

# Finance Data Science

## Lecture 4: PCA and and Factor Models

Laurent El Ghaoui

MFE 230P, Summer 2017  
MFE Program  
Haas School of Business  
UC Berkeley

6/14/2017

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank  
approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA

## Motivation

### Linear Algebra Recap

- Eigenvalues
- Singular values

### Principal Component Analysis

- Overview
- Deflation: iterated variance maximization
- Example

### Low-rank approximations

- Problem
- Link with PCA
- Explained variance and approximation error
- Factor models

### Extensions

- Robust PCA
- Sparse PCA

## Motivation

### Linear Algebra Recap

- Eigenvalues
- Singular values

### PCA

- Overview
- Deflation
- Example

### Low-rank approximations

- Problem
- Link with PCA
- Explained variance
- Factor models

### Extensions

- Robust PCA
- Sparse PCA

## Motivation

### Linear Algebra Recap

- Eigenvalues
- Singular values

### Principal Component Analysis

- Overview
- Deflation: iterated variance maximization
- Example

### Low-rank approximations

- Problem
- Link with PCA
- Explained variance and approximation error
- Factor models

### Extensions

- Robust PCA
- Sparse PCA

## Motivation

### Linear Algebra Recap

- Eigenvalues
- Singular values

### PCA

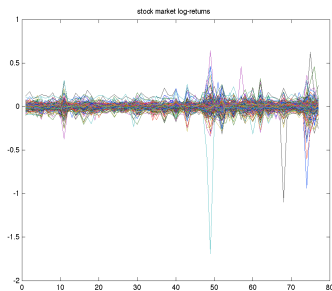
- Overview
- Deflation
- Example

### Low-rank approximations

- Problem
- Link with PCA
- Explained variance
- Factor models

### Extensions

- Robust PCA
- Sparse PCA



Daily log-returns of 77 Fortune 500 companies, 1/2/2007—12/31/2008.

High-dimensional data does not make any sense! (Other than tell us: returns are approximately zero ...)

*In this lecture:*

- ▶ start with a classical unsupervised learning to obtain insights
- ▶ examine a newer method that improves interpretability

## Motivation

### Linear Algebra Recap

Eigenvalues  
Singular values

### PCA

Overview  
Deflation  
Example

### Low-rank approximations

Problem  
Link with PCA  
Explained variance  
Factor models

### Extensions

Robust PCA  
Sparse PCA

## Motivation

## Linear Algebra Recap

- Eigenvalues
- Singular values

## Principal Component Analysis

- Overview
- Deflation: iterated variance maximization
- Example

## Low-rank approximations

- Problem
- Link with PCA
- Explained variance and approximation error
- Factor models

## Extensions

- Robust PCA
- Sparse PCA

## Motivation

## Linear Algebra Recap

- Eigenvalues
- Singular values

## PCA

- Overview
- Deflation
- Example

## Low-rank approximations

- Problem
- Link with PCA
- Explained variance
- Factor models

## Extensions

- Robust PCA
- Sparse PCA

## Theorem (EVD of symmetric matrices)

We can decompose any symmetric  $p \times p$  matrix  $S$  as

$$S = U \Lambda U^T = \sum_{i=1}^p \lambda_i u_i u_i^T,$$

where  $\Lambda = \mathbf{diag}(\lambda_1, \dots, \lambda_p)$ , with  $\lambda_1 \geq \dots \geq \lambda_p$  the eigenvalues, and  $U = [u_1, \dots, u_p]$  is a  $p \times p$  orthogonal matrix ( $U^T U = I_p$ ) that contains the eigenvectors  $u_i$  of  $S$ , that is:

$$S u_i = \lambda_i u_i, \quad i = 1, \dots, p.$$

*Corollary:* If  $S$  is square, symmetric:

$$\lambda_{\max}(S) = \max_{x : \|x\|_2=1} x^T S x. \quad (1)$$

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank  
approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA

# Positive semi-definite (PSD) matrices

A (square) symmetric matrix  $S$  is said to be *positive semi-definite* (PSD) if

$$\forall x, \quad x^T S x \geq 0.$$

In this case, we write  $S \succeq 0$ .

