
Stochastic Variance-Reduced Algorithms for PCA with Arbitrary Mini-Batch Sizes

Cheolmin Kim
Northwestern University

Diego Klabjan
Northwestern University

Abstract

We present two stochastic variance-reduced PCA algorithms and their convergence analyses. By deriving explicit forms of step size, epoch length and batch size to ensure the optimal runtime, we show that the proposed algorithms can attain the optimal runtime with any batch sizes. Also, we establish global convergence of the algorithms based on a novel approach, which studies the optimality gap as a ratio of two expectation terms. The framework in our analysis is general and can be used to analyze other stochastic variance-reduced PCA algorithms and improve their analyses. Moreover, we introduce practical implementations of the algorithms which do not require hyper-parameters. The experimental results show that the proposed methods outperform other stochastic variance-reduced PCA algorithms regardless of the batch size.

1 Introduction

Principal component analysis (PCA) (Jolliffe, 2011) is a fundamental tool for dimensionality reduction in machine learning and statistics. Given a data matrix $A = [a_1 a_2 \dots a_n] \in \mathbb{R}^{d \times n}$ consisting of n data vectors a_1, a_2, \dots, a_n in \mathbb{R}^d , PCA finds a direction w onto which the projections of the data vectors have the largest variance. Assuming that the data vectors are standardized with a mean of zero and standard deviation of one, the PCA problem can be formulated as

$$\begin{aligned} \text{maximize} \quad & f(w) = \frac{1}{2n} \sum_{i=1}^n (a_i^T w)^2 = \frac{1}{2} w^T C w \\ \text{subject to} \quad & \|w\|_2 = 1 \end{aligned} \quad (1)$$

where $C = \frac{1}{n} A A^T \in \mathbb{R}^{d \times d}$ is the covariance matrix of data matrix A . Since the largest eigenvector u_1 of C maximizes $f(w)$, (1) can be solved by the singular value decomposition (SVD) of A . However, the runtime of SVD is $\mathcal{O}(\min\{nd^2, n^2d\})$, which can be expensive in a large-scale setting. An alternative way to solve (1) is to use power iteration (Golub and Van Loan, 2012) which repeatedly applies $w_{t+1} = C w_t / \|C w_t\|$ at each iteration. The sequence of iterates $\{w_t\}$ generated by power iteration is guaranteed to obtain an ϵ -optimal solution after $\mathcal{O}(\frac{1}{\Delta} \log \frac{1}{\epsilon})$ iterations where $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_d \geq 0$ are the eigenvalues of C and $\Delta = 1 - \lambda_2/\lambda_1$. Since each iteration involves multiplying vector w_t with the matrix C , the runtime becomes $\mathcal{O}(nd \frac{1}{\Delta} \log \frac{1}{\epsilon})$. When n and d are both large, the runtime of power iteration is better than that of SVD. Nonetheless, it still largely depends on n and can be prohibitive when Δ is small.

In order to reduce the dependence on Δ or n , many variants have been developed. To improve the dependence on Δ , Xu et al. (2018) propose power iteration with momentum (Power+M) based on the momentum idea of Polyak (1964). With the optimal choice of the momentum parameter $\beta = \lambda_2^2/4$, the total runtime reduces to $\mathcal{O}(nd \frac{1}{\sqrt{\Delta}} \log \frac{1}{\epsilon})$. Also, a stochastic algorithm utilizing a stochastic gradient rather than a full gradient $C w_t$ is introduced in Oja (1982). Since it requires just one data vector at a time, the computational cost per iteration is significantly reduced. However, due to the variance of stochastic gradients, a sequence of diminishing step sizes needs to be adopted, making its progress slow near the optimum.

Built on the recent stochastic variance-reduced gradient (SVRG) technique (Johnson and Zhang, 2013), Shamir (2015, 2016) propose a stochastic variance-reduced version of Oja's algorithm (VR-PCA) and its extension for finding $k \geq 1$ principal components. Utilizing stochastic variance-reduced gradients, VR-PCA works with a constant step size and converges at an exponential rate, reducing the total runtime to $\mathcal{O}(d(n + \frac{1}{\Delta^2}) \log \frac{1}{\epsilon})$. The analysis of VR-PCA considers a batch of size one. While this implies that it works with any batch size,

Table 1: Comparison of stochastic variance-reduced methods for PCA and their convergence analyses. Types of convergence and complexity results are summarized. “Local” means that there is a restriction on the angle between an initial iterate and the first eigenvector u_1 and “global” implies no such restriction. For VR Power and VR HB Power, $\mu \geq 0$ is a parameter that controls the progress of the algorithms through step size $\eta = \Delta^\mu$ depending on batch size.

Algorithm	Convergence	Iteration	Batch Size	Total Runtime	Reference
VR-PCA	Local	$\mathcal{O}\left(\frac{1}{\Delta^2} \log \frac{1}{\epsilon}\right)$	$\mathcal{O}(1)$	$\mathcal{O}\left(d\left(n + \frac{1}{\Delta^2}\right) \log \frac{1}{\epsilon}\right)$	(Shamir, 2015)
VR Power+M	Local	$\mathcal{O}\left(\frac{1}{\Delta^{1/2}} \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\frac{\sqrt{d}}{\Delta^{3/2}}\right)$	$\mathcal{O}\left(d\left(n + \frac{\sqrt{d}}{\Delta^2}\right) \log \frac{1}{\epsilon}\right)$	(Xu et al., 2018)
Fast PCA	Global	$\mathcal{O}\left(\frac{1}{\Delta^2} \text{poly}\left(\log \frac{1}{\epsilon}\right)\right)$	$\mathcal{O}(1)$	$\mathcal{O}\left(d\left(n + \frac{1}{\Delta^2}\right) \text{poly}\left(\log \frac{1}{\epsilon}\right)\right)$	(Garber and Hazan, 2015)
VR Power	Global	$\mathcal{O}\left(\frac{1}{\Delta^{1+\mu}} \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\frac{1}{\Delta^{1-\mu}}\right)$	$\mathcal{O}\left(d\left(n + \frac{1}{\Delta^2}\right) \log \frac{1}{\epsilon}\right)$	[This Paper]
VR HB Power	Global	$\mathcal{O}\left(\frac{1}{\Delta^{1/2+\mu}} \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\frac{1}{\Delta^{3/2-\mu}}\right)$	$\mathcal{O}\left(d\left(n + \frac{1}{\Delta^2}\right) \log \frac{1}{\epsilon}\right)$	[This Paper]

conditions for the step size and the epoch size are not precisely given, making it hard to attain the theoretically guaranteed performance in practice.

