

Subsampling Algorithms for Semidefinite Programming.

Alexandre d'Aspremont*

April 20, 2009

Abstract

We derive a stochastic gradient algorithm for semidefinite optimization using randomization techniques. The algorithm uses subsampling to reduce the computational cost of each iteration and the subsampling ratio explicitly controls the algorithm's granularity, i.e. the tradeoff between cost per iteration and total number of iterations. Furthermore, the total computational cost is directly proportional to the complexity (i.e. rank) of the solution. We study numerical performance on some large-scale problems arising in statistical learning.

1 Introduction

Beyond classic combinatorial relaxations [GW95], semidefinite programming has recently found a new stream of applications in machine learning [LCB⁺02], geometry [WS06], statistics [dBEG06] or graph theory [SBXD05]. All these problems have a common characteristic: they have relatively low precision targets but form very large semidefinite programs for which obtaining second order models is numerically hopeless, which means that Newton based interior point solvers typically fail before completing even a single iteration. Early efforts focused on exploiting structural properties of the problem (sparsity, block patterns, etc), but this has proven particularly hard for semidefinite programs. For very large problem instances, first-order methods remain at this point the only credible alternative. This follows a more general trend in optimization which seeks to significantly reduce the *granularity* of solvers, i.e. reduce the per iteration complexity of optimization algorithms rather than their total numerical complexity, thus allowing at least some progress to be made on problems that are beyond the reach of current algorithms.

In this work, we focus on the following spectral norm minimization problem:

$$\begin{aligned} & \text{minimize} && \left\| \sum_{j=1}^p y_j A_j + C \right\|_2 - b^T y \\ & \text{subject to} && y \in Q, \end{aligned} \tag{1}$$

in the variable $y \in \mathbf{R}^p$, with parameters $A_j \in \mathbf{S}_n$, for $j = 1, \dots, p$, $b \in \mathbf{R}^p$ and $C \in \mathbf{S}_n$, where Q is a compact convex set. Throughout the paper, we also implicitly assume that the set $Q \subset \mathbf{R}^p$ is simple enough so that the complexity of projecting y on Q is relatively low compared to the other steps in the algorithm.

The idea behind this paper stems from a recent result by [JLNS09], who used a mirror descent stochastic approximation algorithm for solving bilinear matrix games (see [Nes09], [PJ92] or [NY83]

*ORFE Department, Princeton University, Princeton, NJ 08544. aspremon@princeton.edu

for more background), where subsampling is used to perform matrix vector products and produce an approximate gradient. Strikingly, the algorithm has a total complexity of $O(n \log n / \epsilon^2)$, when the problem matrix is $n \times n$, hence only requires access to a negligible proportion of the matrix coefficients as the dimension n tends to infinity.

In parallel, recent advances in large deviations and random matrix theory have produced a stream of new randomization results for high dimensional linear algebra (see [FKV04, DKM06, AM07] among many others), motivated by the need to perform these operations on very large scale, sometimes streaming, data sets in applications such as machine learning, signal processing, etc. Similar subsampling techniques have been successfully applied to support vector machine classification [KBH08] or Fourier decomposition. Randomization results were used in [AK07] to solve certain semidefinite programs arising in combinatorial relaxations of graph problems. Randomization was also used in [BLO02] and [BLO05] to approximate subdifferentials of functions that are only differentiable almost everywhere.

Our contribution here is to further reduce the granularity of first-order semidefinite programming solvers by combining subsampling procedures with stochastic approximation algorithms to derive stochastic gradient methods for spectral norm minimization with very low complexity per iteration. In practice, significantly larger per iteration complexity and memory requirements mean that interior point techniques often fail to complete a single iteration on very large problem instances. CPU clock also runs much faster than RAM, so operations small enough to be performed entirely in cache (which runs at full speed) are much faster than those requiring larger data sets. Solver performance on very large problem instances is then often more constrained by memory bandwidth than clock speed, hence everything else being equal, algorithms running many cheap iterations will be much faster than those requiring fewer, more complex ones. Here, subsampling techniques allow us to produce semidefinite optimization algorithms with very low cost per iteration, where all remaining $O(n^2)$ operations have a small constant and can be performed in a single pass over the data.

We also observe that the relative approximation error in computing the spectral norm (or trace norm) of a matrix using subsampling is directly proportional to the numerical rank of that matrix, hence another important consequence of using subsampling techniques to solve large-scale semidefinite programs is that the total complexity of running the algorithm becomes explicitly dependent on the complexity (i.e. rank) of its solution.

The paper is organized as follows. Section 2 surveys some key results on randomized linear algebra and spectral norm approximations. In Section 3 we then derive a stochastic approximation algorithm for spectral norm minimization with very low cost per iteration. In Section 4, we discuss some extensions to statistical learning problems. Finally, we present some numerical experiments in Section 5.

Notation

We write \mathbf{S}_n the set of symmetric matrices of dimension n . For a matrix $X \in \mathbf{R}^{m \times n}$, we write $\|X\|_F$ its Frobenius norm, $\|X\|_{\text{tr}}$ its trace norm, $\|X\|_2$ its spectral norm, $\sigma_i(X)$ its i -th largest singular value and let $\|X\|_\infty = \max_{ij} |X_{ij}|$, while $X^{(i)}$ is the i -th column of the matrix X and $X_{(i)}$ its i -th row. We write $\text{vec}(X)$ the vector of \mathbf{R}^{n^2} obtained by stacking up the columns of the matrix X and $\text{NumRank}(X)$ the numerical rank of the matrix X , where $\text{NumRank}(X) = \|X\|_F^2 / \|X\|_2^2$. Finally, when $x \in \mathbf{R}^n$ is a vector, we write $\|x\|_2$ its Euclidean norm, $\|\cdot\|$ is a general norm on \mathbf{R}^m and $\|\cdot\|_*$ its dual norm.

2 Randomized linear algebra

In this section, we survey several results by [DKM06] which, after a single pass on the data, sample matrix columns to approximate matrix products and produce low rank matrix approximations with a complexity of $O(sn)$ where s only depends on the target precision and not on n .

2.1 Randomized matrix multiplication

Algorithm 1 Matrix multiplication

Input: $A \in \mathbf{R}^{m \times n}$, $B \in \mathbf{R}^{n \times p}$ and s such that $1 \leq s \leq n$.

1: Define a probability vector $p \in \mathbf{R}^n$ such that

$$p_i = \frac{\|A^{(i)}\|_2 \|B_{(i)}\|_2}{\sum_{j=1}^n \|A^{(j)}\|_2 \|B_{(j)}\|_2}, \quad i = 1, \dots, n.$$

2: Define subsampled matrices $C \in \mathbf{R}^{m \times s}$ and $R \in \mathbf{R}^{s \times p}$ as follows.

3: **for** $i = 1$ to s **do**

4: Pick $j \in [1, n]$ with $\mathbf{P}(j = l) = p_l$.

5: Set $C^{(i)} = A^{(j)}/\sqrt{sp_j}$ and $R_{(i)} = B_{(j)}/\sqrt{sp_j}$.

6: **end for**

Output: Matrix product CR approximating AB .

By construction, we have $\mathbf{E}[CR] = AB$, and the following randomization result from [DKM07] controls the precision of the approximations in algorithm 1.

Lemma 1 *Let $A \in \mathbf{R}^{m \times n}$, $B \in \mathbf{R}^{n \times p}$, given a subsampling rate s such that $1 \leq s \leq n$, suppose that $C \in \mathbf{R}^{m \times s}$ and $R \in \mathbf{R}^{s \times p}$ are computed according to algorithm 1 above, then*

$$\mathbf{E}[\|AB - CR\|_F^2] \leq \frac{1}{s} \|A\|_F^2 \|B\|_F^2$$

and if $\beta \in [0, 1]$ with $\eta = 1 + \sqrt{8 \log(1/\beta)}$ then

$$\|AB - CR\|_F^2 \leq \frac{\eta^2}{s} \|A\|_F^2 \|B\|_F^2$$

with probability at least $1 - \beta$.

