# An Implicit Form of Krasulina's k-PCA Update without the Orthonormality Constraint

Ehsan Amid and Manfred K. Warmuth
University of California, Santa Cruz
Google Brain, Mountain View
{eamid, manfred}@google.com

## Implicit vs. Explicit Update

Gradient descent is motivated by

$$\theta_{t+1} = \underset{\theta}{\arg\min} \left( \underbrace{1/2\eta \, \|\theta - \theta_t\|^2}_{\text{inertia}} + \text{loss}(\theta) \right)$$

The actual minimizer of **inertia + loss**

$$\theta_{t+1} = \theta_t - \eta \, \nabla \text{loss}(\theta_{t+1}) \qquad \text{(implicit update)}$$

Commonly approximated by the old gradient

$$\theta_{t+1} \approx \theta_t - \eta \, \nabla \text{loss}(\theta_t) \qquad \text{(explicit update)}$$

**Most of Machine Learning is explicit!**

## k-PCA Loss

**k-PCA:** Given zero mean random variable $y \in \mathbb{R}^d$
Find the (d x d) projection matrix $P$ of rank-k s.t.

$$\ell_{\text{comp}}(P) = \mathbb{E}\left[\|P\,y - y\|^2\right] = \text{tr}\left((I_d - P)\,\mathbb{E}[y\,y^\top]\right)$$

or

$$\ell_{\text{var}}(P) = -P\,\text{tr}(\mathbb{E}[yy^\top])$$

is minimized.

**Online k-PCA:** Observe one example $y_t$ at a time
**Decomposition:** $P = C\,(C^\top C)^{-1} C^\top := CC^\dagger$

Solve for $C$ instead of $P$ !

## GD on the Stiefel Manifold

Manifold of (d x k) orthonormal matrices

$$\text{St}_{(d,k)} = \{C \in \mathbb{R}^{d \times k} \,|\, C^\top C = I_k\}$$

**Krasulina's update:** uses <u>projected</u> gradient

$$\widetilde{C} = C - \eta\,\overline{\nabla}\,\hat{\ell}_{\text{comp}}(C) = C - \eta\,(C\,x_t - y_t)\,x_t^\top$$

where $x_t = C^\top y_t$ and $C^{\text{new}} = \text{QR}(\widetilde{C})$

**Oja's update:** uses <u>unprojected</u> gradient

$$\widetilde{C} = C - \eta\nabla\hat{\ell}_{\text{var}}(C) = C + \eta\,y_t y_t^\top C$$

and $C^{\text{new}} = \text{QR}(\widetilde{C})$
**Both updates are explicit!**

## Implicit Krasulina Update

Without the orthonormality constraint

$$C^{\text{new}} = \underset{\widetilde{C}}{\arg\min} \; 1/2 \left( 1/\eta \, \left\|\widetilde{C} - C\right\|_F^2 + \mathbb{E}\left[\|\widetilde{C}\,\widetilde{C}^\dagger y - y\|^2\right] \right)$$

Approximating $\widetilde{C}^\dagger y$ with $C^\dagger y$ yields the (partially)
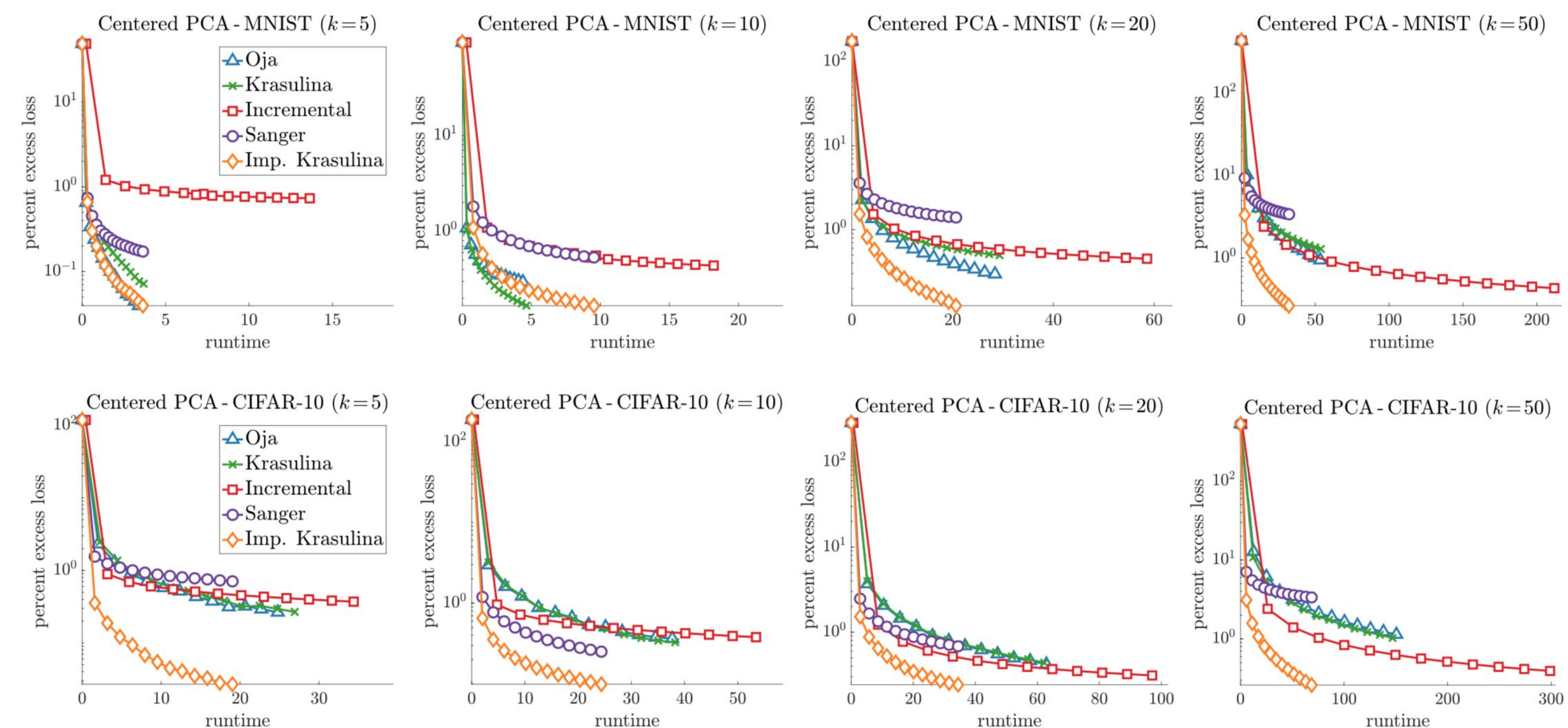**Implicit Krasulina (a.k.a. Sanger) update:**