A stochastic variance-reduced version of Power+M (VR Power+M) is introduced by Xu et al. (2018). Due to the momentum term, the iteration complexity is improved to $\mathcal{O}\left(\frac{1}{\Delta^{1/2}} \log \frac{1}{\epsilon}\right)$. However, a batch size of $\mathcal{O}\left(\frac{\sqrt{d}}{\Delta^{3/2}}\right)$ is required to achieve such iteration complexity, leading to the total runtime of $\mathcal{O}\left(d\left(n + \frac{\sqrt{d}}{\Delta^2}\right) \log \frac{1}{\epsilon}\right)$. This runtime is worse than that of VR-PCA due to the extra dependency on \sqrt{d} . Moreover, if the batch size is not sufficiently large, VR Power+M may diverge, which makes it hard to use.

On other other hand, Garber and Hazan (2015) reduce the PCA problem into a sequence of convex optimization problems. Each convex optimization problem has the form of the least square problem and amounts to one step of inverse power iteration (Golub and Van Loan, 2012). Due to the finite sum structure of the objective function, the SVRG algorithm (Johnson and Zhang, 2013) can be used to solve the least square problem. However, solving this strongly convex optimization problem can be as hard as the original PCA problem since the objective function is $(\lambda_1 - \lambda_2)$ -strongly convex and $(2\lambda_1 - \lambda_2 - \lambda_d)$ -smooth in the accurate regime. Through inexactly solving these problems, an ϵ -optimal solution can be obtained after a poly-logarithmic number of iterations.

The shifted-and-inverted approach is also introduced for the leading eigenvector problem (Garber et al., 2016) and numerous solvers such as coordinate-descent (Wang et al., 2018), SVRG (Garber et al., 2016), accelerated gradient descent, accelerated SVRG (Allen-Zhu and Li, 2016) and Riemannian gradient descent (Xu, 2018) have been developed to solve the least square problem. Other works on power iteration include the noisy (Hardt and Price, 2014) and coordinate-wise (Lei et al., 2016) power methods. The noisy power method considers the power method in a noise setting, which Balcan et al. (2016) extend to provide an improved gap-dependency

analysis. Moreover, power iteration has been analyzed for incremental or online PCA in many works (Allen-Zhu and Li, 2017; Li et al., 2018; Balsubramani et al., 2013; Arora et al., 2012; Boutsidis et al., 2015; Jain et al., 2016; Mitliagkas et al., 2013).

In this paper, we introduce two mini-batch stochastic variance-reduced PCA algorithms (VR Power, VR HB Power) and provide their convergence analyses. They are mini-batch stochastic variance-reduced variants of power iteration (Golub and Van Loan, 2012) and power with momentum method (Xu et al., 2018). While VR-PCA (Shamir, 2015) takes a data vector at a time, VR Power works with any batch sizes and the accompanying analysis reveals that whatever the batch size is, VR Power attains the optimal runtime by appropriately choosing the step size and epoch length. Explicit conditions of the step size, epoch length and batch size to ensure the optimal runtime of VR Power are derived. On the other hand, VR HB Power is an enhanced algorithm of VR Power+M. By adding the step size, VR HB Power works with any batch sizes while VR Power+M can fail unless the batch size is sufficiently large. For any batch sizes, VR HB Power can achieve the optimal runtime if we appropriately choose the step size, epoch length and momentum parameter. We derive explicit expressions for these parameters. Our analysis improves the analysis of VR Power+M by removing the dependency on \sqrt{d} for the batch size. For the comparison of stochastic variance-reduced PCA algorithms and their convergence analyses see Table 1.

In the convergence analyses, we introduce a novel framework of analyzing stochastic variance-reduced algorithms for PCA. For an inner-loop iterate w_t , we decompose $E[(u_k^T w_t)^2]$ with u_k an eigenvector with respect to λ_k into two parts where the first one is the expectation term and the second one is the variance term. To obtain a tight bound for the variance term, we analyze its growth over an epoch rather than focusing on iteration-by-iteration behavior. Based on the Binomial expansion of matrices, we come up with a compact bound

of the variance term, which is used to establish an upper bound of $(E[\|w_t\|^2] - E[(u_1^T w_t)^2])/E[(u_1^T w_t)^2] = \sum_{k=2}^d E[(u_k^T w_t)^2]/E[(u_1^T w_t)^2]$ and derive conditions for the step size, epoch length and batch size to ensure its sufficient decrease.

The concept of representing the optimality gap as the ratio of two expectations has been never used for analyzing stochastic PCA algorithms. However, it results in much simpler convergence statements than probabilistic statements in Shamir (2015) and Xu et al. (2018). Note that we are able to obtain probabilistic statements from the expectation bounds using the Chebyshev inequality. Using the expectation bounds, we can establish global convergence of stochastic PCA algorithms. Although stochastic PCA algorithms have been observed to work well with random initialization (Shamir, 2015), an initial condition such as $|u_1^T \tilde{w}_0| \geq 1/2$ is required in previous probabilistic analyses. In our framework, such condition is not necessary and the rate of convergence does not depend on how far an iterate is from u_1 but is kept the same across iterations, as in the case of deterministic power iteration. The framework introduced in this work is not specific to the proposed algorithms; it can be applied to analyze other stochastic variance-reduced PCA algorithms such as VR-PCA or VR Power+M, deriving expectation bounds for them and resolving their initialization issues.

Our work has the following contributions.

1. We introduce two mini-batch stochastic variance-reduced PCA algorithms. Regardless of the batch size, our algorithms can attain the optimal runtime by appropriately choosing algorithm parameters. Explicit expressions for these parameters are provided.
2. We provide novel convergence analyses for the algorithms where we establish global convergence by deriving a bound for the ratio of two expectation terms. The framework in our convergence analyses is general and can be used to analyze other stochastic variance-reduced PCA algorithms. To this end, we are the first to establish convergence of VR-PCA and VR Power+M for any initial vector and in expectation.
3. We present practical implementations of the algorithms and report numerical experiments. The experimental results on real-world datasets show that our algorithms outperform other stochastic variance-reduced algorithms for any batch size.

The paper is organized as follows. We introduce the algorithms in Section 2 and the convergence analyses in Section 3. Some practical considerations regarding

the implementations of the algorithms are discussed in Section 4 and the experimental results are followed in Section 5.