Proof. See Theorem 1 in [DKM07]. ■

Note that using the adaptive probabilities p_i is crucial here. The error bounds increase by a factor n when $p_i = 1/n$ for example.

2.2 Randomized low-rank approximation

Algorithm 2 below computes the leading singular vectors of the smaller matrix S , which is a subsampled and rescaled version of X . Here, the computational savings come from the fact that we only need to compute singular values of a matrix of dimension $m \times s$ with $s \leq n$. Recall that

Algorithm 2 Low-rank approximation

Input: $X \in \mathbf{R}^{m \times n}$ and k, s such that $1 \leq k \leq s < n$.

- 1: Define a probability vector $p \in \mathbf{R}^n$ such that $p_i = \|X^{(i)}\|_2^2 / \|X\|_F^2$, for $i = 1, \dots, n$.
- 2: Define a subsampled matrix $S \in \mathbf{R}^{m \times s}$ as follows.
- 3: **for** $i = 1$ to s **do**
- 4: Pick an index $j \in [1, n]$ with $\mathbf{P}(j = l) = p_l$.
- 5: Set $S^{(i)} = X^{(j)} / \sqrt{sp_j}$.
- 6: **end for**
- 7: Form the eigenvalue decomposition $S^T S = Y \mathbf{diag}(\sigma) Y^T$ where $Y \in \mathbf{R}^{s \times s}$ and $\sigma \in \mathbf{R}^s$.
- 8: Form a matrix $H \in \mathbf{R}^{m \times k}$ with $H^{(i)} = SY^{(i)} / \sigma_i^{1/2}$.

Output: Approximate singular vectors $H^{(i)}$, $i = 1, \dots, k$.

computing k leading eigenvectors of a symmetric matrix of dimension s only requires matrix vector products, hence can be performed in $O(ks^2)$ operations using iterative algorithms such as the power method or Lanczos method (see [GVL90, Chap. 8-9], as usual we omit the precision target in linear algebra operations, implicitly assuming that it is much finer than ϵ), so that the cost of computing k leading singular vectors of a matrix of size $m \times s$ is $O(kms)$.

This means that, given the probabilities p_i , the total cost of obtaining k approximate singular vectors using algorithm 2 is $O(ksm)$ instead of $O(knm)$ for exact singular vectors. Of course, computing p_i requires mn operations, but can be done very efficiently in a single pass over the data. We now recall the following result from [DKM06] which controls the precision of the approximations in algorithm 2.

Lemma 2 *Let $X \in \mathbf{R}^{m \times n}$ and $1 \leq k \leq s < n$. Given a precision target $\epsilon > 0$, if $s \geq 4/\epsilon^2$ and $H \in \mathbf{R}^{m \times k}$ is computed as in algorithm 2, we have*

$$\mathbf{E}[\|X - H_k H_k^T X\|_2^2] \leq \|X - X_k\|_2^2 + \epsilon \|X\|_F^2$$

and if in addition $s > 4\eta^2/\epsilon^2$ where $\eta = 1 + \sqrt{8 \log(1/\beta)}$ for $\beta \in [0, 1]$, then

$$\|X - H_k H_k^T X\|_2^2 \leq \|X - X_k\|_2^2 + \epsilon \|X\|_F^2$$

with probability at least $1 - \beta$, where X_k is the best rank k approximation of X .

Proof. See Theorem 4 in [DKM06]. ■

An identical precision bound holds in the Frobenius norm when $s \geq 4k/\epsilon^2$. We then show the following lemma on approximating the spectral radius of a symmetric matrix X using algorithm 2.

Lemma 3 *Let $X \in \mathbf{R}^{m \times n}$ and $\beta \in [0, 1]$. Given a precision target $\epsilon > 0$, construct a matrix $S \in \mathbf{R}^{m \times s}$ by subsampling the columns of X as in algorithm 2. Let $\eta = 1 + \sqrt{8 \log(1/\beta)}$ and*

$$s = \eta^2 \frac{\|X\|_2^2}{\epsilon^2} \mathbf{NumRank}(X)^2 \tag{2}$$

we have

$$\mathbf{E}[|\|S\|_2 - \|X\|_2|] \leq \epsilon$$

and

$$\|S\|_2 - \|X\|_2 \leq \epsilon$$

with probability at least $1 - \beta$.

Proof. Using the Hoffman-Wielandt inequality (see [SS90, Th. 3.1] or the proof of [DKM06, Th.2] for example) we get

$$\|S\|_2^2 - \|X\|_2^2 \leq \|SS^T - XX^T\|_F$$

and Jensen's inequality together with the matrix multiplication result in Lemma 1 yields

$$\mathbf{E}[\|SS^T - XX^T\|_F] \leq \frac{\|X\|_F^2}{\sqrt{s}}$$

and

$$\|SS^T - XX^T\|_F \leq \frac{\eta\|X\|_F^2}{\sqrt{s}}$$

with probability at least $1 - \beta$. Combining these two inequalities with the sampling rate in (2)

$$s = \eta^2 \frac{\|X\|_F^4}{\epsilon^2 \|X\|_2^2}$$

yields the desired result. ■

The subsampling rate required to achieve a precision target ϵ has a natural interpretation. Indeed

$$s = \eta^2 \frac{\|X\|_2^2}{\epsilon^2} \mathbf{NumRank}(X)^2$$

is simply the squared ratio of the numerical rank of the matrix X over the relative precision target $\epsilon/\|X\|_2$, times a factor η^2 controlling the confidence level. The numerical rank $\mathbf{NumRank}(X)$ always satisfies $\mathbf{NumRank}(X) = \|X\|_F^2/\|X\|_2^2 \leq \mathbf{Rank}(X)$ and can be seen as a stable relaxation of the rank of the matrix X (see [RV07] for a discussion). Note that, by construction, the subsampled matrix always has lower rank than the matrix X . The expectation bound is still valid if we drop the factor η .

3 Stochastic approximation algorithm

Below, we will use a stochastic approximation algorithm to solve problem (1) when the gradient is approximated using the subsampling algorithms detailed above. We focus on a stochastic approximation of problem (1), written:

$$\min_{y \in Q} f(y) \equiv \mathbf{E} \left[\left\| \pi^{(s)} \left(\sum_{j=1}^p y_j A_j + C \right) \right\|_2 \right] - b^T y \quad (3)$$

in the variable $y \in \mathbf{R}^p$ and parameters $A_j \in \mathbf{S}_n$, for $j = 1, \dots, p$, $b \in \mathbf{R}^p$ and $C \in \mathbf{S}_n$, with $1 \leq s \leq n$ controlling the sampling rate, where the function $\|\pi^{(s)}(\sum_{j=1}^p y_j A_j + C)\|_2$ and a subgradient with respect to y are computed using algorithms 1 and 2. For $X \in \mathbf{S}_n$, we have written $\pi^{(s)}(X)$ the subsampling/scaling operation used in algorithms 1 and 2 with

$$\pi^{(s)}(X) = S, \quad (4)$$

where $0 < s < n$ controls the sampling rate and $S \in \mathbf{R}^{n \times s}$ is the random matrix defined in algorithm 2 whose columns are a scaled sample of the columns of X . We will write $S = \pi_{(s)}(X)$ the matrix obtained by subsampling rows as in algorithm 1. We also define $\mathcal{A} \in \mathbf{R}^{n^2 \times p}$ as the matrix whose columns are given by $\mathcal{A}^{(j)} = \mathbf{vec}(A_j)$, $j = 1, \dots, p$.

3.1 Stochastic approximation algorithm

We show the following lemma approximating the gradient of the function $\|\pi^{(s)}(\sum_{j=1}^p y_j A_j + C)\|_2$ with respect to y and bounding its quadratic variation.

