# Efficient Convex Relaxations for Streaming PCA

Teodor V. Marinov

Johns Hopkins University
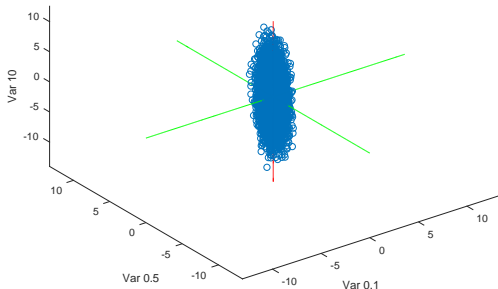
Joint work with Dr. Raman Arora
To appear at NeurIPS 2019

# Outline

# PCA as a geometric problem

- Given a data matrix $X \in \mathbb{R}^{d \times T}$
- Output a $U \in \mathbb{R}^{d \times k}, U^\top U = I_k$ to minimize reconstruction error: $\frac{1}{T}\|X - UU^\top X\|_F^2$

# PCA as a geometric problem

## Optimization problem

$$\underset{U \in \mathbb{R}^{d \times k}}{\text{minimize:}} \quad \frac{1}{T} \|X - UU^\top X\|_F^2$$

$$\text{subject to:} \quad U^\top U = I_k$$

# PCA as a geometric problem

## Optimization problem

$$\underset{U \in \mathbb{R}^{d \times k}}{\text{minimize:}} \quad \frac{1}{T} \|X - UU^\top X\|_F^2$$

$$\text{subject to:} \quad U^\top U = I_k$$

Quickly rewriting the objective:

$$
\begin{aligned}
\frac{1}{T} \|X - UU^\top X\|_F^2 &= \frac{1}{T} \left\{ \|X\|_F^2 + \|UU^\top X\|_F^2 - 2\text{Tr}\left(X^\top UU^\top X\right) \right\} \\
&= \frac{1}{T} \left\{ \|X\|_F^2 + \text{Tr}\left((UU^\top X)^\top UU^\top X\right) - 2\text{Tr}\left(X^\top UU^\top X\right) \right\} \\
&= \frac{1}{T} \left\{ \|X\|_F^2 + \text{Tr}\left(X^\top UU^\top UU^\top X\right) - 2\text{Tr}\left(X^\top UU^\top X\right) \right\} \\
&= \frac{1}{T} \left\{ \|X\|_F^2 - \text{Tr}\left(U^\top XX^\top U\right) \right\}
\end{aligned}
$$

# PCA as variance maximization

## Equivalent optimization problem

$$\underset{U \in \mathbb{R}^{d \times k}}{\text{maximize}} \quad \frac{1}{T} \text{Tr} \left( U^\top X X^\top U \right)$$

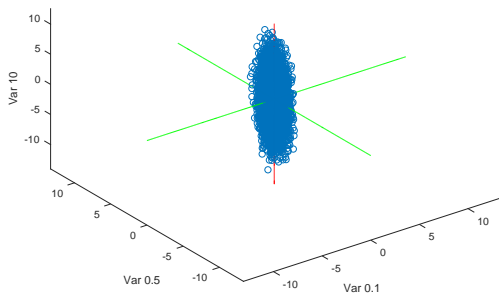$$\text{subject to} \quad U^\top U = I_k$$

# PCA as variance maximization

## Equivalent optimization problem

$$\underset{U \in \mathbb{R}^{d \times k}}{\text{maximize}} \quad \frac{1}{T} \text{Tr} \left( U^\top X X^\top U \right)$$

$$\text{subject to} \quad U^\top U = I_k$$

## Solution to optimization problem

Optimal solution is given by a set of $k$ eigenvectors associated with the top $k$ eigenvalues of $\frac{1}{T} X X^\top$

# The stochastic optimization perspective



Each $\mathrm{x}_t \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} 0.1 & 0 & 0 \\ 0 & 0.5 & 0 \\ 0 & 0 & 10 \end{pmatrix}\right)$ independently.