*From EVD theorem:* for any square, symmetric matrix  $S$ :

$$S \succeq 0 \iff \text{every eigenvalue of } S \text{ is non-negative.}$$

Hence we can numerically (via EVD) check positive semi-definiteness.

## Motivation

### Linear Algebra Recap

Eigenvalues

Singular values

### PCA

Overview

Deflation

Example

### Low-rank approximations

Problem

Link with PCA

Explained variance

Factor models

### Extensions

Robust PCA

Sparse PCA

## Theorem (SVD of general matrices)

We can decompose any non-zero  $p \times m$  matrix  $A$  as

$$A = \sum_{i=1}^r \sigma_i u_i v_i^T = U \Sigma V^T, \quad \Sigma = \mathbf{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0) \in \mathbf{R}^{p \times m}$$

where  $\sigma_1 \geq \dots \geq \sigma_r > 0$  are the singular values, and

$$U = [u_1, \dots, u_m], \quad V = [v_1, \dots, v_p]$$

are square, orthogonal matrices ( $U^T U = I_p$ ,  $V^T V = I_m$ ). The number  $r \leq \min(p, m)$  (the number of non-zero singular values) is called the **rank** of  $A$ . The first  $r$  columns of  $U$ ,  $V$  contains the left- and right singular vectors of  $A$ , respectively, that is:

$$A v_i = \sigma_i u_i, \quad A^T u_i = \sigma_i v_i, \quad i = 1, \dots, r.$$



The SVD of a  $p \times m$  matrix  $A$  is related to the EVD of a (PSD) matrix related to  $A$ .

If  $A = U\Sigma V^T$  is the SVD of  $A$ , then

- ▶ The EVD of  $AA^T$  is  $U\Lambda U^T$ , with  $\Lambda = \Sigma^2$ .
- ▶ The EVD of  $A^T A$  is  $V\Lambda V^T$ .

Hence the left (resp. right) singular vectors of  $A$  are the eigenvectors of the PSD matrix  $AA^T$  (resp.  $A^T A$ ).

[Motivation](#)[Linear Algebra Recap](#)[Eigenvalues](#)[Singular values](#)[PCA](#)[Overview](#)[Deflation](#)[Example](#)[Low-rank approximations](#)[Problem](#)[Link with PCA](#)[Explained variance](#)[Factor models](#)[Extensions](#)[Robust PCA](#)[Sparse PCA](#)

# Computing SVD

## Power iteration algorithm

For a large, sparse matrix  $M$ , we can find left and right singular vectors corresponding to the largest singular value of  $M$  with the *power iteration* algorithm:

$$u \rightarrow \frac{Mv}{\|Mv\|_2}, \quad v \rightarrow \frac{M^T u}{\|M^T u\|_2}.$$

This converges (for arbitrary initial  $u, v$ ) under mild conditions on  $M$ .

Similar efficient algorithm when  $M$  is centered (thus, not necessarily sparse, even if data is).

Google's page rank is based on this kind of algorithm ...

### Motivation

### Linear Algebra Recap

Eigenvalues

Singular values

### PCA

Overview

Deflation

Example

### Low-rank approximations

Problem

Link with PCA

Explained variance

Factor models

### Extensions

Robust PCA

Sparse PCA

## Motivation

## Linear Algebra Recap

- Eigenvalues
- Singular values

## Principal Component Analysis

- Overview
- Deflation: iterated variance maximization
- Example

## Low-rank approximations

- Problem
- Link with PCA
- Explained variance and approximation error
- Factor models

## Extensions

- Robust PCA
- Sparse PCA

## Motivation

## Linear Algebra Recap

- Eigenvalues
- Singular values

## PCA

- Overview
- Deflation
- Example

## Low-rank approximations

- Problem
- Link with PCA
- Explained variance
- Factor models

## Extensions

- Robust PCA
- Sparse PCA

# Principal Component Analysis

## Overview

Principal Component Analysis (PCA) originated in psychometrics in the 1930's. It is now widely used in

- ▶ Exploratory data analysis.
- ▶ Simulation.
- ▶ Visualization.