Lemma 4 *Given $A_j \in \mathbf{S}_n$, for $j = 1, \dots, p$, $b \in \mathbf{R}^p$, $C \in \mathbf{S}_n$ and sampling rates s_1 and s_2 , a (stochastic) subgradient of the function $\|\pi^{(s_1)}(\sum_{j=1}^p y_j A_j + C)\|_2 - b^T y$ with respect to y is given by the vector $w \in \mathbf{R}^p$ with*

$$w = \mathcal{A}^T \mathbf{vec}(vv^T) - b$$

where $v \in \mathbf{R}^n$ is a leading singular vector of the subsampled matrix $S = \pi^{(s_1)}(X)$ formed in algorithm 2. Furthermore, the product $\mathcal{A}^T \mathbf{vec}(vv^T)$ can be approximated using algorithm 1 to form an approximate gradient

$$g = \pi^{(s_2)}(\mathcal{A}^T) \pi_{(s_2)}(\mathbf{vec}(vv^T)) - b,$$

which satisfies

$$\mathbf{E}[g] \in \partial f(y) \quad \text{and} \quad \mathbf{E}[\|g\|_2^2] \leq M^* \equiv 2 \frac{\|\mathcal{A}\|_F^2}{s_2} + 2\|b\|_2^2. \quad (5)$$

Proof. Iterated expectations give $\mathbf{E}[g] = \mathbf{E}[w] \in \partial f(y)$. The sampling probabilities p_i used in approximating the matrix vector product $\mathcal{A}^T \mathbf{vec}(vv^T)$ following algorithm 1 are defined as

$$p_i = \frac{\|\mathcal{A}_{(i)}\|_2 |\mathbf{vec}(vv^T)_i|}{\sum_{j=1}^n \|\mathcal{A}_{(j)}\|_2 |\mathbf{vec}(vv^T)_j|}, \quad i = 1, \dots, n.$$

As in [DKM07, Lemma 3], the quadratic variation of the approximate product $\pi^{(s_2)}(\mathcal{A}^T) \pi_{(s_2)}(\mathbf{vec}(vv^T))$ is then given by

$$\mathbf{E}[\|\pi^{(s_2)}(\mathcal{A}^T) \pi_{(s_2)}(\mathbf{vec}(vv^T))\|_F^2] = \sum_{i=1}^{n^2} \frac{\|\mathcal{A}_{(i)}\|_2^2 \mathbf{vec}(vv^T)_i^2}{s_2 p_i}.$$

With p_i defined as above, we get

$$\begin{aligned} \sum_{i=1}^{n^2} \frac{\|\mathcal{A}_{(i)}\|_2^2 \mathbf{vec}(vv^T)_i^2}{s_2 p_i} &\leq \frac{\left(\sum_{i=1}^{n^2} \|\mathcal{A}_{(i)}\|_2 \mathbf{vec}(vv^T)_i\right)^2}{s_2} \\ &\leq \frac{\|\mathcal{A}\|_F^2}{s_2} \end{aligned}$$

by Cauchy-Schwarz, because $\|\mathbf{vec}(vv^T)\|_2^2 = \|vv^T\|_F^2 = \|v\|_2^4 = 1$, hence the desired result. \blacksquare

We now use this result to produce an explicit bound on the complexity of solving problems (3) and (1) by subsampling using a stochastic approximation algorithm. In this section, we let $\|\cdot\|$ be a general norm on \mathbf{R}^p , we write $\|\cdot\|_*$ its dual norm and define $\delta^*(p)$ as the smallest number such

that $\|y\|_2 \leq \delta^*(p)^{1/2} \|y\|_*$ for all $y \in \mathbf{R}^p$. Following the notation in [JLNS09, §2.3], we let $\omega(x)$ be a distance generating function, i.e. a function such that

$$Q^o = \left\{ x \in Q : \exists y \in \mathbf{R}^m, x \in \underset{u \in Q}{\operatorname{argmin}} [y^T u + \omega(u)] \right\}$$

is a convex set. We assume that $\omega(x)$ is strongly convex on Q^o with modulus α with respect to the norm $\|\cdot\|$, which means

$$(y - x)^T (\nabla \omega(y) - \nabla \omega(x)) \geq \alpha \|y - x\|^2, \quad x, y \in Q^o.$$

We then define a prox-function $V(x, y)$ on $Q^o \times Q$ as follows:

$$V(x, y) \equiv \omega(y) - [\omega(x) + \nabla \omega(x)^T (y - x)],$$

which is nonnegative and strongly convex with modulus α with respect to the norm $\|\cdot\|$. The prox-mapping associated to V is then defined as

$$P_x^{Q, \omega}(y) \equiv \underset{z \in Q}{\operatorname{argmin}} \{y^T (z - x) + V(x, z)\}. \quad (6)$$

Finally, we define the ω diameter of the set Q as:

$$D_{\omega, Q} \equiv (\max_{z \in Q} \omega(z) - \min_{z \in Q} \omega(z))^{1/2} \quad (7)$$

and we let γ_l for $l = 0, \dots, N$ be a step size strategy.

Algorithm 3 Spectral norm minimization using subsampling

Input: Matrices $A_j \in \mathbf{S}_n$, for $j = 1, \dots, p$, $b \in \mathbf{R}^p$ and $C \in \mathbf{S}_n$, sampling rates s_1 and s_2 .

- 1: Pick initial $y_0 \in Q$
- 2: **for** $l = 1$ to N **do**
- 3: Compute $v \in \mathbf{R}^n$, the leading singular vector of the matrix $\pi^{(s_1)}(\sum_{j=1}^p y_{l,j} A_j + C)$, subsampled according to algorithm 2 with $k = 1$ and $s = s_1$.
- 4: Compute the approximate subgradient $g_l = \pi^{(s_2)}(\mathcal{A}^T) \pi_{(s_2)}(\mathbf{vec}(v v^T)) - b$, by subsampling the matrix product using algorithm 1 and $s = s_2$.
- 5: Set $y_{l+1} = P_{y_l}^{Q, \omega}(\gamma_l g_l)$.
- 6: Update the running average $\tilde{y}_N = \sum_{k=0}^N \gamma_k y_k / \sum_{k=0}^N \gamma_k$.
- 7: **end for**

Output: An approximate solution $\tilde{y}_N \in \mathbf{R}^p$ of problem (3) with high probability.

The following results control the convergence of the robust stochastic approximation algorithm 3 (see [JLNS09], [Nes09], [PJ92] or [NY83] for further details). We call \bar{y} the optimal solution of problem (3), the lemma below characterizes convergence speed in expectation.

Lemma 5 Given $N > 0$, let M^* be defined as in (5) by

$$M^* = 2 \frac{\|\mathcal{A}\|_F^2}{s_2} + 2 \|b\|_2^2,$$

using a fixed step size strategy with

$$\gamma_l = \frac{D_{\omega, Q}}{\delta^*(p)M^*} \sqrt{\frac{2}{\alpha N}}, \quad l = 1, \dots, N$$

we have, after N iterations of algorithm 3

$$\mathbf{E}[f(\tilde{y}_N) - f(\bar{y})] \leq D_{\omega, Q} \delta^*(p) M^* \sqrt{\frac{2}{\alpha N}}$$

and

$$f(\tilde{y}_N) - f(\bar{y}) \geq \epsilon$$

with probability less than $\frac{D_{\omega, Q} \delta^*(p) M^*}{\epsilon} \sqrt{\frac{2}{\alpha N}}$.

Proof. By construction $\mathbf{E}[\|g\|_*^2] \leq \delta^*(p)M^*$, the rest follows from [JLNS09, §2.3] for example. ■

Lemma 5 means that we need at most

$$N = \frac{2D_{\omega, Q}^2 (\delta^*(p)M^*)^2}{\alpha \epsilon^2 \beta^2}$$

iterations to get an ϵ solution to problem (3) with confidence at least $1 - \beta$. We call y^* the solution to the deterministic spectral norm minimization problem (1). We can now bound the suboptimality of \tilde{y}_N in problem (1) with high probability.

Theorem 1 *If the sampling rate s_1 is set to*

$$s_1 = \frac{\left\| \sum_{j=1}^p y_j^* A_j + C \right\|_2^2}{\epsilon^2} \mathbf{NumRank} \left(\sum_{j=1}^p y_j^* A_j + C \right)^2 \quad (8)$$

then after

$$N = \frac{2D_{\omega, Q}^2 (\delta^*(p)M^*)^2}{\alpha \epsilon^2 \beta^2} \quad (9)$$

iterations of algorithm 3, we have

$$\left\| \sum_{j=1}^p \tilde{y}_{N,j} A_j + C \right\|_2 - b^T \tilde{y}_N - \left\| \sum_{j=1}^p y_j^* A_j + C \right\|_2 + b^T y^* \leq \epsilon$$

with probability at least $1 - \beta$.