2 Stochastic Variance-Reduced Algorithms for PCA

We consider two mini-batch stochastic variance-reduced algorithms for PCA. The first one is a mini-batch version of VR-PCA (Shamir, 2015) and the second one is an enhanced version of VR Power+M (Xu et al., 2018) with a step size incorporated. For eigenpairs (λ_k, u_k) of $C = \frac{1}{n} \sum_{i=1}^n a_i a_i^T$, we assume that the eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_d$ satisfy $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_d \geq 0$ and the eigenvectors u_1, u_2, \dots, u_d form an orthonormal basis. Since a symmetric matrix is orthogonally diagonalizable, we can assume that such eigenvectors exist without loss of generality. We assume that all norms are L_2 for vectors and spectral for matrices.

Variance reduction algorithms have an outer loop and an inner loop. They periodically compute exact gradients at each outer iteration and use it in inner iterations to reduce the variance of stochastic gradients. Let \tilde{w}_s and w_t denote an outer-loop and inner-loop iterate, respectively. To get a stochastic variance-reduced gradient of an inner loop iterate w_t , we first decompose the inner loop iterate w_t into two parts as

$$w_t = \frac{(\tilde{w}_s^T w_t)}{\|\tilde{w}_s\|^2} \tilde{w}_s + \left(I - \frac{\tilde{w}_s \tilde{w}_s^T}{\|\tilde{w}_s\|^2} \right) w_t$$

using the outer loop iterate \tilde{w}_s . In the above decomposition, the former term represents the projection of w_t on \tilde{w}_s while the latter term represents the remaining vector. Utilizing the exact gradient \tilde{g}_s at \tilde{w}_s , the exact gradient at the first term can be computed as

$$\nabla f \left(\frac{(\tilde{w}_s^T w_t)}{\|\tilde{w}_s\|^2} \tilde{w}_s \right) = \frac{(\tilde{w}_s^T w_t)}{\|\tilde{w}_s\|^2} C \tilde{w}_s = \frac{(\tilde{w}_s^T w_t)}{\|\tilde{w}_s\|^2} \tilde{g}_s.$$

On the other hand, a stochastic sample S_t is used to compute a stochastic gradient at the second term as

$$\frac{1}{|S_t|} \sum_{l \in S_t} a_l a_l^T \left(I - \frac{\tilde{w}_s \tilde{w}_s^T}{\|\tilde{w}_s\|^2} \right) w_t.$$

This results in the following stochastic variance-reduced gradient g_t at w_t as

$$g_t = \frac{(\tilde{w}_s^T w_t)}{\|\tilde{w}_s\|^2} \tilde{g}_s + \frac{1}{|S_t|} \sum_{l \in S_t} a_l a_l^T \left(I - \frac{\tilde{w}_s \tilde{w}_s^T}{\|\tilde{w}_s\|^2} \right) w_t. \quad (2)$$

2.1 VR Power

Using the stochastic variance-reduced gradient g_t , we obtain a stochastic variance reduced version of Power

iteration as

$$w_{t+1} \leftarrow (1 - \eta)w_t + \eta g_t. \quad (3)$$

This update rule has a similar form as the one in VR-PCA, which repeats

$$w_{t+1} \leftarrow w_t + \bar{\eta} (a_{i_t}^T w_t - a_{i_t}^T \tilde{w}_s) + \tilde{g}_s. \quad (4)$$

Note that (3) generalizes (4) in the following two senses. First, we can obtain an update rule of (4) by letting $\eta = (1 + \bar{\eta})/\bar{\eta}$ in (3). Second, with the choice of $\eta = 1$, we can recover deterministic power iteration from (3) while (4) does not. Using update rule (3), we have VR Power exhibited in Algorithm 1.

Algorithm 1 VR Power

Parameters: step size η , mini-batch size $|S|$, epoch length m

Input: data vectors $a_i, i = 1, \dots, n$
randomly initialize outer iterate \tilde{w}_0

for $s = 0, 1, \dots$ **do**

$$\tilde{g} \leftarrow C\tilde{w}_s$$

$$w_0 \leftarrow \tilde{w}_s$$

$$w_1 \leftarrow (1 - \eta)w_0 + \eta\tilde{g}$$

for $t = 1, 2, \dots, m - 1$ **do**

sample a mini-batch $S_t \subset \{1, \dots, n\}$ of size $|S|$
uniformly at random

$$g_t \leftarrow \frac{1}{|S_t|} \sum_{l \in S_t} a_l a_l^T \left(I - \frac{w_0 w_0^T}{\|w_0\|^2} \right) w_t + \frac{(w_t^T w_0)}{\|w_0\|^2} \tilde{g}$$

$$w_{t+1} \leftarrow (1 - \eta)w_t + \eta g_t$$

end for

$$\tilde{w}_{s+1} \leftarrow w_m$$

end for

When per sample cost is as expensive as per iteration cost, VR Power is an efficient algorithm since it attains the optimal sample complexity. However, if per sample cost is cheap, it might not be effective since its iteration complexity does not improve beyond $\mathcal{O}(\frac{1}{\Delta} \log(\frac{1}{\epsilon}))$. For this reason, we introduce VR HB Power which works better in the latter setting.

2.2 VR HB Power

Using g_t , we obtain a stochastic variance-reduced heavy ball power iteration as

$$w_{t+1} \leftarrow 2((1 - \eta)w_t + \eta g_t) - \beta w_{t-1} \quad (5)$$

where $\eta \in (0, 1]$ is the step size and β is the momentum parameter. Note that we can recover the deterministic heavy ball power iteration from (5) when the step size η is set to 1 and the exact gradient $g_t = Cw_t$ is used. The mechanism of controlling the progress of the algorithm using the step size η is not present in VR Power+M (Xu et al., 2018). As a result, it fails to converge unless

the mini-batch size $|S|$ is sufficiently large. To the contrary, our algorithm works with any mini-batch size $|S|$ due to the presence of the step size η . By selecting an appropriate value of η depending on the size of $|S|$ and m , we can always ensure that the variance terms do not grow faster than expectation terms. Having update rule (5), VR HB Power is described in Algorithm 2.