Application fields include

- ▶ Finance, marketing, economics.
- ▶ Biology, medicine.
- ▶ Engineering design, signal compression and image processing.
- ▶ Search engines, data mining.

### Motivation

### Linear Algebra Recap

Eigenvalues

Singular values

### PCA

Overview

Deflation

Example

### Low-rank approximations

Problem

Link with PCA

Explained variance

Factor models

### Extensions

Robust PCA

Sparse PCA

PCA finds “principal components” (PCs), *i.e.* **orthogonal** directions of *maximal variance*.

- ▶ PCs are computed via EVD of covariance matrix.
- ▶ Alternatively, PCs can be found directly via SVD of (centered) data matrix.
- ▶ Can be interpreted as a “factor model” of original data matrix.

## *Applications in finance:*

- ▶ General understanding of market data.
- ▶ Underlies many theoretical models (such as CAPM).
- ▶ Modeling of term structure of interest rates.
- ▶ Portfolio hedging and immunization.
- ▶ Risk analysis, scenario generation.
- ▶ Obtain speed-ups in some portfolio optimization problems.

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank  
approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA

# Variance maximization problem

Let  $S$  be the (empirical) covariance matrix. *Variance maximization problem:*

$$\max_x x^T S x : \|x\|_2 = 1.$$

Assume the EVD of  $S$  is given:

$$S = \sum_{i=1}^p \lambda_i u_i u_i^T,$$

with  $\lambda_1 \geq \dots \geq \lambda_p$ , and  $U = [u_1, \dots, u_p]$  is orthogonal ( $U^T U = I$ ). Then a solution to

$$\max_{x : \|x\|_2=1} x^T S x$$

is  $x^* = u_1$ , with  $u_1$  an eigenvector of  $S$  that corresponds to its largest eigenvalue  $\lambda_1$ .

Alternatively,  $u_1$  can be found directly via SVD of (centered) data matrix (see later).

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank  
approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA

# Finding orthogonal directions

## A deflation method

Once we've found a direction with high variance, can we repeat the process and find other ones?

### *Deflation method:*

- ▶ Project data points on the subspace orthogonal to the direction we found.
- ▶ Find a direction of maximal variance for projected data.

The process stops after  $p$  steps ( $p$  is the dimension of the whole space), but can be stopped earlier (to find only  $k$  directions, with  $k \ll p$ ).

### Motivation

### Linear Algebra Recap

Eigenvalues  
Singular values

### PCA

Overview  
**Deflation**  
Example

### Low-rank approximations

Problem  
Link with PCA  
Explained variance  
Factor models

### Extensions

Robust PCA  
Sparse PCA

# Finding orthogonal directions

## Result

It turns out that the direction that solves

$$\max_x \mathbf{var}(x) : x^T u_1 = 0$$

is  $u_2$ , an eigenvector corresponding to the second-to-largest eigenvalue.

After  $k$  steps of the deflation process, the directions returned are  $u_1, \dots, u_k$ . Thus we can compute  $k$  directions of largest variance in *one* eigenvalue decomposition of the covariance matrix.

### Motivation

### Linear Algebra Recap

Eigenvalues

Singular values

### PCA

Overview

**Deflation**

Example

### Low-rank approximations

Problem

Link with PCA

Explained variance

Factor models

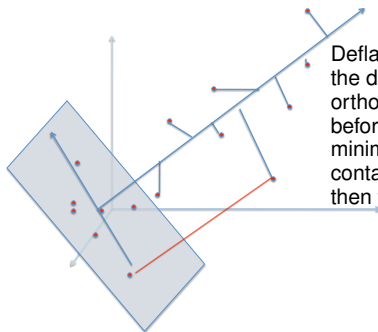
### Extensions

Robust PCA

Sparse PCA



# Geometry of deflation



Deflation consists in projecting the data on a hyperplane orthogonal to the line found before; a new minimum-distance line contained in the hyperplane is then found.