# The stochastic optimization perspective

- Assume each column $x_t$ of $X$ is i.i.d as some probability law $\mathcal{D}$.
- Empirical problem is just a proxy for the stochastic optimization problem:

## Stochastic optimization problem

$$\underset{U \in \mathbb{R}^{d \times k}}{\text{minimize:}} \quad \mathbb{E}_{x \sim \mathcal{D}} \|x - UU^{\top}x\|_F^2$$

$$\text{subject to:} \quad U^{\top}U = I_k$$

# The stochastic optimization perspective

## Minimizing reconstruction

$$\underset{U \in \mathbb{R}^{d \times k}}{\text{minimize:}} \quad \mathbb{E}_{x \sim \mathcal{D}} \| x - UU^\top x \|_F^2$$

$$\text{subject to:} \quad U^\top U = I_k$$

$$\iff$$

## Variance maximization

$$\underset{U \in \mathbb{R}^{d \times k}}{\text{maximize}} \quad \mathbb{E}_{x \sim \mathcal{D}} \text{Tr}\left( U^\top x x^\top U \right)$$

$$\text{subject to} \quad U^\top U = I_k \tag{1}$$

# Oja's algorithm

- One possible way to solve Problem 1 is to use Stochastic Gradient Descent (SGD)
- Gradient for objective in Problem 1 is $2\mathbb{E}_{x \sim \mathcal{D}}[xx^\top]U$
- Since we do not have direct access to distribution $\mathcal{D}$, we use an unbiased estimator of the gradient based on a sample $x_t \sim \mathcal{D}$ given by $x_t x_t^\top U$

# Oja's algorithm

SGD on Problem 1 is also known as Stochastic Power Method (Oja's algorithm) Allen-Zhu and Li [2017]

$$U_t \leftarrow (I + \eta_t x_t x_t^\top) U_{t-1}, U_t = \mathsf{Orth}(U_t); \qquad (2)$$

### Convergence guarantee (informal)

After $T$ iterations of Oja's algorithm w.p. $1 - \delta$ it holds that

$$\|U_*^\top U_T\|_F^2 \le k - \tilde{\mathcal{O}}\left(\frac{1}{\Delta(C)^2 T}\right),$$

where $C = \mathbb{E}_{x \sim \mathcal{D}}[xx^\top]$ and $\Delta(C)$ is the eigengap at the $k$-th eigenvalue.

$$\Delta(C) := \lambda_k(C) - \lambda_{k+1}(C)$$

# Convex relaxations to PCA

## Maximizing variance formulation

$$\underset{U \in \mathbb{R}^{d \times k}}{\text{maximize}} \quad \mathbb{E}_{x \sim \mathcal{D}} \text{Tr}\left(UU^\top xx^\top\right)$$

$$\text{subject to} \quad U^\top U = I_k$$

The above is equivalent to:

## Maximizing variance formulation

$$\underset{P \in \mathbb{R}^{d \times d}}{\text{maximize}} \quad \mathbb{E}_{x \sim \mathcal{D}} \text{Tr}\left(Pxx^\top\right)$$

$$\text{subject to} \quad P^2 = P, P^\top = P, \text{rank}(P) = k$$

# Convex relaxations to PCA (continued)

The convex hull of $\{P \in \mathbb{R}^{d \times d} : P^2 = P, P^\top = P, \mathsf{rank}(P) = k\}$ is
$\{P \in \mathbb{R}^{d \times d} : \mathrm{Tr}\,(P) \leq k, 0 \preceq P \preceq I, P^\top = P\}$

## Convex relaxation [Arora et al., 2013]

$$\begin{aligned}
\underset{P \in \mathbb{R}^{d \times d}}{\text{maximize}} \quad & \mathrm{Tr}\,(PC) \\
\text{subject to} \quad & \mathrm{Tr}\,(P) \leq k, 0 \preceq P \preceq I, P^\top = P
\end{aligned}, \tag{3}$$

## Convex relaxation with regularization [Mianjy and Arora, 2018]

$$\begin{aligned}
\underset{P \in \mathbb{R}^{d \times d}}{\text{maximize}} \quad & \mathrm{Tr}\,(PC) - \frac{\lambda}{2}\|P\|_F^2 \\
\text{subject to} \quad & \mathrm{Tr}\,(P) \leq k, 0 \preceq P \preceq I, P^\top = P
\end{aligned}. \tag{4}$$

# Convex relaxations to PCA (continued)

## Projected SGD for Problem 3 (MSG)

$$P_t \leftarrow \mathcal{P}\left(P_{t-1} + \eta_t x_t x_t^\top\right)$$

## Projected SGD for Problem 4 (RMSG)

$$P_t \leftarrow \mathcal{P}\left((1 - \lambda\eta_t)P_{t-1} + \eta_t x_t x_t^\top\right)$$

In the above $\mathcal{P}(\cdot)$ is the projection onto the convex set of constraints $\{P : \mathrm{Tr}(P) \leq k, 0 \preceq P \preceq I, P^\top = P\}$.

