouzias, M. W. Mahoney, and P. Drineas. Stochastic dimensionality reduction for k-means clustering. *arXiv preprint arXiv:1110.2897*, 2011. [13] C. Boutsidis, P. Drineas, and M. Magdon-Ismail. Near-optimal column-based matrix reconstruction. *SIAM Journal on Computing*, pages 687?717, 2014. [14] S. Chen, R. Varma, A. Sandryhaila, and J. Kova?cevi?c. Discrete signal processing on graphs: Sampling theory. *IEEE Transactions on Signal Processing*, 63(24):6510?6523, 2015. [15] A. ?ivril and M. Magdon-Ismail. On selecting a maximum volume sub-matrix of a matrix and related problems. *Theoretical Computer Science*, pages 4801?4811, 2009. [16] M. Derezhinski and M. K. Warmuth. Unbiased estimates for linear regression via volume sampling. *Advances in Neural Information Processing Systems (NIPS)*, 2017. [17] A. Deshpande and L. Rademacher. Efficient volume sampling for row/column subset selection. In *Foundations of Computer Science (FOCS)*, 2010 51st Annual IEEE Symposium on, pages 329?338. IEEE, 2010. [18] M. Elkin, Y. Emek, D. A. Spielman, and S.-H. Teng. Lower-stretch spanning trees. *SIAM Journal on Computing*, 2008. [19] T. Feder and M. Mihail. Balanced matroids. In *Symposium on Theory of Computing (STOC)*, pages 26?38, 1992. [20] V. Fedorov. Theory of optimal experiments. Preprint 7 lsm, Moscow State University, 1969. [21] V. Fedorov. Theory of optimal experiments. Academic Press, 1972. [22] A. Frieze, N. Goyal, L. Rademacher, and S. Vempala. Expanders via random spanning trees. *SIAM Journal on Computing*, 43(2):497?513, 2014. [23] S. O. Gharan, A. Saberi, and M. Singh. A randomized rounding approach to the traveling salesman problem. In *IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 550?559, 2011. [24] N. Halko, P.-G. Martinsson, and J. A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*, 53(2):217?288, 2011. [25] S. Joshi and S. Boyd. Sensor selection via convex optimization. *IEEE Transactions on Signal Processing*, pages 451?462, 2009. [26] A. Kulesza and B. Taskar. Determinantal Point Processes for machine learning, volume 5. *Foundations and Trends in Machine Learning*, 2012. [27] C. Li, S. Jegelka, and

S. Sra. Fast mixing markov chains for strongly Rayleigh measures, DPPs, and constrained sampling. In *Advances in Neural Information Processing Systems (NIPS)*, 2016. [28] C. Li, S. Sra, and S. Jegelka. Gaussian quadrature for matrix inverse forms with applications. In *ICML*, pages 1766?1775, 2016. [29] R. Lyons. Determinantal probability measures. *Publications Mathématiques de l’Institut des Hautes Études Scientifiques*, 98(1):167?212, 2003.

9

[30] P. Ma, M. Mahoney, and B. Yu. A statistical perspective on algorithmic leveraging. In *Journal of Machine Learning Research (JMLR)*, 2015. [31] O. Macchi. The coincidence approach to stochastic point processes. *Advances in Applied Probability*, 7(1), 1975. [32] A. Magen and A. Zouzias. Near optimal dimensionality reductions that preserve volumes. In *Approximation, Randomization and Combinatorial Optimization. Algorithms and Techniques*, pages 523?534. Springer, 2008. [33] A. J. Miller and N.-K. Nguyen. A fedorov exchange algorithm for d-optimal design. *Journal of the royal statistical society*, 1994. [34] M. D. Morris and T. J. Mitchell. Exploratory designs for computational experiments. *Journal of Statistical Planning and Inference*, 43:381?402, 1995. [35] N.-K. Nguyen and A. J. Miller. A review of some exchange algorithms for constructing discrete optimal designs. *Computational Statistics and Data Analysis*, 14:489?498, 1992. [36] R. Pemantle. Towards a theory of negative dependence. *Journal of Mathematical Physics*, 41: 1371?1390, 2000. [37] R. Pemantle and Y. Peres. Concentration of Lipschitz functionals of determinantal and other strong Rayleigh measures. *Combinatorics, Probability and Computing*, 23:140?160, 2014. [38] F. Pukelsheim. Optimal design of experiments. *SIAM*, 2006. [39] D. Spielman and N. Srivastava. Graph sparsification by effective resistances. *SIAM J. Comput.*, 40(6):1913?1926, 2011. [40] D. A. Spielman and S.-H. Teng. Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems. In *STOC*, 2004. [41] M. Tsitsvero, S. Barbarossa, and P. D. Lorenzo. Signals on graphs: Uncertainty principle and sampling. *IEEE Transactions on Signal Processing*, 64(18):4845?4860, 2016. [42] D. Wagner. Multivariate stable polynomials: theory and applications. *Bulletin of the American Mathematical Society*, 48(1):53?84, 2011. [43] Y. Wang, A. W. Yu, and A. Singh. On Computationally Tractable Selection of Experiments in Regression Models. *ArXiv e-prints*, 2016. [44] Y. Zhao, F. Pasqualetti, and J. Cortés. Scheduling of control nodes for improved network controllability. In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pages 1859? 1864, 2016. [45] R. Zhu, P. Ma, M. W. Mahoney, and B. Yu. Optimal subsampling approaches for large sample linear regression. *arXiv preprint arXiv:1509.05111*, 2015.

10