Efficient Convex Relaxations for Streaming PCA

Raman Arora, Teodor V. Marinov Johns Hopkins University, Baltimore, MD 21204

1. Streaming PCA in Stochastic Setting

- Given a sequence of vectors $(\mathbf{x}_t)_{t=1}^{\infty}$ i.i.d. $\mathbf{x}_t \sim \mathcal{D}$.
- Minimize reconstruction error $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \|\mathbf{x} P\mathbf{x}\|^2$, where $P \in$ $\mathcal{P}_k = \{ \mathbf{P} : \mathbf{P}^2 = \mathbf{P}, \mathbf{P}^\top = \mathbf{P}, rank(\mathbf{P}) = k \}.$
- Equivalently solve:

• Usually solved by Oja's algorithm:

$$\overline{\mathbf{U}_{t+1} \leftarrow (\mathbf{I} + \eta_t \mathbf{x}_t \mathbf{x}_t^\top) \mathbf{U}_t, \mathbf{U}_{t+1}} \leftarrow \text{Orth}(\mathbf{U}_{t+1})$$

$$P_{t+1} \leftarrow \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top$$

2. A Convex Relaxation

- Convexify $\mathscr{P}_k \to \mathscr{C} = \{P : Tr(P) = k, 0 \leq P \leq I, P = P^{\top}\}.$
- Can now solve the following convex optimization problem:

 $\underset{P \in \mathbb{R}^{d \times d}}{\mathsf{maximize}} \quad \mathbb{E}_{\mathbf{x} \sim \mathfrak{D}} \langle P, \mathbf{x} \mathbf{x}^{\top} \rangle$ subject to $P \in \mathscr{C}$

• Can also add regularization to achieve distribution dependent guarantees:

$$\begin{bmatrix} \underset{P \in \mathbb{R}^{d \times d}}{\text{maximize}} & \mathbb{E}_{\mathbf{x} \sim \mathfrak{D}} \langle P, \mathbf{x} \mathbf{x}^{\top} \rangle - \frac{\lambda}{2} \|P\|_F^2 \\ \text{subject to} & P \in \mathscr{C} \end{bmatrix}$$
(3)

3. Prior Work Guarantees

- Let P^* be the optimal solution to Problem 1, $C = \mathbb{E}[\mathbf{x}\mathbf{x}^{\top}]$, and $\Delta(\mathbf{C}) = \lambda_k(\mathbf{C}) - \lambda_{k+1}(\mathbf{C})$.
- Projected Stochastic Gradient Descent on Problem 2 (MSG) is: $P_{t+1} = \Pi_{\mathscr{C}}(P_t + \eta_t \mathbf{x}_t \mathbf{x}_t^{\top})$
- Projected Stochastic Gradient Descent on Problem 3 (RMSG) is:

 $P_{t+1} = \Pi_{\mathscr{C}}((1 - \lambda \eta_t)P_t + \eta_t \mathbf{x}_t \mathbf{x}_t^{\top})$ • MSG and RMSG guarantee that $\mathbb{E}\langle P^* - P_t, C \rangle$ is in $\tilde{O}(1/\sqrt{t})$ and $\tilde{O}(1/(\Delta(C)t))$ [2].

• RMSG is statistically optimal, however, MSG and RMSG can take up to $\Omega(d^3)$ computations per iteration!

4. Meta-algorithm and Intuition

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

Output: P_T

- 1: Initialize P₁ from a warm start
- Form unbiased mini-batched estimator of gradient:
- 6: $P_{t+1} = \Pi(P_{t+1/2})$ (gradient descent upgrade) or equivalently
- 8: $P_{t+1} = U_{t+1}U_{t+1}^{\top}$ (efficient update)

Lemma 1 Let P_t be rank k and suppose $\|C - C_t\| \le \epsilon$. Then a sufficient condition for P_{t+1} to be rank k is

Intuition behind algorithm: Lemma 1 gives a sufficient conditions for P_{t+1} to be a projection matrix, given that P_t is a projection matrix. The

mini-batched stochastic gradients

- 2: **for** t = 1, ..., T 1 **do**
- 3: $\eta_t = \Theta\left(\frac{1}{\Delta(C)t}\right)$
- $\mathsf{C}_t = rac{1}{n} \sum_{l=1}^n \mathbf{x}_{t_l} \mathbf{x}_{t_l}^{ op}$
- 5: $P_{t+1/2} \leftarrow (1 \frac{\Delta(C)}{2} \eta_t) P_t + \eta_t C_t$
- 7: $U_{t+1} = \text{Top-k}\left(\left[\sqrt{1 \Delta(C)\eta_t/2}U_t, \sqrt{\eta_t}\mathbf{X}_t\right]\right)$
- 9: **end for**

$$\langle P_t, C \rangle \ge \langle P^*, C \rangle - \frac{\Delta(C)}{4} + \epsilon(k+1).$$
 (4)

sufficient condition translates to:

- P_t is close enough to P^* in objective and hence the warm start of the algorithm
- C_t is close enough to C with high probability which results in the

7. REFERENCES

- Allen-Zhu, Z. and Li. Y. First efficient convergence for streaming k-pca: a global, gap-free, and near-optimal rate. FOCS, 2017.
- [2] Mianjy, P. and Arora, R. Stochastic PCA with ℓ_2 and ℓ_1 Regularization. ICML, 2018.