# Convex relaxations to PCA (continued)

Running projected SGD on the above problems, comes with the following guarantees:

### Convergence guarantee for MSG (informal)

After $T$ iterations of MSG, it holds that

$$\mathbb{E}[\langle P_* - P_T, C \rangle] \leq \tilde{\mathcal{O}}\left(\frac{1}{\sqrt{T}}\right).$$

### Convergence guarantee for RMSG (informal)

After $T$ iterations of RMSG, it holds that

$$\mathbb{E}[\langle P_* - P_T, C \rangle] \leq \tilde{\mathcal{O}}\left(\frac{1}{\Delta(C)^2 T}\right).$$

# Angle between subspaces and suboptimality in objective

- Oja's guarantee is of the form $k - \|U_*^\top U_T\|_F^2 \le \epsilon$
- MSG and RMSG guarantees are of the form $\langle U_* U_*^\top - U_T U_T^\top, C \rangle \le \epsilon$
- We have [Mianjy and Arora, 2018]

$$\langle U_* U_*^\top - U_T U_T^\top, C \rangle \le \lambda_1(C)(k - \|U_*^\top U_T\|_F^2)$$

- No known relation in opposite direction

- The computational complexity of Oja's algorithm per iteration is $O(dk^2)$ (can be reduced to $O(dk)$ if we do not call the ($Orth$) procedure)
- The computational complexity of MSG and RMSG per iteration is $O(d\text{rank}(\text{P}_t)^2)$
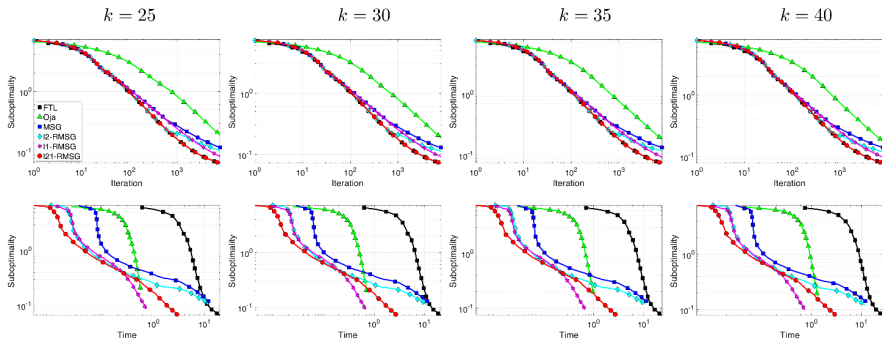- Worst case for MSG and RMSG is $O(d^3)$

Figure 3: Experiment on MNIST by Mianjy and Arora [2018]

# Controlling the rank of $P_t$

- If distribution is well-behaved rank of $P_t$ stays in $O(k)$
- We get a similar computational complexity to Oja's algorithm but experiments suggest MSG and RMSG have better convergence properties
- Can we formalize the above statements through theoretical results?

- If distribution is well-behaved rank of $P_t$ stays in $O(k)$
- We get a similar computational complexity to Oja's algorithm but experiments suggest MSG and RMSG have better convergence properties
- Can we formalize the above statements through theoretical results?

Yes, if we make slight tweaks to MSG and RMSG!

# MB-RMSG

**Input:** Stream of data $\{x_{t_n}\}$, parameters $\Delta(C)$, probability of failure $\delta$, number of components $k$

**Output:** $P_T$

1: $n = \log(d/\delta) \frac{128k \log(1/\delta)}{\Delta(C)^5}$

2: $P_1 = \text{Top-k}(\frac{1}{n} \sum_{l=1}^{n} x_{0_l} x_{0_l}^\top)$

3:                                     %% $\{x_{0_l}\}_{l=1}^{n}$ is the warm-start mini-batch

4: $n = \log\left(\frac{Td}{\delta}\right) \frac{8(k+1)}{\Delta(C)^2}$

5: **for** $t = 1, \ldots, T-1$ **do**

6:      $\eta_t = \frac{1}{\frac{\Delta(C)}{2}\left(t + \frac{128\log\left(\frac{1}{\delta}\right)}{\Delta(C)^3}\right)}$