Proof. Recall that we have written y^* the solution to the deterministic problem (1), \bar{y} the solution to the stochastic problem (3) and \tilde{y}_N the N -th iterate of algorithm 3. With s_1 defined as above, Lemma 3 means that

$$\mathbf{E} \left[\left| \left\| \sum_{j=1}^p y_j^* A_j + C \right\|_2 - \left\| \pi^{(s)} \left(\sum_{j=1}^p y_j^* A_j + C \right) \right\|_2 \right| \right] \leq \epsilon$$

and Jensen's inequality yields

$$\left| \left\| \left(\sum_{j=1}^p y_j^* A_j + C \right) \right\|_2 - b^T y^* - f(y^*) \right| \leq \epsilon.$$

which bounds the difference between the minimum of problem (1) and the value $f(y^*)$ of its stochastic approximation in (3). Lemma 5 then shows that

$$f(\tilde{y}_N) - f(\bar{y}) \geq \epsilon$$

with probability less than β . We conclude using the fact that $f(\bar{y}) \leq f(y^*)$ by definition. ■

This result allows us to bound the *oracle* complexity of solving (1) by subsampling. In practice of course, both the spectral norm and the numerical rank of the solution matrix $\sum_{j=1}^p y_j^* A_j + C$ are unknown and one needs to check for convergence using the duality gap computed in what follows, then adjust the precision target if necessary (in at most $\log_2 n$ binary search steps on s_1 for example). This means that, everything else being fixed, subsampling decreases the complexity of the main step in the algorithm by a factor $\log_2 n/n$ in the worst case. The next section provides a detailed analysis of the complexity of algorithm 3 as a function of ϵ, s_1 and s_2 .

3.2 Complexity

We now study in detail the complexity of algorithm 3. Suppose we are given a precision target ϵ and fix the sampling rate s_2 arbitrarily between 1 and n^2 , with the sampling rate s_1 set as in Theorem 1. The cost of each iteration in algorithm 3 breaks down as follows.

- On line 3: Computing the leading singular vector v , using algorithm 2 with $k = 1$. This means first computing the probabilities p_i at a cost of $O(n^2)$ operations. Forming the matrix $S = \pi^{(s_1)}(\sum_{j=1}^p y_{l,j} A_j + C)$ costs $O(ns_1)$ operations. It remains to compute the leading singular vector of S using the Lanczos method (see [GVL90, Chap. 8-9]), at a cost of $O(ns_1)$. The total numerical cost of this step is then bounded by $c_1 n^2 + c_2 ns_1$ where c_1 and c_2 are absolute constants. Here, c_1 is always less than ten while c_2 is the number of iterations required by the Lanczos method to reach a fixed precision target (typically 1e-8 or better here) hence we have $c_1 \ll c_2$.
- On line 4: Computing the approximate subgradient $g_l = \pi^{(s_2)}(\mathcal{A}^T) \pi_{(s_2)}(\text{vec}(vv^T)) - b$, by subsampling the matrix product using algorithm 1. This means again forming the vector p at a cost of $O(n^2)$ (the row norms of \mathcal{A} can be precomputed). Computing the subsampled matrix vector product then costs $O(ps_2)$. Both of these complexity bounds have low constants.
- On line 5: Computing the projection $y_{l+1} = P_{y_l}^{Q,\omega}(\gamma_l g_l)$, whose numerical cost will be denoted by $c(p)$.

Let us remark in particular that all $O(n^2)$ operations above only require one pass over the data, which means that the entire data set does not need to fit in memory. Using the bound on the quadratic variation of the gradient computed in Lemma 4, we can then bound the number of iterations required by algorithm 3 to produce a ϵ -solution to problem (1) with probability at least $1 - \beta$. Let us call $Y^* = \sum_{j=1}^p y_j^* A_j + C$, and recall that $\eta = 1 + \sqrt{8 \log(1/\beta)}$, Table 1 summarizes these complexity bounds and compares them with complexity bounds for a stochastic approximation algorithm without subsampling.

Complexity	Stoch. Approx.	Stoch. Approx. with Subsampling
Per Iter.	$c_4 n^2 p + c(p)$	$c_1 n^2 + c_3 p s_2 + c_2 n \eta^2 \frac{\ Y^*\ _2^2}{\epsilon^2} \mathbf{NumRank}(Y^*)^2 + c(p)$
Num. Iter.	$\frac{2D_{\omega, Q}^2 \delta^*(p)^2 (\ A\ _F^2 + \ b\ _2^2)}{\alpha \epsilon^2 \beta^2}$	$\frac{2D_{\omega, Q}^2 \delta^*(p)^2 \left(\frac{\ A\ _F^2}{s_2} + \ b\ _2^2 \right)}{\alpha \epsilon^2 \beta^2}$

Table 1: Complexity of solving problem (1) using subsampled stochastic approximation method versus original algorithm. Here c_1, \dots, c_4 are absolute constants with $c_1, c_3 \ll c_2, c_4$.

We observe that subsampling affects the complexity of solving problem (1) in two ways. Decreasing the (matrix product) subsampling rate $s_2 \in [1, n^2]$ decreases the cost of each iterations but increases the number of iterations in the same proportion, hence has no explicit effect on the total complexity bound. In practice of course, because of higher cache memory speed and better bandwidth on smaller problems, cheaper iterations tend to run more efficiently than more complex ones. The impact of the (singular vector) subsampling rate $s_1 \in [1, n]$ is much more important however, since computing the leading eigenvector of the current iterate is the most complex step in the algorithm when solving problem (1) using stochastic approximation. Because $c_1, c_3 \ll c_2$, the per iteration complexity solving large-scale problems essentially follows

$$n \eta^2 \frac{\|Y^*\|_2^2}{\epsilon^2} \mathbf{NumRank}(Y^*)^2$$

hence explicitly depends on both the numerical rank of the solution matrix $Y^* = \sum_{j=1}^p y_j^* A_j + C$ and on the relative precision target $\epsilon/\|Y^*\|_2$. This means that problems with simpler solutions will be solved more efficiently than problems whose solutions has a high rank.

The choice of norm $\|\cdot\|$ and distance generating function also has a direct impact on complexity through $c(p)$ and $\delta^*(p)M^*$. Unfortunately here, subsampling error bounds are only available in the Frobenius and spectral norms hence part of the benefit of choosing optimal norm/distance generating function combinations is sometimes lost in the norm ratio bound $\delta^*(p)$.

Finally, subsampling can have a more subtle effect on complexity. By construction, solutions to problem (1) tend to have multiple leading singular values which coalesce near the optimum. Introducing noise by subsampling can potentially break this degeneracy and increase the gap between leading eigenvalues. Since the complexity of the algorithm depends in great part on the complexity of computing a leading singular vector using iterative methods such as the power method or the Lanczos method, and the complexity of these methods decreases as the gap between the two leading singular values increases, subsampling can also improve the efficiency of iterative singular value computations.

3.3 Surrogate Duality Gap

In practice, we often have no a priori knowledge of $\mathbf{NumRank}(Y^*)^2$ and if the sampling rate s is set too low, it's possible for the algorithm to terminate at a suboptimal point Y where the subsampling error is less than ϵ (if the error at the true optimal point Y^* is much larger than ϵ). In order to find the optimal sampling rate s by bisection, we first need to check for suboptimality in (1), so we show how to track convergence in algorithm 3 by computing a surrogate duality gap.