5. Main results

High Probability Guarantees

The following holds for Algorithm MB-MSG: with probability at least $1 - \delta$, for all $t \leq T$

$$\langle \mathbf{P}^* - \mathbf{P}_t, \mathbf{C} \rangle \le O\left(\frac{k^4 \log(1/\delta)(\log(T))^2}{\sqrt{t + \frac{1}{\gamma}}}\right),$$

where $\gamma = O\left(\frac{\Delta(C)^2}{(k\log(1/\delta))^2}\right)$. Further, it holds that P_t is a rank-k projection matrix.

The following holds for Algorithm MB-RMSG: with probability at least $1 - \delta$, for all $t \leq T$

$$\langle \mathbf{P}^* - \mathbf{P}_t, \mathbf{C} \rangle \le \frac{32 \log (3e/\delta)}{\Delta(\mathbf{C})^2 \left(t + \frac{1}{\gamma} - 1\right)},$$

where $\gamma = \frac{\Delta(C)^3}{128\log(1/\delta)}$. Further, it holds that P_t is a rank-k projection matrix.

Guarantees in Expectation

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

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

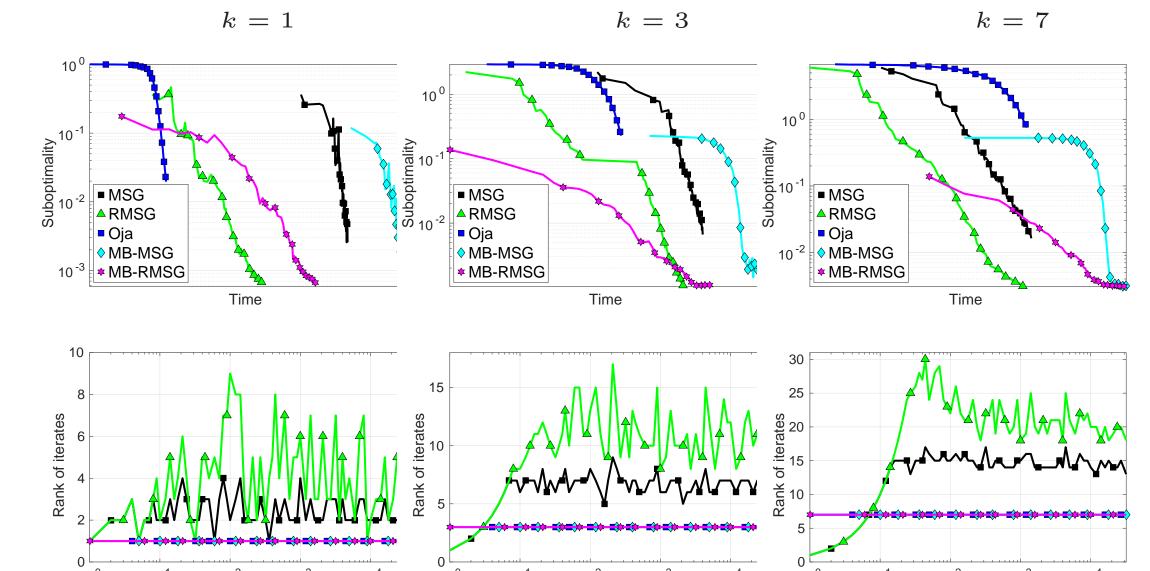
Comparison to Oja's Algorithm

(Informal [1]) The following holds for Algorithm Oja's **Algorithm**: with probability at least $1 - \delta$, for all $t \leq T$

$$\langle \mathbf{P}^* - \mathbf{P}_t, \mathbf{C} \rangle \le \tilde{O}\left(\frac{1}{t\Delta(\mathbf{C})^2}\right)$$

- Total computational complexity for ϵ -suboptimality for **MB-RMSG** is $\tilde{O}\left(dk^2/(\epsilon\Delta(\mathbf{C})^2)\times\min(d\Delta(\mathbf{C})),1\right)$
- Total computational complexity for ϵ -suboptimality for Oja's Algorithm is $\tilde{O}\left(dk/(\epsilon\Delta(C)^2)\right)$

6. Experimental Results



Comparisons of Oja, MSG, RMSG, MB-MSG and MB-RMSG on synthetic data on the left and MNIST, on the right, in terms of runtime (top) and rank of iterates