## Motivation

### Linear Algebra Recap

Eigenvalues  
Singular values

### PCA

Overview  
**Deflation**  
Example

### Low-rank approximations

Problem  
Link with PCA  
Explained variance  
Factor models

### Extensions

Robust PCA  
Sparse PCA

# Measuring quality

How well is data approximated by its projections on the successive subspaces?

*Approach:* compare sum of variances contained in the  $k$  directions found, with total variance.

*Explained variance:* measured by the ratio

$$\frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_p} = \frac{\sigma_1^2 + \dots + \sigma_k^2}{\sigma_1^2 + \dots + \sigma_p^2},$$

where  $\lambda_1 \geq \dots \geq \lambda_p$  are the eigenvalues of the covariance matrix, and  $\sigma_1 \geq \dots \geq \sigma_p$  are the singular values of the (centered) data matrix.

## Motivation

### Linear Algebra Recap

Eigenvalues

Singular values

## PCA

Overview

Deflation

Example

## Low-rank approximations

Problem

Link with PCA

Explained variance

Factor models

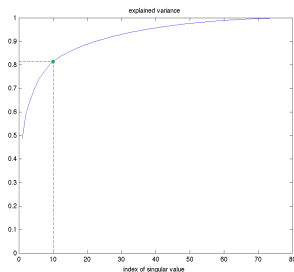
## Extensions

Robust PCA

Sparse PCA

# Example

## PCA of market data



Data: Daily log-returns of 77 Fortune 500 companies,  
1/2/2007—12/31/2008.

- ▶ Plot shows the eigenvalues of covariance matrix in decreasing order.
- ▶ First ten components explain 80% of the variance.
- ▶ Largest magnitude of eigenvector for 1st component correspond to financial sector (FABC, FTU, MER, AIG, MS).

### Motivation

### Linear Algebra Recap

Eigenvalues

Singular values

### PCA

Overview

Deflation

Example

### Low-rank approximations

Problem

Link with PCA

Explained variance

Factor models

### Extensions

Robust PCA

Sparse PCA

## Motivation

## Linear Algebra Recap

- Eigenvalues
- Singular values

## Principal Component Analysis

- Overview
- Deflation: iterated variance maximization
- Example

## Low-rank approximations

- Problem
- Link with PCA
- Explained variance and approximation error
- Factor models

## Extensions

- Robust PCA
- Sparse PCA

## Motivation

## Linear Algebra Recap

- Eigenvalues
- Singular values

## PCA

- Overview
- Deflation
- Example

## Low-rank approximations

- Problem
- Link with PCA
- Explained variance
- Factor models

## Extensions

- Robust PCA
- Sparse PCA

# Low-rank approximation of a matrix

For a given  $p \times m$  matrix  $A$ , and integer  $k \leq m, p$ , the  *$k$ -rank approximation* problem is

$$A^{(k)} := \arg \min_X \|X - A\|_F : \mathbf{Rank}(X) \leq k,$$

where  $\|\cdot\|_F$  is the Frobenius norm (Euclidean norm of the vector formed with all the entries of the matrix). The solution is

$$A^{(k)} = \sum_{i=1}^k \sigma_i u_i v_i^T,$$

where

$$A = U \Sigma V^T = \sum_{i=1}^r \sigma_i u_i v_i^T$$

is an SVD of the matrix  $A$ .

## Motivation

### Linear Algebra Recap

- Eigenvalues
- Singular values

### PCA

- Overview
- Deflation
- Example

### Low-rank approximations

- Problem**
- Link with PCA
- Explained variance
- Factor models

### Extensions

- Robust PCA
- Sparse PCA

# Low-rank approximation

Interpretation: rank-one case

Assume data matrix  $A \in \mathbf{R}^{p \times m}$  represents time-series data (each row is a time-series). Assume also that  $A$  is rank-one, that is,  $A = uv^T \in \mathbf{R}^{p \times m}$ , where  $u, v$  are vectors. Then

$$A = \begin{pmatrix} a_1^T \\ \vdots \\ a_m^T \end{pmatrix}, \quad a_j(t) = \sigma_1 u(j)v(t), \quad 1 \leq j \leq p, \quad 1 \leq t \leq m.$$

Thus, each time-series is a “scaled” copy of the time-series represented by  $v$ , with scaling factors given in  $u$ . We can think of  $v$  as a “factor” that drives all the time-series.