Algorithm 2 VR HB Power

Parameters: step size η , momentum β , mini-batch size $|S|$, epoch length m

Input: data vectors $a_i, i = 1, \dots, n$
randomly initialize outer iterate \tilde{w}_0

for $s = 0, 1, \dots$ **do**

$$\tilde{g} \leftarrow C\tilde{w}_s$$

$$w_0 \leftarrow \tilde{w}_s$$

$$w_1 \leftarrow (1 - \eta)w_0 + \eta\tilde{g}$$

for $t = 1, 2, \dots, m - 1$ **do**

sample a mini-batch $S_t \subset \{1, \dots, n\}$ of size $|S|$
uniformly at random

$$g_t \leftarrow \frac{1}{|S_t|} \sum_{l \in S_t} a_l a_l^T \left(I - \frac{w_0 w_0^T}{\|w_0\|^2} \right) w_t + \frac{(w_t^T w_0)}{\|w_0\|^2} \tilde{g}$$

$$w_{t+1} \leftarrow 2((1 - \eta)w_t + \eta g_t) - \beta w_{t-1}$$

end for

$$\tilde{w}_{s+1} \leftarrow w_m$$

end for

3 Convergence Analyses

In this section, we provide convergence analyses for VR Power and VR HB Power. Before presenting the convergence analyses, we first introduce some notation.

3.1 Notation

Let C_t and P be the sample covariance matrix at inner iteration t and the projection matrix to the space orthogonal to the outer iterate $w_0 = \tilde{w}_s$ as

$$C_t = \frac{1}{|S_t|} \sum_{l \in S_t} a_l a_l^T, \quad P = I - \frac{w_0 w_0^T}{\|w_0\|^2}. \quad (6)$$

Using (6), we can write g_t as $g_t = \eta C w_t + \eta(C_t - C)P w_t$. Next, we characterize the variance of sample covariance matrix C_t as

$$K = E[\|(C_t - C)^2\|], \quad \sigma^2 = E[\|a_{i_t} a_{i_t}^T - C\|^2].$$

Then, for $M_k = E[(C_t - C)u_k u_k^T (C_t - C)]$, we have

$$\|M_k\| \leq K = \frac{\sigma^2}{|S|}. \quad (7)$$

For the analysis of VR HB Power, we define

$$\alpha_k(\eta) = 4(1 - \eta + \eta\lambda_k)^2, \quad \beta(\eta) = (1 - \eta + \eta\lambda_2)^2. \quad (8)$$

Also, we let $p_t(\alpha, \beta)$ and $q_t(\alpha, \beta)$ be the Chebyshev polynomials of the first and the second kind (Mason and Handscomb, 2002) respectively such that

$$p_t(\alpha, \beta) = (\alpha - \beta)p_{t-1}(\alpha, \beta) - \beta(\alpha - \beta)p_{t-2}(\alpha, \beta) + \beta^3 p_{t-3}(\alpha, \beta), \quad (9)$$

$$q_t(\alpha, \beta) = (\alpha - \beta)q_{t-1}(\alpha, \beta) - \beta(\alpha - \beta)q_{t-2}(\alpha, \beta) + \beta^3 q_{t-3}(\alpha, \beta) \quad (10)$$

for $t \geq 3$ and

$$p_0(\alpha, \beta) = 1, p_1(\alpha, \beta) = \frac{\alpha}{4}, p_2(\alpha, \beta) = \left(\frac{\alpha}{2} - \beta\right)^2, \quad (11)$$

$$q_0(\alpha, \beta) = 1, q_1(\alpha, \beta) = \alpha, q_2(\alpha, \beta) = (\alpha - \beta)^2. \quad (12)$$

Since the first eigenvector u_1 of the covariance matrix C is an optimal solution to (1), the optimality gap is measured as $\sum_{k=2}^d (u_k^T w_t)^2 / (u_1^T w_t)^2$, representing how closely w_t is aligned with u_1 . Note that this ratio is zero if $w_t = u_1$. Our analysis studies it in expectation, providing a bound for $\theta_t = \sum_{k=2}^d E[(u_k^T w_t)^2] / E[(u_1^T w_t)^2]$ given fixed s and $\tilde{\theta}_s = \sum_{k=2}^d E[(u_k^T \tilde{w}_s)^2] / E[(u_1^T \tilde{w}_s)^2]$ for an inner loop iterate w_t and an outer loop iterate \tilde{w}_s , respectively.

3.2 VR Power

In Lemmas 3.1, 3.2 and 3.3, we consider a single epoch, which corresponds to one inner loop iteration starting with w_0 .

Lemma 3.1. *For any $\eta \in (0, 1]$, $1 \leq k \leq d$ and $1 \leq t \leq m$, we have*

$$E[(u_k^T w_t)^2] = (1 - \eta + \eta\lambda_k)^{2t} E[(u_k^T w_0)^2] + \eta^2 \sum_{i=1}^{t-1} (1 - \eta + \eta\lambda_k)^{2(t-i-1)} E[w_i^T P M_k P w_i].$$

Lemma 3.1 decomposes $E[(u_k^T w_t)^2]$ into two parts. The first part represents the expectation term which grows at a rate of $(1 - \eta + \eta\lambda_k)^2$ and the second part is the variance term which increases as w_t strides away from w_0 as captured by $E[w_t^T P M_k P w_t]$.

Lemma 3.2. *For any $\eta \in (0, 1]$, $1 \leq k \leq d$ and $1 \leq t \leq m$, we have*

$$\sum_{k=2}^d E[w_t^T P M_k P w_t] \leq 2K \cdot \sum_{k=2}^d E[(u_k^T w_0)^2] \cdot ((1 - \eta + \eta\lambda_1)^2 + \eta^2 K)^t.$$

Moreover, if $0 < \frac{\eta^2 K m}{(1 - \eta + \eta\lambda_1)^2} < 1$, then we have

$$\theta_m \leq \left[\left(\frac{1 - \eta + \eta\lambda_2}{1 - \eta + \eta\lambda_1} \right)^{2m} + \frac{4\eta^2 K m}{(1 - \eta + \eta\lambda_1)^2} \right] \cdot \theta_0.$$

Lemma 3.2 provides a bound for $\sum_{k=2}^d E[w_t^T P M_k P w_t]$, which grows at a rate not greater than $(1 - \eta + \eta\lambda_1)^2 + \eta^2 K$. Using this bound and assuming some condition on η , K , and m , a bound on θ_m is derived as a function of θ_0 , η , m , and K . In Lemma 3.3, we present explicit conditions for η , m , and $|S|$ to ensure a sufficient decrease of θ_m .