7:      $C_t \leftarrow \frac{1}{n} \sum_{l=1}^{n} x_{t_l} x_{t_l}^\top$

8:                                 %% $\{x_{t_l}\}_{l=1}^{n}$ is the mini-batch for the $t^{th}$ epoch

9:      $P_{t+1/2} \leftarrow (1 - \frac{\Delta(C)}{2}\eta_t)P_t + \eta_t C_t$

10:      $P_{t+1} = \mathcal{P}(P_{t+1/2})$

11: **end for**

# Formal guarantee 1

## Theorem

*There exists an algorithm, solving Problem 4, which after $T$ iterations, with probability at least $1 - 3e\delta$, returns a sequence of iterates $\{P_t\}_{t=1}^{T}$, such that for all $t \leq T$*

$$\langle P^* - P_t, C \rangle \leq \frac{32 \log (1/\delta)}{\Delta(C)^2 \left( t + \frac{1}{\gamma} - 1 \right)},$$

*where $\gamma = \frac{\Delta(C)^3}{128 \log(1/\delta)}$. Further, for all $t \leq T$ it holds that $P_t$ is a rank-$k$ projection matrix and the per-iteration computational complexity of the algorithm is bounded by $\tilde{O} \left( \frac{dk^2}{\Delta(C)^2} + dk^2 \right)$.*

- Total computational complexity for $\epsilon$-suboptimality of MB-RMSG is $\tilde{O}(\frac{dk^2}{\Delta(\mathrm{C})^4\epsilon})$
- Total computational complexity for Oja's algorithm to reach $\epsilon$-suboptimality is $\tilde{O}(\frac{dk}{\Delta(\mathrm{C})^2\epsilon})$

<div align="center">We can do better in expectation!</div>

# Formal guarantee 2

## Theorem

*Let $\mathcal{A}$ be the event that for all $t \in [T]$ it holds that $\|C_t - C\| \leq \frac{\Delta(C)}{8(k+1)}$ and $P_t$ is a rank-$k$ projection matrix. Then Algorithm 22 guarantees that $\mathcal{A}$ occurs with probability at least $1 - \delta$ and that*

$$\mathbb{E}\left[\langle P^* - P_T, C \rangle | \mathcal{A}\right] \leq \tilde{O}\left(\frac{\Delta(C)}{T} + \min(\Delta(C) \times d, 1)\frac{1}{kT}\right).$$

Above theorem implies that the total computational complexity for achieving $\epsilon$-suboptimality is $\tilde{O}\left(\frac{dk^2}{\epsilon\Delta(C)^2} \times \min(d\Delta(C)), 1)\right)$, which is only a factor of $k$ away from Oja's algorithm whenever the gap is large, and actually improves by a factor of $1/\Delta(C)$ over Oja's in the case when $\Delta(C) \in o(1/kd)$!

# Key lemma

## Lemma

*Let $P_t$ be rank k and suppose $\|C - C_t\| \leq \beta$. Then, a sufficient condition for $P_{t+1}$ to be rank k is*

$$\langle P_t, C \rangle \geq \langle P^*, C \rangle - \frac{\Delta(C)}{2} + \frac{\lambda}{2} + \beta(k+1). \tag{5}$$

- Projection works by shifting all eigenvalues of $P_t + \eta_t C_t$ and then clipping them between 0 and 1
- Let $\lambda_k(P_t + \eta_t C_t) = 1 + \lambda_k$ and $\lambda_{k+1}(P_t + \eta_t C_t) = \lambda_{k+1}$ for some $\lambda_k$ and $\lambda_{k+1}$
- If $\lambda_k > \lambda_{k+1}$, then the shift from the projection is larger than $\lambda_{k+1}$ and thus projection will clip $\lambda_{k+1}$ to 0.
- We can guarantee that this happens with high probability if $P_t$ is close enough to $P^*$ and $C_t$ is close enough to $C$.