**Geometry:** if a data matrix is rank-one, then all the data points are on a single line.

## Motivation

### Linear Algebra Recap

Eigenvalues  
Singular values

### PCA

Overview  
Deflation  
Example

### Low-rank approximations

#### Problem

Link with PCA  
Explained variance  
Factor models

### Extensions

Robust PCA  
Sparse PCA

# Low-rank approximation

Interpretation: low-rank case

When  $A$  is rank  $k$ , that is,

$$A = USV^T, \quad U \in \mathbf{R}^{p \times k}, \quad S = \text{diag}(\sigma_1, \dots, \sigma_k) \in \mathbf{R}^{k \times k}, \quad V \in \mathbf{R}^{m \times k},$$

we can express the  $j$ -th row of  $A$  as

$$a_j(t) = \sum_{i=1}^k \sigma_i u_i(j) v_i(t), \quad 1 \leq j \leq p, \quad 1 \leq t \leq m.$$

Thus, each time-series is the **sum** of scaled copies of  $k$  time-series represented by  $v_1, \dots, v_k$ , with scaling factors given in  $u_1, \dots, u_k$ .

We can think of  $v_i$ 's as the few “factors” that drive all the time-series.

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank  
approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA

PCA can be obtained directly (without forming the covariance matrix) via a *low-rank approximation* to the centered data matrix  $A_c$ :

$$A_c \approx \hat{A}_c^{(k)} := \sum_{i=1}^k \sigma_i u_i v_i^T$$

Each  $v_i$  is a particular factor, and  $u_i$ 's contain scalings.

That is,  $\hat{A}_c^{(k)}$  solves the problem

$$\arg \min_X \|X - A_c\|_F : \mathbf{Rank}(X) \leq k,$$

where  $\|X\|_F^2$  is the sum of the squares of the entries of  $X$ .

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank  
approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA



## Link with power iteration algorithm

For a large, sparse matrix  $M$ , we can find left and right singular vectors corresponding to the largest singular value of  $M$  with the *power iteration* algorithm:

$$u \rightarrow \frac{Mv}{\|Mv\|_2}, \quad v \rightarrow \frac{M^T u}{\|M^T u\|_2}.$$

This converges (for arbitrary initial  $u, v$ ) under mild conditions on  $M$ .

*Interpretation:* power iteration can be obtained by solving

$$\min_{p, q} \|M - pq^T\|_F$$

alternatively over  $p, q$  (in the P.I. algorithm above,  $u, v$  are just normalized versions of  $p, q$ ).

[Motivation](#)[Linear Algebra Recap](#)[Eigenvalues](#)[Singular values](#)[PCA](#)[Overview](#)[Deflation](#)[Example](#)[Low-rank  
approximations](#)[Problem](#)[Link with PCA](#)[Explained variance](#)[Factor models](#)[Extensions](#)[Robust PCA](#)[Sparse PCA](#)

# Low-rank approximation of covariance matrix

PCA also (implicitly) forms a low-rank approximation of the empirical covariance matrix.

The corresponding approximate covariance matrix is rank  $k$ :

$$S \approx S^{(k)} := U\Lambda^{(k)}U^T = FF^T$$

where  $\Lambda^{(k)} = \mathbf{diag}(\lambda_1, \dots, \lambda_k, 0, \dots, 0)$ , with  $\lambda_i = \sigma_i^2$ , and

$$F = U^{(k)} \mathbf{diag}(\sigma_1, \dots, \sigma_k),$$

where  $U^{(k)}$  contains the first  $k$  columns of  $U$ .