The dual of problem (1) is given by

$$\begin{aligned} & \text{maximize} && \mathbf{Tr}(CX) - S_Q(w) \\ & \text{subject to} && w_j = b_j - \mathbf{Tr}(A_j X), \quad j = 1, \dots, p \\ & && \|X\|_{\text{tr}} \leq 1, \end{aligned} \tag{10}$$

in the variables $X \in \mathbf{S}_n$ and $w \in \mathbf{R}^p$, where $S_Q(v)$ is the support function of the set Q , defined as

$$S_Q(w) \equiv \max_{y \in Q} w^T y.$$

For instance, when Q is an Euclidean ball of radius B , problem (10) becomes

$$\begin{aligned} & \text{maximize} && \mathbf{Tr}(CX) - B\|w\|_2 \\ & \text{subject to} && w_j = b_j - \mathbf{Tr}(A_j X), \quad j = 1, \dots, p \\ & && \|X\|_{\text{tr}} \leq 1, \end{aligned} \tag{11}$$

in the variables $X \in \mathbf{S}_n$ and $w \in \mathbf{R}^p$. The vector v in algorithm 3 always satisfies $\|v\|_{\text{tr}} \leq 1$, hence we can track convergence in solving (1) by computing the following surrogate duality gap

$$\left\| \sum_{j=1}^p y_j A_j + C \right\|_2 - b^T y - v^T C v + S_Q(w) \tag{12}$$

where $w_i = b_i - v^T A_i v$ for $i = 1, \dots, p$.

4 Extensions

In this section, motivated by applications in statistical learning, we discuss direct extensions of the results detailed above to the problem of minimizing the sum of the k largest singular values of an affine combination of matrices. We then briefly discuss other matrix subsampling techniques.

4.1 Minimizing the sum of the k largest singular values

Here, we focus on the problem of minimizing the sum of the k largest singular values of an affine combination of matrices

$$\min_{y \in Q} \sum_{i=1}^k \sigma_i \left(\sum_{j=1}^p y_j A_j + C \right) - b^T y \tag{13}$$

in the variable $y \in \mathbf{R}^p$, with parameters $A_j \in \mathbf{S}_n$, for $j = 1, \dots, p$, $b \in \mathbf{R}^p$ and $C \in \mathbf{S}_n$. As in the previous section, we also form its stochastic approximation

$$\min_{y \in Q} f(y) \equiv \mathbf{E} \left[\sum_{i=1}^k \sigma_i \left(\pi^{(s)} \left(\sum_{j=1}^p y_j A_j + C \right) \right) \right] - b^T y \tag{14}$$

in the variable $y \in \mathbf{R}^p$, with $1 \leq s \leq n$ controlling the sampling rate. We now prove an analog of Lemma 3 for this new objective function.

Lemma 6 Let $X \in \mathbf{R}^{m \times n}$ and $\beta \in [0, 1]$. Given a precision target $\epsilon > 0$, $k \geq 1$ and a matrix $S \in \mathbf{R}^{m \times s}$ constructed by subsampling the columns of X as in algorithm 2, let $\eta = 1 + \sqrt{8 \log(1/\beta)}$ and

$$s = \eta^2 \frac{(\sum_{i=1}^k \sigma_i(X))^2 \mathbf{NumRank}(X)^2}{\epsilon^2 k^2} \kappa(X)^4 \mathbf{Rank}(X) \quad (15)$$

where $\kappa(X) = \sigma_1(X)/\sigma_r(X)$ with $r = \min\{k, \mathbf{Rank}(X)\}$, we have

$$\mathbf{E} \left[\sum_{i=1}^k |\sigma_i(X) - \sigma_i(S)| \right] \leq \epsilon$$

and

$$\sum_{i=1}^k |\sigma_i(X) - \sigma_i(S)| \leq \epsilon$$

with probability at least $1 - \beta$.

Proof. Because $\mathbf{Rank}(SS^T) \leq \mathbf{Rank}(XX^T)$ by construction, we always have

$$\begin{aligned} \sum_{i=1}^k |\sigma_i^2(X) - \sigma_i^2(S)| &= \sum_{i=1}^k |\sigma_i(X) - \sigma_i(S)| (\sigma_i(X) + \sigma_i(S)) \\ &\geq \sigma_r(X) \sum_{i=1}^k |\sigma_i(X) - \sigma_i(S)| \end{aligned}$$

where $r = \min\{k, \mathbf{Rank}(X)\}$. Because the sum of the k largest singular values is a unitarily invariant norm on \mathbf{S}_n (see [HJ91, §3.4]), Mirsky's theorem (see [SS90, Th. 4.11] for example) shows that

$$\begin{aligned} \sum_{i=1}^k |\sigma_i^2(X) - \sigma_i^2(S)| &= \sum_{i=1}^k |\sigma_i(XX^T) - \sigma_i(SS^T)| \\ &\leq \sum_{i=1}^k \sigma_i(XX^T - SS^T) \end{aligned}$$

and because, by construction, the range of SS^T is included in the range of XX^T , we must have $\mathbf{Rank}(XX^T - SS^T) \leq \mathbf{Rank}(XX^T)$ and

$$\sum_{i=1}^k \sigma_i(XX^T - SS^T) \leq \sqrt{\mathbf{Rank}(X)} \|XX^T - SS^T\|_F$$

Jensen's inequality together with the matrix multiplication result in Lemma 1 yield

$$\mathbf{E}[\|SS^T - XX^T\|_F] \leq \frac{\|X\|_F^2}{\sqrt{s}}$$

and

$$\|SS^T - XX^T\|_F \leq \frac{\eta \|X\|_F^2}{\sqrt{s}}$$

with probability at least $1 - \beta$. Combining these inequalities with the sampling rate in (15)

$$s = \eta^2 \frac{\|X\|_F^4 \mathbf{Rank}(X)}{\epsilon^2 \sigma_r(X)^2}$$

and using

$$\frac{\|X\|_F^4}{(\sum_{i=1}^k \sigma_i(X))^2 \sigma_r(X)^2} \leq \frac{\mathbf{NumRank}(X)^2}{k^2} \kappa(X)^4$$

yields the desired result. ■

Once again, the subsampling rate in the above lemma has a clear interpretation, indeed

$$\eta^2 \frac{(\sum_{i=1}^k \sigma_i(X))^2}{\epsilon^2} \frac{\mathbf{NumRank}(X)^2}{k^2} \kappa(X)^4 \mathbf{Rank}(X)$$

is the product of a term representing relative precision, a term reflecting the rank of X and a term in $\kappa(X)$ representing its (pseudo) condition number. Note that the bound can be further refined when $\sigma_r \leq \epsilon$. Lemma 6 allows us to compute the gradient by subsampling when using algorithm 3 to solve problem (13). The remaining steps in the algorithm are identical, except that the matrix vv^T is replaced by a combination of matrices formed using the k leading singular vectors.

4.2 Other sampling techniques

For completeness, we recall below the subsampling procedure in [AM07]. The key idea behind this result is that, as the matrix dimension n grows and given a fixed, scale invariant precision target $\|X\|_F/\epsilon$, the norm $\|X\|_\infty$ of individual coefficients in X typically becomes negligible and we can randomly discard the majority of them while keeping important spectral features of X mostly intact.

Lemma 7 *Given $X \in \mathbf{S}_n$ and $\epsilon > 0$, we define a subsampled matrix S whose coefficients are independently distributed as:*

$$S_{ij} = \begin{cases} X_{ij}/p & \text{with probability } p, \\ 0 & \text{otherwise.} \end{cases} \quad (16)$$

when $i \geq j$, and $S_{ij} = S_{ji}$ otherwise. Assume that $1 \geq p \geq (8 \log n)^4/n$, then

$$\|X - S\|_2 \leq 4\|X\|_\infty \sqrt{n/p}.$$

with probability at least $1 - \exp(-19(\log n)^4)$.

Proof. See [AM07, Th. 1.4]. ■

At first sight here, bounding the approximation error means letting the probability p grow relatively fast as n tends to infinity. However, because $\|X\|_\infty/\epsilon$ is typically much smaller than $\|X\|_F/\epsilon$, this subsampling ratio p can often be controlled. Adaptive subsampling, i.e. letting p vary with the magnitude of the coefficients in X , can further improve these results (see [AM07, §4] for details). The average number of nonzero coefficients in the subsampled matrix can be bounded using the structure of X . Note that the constants in this result are all very large (in particular, $1 \geq p \geq (8 \log n)^4/n$ implies $n \geq 10^9$) so despite its good empirical performance in low dimensions, the result presented above has to be understood in an asymptotic sense.

