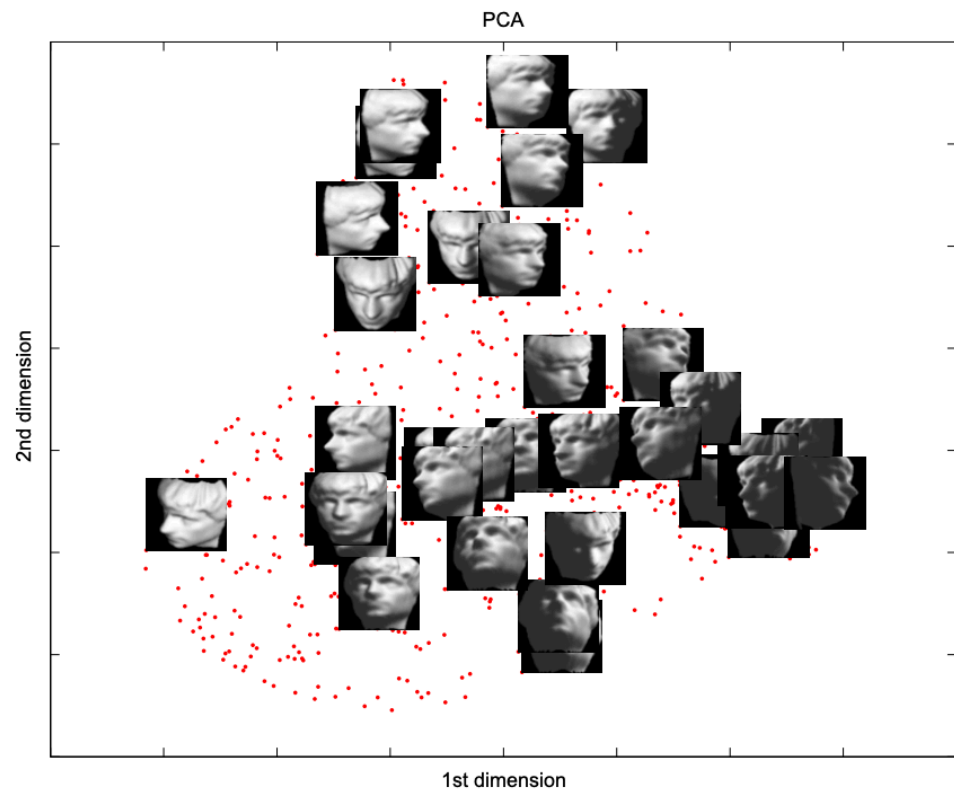
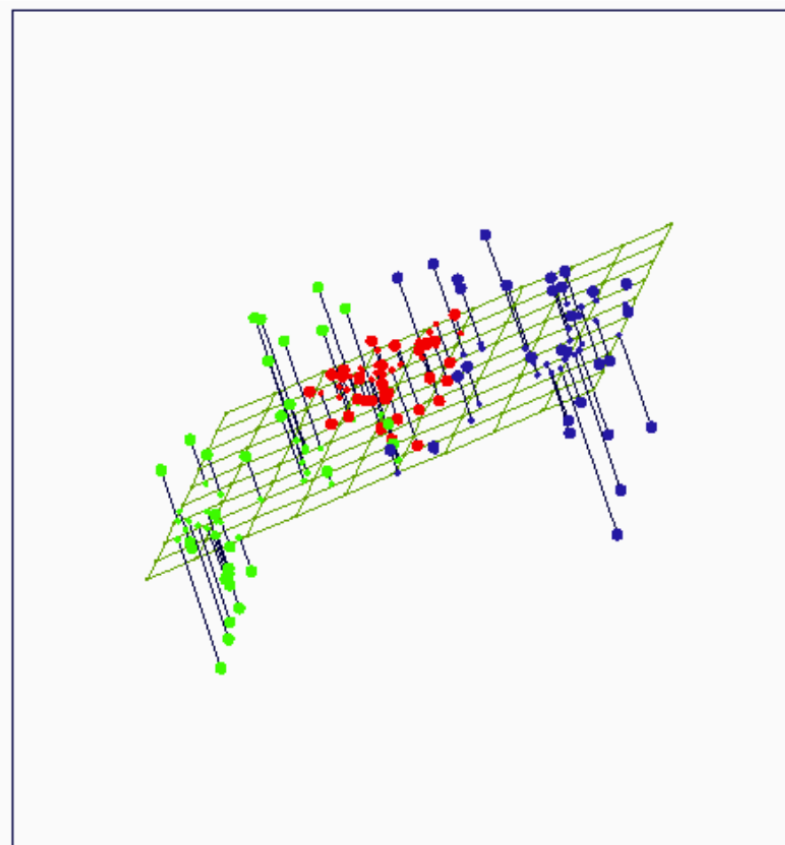