# Rest of analysis

- The rest of the analysis requires last iterate SGD guarantees with high probability
- For results in expectation we need to adapt the analysis for smooth SGD
- Both of these are non-trivial to do!
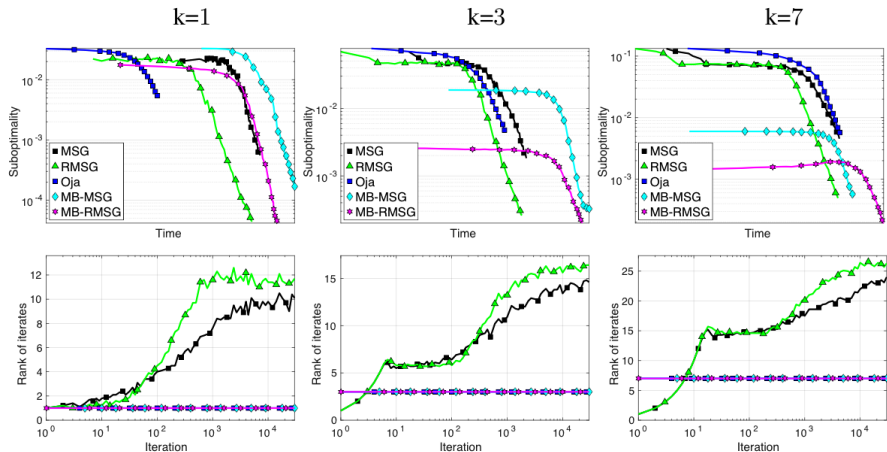
# But does it work?



Figure 4: Experiments on MNIST

# Open problems

- We gave an algorithm based on the convex relaxation to PCA which achieves (almost) optimal rates
- But goal was to study MSG and RMSG directly
- We need better tools to control rank of MSG and RMSG iterates

# Other related work

- There are many other works on Streaming PCA, mainly focused on studying Oja's algorithm [De Sa et al., 2014, Hardt and Price, 2014, Balcan et al., 2016, Jain et al., 2016, Shamir, 2016a,b, Allen-Zhu and Li, 2017, Li et al., 2018]

- Other important work in the Online Learning setting is by Warmuth and Kuzmin [2008], Grabowska and Kotłowski [2018], Garber [2018]

Zeyuan Allen-Zhu and Yuanzhi Li. 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 (FOCS)*, pages 487–492. IEEE, 2017.

Raman Arora, Andy Cotter, and Nati Srebro. Stochastic optimization of pca with capped msg. In *Advances in Neural Information Processing Systems*, pages 1815–1823, 2013.

Maria-Florina Balcan, Simon Shaolei Du, Yining Wang, and Adams Wei Yu. An improved gap-dependency analysis of the noisy power method. In *Conference on Learning Theory*, pages 284–309, 2016.

Christopher De Sa, Kunle Olukotun, and Christopher Ré. Global convergence of stochastic gradient descent for some non-convex matrix problems. *arXiv preprint arXiv:1411.1134*, 2014.

Dan Garber. On the regret minimization of nonconvex online gradient ascent for online pca. *arXiv preprint arXiv:1809.10491*, 2018.

Monika Grabowska and Wojciech Kotłowski. Online principal component analysis for evolving data streams. In *International Symposium on Computer and Information Sciences*, pages 130–137. Springer, 2018.

Moritz Hardt and Eric Price. The noisy power method: A meta algorithm with applications. In *Advances in Neural Information Processing Systems*, pages 2861–2869, 2014.

Prateek Jain, Chi Jin, Sham M Kakade, Praneeth Netrapalli, and Aaron Sidford. Streaming pca: Matching matrix bernstein and near-optimal finite sample guarantees for oja's algorithm. In *Conference on Learning Theory*, pages 1147–1164, 2016.

Chris Junchi Li, Mengdi Wang, Han Liu, and Tong Zhang. Near-optimal stochastic approximation for online principal component estimation. *Mathematical Programming*, 167(1):75–97, 2018.

Poorya Mianjy and Raman Arora. Stochastic pca with $\ell_2$ and $\ell_1$ regularization. In *International Conference on Machine Learning*, pages 3531–3539, 2018.

Ohad Shamir. Convergence of stochastic gradient descent for pca. In *International Conference on Machine Learning*, pages 257–265, 2016a.

Ohad Shamir. Fast stochastic algorithms for svd and pca: Convergence properties and convexity. In *International Conference on Machine Learning*, pages 248–256, 2016b.

Manfred K Warmuth and Dima Kuzmin. Randomized online pca algorithms with regret bounds that are logarithmic in the dimension. *Journal of Machine Learning Research*, 9(Oct):2287–2320, 2008.