5 Applications & numerical results

In this section, we first detail a few instances of problem (1) arising in statistical learning. We then study the numerical performance of the methods detailed here on large scale problems.

5.1 Matrix factorization and collaborative filtering

Matrix factorization methods have been heavily used to solve collaborative filtering problems (e.g. the *Netflix* problem) and we refer the reader to [Sre04], [Bac07], [RFP07] or [CR08] for details. All these references form a particular instance of problem (13), written

$$\begin{aligned} & \text{minimize} && \left\| \sum_{j=1}^p y_j A_j + C \right\|_{\text{tr}} - b^T y \\ & \text{subject to} && y \in Q, \end{aligned} \tag{17}$$

in the variable $y \in \mathbf{R}^p$ where Q is a low dimension norm ball for example and the matrices A_j have a block format with only a few nonzero coefficients. Here, the trace norm can be understood as a convex lower bound on the rank function (as in [FHB01]) but sometimes also has a direct interpretation in terms of learning (see [Sre04]). In this particular case, the complexity of the main step in the algorithm (i.e. computing the gradient) is controlled by the sampling rate in Lemma 6, which can be simplified here to

$$s = \eta^2 \frac{\|Y^*\|_{\text{tr}}^2}{\epsilon^2} \kappa(Y^*)^2 \mathbf{Rank}(Y^*)$$

where $Y^* = \sum_{j=1}^p y_j^* A_j + C$ and $\kappa(Y^*) = \sigma_1(Y^*)/\sigma_r(Y^*)$ with $r = \mathbf{Rank}(Y^*)$. The bound can be further refined when $\sigma_r \leq \epsilon$. In practice, the complexity of solving problem (17) can often be further reduced using the simple observation that an optimal solution of (13) will also be optimal in (17) whenever $\mathbf{Rank}(Y_k^*) < k$, where Y_k^* is the optimal solution to (13) here. Once again, the sampling rate s has a natural interpretation as the product of a relative precision term, a term reflecting the condition number of the solution and the rank of the optimal solution. It means in particular that problems whose solutions have a lower rank are explicitly easier to solve than problems with more complex solutions.

5.2 LASSO

Consider a particular instance of problem (13) written

$$\begin{aligned} & \text{minimize} && \|y\|_1 \\ & \text{subject to} && \|Ax - b\|_2 \leq \sigma \end{aligned} \tag{18}$$

in the variable $y \in \mathbf{R}^n$, with $A \in \mathbf{R}^{m \times n}$, $b \in \mathbf{R}^m$ and $\sigma > 0$. This is a (somewhat trivial) version of problem (13), where the matrices are diagonal and Q is an ellipsoid. This problem is directly related to LASSO, i.e. ℓ_1 -penalized regression (see [Tib96]). In the diagonal case, the low rank matrix approximation produced by algorithm 2 simply picks s coefficients of y with probability proportional to their magnitude. Here, computing the gradient is trivial, but computational savings come from the fact that the bound on the quadratic variation of the gradient in (5) is now equal to the sampling rate, so $M^* = s$ in the complexity estimate (9). Since s is chosen as above, with

$$s = \eta^2 \frac{\|y^*\|_1}{\epsilon^2} \kappa(y^*)^2 \mathbf{Card}(y^*)$$

where $\kappa(y^*) = y_{[1]}/y_{[r]}$ with $r = \mathbf{Card}(y^*)$, this means that the complexity of solving problem (18) is (explicitly) proportional to the cardinality of the solution $\mathbf{Card}(y^*)$. Of course, this algorithm is not competitive with specialized algorithms for solving (18), but this subsampling bound provides some theoretical support for the empirical observations made in [DT06] using homotopy methods.

5.3 Fastest mixing Markov chain on a graph

As in [BDX04], suppose we are given a connected graph with vertex set $\mathcal{V} = \{1, \dots, n\}$ and edge set $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$, with $(i, j) \in \mathcal{E} \Leftrightarrow (j, i) \in \mathcal{E}$, where all vertices have a self-loop. We define a Markov chain on this graph with transition probability matrix $P \in \mathbf{R}^{n \times n}$, where

$$P_{ij} = \mathbf{Prob}(X_{t+1} = j | X_t = i), \quad i, j = 1, \dots, n$$

This matrix satisfies $P = P^T$ and $P\mathbf{1} = \mathbf{1}$, which means that the equilibrium distribution of this Markov chain is uniform and the largest singular value of P is equal to one. The asymptotic rate of convergence of this Markov chain to its equilibrium distribution is controlled by the second singular value of P , with smaller values of $\sigma_2(P)$ producing faster convergence. [BDX04] exploited this property to show that the fastest mixing Markov chain on the graph $(\mathcal{V}, \mathcal{E})$ could be computed by minimizing $\sigma_2(P)$ over all possible transition matrices on the graph, i.e. by solving

$$\begin{aligned} & \text{minimize} && \sigma_2(P) \\ & \text{subject to} && P \geq 0, P\mathbf{1} = \mathbf{1}, P = P^T, \\ & && P_{ij} = 0, \quad (i, j) \notin \mathcal{E} \end{aligned} \tag{19}$$

in the variable $P \in \mathbf{R}^{n \times n}$. The optimal mixing rate is often significantly faster than the rate provided by classical chains such as maximum degree or Metropolis-Hastings. Because $\sigma_1(P) = 1$ here, this is a particular instance of problem (13) where $k = 2$ and the matrices A_j are sparse. Projections on Q can be handled as in [BDX04]. Once again, the complexity of the main step in the algorithm (i.e. computing the gradient) is controlled by the sampling rate in Lemma 6, which simplifies to

$$s = \eta^2 \frac{1}{\epsilon^2 \sigma_2(P^*)^2} \mathbf{NumRank}(P^*)^2 \mathbf{Rank}(P^*)$$

where P^* is the transition matrix of the fastest mixing Markov chain. Once again, simpler transition matrices mean faster convergence.

5.4 Numerical experiments

In this section, we test the quality of the subsampling approximations detailed in Section 2 on various matrices. We also evaluate the performance of the algorithms detailed above on large scale problem instances. Numerical code reproducing these experiments is available from the author's webpage.

Randomized low-rank approximations. Here, we first measure the quality of the randomized low-rank matrix approximation on both randomly generated matrices and on covariance matrices formed using gene expression data. Because the spectrum of naive large scale random matrices is very structured, these examples are too simple to appropriately benchmark numerical error in

algorithm 2. Fortunately, as we will see below, generating random symmetric matrices with a given spectral measure is straightforward.

Suppose $X \in \mathbf{S}_n$ is a matrix with normally distributed coefficients, $X_{ij} \sim \mathcal{N}(0, 1)$, $i, j = 1, \dots, n$. If we write its QR decomposition, $X = QR$ with $Q, R \in \mathbf{R}^{n \times n}$, then the orthogonal matrix Q is Haar distributed on the orthogonal group \mathcal{O}_n (see [Dia03] for example). This means that to generate a random matrix with given spectrum $\mu \in \mathbf{R}^n$, we generate a normally distributed matrix X , compute its QR decomposition and the matrix $Q \mathbf{diag}(\mu) Q^T$ will be uniformly distributed on the set of symmetric matrices with spectrum μ . Because the spectral measure of “natural” covariance matrices often follows a power law (Tracy-Widom in the Gaussian case, see [Joh01] and [EK07] for a discussion), we sample the spectrum μ from a beta distribution with various exponents to get realistic random matrices with a broad range of numerical ranks. We also use a covariance matrix formed using the gene expression data set in [ABN⁺99].