Exponentially convergent stochastic k-PCA *without* variance reduction

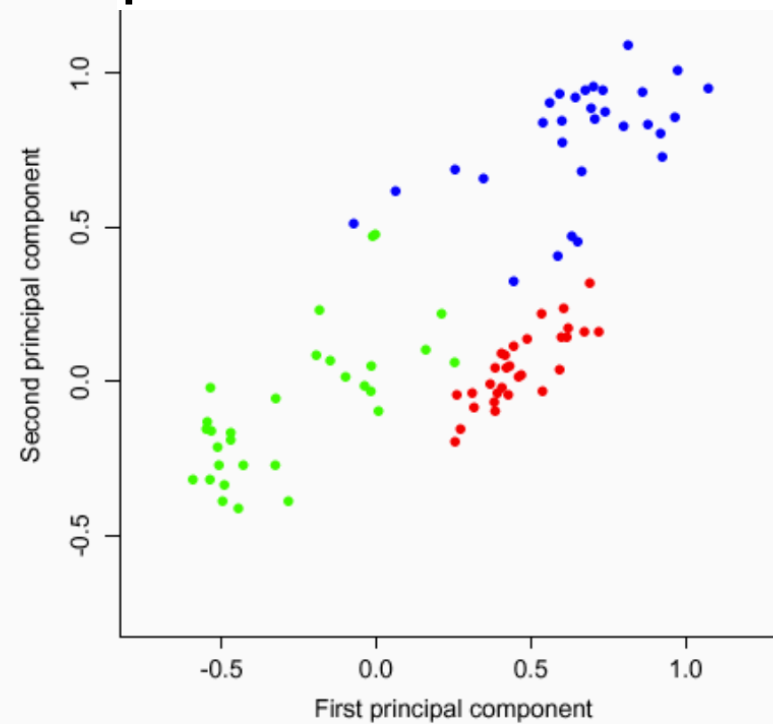
Cheng Tang
Amazon (AWS AI Labs)



k-PCA: popular dimension reduction method with many applications



picture credit to Ali Ghodsi



The best rank-two linear approximation

Batch vs Online k-PCA

Batch k-pca

- ★ Apply k-PCA on the entire dataset
- ★ Usually done by SVD

Online k-PCA

- ★ Approximate k-PCA without observing all data
- ★ Data points come one-by-one

Two popular online 1-PCA algorithms & their convergence rates

Oja's rule

- ☆ Online version of power method (has a *length normalization* step)
- ☆ Converges in order $O(1/t)$ for 1-PCA
- ☆ Connection to neural networks

$$w^t \leftarrow \frac{w^{t-1} + \eta^t x x^T w^{t-1}}{\|w^{t-1} + \eta^t x x^T w^{t-1}\|}$$

Balsubramani-Dasgupta-Freud 2013

Krasulina's method

- ☆ Stochastic gradient ascent on Rayleigh quotient for 1-PCA
- ☆ Converges in order $O(1/t)$ for 1-PCA

$$w^t \leftarrow w^{t-1} + \eta^t \left(x x^T - \frac{(x^T w^{t-1})^2}{\|w^{t-1}\|^2} I_d \right) w^{t-1}$$

Balsubramani-Dasgupta-Freud 2013

Two popular online 1-PCA algorithms & their convergence rates

Oja's rule

- ★ Easy to generalize to k-PCA
- ★ Converges in order $O(1/t)$ for k-PCA by recent analysis

Allen-Zhu 2017

Krasulina's method

- ★ Not obvious how to generalize to k-PCA

VR-PCA

- ★ Mixed Oja's rule (online) & power iteration (batch)
- ★ Exponential convergence rate

Shamir 2016

Algorithm 1 VR-PCA: Vector version ($k = 1$)