$$x_t = C^\dagger y_t \qquad\qquad\qquad \text{(E-Step)}$$

$$C^{\text{new}} = C - \frac{\eta}{1 + \eta\,\|x_t\|^2}\,(Cx_t - y_t)\,x_t^\top \qquad \text{(M-Step)}$$

- No QR-step, needs to keep track of $C^\dagger$ instead
- Has an **online EM** interpretation! (see [1])

## Results

### Online k-PCA



### Sensitivity to Learning Rate

| Method | $\eta_0$-Scale | MNIST | | | CIFAR10 | | |
|---|---|---|---|---|---|---|---|
| | | $k = 5$ | $k = 10$ | $k = 20$ | $k = 5$ | $k = 10$ | $k = 20$ |
| Batch PCA | – | 35.16 | 26.95 | 18.74 | 86.98 | 65.70 | 48.65 |
| Oja | $0.1\times$ | $35.79 \pm 0.38$ | $29.78 \pm 0.38$ | $21.55 \pm 0.25$ | $116.26 \pm 2.58$ | $66.81 \pm 0.58$ | $58.28 \pm 0.65$ |
| | $1\times$ | $35.19 \pm 0.05$ | $26.98 \pm 0.02$ | $18.80 \pm 0.03$ | $87.28 \pm 0.43$ | $65.90 \pm 0.01$ | $48.86 \pm 0.05$ |
| | $10\times$ | $35.29 \pm 0.01$ | $27.08 \pm 0.01$ | $19.01 \pm 0.01$ | $87.22 \pm 0.04$ | $67.68 \pm 0.13$ | $50.36 \pm 0.13$ |
| Krasulina | $0.1\times$ | $35.78 \pm 0.37$ | $29.92 \pm 0.37$ | $21.44 \pm 0.24$ | $116.86 \pm 2.58$ | $66.76 \pm 0.65$ | $57.62 \pm 0.65$ |
| | $1\times$ | $35.17 \pm 0.00$ | $26.96 \pm 0.01$ | $18.79 \pm 0.02$ | $87.04 \pm 0.52$ | $65.90 \pm 0.01$ | $48.87 \pm 0.05$ |
| | $10\times$ | $35.30 \pm 0.01$ | $27.09 \pm 0.01$ | $19.02 \pm 0.01$ | $87.17 \pm 0.04$ | $67.87 \pm 0.16$ | $50.56 \pm 0.15$ |
| Incremental | – | $35.26 \pm 0.10$ | $27.01 \pm 0.04$ | $18.83 \pm 0.05$ | $87.50 \pm 0.46$ | $65.75 \pm 0.05$ | $48.82 \pm 0.10$ |
| Sanger | $0.1\times$ | $39.09 \pm 0.78$ | $30.61 \pm 0.48$ | $22.03 \pm 0.23$ | $90.74 \pm 1.84$ | $70.70 \pm 2.18$ | $54.04 \pm 0.95$ |
| | $1\times$ | $36.27 \pm 0.93$ | $28.00 \pm 0.38$ | $19.82 \pm 0.18$ | $88.21 \pm 0.81$ | $66.84 \pm 0.88$ | $50.73 \pm 0.62$ |
| | $10\times$ | $42.36 \pm 1.62$ | $34.48 \pm 0.95$ | $25.50 \pm 0.46$ | $99.78 \pm 3.69$ | $77.88 \pm 2.62$ | $60.21 \pm 1.08$ |
| Imp. Krasulina | $0.1\times$ | $35.17 \pm 0.01$ | $26.96 \pm 0.01$ | $18.78 \pm 0.02$ | $87.00 \pm 0.00$ | $65.75 \pm 0.05$ | $48.75 \pm 0.03$ |
| | $1\times$ | $35.17 \pm 0.01$ | $26.97 \pm 0.03$ | $18.77 \pm 0.01$ | $87.01 \pm 0.00$ | $65.78 \pm 0.09$ | $48.74 \pm 0.04$ |
| | $10\times$ | $35.17 \pm 0.01$ | $26.98 \pm 0.05$ | $18.77 \pm 0.01$ | $87.02 \pm 0.04$ | $65.74 \pm 0.01$ | $48.76 \pm 0.03$ |

### Distributed Setting

**Reference:** [1] E. Amid and M. K. Warmuth, Divergence-based Motivation for Online EM and Combining Hidden Variable Models, *arXiv preprint arXiv:1902.04107, 2019.*