In Figure 1, we plot relative error $\epsilon/\|X\|_2$ versus numerical rank $\mathbf{NumRank}(X)$ in loglog scale with 20% subsampling and $n = 500$ on random matrices generated as above and on the gene expression covariance from [ABN⁺99]. We notice that, on these experiments, the relative error grows at most linearly with the numerical rank of the matrix, as predicted by Lemma 3. We then plot the histogram in semilog scale of relative error $\epsilon/\|X\|_2$ over theoretical bound $\eta \mathbf{NumRank}(X)/\sqrt{s}$ for random matrices with $n = 500$. In Figure 2, we plot relative error $\epsilon/\|X\|_2$ versus sampling rate s , in loglog scale, for a gene expression covariance with $n = 500$. Once gain, the error decreases as $1/\sqrt{s}$ as predicted by Lemma 3. We also plot the median speedup factor (over ten runs) in computing largest magnitude eigenvalues using ARPACK with and without subsampling on a gene expression covariance matrix with $n = 2000$, for various values of the sampling ratio s/n . Note that both exact and subsampled eigenvalues are computed using direct MEX calls to ARPACK by [LSY98], as `eigs` (MATLAB’s interface to ARPACK) carries a massive overhead. In all the experiments above, the confidence level used in computing η was set to 99%.

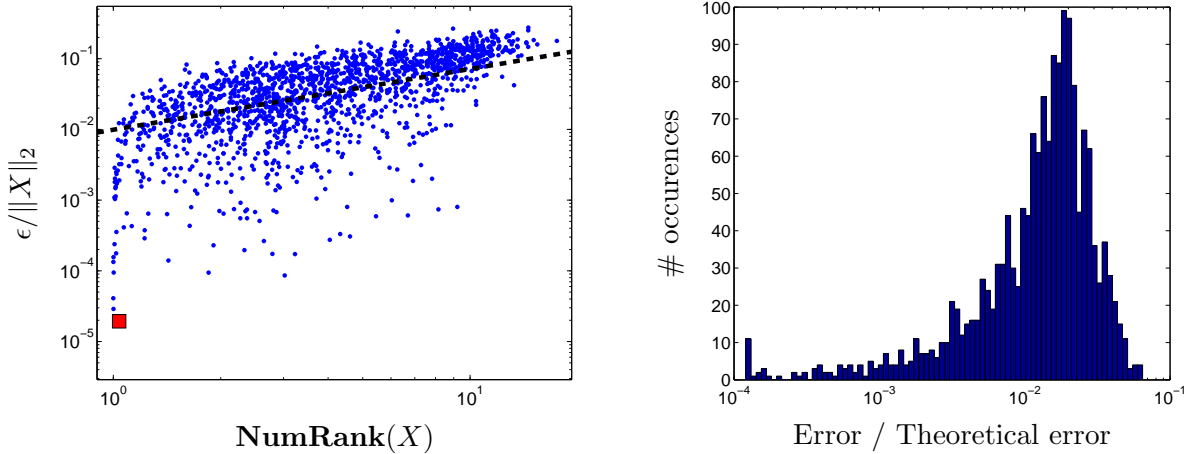


Figure 1: *Left:* Loglog plot of relative error $\epsilon/\|X\|_2$ versus numerical rank $\mathbf{NumRank}(X)$ with 20% subsampling and $n = 500$ on random matrices (blue dots) and gene expression covariance (red square). The dashed line has slope one in loglog scale. *Right:* Histogram plot in semilog scale of relative error $\epsilon/\|X\|_2$ over theoretical bound $\eta \mathbf{NumRank}(X)/\sqrt{s}$ for random matrices with $n = 500$.

Stochastic approximation with subsampling. In Figure 3, we generate a sample ratings matrix $X = VV^T$ for the collaborative filtering problem in §5.1, where V is a discrete feature matrix $V \in \{0, 1, 2\}^{100 \times 3}$. We “observe” only 30% of the ratings and solve problem (13) with $k = 4$ to approximately reconstruct the full ratings matrix. We plot objective value versus CPU time in seconds for this sample matrix factorization problem, using a stochastic approximation algorithm with deterministic gradient or the subsampled gradient algorithm 3 with subsampling ratio s_1/n set at 20%. We also plot surrogate duality gap versus CPU time on the same example. We notice that while the subsampled algorithm converges much faster than the deterministic one, the quality of the surrogate dual points and duality gap produced using subsampled gradients as in §3.3 is worst than in the deterministic case.

In Table 2, using the same 20% sampling rate we compare CPU time versus problem dimension n for subsampled and deterministic algorithms when solving the following instance of problem (1)

$$\begin{aligned} & \text{minimize} && \|C + X\|_2 \\ & \text{subject to} && \|X\|_\infty \leq \rho \end{aligned}$$

in the variable $X \in \mathbf{S}_n$ where C is a covariance matrix constructed using a subset of size n of the variables in [ABN⁺99], for various values of n . Finally, we generate and solve sample collaborative filtering problems as in (13) for ratings matrix of various dimensions n . We report median CPU time over ten sample problems in Table 3. Here, subsampling speeds up the algorithm by an order of magnitude, however the stochastic approximation algorithm is still not competitive with (non convex) local minimization techniques over low rank matrices.

n	Deterministic	Subsampling	Speedup factor
500	5	5	0.92
750	19	13	1.40
1000	32	24	1.31
1500	107	58	1.84
2000	281	120	2.34

Table 2: CPU time (in seconds) versus problem dimension n for deterministic and subsampled stochastic approximation algorithms on spectral norm minimization problems.

n	Deterministic	Subsampling	Speedup factor
100	154	23	6.67
200	766	63	12.2
500	4290	338	12.7

Table 3: Median CPU time (in seconds) versus problem dimension n for deterministic and subsampled stochastic approximation algorithms on collaborative filtering problems.

Acknowledgements

The author would like to acknowledge support from NSF DMS-0625352, SES-0835550 (CDI) and CAREER awards, a Peek junior faculty fellowship and a Howard B. Wentz Jr. award.

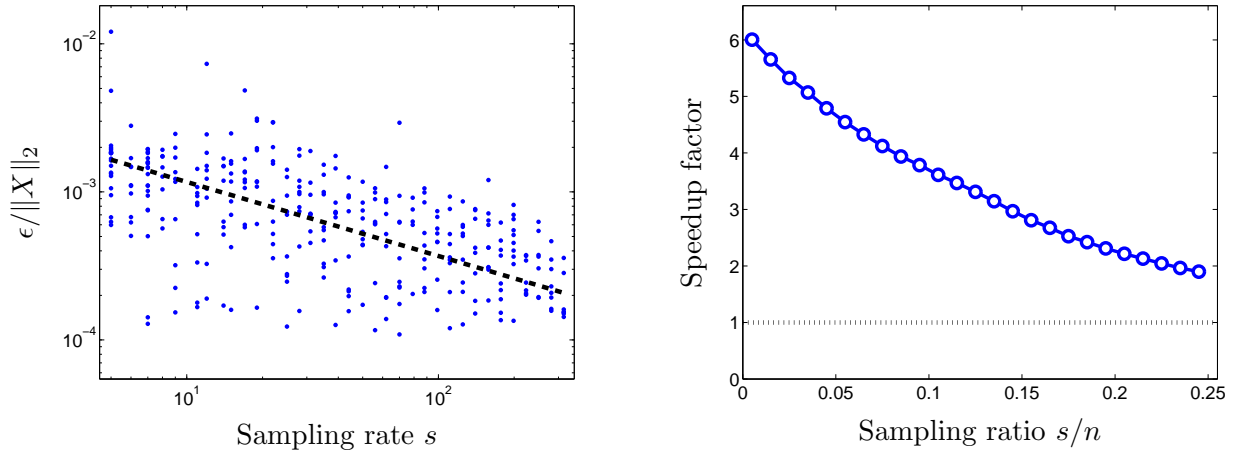


Figure 2: *Left:* Loglog plot of relative error $\epsilon/\|X\|_2$ versus sampling rate s for a gene expression covariance with $n = 500$. The dashed line has slope $-1/2$ in loglog scale. *Right:* Plot of median speedup factor in computing largest magnitude eigenvalue, using ARPACK with and without subsampling on a gene expression covariance matrix with $n = 2000$, for various values of the sampling ratio s/n .