Parameters: Step size η , epoch length m

Input: Data matrix $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)$; Initial unit vector $\tilde{\mathbf{w}}_0$

for $s = 1, 2, \dots$ **do**

$$\tilde{\mathbf{u}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i (\mathbf{x}_i^\top \tilde{\mathbf{w}}_{s-1})$$

$$\mathbf{w}_0 = \tilde{\mathbf{w}}_{s-1}$$

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

Pick $i_t \in \{1, \dots, n\}$ uniformly at random

$$\mathbf{w}'_t = \mathbf{w}_{t-1} + \eta (\mathbf{x}_{i_t} (\mathbf{x}_{i_t}^\top \mathbf{w}_{t-1} - \mathbf{x}_{i_t}^\top \tilde{\mathbf{w}}_{s-1}) + \tilde{\mathbf{u}})$$

$$\mathbf{w}_t = \frac{1}{\|\mathbf{w}'_t\|} \mathbf{w}'_t$$

end for

$$\tilde{\mathbf{w}}_s = \mathbf{w}_m$$

end for

Algorithm 2 VR-PCA: Block version

Parameters: Rank k , Step size η , epoch length m

Input: Data matrix $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)$; Initial $d \times k$ matrix \tilde{W}_0 with orthonormal columns

for $s = 1, 2, \dots$ **do**

$$\tilde{U} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i (\mathbf{x}_i^\top \tilde{W}_{s-1})$$

$$W_0 = \tilde{W}_{s-1}$$

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

$B_{t-1} = VU^\top$, where USV^\top is an SVD decomposition of $W_{t-1}^\top \tilde{W}_{s-1}$

▷ Equivalent to

$$B_{t-1} = \arg \min_{B^\top B = I} \|W_{t-1} - \tilde{W}_{s-1} B\|_F^2$$

Pick $i_t \in \{1, \dots, n\}$ uniformly at random

$$W'_t = W_{t-1} +$$

$$\eta \left(\mathbf{x}_{i_t} (\mathbf{x}_{i_t}^\top W_{t-1} - \mathbf{x}_{i_t}^\top \tilde{W}_{s-1} B_{t-1}) + \tilde{U} B_{t-1} \right)$$

$$W_t = W'_t (W_t'^\top W_t')^{-1/2}$$

end for

$$\tilde{W}_s = W_m$$

end for

Information-theoretic lower bound

- This implies a $O(1/t)$ upper bound on convergence rate of *ALL* online k-PCA algorithms

$$\mathbb{E}[\Delta^n] \geq \Omega\left(\frac{\sigma^2}{n}\right)$$

for $\frac{\lambda_1 \lambda_{k+1}}{(\lambda_k - \lambda_{k+1})^2}$,

Vu and Lei, AISTATS, 2012

Our observation

When data is of rank 1, Krasulina's method for 1-PCA can be viewed as

$$W^t \leftarrow W^{t-1} + \eta^t s^t (r^t)^T$$

$$Var(s^t r^t) \propto \text{loss}$$

Krasulina's method is already doing variance-reduction on low-rank data!

Information-theoretic lower bound — a closer look

- The lower bound to-the-right implies a $1/t$ upper bound on convergence rate of online k-PCA algorithms **in the general case**
- Our result says otherwise when data is **low-rank**

$$\mathbb{E}[\Delta^n] \geq \Omega\left(\frac{\sigma^2}{n}\right)$$

for $\frac{\lambda_1 \lambda_{k+1}}{(\lambda_k - \lambda_{k+1})^2},$

Vu and Lei, AISTATS, 2012

Matrix Krasulina

$$W^t \leftarrow W^{t-1} + \eta^t \underbrace{s^t}_{\text{Projection}} \underbrace{(r^t)^T}_{\text{Residual}}$$