## Motivation

### Linear Algebra Recap

Eigenvalues  
Singular values

### PCA

Overview  
Deflation  
Example

### Low-rank approximations

Problem  
Link with PCA  
Explained variance  
Factor models

### Extensions

Robust PCA  
Sparse PCA

# Explained variance and approximation error

Recall the *explained variance* ratio

$$\frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_p} = \frac{\sigma_1^2 + \dots + \sigma_k^2}{\sigma_1^2 + \dots + \sigma_p^2},$$

The (square-root) of the explained variance ratio is the relative approximation error:

$$\frac{\|\hat{A}_c^{(k)} - A_c\|_F}{\|A_c\|_F} = \frac{\sigma_1^2 + \dots + \sigma_k^2}{\sigma_1^2 + \dots + \sigma_p^2}.$$

We can also express the above as

$$\frac{\text{Tr}(S - S^{(k)})}{\text{Tr } S} = \frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_p}.$$

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank  
approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA

# Stochastic interpretation of low-rank approximations

Assume that the (e.g., price) observations  $y$  are the generated by stochastic model (here  $F \in \mathbf{R}^{n \times k}$ )

$$y = F\xi$$

where  $\xi \in \mathbf{R}^k$  are independent random variables ( $k \leq p$ ), with zero mean and identity covariance matrix. Here,  $F \in \mathbf{R}^{p \times k}$  is the *loading* matrix.

Then the covariance matrix of  $y$  is

$$S = \mathbf{E}(F\xi\xi^T F^T) = FF^T.$$

In effect we are postulating that a few factors drive the market observations. This is the *same* as the PCA seen before!

## Motivation

### Linear Algebra Recap

Eigenvalues  
Singular values

### PCA

Overview  
Deflation  
Example

### Low-rank approximations

Problem  
Link with PCA  
Explained variance

### Factor models

### Extensions

Robust PCA  
Sparse PCA

More generally we can assume that market observations are of the form

$$y = F\xi + \sigma e,$$

with  $(\xi, e)$  a zero-mean, random variable with identity covariance matrix, and  $\sigma > 0$  is a parameter.

- ▶  $\xi$  contains the market factors;
- ▶  $e$  is a noise term that affect each observation independently (“idiosyncratic noise”).

The covariance matrix of  $y$  is of the form

$$S = FF^T + D^2$$

with  $D^2 = \sigma^2 I$ . We can fit this model (*i.e.*, find  $F, \sigma$ ) via SVD.

More general factor models allow for idiosyncratic noises with different variances (see lecture 4). However SVD cannot be used directly and the fitting problem is more challenging.

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank  
approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA

## Motivation

## Linear Algebra Recap

- Eigenvalues
- Singular values

## Principal Component Analysis

- Overview
- Deflation: iterated variance maximization
- Example

## Low-rank approximations

- Problem
- Link with PCA
- Explained variance and approximation error
- Factor models

## Extensions

- Robust PCA
- Sparse PCA

## Motivation

## Linear Algebra Recap

- Eigenvalues
- Singular values

## PCA

- Overview
- Deflation
- Example

## Low-rank approximations

- Problem
- Link with PCA
- Explained variance
- Factor models

## Extensions

- Robust PCA
- Sparse PCA

PCA is based on the assumption that the data matrix can be (approximately) written as a low-rank matrix:

$$A = LR^T,$$

with  $L \in \mathbf{R}^{p \times k}$ ,  $R \in \mathbf{R}^{m \times k}$ , with  $k \ll m, p$ .

*Robust PCA* [2] assumes that  $A$  has a “low-rank plus sparse” structure:

$$A = N + LR^T$$

where “noise” matrix  $N$  is sparse (has many zero entries).

How do we discover  $N, L, R$  based on  $A$ ?

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank  
approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA

## Motivation

### Linear Algebra Recap

Eigenvalues  
Singular values

### PCA

Overview  
Deflation  
Example

### Low-rank approximations

Problem  
Link with PCA  
Explained variance  
Factor models

### Extensions

Robust PCA  
Sparse PCA

In robust PCA, we solve the **convex** problem

$$\min_N \|A - N\|_* + \lambda \|N\|_1$$

where  $\|\cdot\|_*$  is the so-called nuclear norm (sum of singular values) of its matrix argument. At optimum,  $A - N$  has usually low-rank.