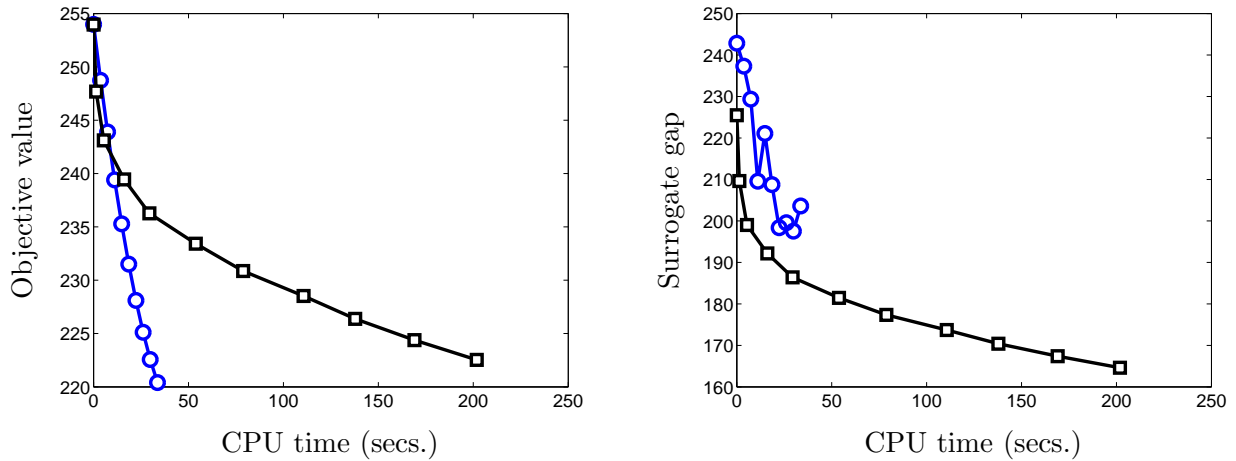


Figure 3: *Left:* Objective value versus CPU for a sample matrix factorization problem in dimension 100, using a deterministic gradient (squares) or a subsampled gradient with subsampling rate set at 20% (circles). *Right:* Surrogate duality gap versus CPU time on the same example.

References

- [ABN⁺99] A. 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. *Cell Biology*, 96:6745–6750, 1999.
- [AK07] S. Arora and S. Kale. A combinatorial, primal-dual approach to semidefinite programs. In *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*, pages 227–236, 2007.
- [AM07] D. Achlioptas and F. McSherry. Fast computation of low-rank matrix approximations. *Journal of the ACM*, 54(2), 2007.
- [Bac07] F. Bach. Consistency of trace norm minimization. *To appear in Journal of Machine Learning Research*. *Arxiv preprint arXiv:0710.2848*, 2007.
- [BDX04] S. Boyd, P. Diaconis, and L. Xiao. Fastest mixing Markov chain on a graph. *Siam Review*, 46(4):667–690, 2004.
- [BLO02] J.V. Burke, AS Lewis, and ML Overton. Approximating Subdifferentials by Random Sampling of Gradients. *Mathematics of Operations Research*, 27(3):567–584, 2002.
- [BLO05] J.V. Burke, A.S. Lewis, and M.L. Overton. A Robust Gradient Sampling Algorithm for Nonsmooth, Nonconvex Optimization. *SIAM Journal on Optimization*, 15(3):751–779, 2005.
- [CR08] E.J. Candes and B. Recht. Exact matrix completion via convex optimization. *preprint*, 2008.
- [dBEG06] A. d’Aspremont, O. Banerjee, and L. El Ghaoui. First-order methods for sparse covariance selection. *To appear in SIAM Journal on Matrix Analysis and Applications*, 2006.
- [Dia03] P. Diaconis. Patterns in eigenvalues: The 70th Josiah Willard Gibbs lecture. *Bulletin of the American Mathematical Society*, 40(2):155–178, 2003.
- [DKM06] P. Drineas, R. Kannan, and M.W. Mahoney. Fast Monte Carlo Algorithms for Matrices II: Computing a Low-Rank Approximation to a Matrix. *SIAM Journal on Computing*, 36:158, 2006.
- [DKM07] P. Drineas, R. Kannan, and M.W. Mahoney. Fast Monte Carlo algorithms for matrices I: Approximating matrix multiplication. *SIAM Journal on Computing*, 36(1):132–157, 2007.
- [DT06] D. Donoho and Y. Tsaig. Fast solution of ℓ_1 -norm minimization problems when the solution may be sparse. *Preprint*, 2006.
- [EK07] N. El Karoui. Tracy-Widom limit for the largest eigenvalue of a large class of complex sample covariance matrices. *Annals of Probability*, 35(2):663–714, 2007.
- [FHB01] M. Fazel, H. Hindi, and S. Boyd. A rank minimization heuristic with application to minimum order system approximation. *Proceedings American Control Conference*, 6:4734–4739, 2001.
- [FKV04] A. Frieze, R. Kannan, and S. Vempala. Fast monte-carlo algorithms for finding low-rank approximations. *Journal of the ACM (JACM)*, 51(6):1025–1041, 2004.
- [GVL90] G.H. Golub and C.F. Van Loan. Matrix computation. *North Oxford Academic*, 1990.
- [GW95] M.X. Goemans and D.P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *J. ACM*, 42:1115–1145, 1995.
- [HJ91] R.A. Horn and C.R. Johnson. *Topics in matrix analysis*. Cambridge university press, 1991.
- [JLNS09] A. Juditsky, G. Lan, A. Nemirovski, and A. Shapiro. Stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.

- [Joh01] I.M. Johnstone. On the distribution of the largest eigenvalue in principal components analysis. *Annals of Statistics*, pages 295–327, 2001.
- [KBH08] Krishnan Kumar, Chiru Bhattacharya, and Ramesh Hariharan. A randomized algorithm for large scale support vector learning. In J.C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20*. MIT Press, Cambridge, MA, 2008.
- [LCB⁺02] G. R. G. Lanckriet, N. Cristianini, P. Bartlett, L. El Ghaoui, and M. I. Jordan. Learning the kernel matrix with semi-definite programming. *19th International Conference on Machine Learning*, 2002.
- [LSY98] R.B. Lehoucq, D.C. Sorensen, and C. Yang. *ARPACK: Solution of Large-scale Eigenvalue Problems with Implicitly Restarted Arnoldi Methods*. Society for Industrial & Applied Mathematics, 1998.
- [Nes09] Y. Nesterov. Primal-dual subgradient methods for convex problems. *Mathematical programming Series B*, 120(1):221–259, 2009.
- [NY83] AS Nemirovsky and DB Yudin. *Problem complexity and method efficiency in optimization*. 1983.
- [PJ92] BT Polyak and AB Juditsky. Acceleration of stochastic approximation by averaging. *SIAM Journal on Control and Optimization*, 30:838, 1992.
- [RFP07] B. Recht, M. Fazel, and P.A. Parrilo. Guaranteed Minimum-Rank Solutions of Linear Matrix Equations via Nuclear Norm Minimization. *Arxiv preprint arXiv:0706.4138*, 2007.
- [RV07] Mark Rudelson and Roman Vershynin. Sampling from large matrices: An approach through geometric functional analysis. *J. ACM*, 54(4):21, 2007.
- [SBXD05] J. Sun, S. Boyd, L. Xiao, and P. Diaconis. The fastest mixing Markov process on a graph and a connection to a maximum variance unfolding problem. *SIAM Review*, 2005.
- [Sre04] N. Srebro. *Learning with Matrix Factorization*. PhD thesis, Massachusetts Institute of Technology, 2004.
- [SS90] G.W. Stewart and J. Sun. Matrix perturbation theory. 1990.
- [Tib96] R. Tibshirani. Regression shrinkage and selection via the LASSO. *Journal of the Royal statistical society, series B*, 58(1):267–288, 1996.
- [WS06] K.Q. Weinberger and L.K. Saul. Unsupervised Learning of Image Manifolds by Semidefinite Programming. *International Journal of Computer Vision*, 70(1):77–90, 2006.