Orthonormalize rows of W^t

Makes sure r is orthogonal to the k -subspace

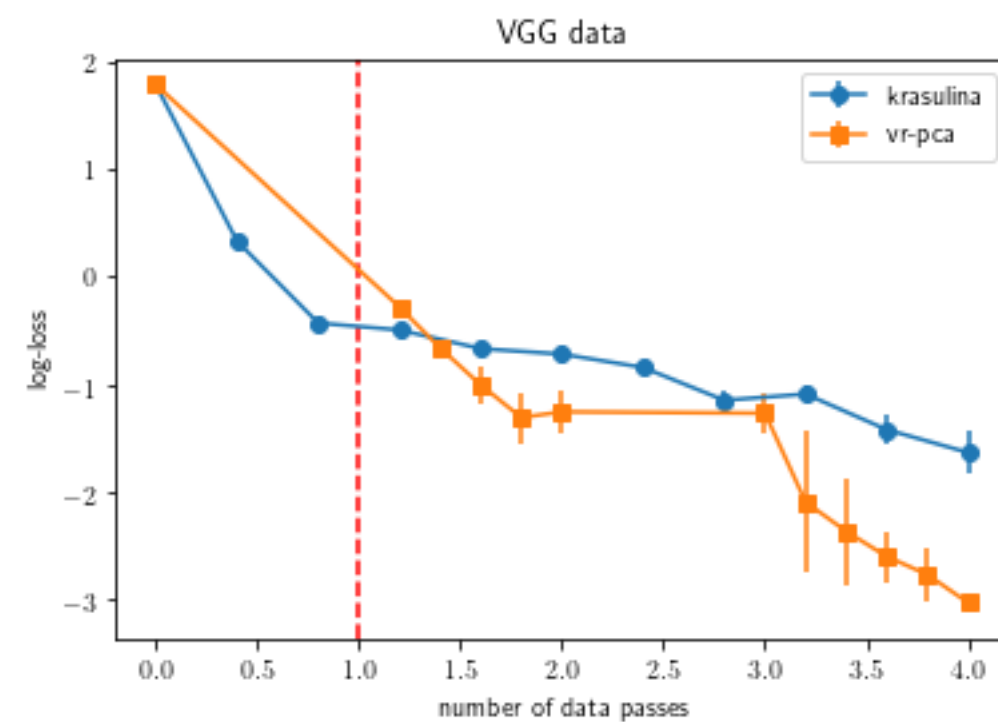
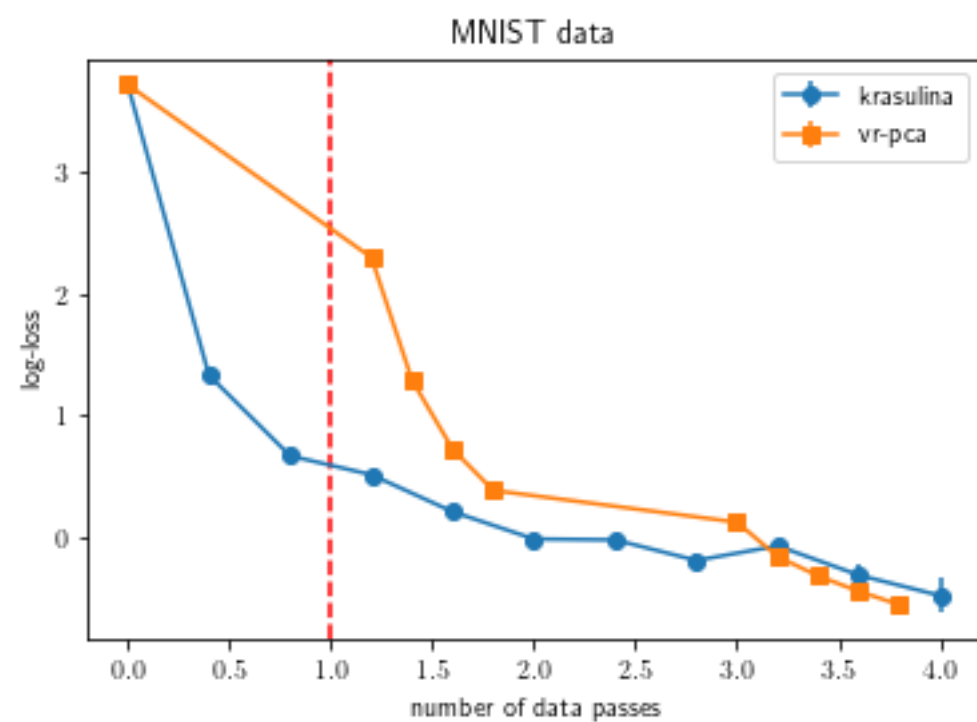
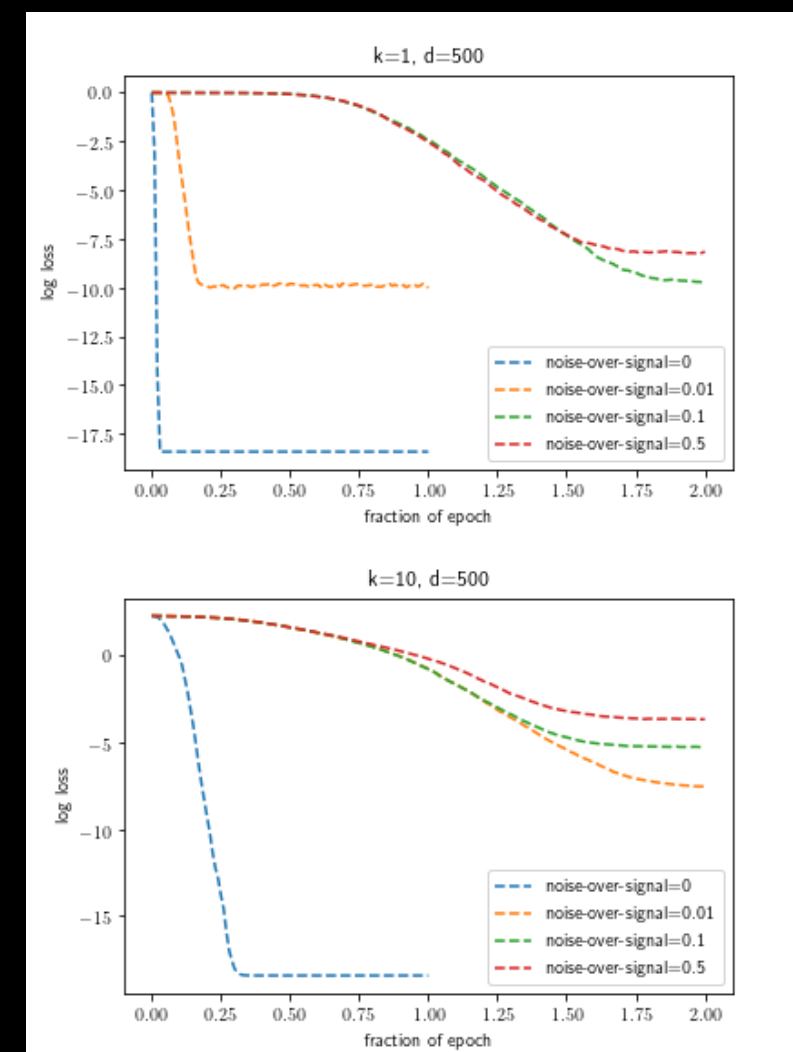
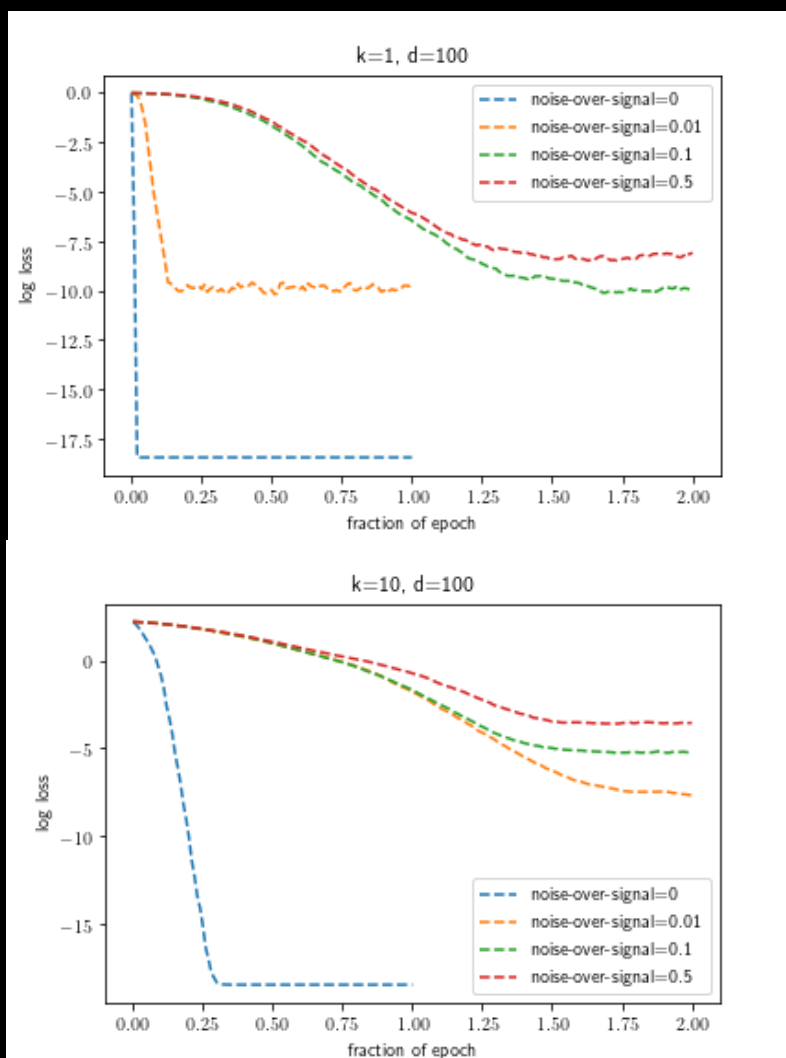
Our main results

If $\text{Cov}(X)$ is of low rank
Matrix Krasulina converges
exponentially for k-PCA

- ★ Set learning rate to be constant

In general case
Matrix Krasulina converges
of order $O(1/t)$ for k-PCA

- ★ Set learning rate to be $\sim 1/t$



Open problem

Can we characterize convergence rate
when data is *nearly low-rank* ?