Lemma 3.3. *Let $\eta = \Delta^\mu$ for some $\mu \geq 0$. If m and $|S|$ satisfy*

$$m = \left\lceil \frac{(1 - \eta + \eta\lambda_1) \log 2}{2\eta\lambda_1 \Delta} \right\rceil \quad (13)$$

and

$$|S| \geq \frac{16\eta^2 \sigma^2 m}{(1 - \eta + \eta\lambda_1)^2}, \quad (14)$$

then we have $\theta_m \leq 3/4 \cdot \theta_0$.

For any $\mu \geq 0$ such that $\eta = \Delta^\mu$, Lemma 3.3 provides explicit values of m and $|S|$ to ensure a sufficient decrease of θ_m . In the analysis of VR-PCA, exact values of η and m to ensure the optimal runtime have not been provided. Instead, only the orders of η and m have been provided such that $\eta = c_1 \Delta$ and $m = c_2 / \Delta^2$, making it hard to obtain the optimal runtime in practice. Contrary to it, our analysis provides explicit expressions for m and $|S|$, being more practical. Moreover, since the term on the right-hand side of (14) goes to zero as μ increases, it can be also stated that for any $|S| \geq 1$, there exists some $\mu \geq 0$ and thus $\eta = \Delta^\mu$ and m (see (14)) such that $\theta_m \leq 3/4 \cdot \theta_0$ holds. This implies that VR Power can always attain a sufficient decrease of θ_m no matter what $|S|$ is used. We next give the main result.

Theorem 3.4. *Suppose that an initial vector \tilde{w}_0 satisfies $u_1^T \tilde{w}_0 \neq 0$ and let $\tilde{\theta}_0 = (1 - (u_1^T \tilde{w}_0)^2) / (u_1^T \tilde{w}_0)^2 \geq \epsilon$ for some $\epsilon > 0$. If $\eta = \Delta^\mu$ and m and $|S|$ satisfy (13) and (14), after $\tau = \lceil \log(\tilde{\theta}_0 / \epsilon) / \log(4/3) \rceil$ epochs of VR Power, we have $\tilde{\theta}_\tau \leq \epsilon$.*

Theorem 3.4 present a convergence result for τ epochs. Note that our result requires only a trivial assumption on $\tilde{\theta}_0$ and thus establishes global convergence. Also, since $\tau = \mathcal{O}(\log(\frac{1}{\epsilon}))$, only a logarithmic number of inner loops is needed to be performed to obtain ϵ -accuracy.

3.3 VR HB Power

The following Lemmas 3.5, 3.6 and 3.7 are counterparts of Lemmas 3.1, 3.2 and 3.3 for VR HB Power. For the momentum parameter β , we let $\beta = \beta(\eta)$ which is defined in (8). As in the analysis of VR Power, we first consider a single epoch with an initial inner loop iterate w_0 .

Lemma 3.5. For any $\eta \in (0, 1]$, $1 \leq k \leq d$ and $1 \leq t \leq m$, we have

$$E[(u_k^T w_t)^2] = p_t(\alpha_k(\eta), \beta(\eta))E[(u_k^T w_0)^2] + 4\eta^2 \sum_{r=1}^{t-1} q_{t-r-1}(\alpha_k(\eta), \beta(\eta))E[w_r^T P M_k P w_r].$$

Lemma 3.5 breaks $E[(u_k^T w_t)^2]$ into the sum of expectation part and variance part. While the expectation term is a function of the Chebyshev polynomial of the first kind, the variance part is a function of the Chebyshev polynomials of the second kind. That being said, the variance term grows faster and thus we need a careful analysis for it.

Lemma 3.6. For any $\eta \in (0, 1]$, $1 \leq k \leq d$, and $1 \leq t \leq m$, we have

$$\sum_{k=2}^d E[w_t^T P M_k P w_t] \leq 4K \cdot \left(1 + \frac{4\eta^2 K}{\alpha_1(\eta) - 4\beta(\eta)}\right)^{t-1} \cdot \left(\frac{\sqrt{\alpha_1(\eta)}}{2} + \frac{\sqrt{\alpha_1(\eta) - 4\beta(\eta)}}{2}\right)^{2t} \cdot \sum_{k=2}^d E[(u_k^T w_0)^2].$$

Moreover, if $0 < \frac{4\eta^2 K m}{\alpha_1(\eta) - 4\beta(\eta)} < 1$, then we have

$$\theta_m \leq \left(\frac{p_m(\alpha_2(\eta), \beta(\eta))}{p_m(\alpha_1(\eta), \beta(\eta))} + \frac{128\eta^2 K m}{\alpha_1(\eta) - 4\beta(\eta)}\right) \cdot \theta_0.$$

Lemma 3.6 provides a bound for $\sum_{k=2}^d E[w_t^T P M_k P w_t]$. Note that it depends on Δ and blows up as Δ goes to zero due to the term involving $1/(\alpha_1(\eta) - 4\beta(\eta))$. Due to this dependency, VR HB Power tends to require a larger batch size than VR Power given the same values of η and m . Lemma 3.6 also establishes a bound for θ_m as a function of θ_0 , η , m and K under some assumption.

Lemma 3.7. For some $\mu \geq 0$, let $\eta = \Delta^\mu$ and

$$m = \left\lceil \left[\frac{1 - \eta + \eta\lambda_1}{\eta\lambda_1\Delta + \sqrt{\eta\lambda_1\Delta(2(1-\eta) + \eta(\lambda_1 + \lambda_2))}} + \frac{\sqrt{\eta\lambda_1\Delta(2(1-\eta) + \eta(\lambda_1 + \lambda_2))}}{\eta\lambda_1\Delta + \sqrt{\eta\lambda_1\Delta(2(1-\eta) + \eta(\lambda_1 + \lambda_2))}} \right] \frac{\log 8}{2} \right\rceil \quad (15)$$

and

$$|S| \geq \frac{128\eta\sigma^2 m}{\lambda_1\Delta [2(1-\eta) + \eta(\lambda_1 + \lambda_2)]}. \quad (16)$$

Then, we have $\theta_m \leq 3/4 \cdot \theta_0$.

Lemma 3.7 provides explicit conditions for m and $|S|$ to ensure a sufficient decrease of θ_m . Note that when $\mu = 0$, we have $|S| \geq \mathcal{O}(\frac{1}{\Delta^{3/2}})$, which improves the analysis of VR Power+M in Xu et al. (2018) by removing the dependency on \sqrt{d} . Also, for any $|S| \geq 1$,

there exists some η and m satisfying the conditions in Lemma 3.7. This implies that VR HB Power works with any batch size while VR Power+M does not. The overall convergence is established next.

