# A Comparison of Stochastic Methods for PCA as Applied to Streaming Facial Recognition[*]

**Corbin Rosset**

Johns Hopkins University

3400 N. Charles St.

Baltimore, MD 21218, USA

`crosset2@jhu.edu`

**Edmund Duhaime**

Johns Hopkins University

3400 N. Charles St.

Baltimore, MD 21218, USA

`eduhaim1@jhu.edu`

## Abstract

This study applies a variety of stochastic and incremental techniques for principle component analysis (PCA) to quickly learn maximally informative linear subspaces on a streaming version of the Yale Face Data Set B. We compare the performance of a K-nearest neighbor classifier and a bagged tree learner on the best subspaces learned by each of: Stochastic Power Method (SPM), Incremental PCA (IPCA), Matrix Stochastic Descent (MSG), and Sparse PCA (SPCA). We compare and contrast the theoretical properties such as correctness, space, and iteration complexity. In all cases, the memory required to store a subspace capable of achieving accuracy similar to that of the baseline classifier was 2-4% of the size of the input data. These algorithms play an important role in improving the scalability of facial recognition in a streaming setting.

## 1 Introduction

The premise of subspace learning is to map a data matrix $X \in \mathbb{R}^{d \times n}$ of $n$ examples each of dimension $d$ to a $k$ dimensional subspace, $k << d$ that preserves maximal "information". The motivation is that most data sets capture redundant, verbose, or noisy features that mask the underlying structure of the data. It is often the case that natural processes are largely controlled by finitely many simpler mechanisms that operate in lower dimensions. For example, latent variables captured by images such as

shadows, reflections, perspective can drastically alter pixel values even though such phenomenon are easily parameterized. PCA seeks to find low dimensional linear representations of these latent variables such that, when transformed back into the original space, the loss of information is minimal. In the next section, we will pose PCA as an optimization problem with multiple equivalent interpretations, and then derive its empirical solution. We will then describe and derive solutions to the state of the art stochastic and incremental approximation algorithms to the PCA objective. We applied each of the algorithms on the Yale Face Data Set B in a streaming fashion and trained a K-NN neighbor classifier on the respective learned subspaces.

## 2 PCA and its Stochastic Variants

Given $n$ data points in $\mathbb{R}^d$, find an orthogonal projection matrix $U \in \mathbb{R}^{d \times k}$ such that the projection $\hat{x} \in \mathbb{R}^k$ of a data vector $x \in \mathbb{R}^d$ given by $\hat{x} = UU^\top x$ minimizes the empirical reconstruction error:

$$\begin{aligned} \text{Error} &= \frac{1}{n} \sum_{i=1}^{n} \|x_i - UU^\top x_i\|_2^2 \\ &= \frac{1}{n} \sum_{i=1}^{n} \left( \|x\|_2^2 - \|U^\top x_i\|_2^2 \right) \end{aligned} \quad (1)$$

In Equation 1, the term $\frac{1}{n} \sum_{i=1}^{n} \|x_i\|_2^2$ is constant given any data set $x$. So in order to minimize the reconstruction, or projection, error, we must maximize $\frac{1}{n} \sum_{i=1}^{n} \|U^\top x_i\|_2^2$. This is equivalent to maximizing the trace of $U^\top \left[ \frac{1}{n} \sum_{i=1}^{n} x_i x_i^\top \right] U$ where $C_{xx} =$

$\frac{1}{n}\sum_{i=1}^{n} x_i x_i^\top$ is the empirical covariance matrix.

Since the $i$'th diagonal entry of $C_{xx}$ is the variance [1] of $x_i$, $var(x_i) = \mathbb{E}[x_i x_i^\top]$ the trace of $C_{xx}$ is the variance of the whole data matrix. Hence it is natural to interpret $\mathbb{E}[\|U^\top x_i\|_2^2] = \mathbb{E}[U^\top x_i x_i^\top U]$ as the variance of $x_i$ captured by $U$. Thus, PCA serves to minimize reconstruction error, or maximize the variance of the data captured by the subspace $U$. Derived in the appendix, it turns out that $U$ is comprised of the top-$k$ eigenvectors of $C_{xx}$ sorted in decreasing order of eigenvalue: $U = V_{:,1:k}$ for $C_{xx} = VSV^\top$.

This solution holds for any distribution $\mathcal{D}$ of the data as long as each of the $n$ samples are drawn i.i.d from it. For a correct empirical covariance matrix, a simple $O(d^2 mk)$ time algorithm known as the power iteration method calculates $U \in \mathbb{R}^{d\times k}$ using $O(m)$ iterations per component [2].

The most pressing challenge is the temporal[3] and spatial dependence of these solutions on powers of $d$, which can easily be intractable for $d \geq 10^5$, which is common for images. Secondly, a batch algorithm may not satisfy modern computational desires, for many modern applications require responses to realtime streaming data arriving at high frequency.