Theorem 3.8. Suppose that an initial vector \tilde{w}_0 satisfies $u_1^T \tilde{w}_0 \neq 0$ and let $\theta_0 = (1 - (u_1^T \tilde{w}_0)^2)/(u_1^T \tilde{w}_0)^2 \geq \epsilon$ for some $\epsilon > 0$. If $\eta = \Delta^\mu$ and m and $|S|$ satisfy (15) and (16), after $\tau = \lceil \log(\tilde{\theta}_0/\epsilon)/\log(4/3) \rceil$ epochs of VR HB Power, we have $\tilde{\theta}_\tau \leq \epsilon$.

The global convergence result in Theorem 3.8 is based on the single epoch result in Lemma 3.7. Since $\tau = \mathcal{O}(\log(\frac{1}{\epsilon}))$, the iteration complexity of VR HB Power is $\tau m = \mathcal{O}(\frac{1}{\Delta^{1/2+\mu/2}} \log(\frac{1}{\epsilon}))$. On the other hand, from $|S| = \mathcal{O}(\frac{1}{\Delta^{3/2-\mu/2}})$, the sample complexity amounts to $\mathcal{O}((n + \frac{1}{\Delta^2}) \log(\frac{1}{\epsilon}))$. Note that VR HB Power has the same sample complexity as VR Power but may have small iteration complexity. Therefore, if per sample cost is cheaper than per iteration cost, VR HB Power can be more efficient than VR Power.

4 Practical Considerations

In this section, we discuss some practical aspects implementing the proposed algorithms. First, to ensure that the algorithms are numerically stable, we consider normalizations as introduced in Shamir (2015) and Xu et al. (2018). After updating w_{t+1} , we normalize w_{t+1} as $w_{t+1} \leftarrow w_{t+1}/\|w_{t+1}\|_2$ in VR Power and update w_t and w_{t+1} as $w_t \leftarrow w_t/\|w_{t+1}\|_2$ and $w_{t+1} \leftarrow w_{t+1}/\|w_{t+1}\|_2$ in VR HB Power. Since these scaling schemes do not impact the sample paths of $w_t/\|w_t\|$, we can obtain the same results with numerical stability.

Another practical issue with the implementations of VR Power and VR HB Power is to estimate λ_1 and λ_2 . As appearing in Lemma 3.3 and Lemma 3.7, accurate values of λ_1 and λ_2 are essential to determine the values of η , m , and β (for VR HB Power). In the experiments, the mini-batch size $|S|$ is given as some percentage of n , so no estimation is required for $|S|$. In order to estimate λ_1 and λ_2 at a regular interval (at the start of each inner-loop), we use the exact gradients of two consecutive outer-loop iterates \tilde{w}_{s-1} and \tilde{w}_s . Since we expect that \tilde{w}_s approaches u_1 as the iterations advance, using the Rayleigh quotient, we estimate λ_1 as

$$\hat{\lambda}_1 = \frac{(\tilde{w}_s)^T C (\tilde{w}_s)}{(\tilde{w}_s)^T \tilde{w}_s}. \quad (17)$$

To estimate λ_2 in the same way, we need an estimate of u_2 . In Power iteration, an iterate first approaches the subspace spanned by u_1 and u_2 before converging to u_1 . That being said, after a number of iterations, we can approximate it by a linear combination of u_1

and u_2 . Based on this observation, we estimate u_2 as

$$\hat{u}_2 = \tilde{w}_{s-1} - (\tilde{w}_{s-1}^T \tilde{w}_s) \tilde{w}_s. \quad (18)$$

The idea of the above estimation is to project \tilde{w}_{s-1} to the space orthogonal to \tilde{w}_s . If $\tilde{w}_s \approx u_1$ and $\tilde{w}_{s-1} \approx \alpha_1 u_1 + \alpha_2 u_2$ for some $\alpha_1, \alpha_2 (\neq 0)$, we have $\hat{u}_2 \approx u_2$. Using the Rayleigh quotient of \hat{u}_2 , we estimate λ_2 as

$$\hat{\lambda}_2 = \frac{\tilde{w}_{s-1}^T C \tilde{w}_{s-1} - 2\theta_s \tilde{w}_s^T C \tilde{w}_{s-1} + \theta_s^2 \tilde{w}_s^T C \tilde{w}_s}{1 - \theta_s^2} \quad (19)$$

where $\theta_s = \tilde{w}_{s-1}^T \tilde{w}_s$. While two matrix-vector multiplications, $C \tilde{w}_{s-1}$ and $C \tilde{w}_s$, are involved in computing (17) and (19), they incur no extra computation since they are the exact gradients of \tilde{w}_{s-1} and \tilde{w}_s , which are computed regardless of the estimation. As a result, we can obtain $\hat{\lambda}_1$ and $\hat{\lambda}_2$ by only computing some inner products. For initial estimation of $\hat{\lambda}_1$ and $\hat{\lambda}_2$, we run Power iteration five times and use the last two iterates. Note that the exact gradient of the last iterate is computed at the start of the very first outer-loop iteration.

Given $|S|$ and estimates of λ_1 and λ_2 , we use bisection search to find $\eta \in (0, 1]$ such that the terms on the right-hand sides of (14) and (16) are almost equal to $|S|$. After η is found, we use (13) and (15) to determine m .

5 Numerical Experiments

In this section, we test the performance of VR Power and VR HB Power with that of (i) VR-PCA (Shamir, 2015), (ii) VR Power+M (Xu et al., 2018) and (iii) Fast PCA (Garber and Hazan, 2015) for finding the first eigenvector u_1 of the covariance matrix C constructed by data vectors $a_i, i = 1, \dots, n$ from real world datasets. Note that all present stochastic variance-reduced PCA algorithms are compared in this experiment.

5.1 Datasets

The datasets include ijcn (Prokhorov, 2001), covertype (Blackard and Dean, 1999), YearPredictionMSD (Bertin-Mahieux et al., 2011) and MNIST (LeCun et al., 1998) as summarized in Table 2. All of them are obtained either from the UCI repository (Dheeru and Karra Taniskidou, 2017) or the LIBSVM library (Chang and Lin, 2011). They are carefully chosen to incorporate a variety of datasets in terms of size and eigen-gap. The first three datasets are standardized with a mean of zero and standard deviation of one while the last one is scaled to the range between 0 and 1 to preserve its sparsity.

Table 2: A summary of datasets

DATASET	n	d	Δ
ICJNN(TEST)	91,701	22	0.0079
COV	581,012	54	0.2106
MSD	463,715	90	0.3224
MNIST	70,000	764	0.8851

5.2 Settings

In order to report a comprehensive comparison of the algorithms, we consider two settings for selecting hyper-parameters. In the first setting, we use hyper-parameter tuning. Specifically, we use a grid search to find the best values of η , m and $|S| = \rho\%$ of each algorithm and dataset where $\eta \in \{0.01, 0.05, 0.1, 0.2, 0.4, 0.6, 0.8, 1.0\}$, $m \in \{25, 50, 100, 200\}$ and $\rho \in \{1, 2, 5, 10\}$.

In the second setting, we use the following theoretically derived or recommended hyper-parameter values.

- VR-PCA: $\eta = \sqrt{n} / \sum_{i=1}^n \|a_i\|^2$, $m = n$, $|S| = 1$.
- VR Power+M: $\beta = \lambda_2^2/4$, $\sigma^2 = \sum_{i=1}^n \|a_i\|^2/n$,

$$|S| = \frac{\lambda_2 \log 16}{\sqrt{\lambda_1^2 - \lambda_2^2}}, \quad T = \frac{512 \log 16 \lambda_2 \sigma^2 \sqrt{d}}{\sqrt{\lambda_1^2 - \lambda_2^2}}.$$

- Fast PCA: $\delta = \lambda_1 - \lambda_2$. We only consider the accurate regime. In order to solve each problem, we use SVRG (Johnson and Zhang, 2013) with $\tilde{\epsilon} = 10^{-3}$,

$$\eta = \frac{\lambda_1 - \lambda_2}{7(2\lambda_1 + \lambda_2)^2}, \quad m = \left\lceil \frac{1}{2\eta^2(2\lambda_1 + \lambda_2)^2} \right\rceil.$$

- VR Power, VR HB Power: $|S| = \rho\% \cdot n$ for $\rho \in \{1, 2\}$ and $\sigma^2 = \sum_{i=1}^n \|a_i\|^2/n$. For η and m , we use bisection search explained in Section 4. Also, the scaling schemes in Section 4 are used to ensure numerical stability. The exact values of λ_1 and λ_2 are used to find η and m .
- PF VR Power, PF VR HB Power: As opposed to VR Power and VR HB Power, adaptive estimates of $\hat{\lambda}_1$ and $\hat{\lambda}_2$ obtained by the procedure in Section 4 are used to find η and m .

5.3 Results

Figure 1 displays the experimental result with hyper-parameter tuning. In the figure, the x-axis represents time in seconds and the y-axis represents the optimality gap, $1 - (\tilde{w}_s^T u_1)^2$, in the log-scale. Since VR-PCA and VR Power are related algorithms, their performances are similar except for cov where the step size of VR-PCA is tuned to the largest possible value of 1.0. If

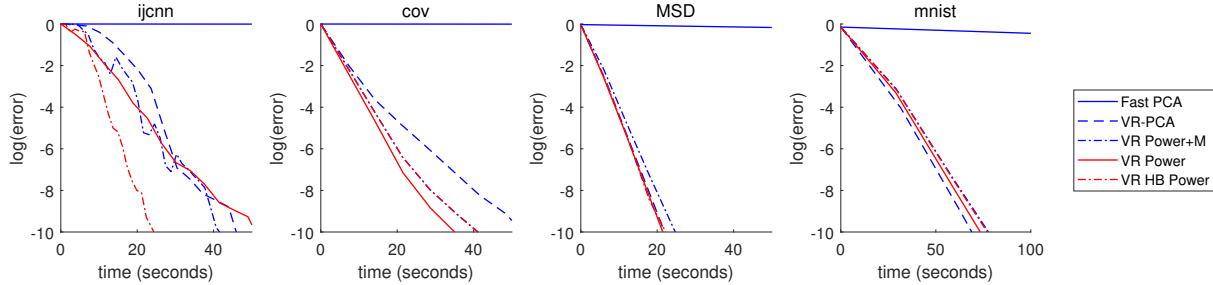


Figure 1: The comparison of stochastic variance-reduced PCA algorithms with hyper-parameters tuned

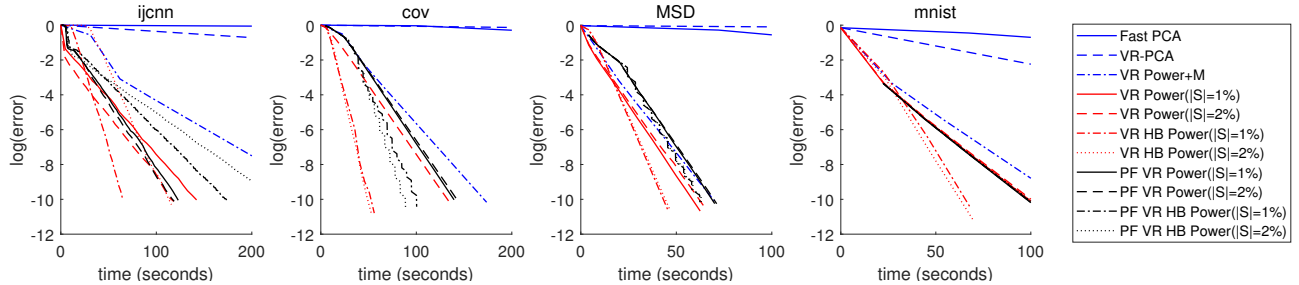


Figure 2: The comparison of stochastic variance-reduced PCA algorithms with recommended hyper-parameters and parameter-free algorithms

some larger values are included in the grid, VR-PCA would have a similar performance to VR Power even for cov. On the other hand, VR HB Power always performs better than VR Power+M due to its additional control through the step size. VR HB Power works particularly well for ijcnn which has the smallest eigen-gap. If the eigen-gap is large, the performance of VR HB Power is not much different from the performances of VR Power+M, VR-PCA and VR Power. We were not able to find good hyperparameters for Fast PCA.