Subspace learning (linear or nonlinear) is ubiquitous in data-driven applications such as compression, de-noising, visualization and matrix completion.

## 2.1 Stochastic Power Method

The optimization problem describing PCA,

$$\max_{U\in\mathbb{R}^{d\times k}} \quad \mathbb{E}_{\mathcal{D}}\left[\text{trace}\left(U^\top xx^\top U\right)\right]$$
$$\text{subject to}\quad U^\top U = I \tag{2}$$

is in fact nonconvex due to the constraint and the objective function being a maximization. A convex relaxation is given by the constraint that $U^\top U \preceq I$ meaning all eigenvalues of $C_{xx}$ are at most one, and

optimally equally to one. Assuming $\mathcal{D}$ is unknown but unchanging, the goal is to update $U$ directly as samples $x_i$ arrive sequentially and independently. Since we have access to neither $\mathcal{D}$ nor the population of samples ahead of time, neither $C_{xx}$ nor its empirical estimate can be computed. However, the i.i.d assumption admits $\mathbb{E}[x_t x_t^\top] = C_{xx}$.

In this setting, gradient descent is a viable algorithm, with updates of the form

$$U^{(t+1)} = \mathcal{P}_{orth}\left(U^{(t)} + \eta_t x_t x_t^\top U^{(t)}\right) \tag{3}$$

The projection of the updated $U$ onto the set of orthonormal matrices, $\mathcal{P}_{orth}$ can be accomplished in $O(k^2 d)$ time using Gram-Schmidt or QR factorization. Notice that instead of $O(d^2)$ memory, the updates require only $O(kd)$. The time complexity is $O(Tkd)$ for $T$ iterations. While this algorithm converges with probability 1, the rate of convergence is unknown. Obviously, if the true covariance $C_{xx}$ were known, then the objective function would become trace $\left(U^\top C_{xx} U\right)$ with derivative $2C_{xx}U$, which would replace the gradient term in the update above. The

## 2.2 Incremental PCA

notes from class

## 2.3 Matrix Stochastic Gradient

notes from class

## 2.4 Sparse PCA

notes from class

## 3 K-NN for High-Dimensional Data

discuss johnson lindenstrauss lemma popularly applied to proximity problems to preserve pairwise distances in a lower dimensional space

## 4 Experiments

**NED: describe the contents of the data, how we removed bad images, just like proposal...**

## 5 Results

**NED: pictures of 1) reconstructed faces with 5, 10, 15... principle components 2) table of the best**

---

[1] It is assumed $X$ is centered $\mathbb{E}[x_i] = 0$, eliminating the $(\mathbb{E}[x_i])^2$ term in the variance.

[2] the iteration complexity depends on the eigengap between consecutive components; a larger difference in eigenvalues implies faster convergence

[3] the fastest known algorithm for brute force eigendecomposition of $C_{xx}$ is $O(d^3 + d^2 log^2 d)$

| Type of Text | Font Size | Style |
|---|---|---|
| paper title | 15 pt | bold |
| author names | 12 pt | bold |
| author affiliation | 12 pt | |
| the word "Abstract" | 12 pt | bold |
| section titles | 12 pt | bold |
| document text | 11 pt | |
| abstract text | 10 pt | |
| captions | 10 pt | |
| bibliography | 10 pt | |
| footnotes | 9 pt | |

Table 1: Font guide.

**dimension and associated accuracy achieved by each algorithm 3) some examples of "eigenfaces" (the top five principle components shown as images) as well as their thresholded versions. You'll find these in the figures directory. 4) graphs of accuracy versus number of principle components for some of the algorithms (choose the better looking ones)**

**Captions**: Provide a caption for every illustration; number each one sequentially in the form: "Figure 1. Caption of the Figure." "Table 1. Caption of the Table." Type the captions of the figures and tables below the body, using 10 point text.

## 6    Conclusion

all algorithms for pca should learn the same subspace, up to small rotations and scalings...

discuss applications of streaming subspace learning

there are extensions for incremental kernal pca and dealing with the case of missing data...

## Acknowledgments

Do not number the acknowledgment section.

## References

the following are examples

Alfred V. Aho and Jeffrey D. Ullman. 1972. *The Theory of Parsing, Translation and Compiling*, volume 1. Prentice-Hall, Englewood Cliffs, NJ.

American Psychological Association. 1983. *Publications Manual*. American Psychological Association, Washington, DC.

Association for Computing Machinery. 1983. *Computing Reviews*, 24(11):503–512.

Ashok K. Chandra, Dexter C. Kozen, and Larry J. Stockmeyer. 1981. Alternation. *Journal of the Association for Computing Machinery*, 28(1):114–133.

Dan Gusfield. 1997. *Algorithms on Strings, Trees and Sequences*. Cambridge University Press, Cambridge, UK.