Figure 2 shows the experimental result without parameter tuning. In the figure, regardless of the batch size, VR Power and VR HB Power outperform VR-PCA, VR Power+M and Fast PCA. Although VR Power and VR-PCA are similar algorithms, the performance of VR Power is much better than that of VR-PCA due to the choice of η and m . While VR Power precisely choose the values of η and m depending on the values of λ_1, λ_2 and $|S|$, VR-PCA does not utilize such information and let them depend only on n . As a result, the step size is too small and the epoch length is too large, leading to slow convergence. On the other hand, due to the extra dependency on \sqrt{d} , VR Power+M requires too large samples and thus it is slower than VR Power even for ijcnn which has the smallest eigen-gap. The epoch length m of SVRG in Fast PCA is of the order of $1/\Delta^2$. Therefore, Fast PCA takes a significant amount of time to solve each convex problem, which makes its optimality gap not decrease as sharply as other algorithms. On the other hand, PF VR HB Power takes

longer than VR HB Power while the performance of PF VR Power looks very similar to that of VR Power. This is because VR HB Power has the additional momentum parameter β , which makes its performance more affected by estimation errors. Nevertheless, both parameter-free algorithms work very well compared to VR-PCA, VR Power+M and Fast PCA.

6 Conclusion

In this paper, we present two mini-batch stochastic variance-reduced algorithms for PCA and derive exact forms of their parameters to attain the optimal runtime. Our results show that for any batch size, the optimal runtime can be achieved by appropriately choosing the step size and epoch length. We also introduce practical implementations which automatically find such values depending on batch sizes. The framework used in our analysis is not specific to the proposed algorithms but can be applied to analyze other stochastic variance-reduced PCA algorithms and improve their results. In our framework, the optimality gap is measured as the ratio of two expectation terms and this enables us to develop global convergence statements. Experimental results show that the proposed algorithms work well for arbitrary batch sizes.

References

- Allen-Zhu, Z. and Li, Y. (2016). LazySVD: Even faster SVD decomposition yet without agonizing pain. In *Advances in Neural Information Processing Systems*, pages 974–982.
- Allen-Zhu, Z. and Li, Y. (2017). First efficient convergence for streaming k-PCA: a global, gap-free, and near-optimal rate. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science*, pages 487–492. IEEE.
- Arora, R., Cotter, A., Livescu, K., and Srebro, N. (2012). Stochastic optimization for PCA and PLS. In *2012 50th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 861–868. IEEE.
- Balcan, M.-F., Du, S. S., Wang, Y., and Yu, A. W. (2016). An improved gap-dependency analysis of the noisy power method. In *Conference on Learning Theory*, pages 284–309.
- Balsubramani, A., Dasgupta, S., and Freund, Y. (2013). The fast convergence of incremental PCA. In *Advances in Neural Information Processing Systems*, pages 3174–3182.
- Bertin-Mahieux, T., Ellis, D. P., Whitman, B., and Lamere, P. (2011). The million song dataset. In *Proceedings of the 12th International Conference on Music Information Retrieval*.
- Blackard, J. A. and Dean, D. J. (1999). Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables. *Computers and Electronics in Agriculture*, 24(3):131–151.
- Boutsidis, C., Garber, D., Karnin, Z., and Liberty, E. (2015). Online principal components analysis. In *Proceedings of the Twenty-Sixth annual ACM-SIAM Symposium on Discrete Algorithms*, pages 887–901. Society for Industrial and Applied Mathematics.
- Chang, C.-C. and Lin, C.-J. (2011). Libsvm: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2(3):27.
- Dheeru, D. and Karra Taniskidou, E. (2017). UCI machine learning repository.
- Garber, D. and Hazan, E. (2015). Fast and simple PCA via convex optimization. *arXiv preprint arXiv:1509.05647*.
- Garber, D., Hazan, E., Jin, C., Kakade, S. M., Musco, C., Netrapalli, P., and Sidford, A. (2016). Faster eigenvector computation via shift-and-invert preconditioning. In *International Conference on Machine Learning*, pages 2626–2634.
- Golub, G. H. and Van Loan, C. F. (2012). *Matrix computations*, volume 3. JHU Press.
- Hardt, M. and Price, E. (2014). The noisy power method: A meta algorithm with applications. In *Advances in Neural Information Processing Systems*, pages 2861–2869.
- Jain, P., Jin, C., Kakade, S. M., Netrapalli, P., and Sidford, A. (2016). Streaming PCA: Matching matrix Bernstein and near-optimal finite sample guarantees for Ojas algorithm. In *Conference on Learning Theory*, pages 1147–1164.
- Johnson, R. and Zhang, T. (2013). Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems*, pages 315–323.
- Jolliffe, I. (2011). *Principal component analysis*. Springer.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.
- Lei, Q., Zhong, K., and Dhillon, I. S. (2016). Coordinate-wise power method. In *Advances in Neural Information Processing Systems*, pages 2064–2072.
- Li, C. J., Wang, M., Liu, H., and Zhang, T. (2018). Near-optimal stochastic approximation for online principal component estimation. *Mathematical Programming*, 167(1):75–97.
- Mason, J. C. and Handscomb, D. C. (2002). *Chebyshev polynomials*. Chapman and Hall/CRC.
- Mitliagkas, I., Caramanis, C., and Jain, P. (2013). Memory limited, streaming PCA. In *Advances in Neural Information Processing Systems*, pages 2886–2894.
- Oja, E. (1982). Simplified neuron model as a principal component analyzer. *Journal of Mathematical Biology*, 15(3):267–273.
- Polyak, B. T. (1964). Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17.
- Prokhorov, D. (2001). Ijcnv 2001 neural network competition. *Slide Presentation in International Joint Conference on Neural Networks*, 1:97.
- Shamir, O. (2015). A stochastic PCA and SVD algorithm with an exponential convergence rate. In *International Conference on Machine Learning*, pages 144–152.
- Shamir, O. (2016). Fast stochastic algorithms for SVD and PCA: Convergence properties and convexity. In *International Conference on Machine Learning*, pages 248–256.

- Wang, J., Wang, W., Garber, D., and Srebro, N. (2018). Efficient coordinate-wise leading eigenvector computation. In *Algorithmic Learning Theory*, pages 806–820.
- Xu, P., He, B., De Sa, C., Mitliagkas, I., and Re, C. (2018). Accelerated stochastic power iteration. In *International Conference on Artificial Intelligence and Statistics*, pages 58–67.
- Xu, Z. (2018). Gradient descent meets shift-and-invert preconditioning for eigenvector computation. In *Advances in Neural Information Processing Systems*, pages 2830–2839.