**Motivation:** the nuclear norm is akin to the  $l_1$ -norm of the vector of singular values, and  $l_1$ -norm minimization encourages sparsity of its argument.



Here is a matlab snippet that solves a robust PCA problem via CVX, given integers  $n, m$ , a  $n \times m$  matrix  $A$  and non-negative scalar  $\lambda$  exist in the workspace:

```
cvx_begin
variable X(n,m);
minimize( norm_nuc(A-X) + lambda*norm(X(:),1))
cvx_end
```

- ▶ Note the use of `norm_nuc`, which stands for the nuclear norm.
- ▶ In practice, this CVX code does not run on large matrices, for memory limitations.
- ▶ Efficient specialized algorithms exist [2].

Alternatively, we can use a power iteration-like algorithm [6], alternating over  $L, R$

$$\min_{L,R} \|X - LR^T\|_1 : L \in \mathbf{R}^{n \times k}, R \in \mathbf{R}^{m \times k},$$

Each step is a convex problem.

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank  
approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA

## Motivation

### Linear Algebra Recap

Eigenvalues  
Singular values

### PCA

Overview  
Deflation  
Example

### Low-rank approximations

Problem  
Link with PCA  
Explained variance  
Factor models

### Extensions

Robust PCA  
Sparse PCA

One of the issues with PCA is that it does not yield principal directions that are easily **interpretable**:

- ▶ The principal directions are really combinations of all the relevant features (say, assets).
- ▶ Hence we cannot interpret them easily.
- ▶ The previous thresholding approach (select features with large components, zero out the others) can lead to much degraded explained variance.

# Sparse PCA

## Problem definition

Modify the variance maximization problem:

$$\max_x x^T S x - \lambda \mathbf{Card}(x) : \|x\|_2 = 1,$$

where penalty parameter  $\lambda \geq 0$  is given, and  $\mathbf{Card}(x)$  is the cardinality (number of non-zero elements) in  $x$ .

The problem is **hard** but can be approximated via convex relaxation.

### Motivation

### Linear Algebra Recap

Eigenvalues  
Singular values

### PCA

Overview  
Deflation  
Example

### Low-rank approximations

Problem  
Link with PCA  
Explained variance  
Factor models

### Extensions

Robust PCA  
Sparse PCA

Express  $S$  as  $S = R^T R$ , with  $R = [r_1, \dots, r_p]$  (each  $r_i$  corresponds to one feature, e.g. asset).

## Theorem (Safe feature elimination)

We have

$$\begin{aligned} \max_{x : \|x\|_2=1} x^T S x - \lambda \mathbf{Card}(x) = \\ \max_{z : \|z\|_2=1} \sum_{i=1}^p \max(0, (r_i^T z)^2 - \lambda). \end{aligned}$$

- ▶ Reduces to ordinary formula when  $\lambda = 0$ .
- ▶ When  $\lambda > 0$  problem is hard, not amenable to SVD methods.

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank  
approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA

## Corollary

If  $\lambda > \|r_i\|_2^2 = S_{ii}$ , we can safely remove the  $i$ -th feature (row/column of  $S$ ).

## Proof.

If  $\|z\|_2 = 1$ , then  $|r_i^T z| \leq \|r_i\|_2$ . ■

- ▶ The presence of the penalty parameter allows to prune out dimensions in the problem.
- ▶ Criterion simply based on variance of each feature (*i.e.*, directional variance along unit vectors).
- ▶ In practice, we want  $\lambda$  high as to allow better interpretability.
- ▶ Hence, interpretability requirement makes the problem easier in some sense!

### Motivation

### Linear Algebra Recap

Eigenvalues  
Singular values

### PCA

Overview  
Deflation  
Example

### Low-rank approximations

Problem  
Link with PCA  
Explained variance  
Factor models

### Extensions

Robust PCA  
Sparse PCA

## Motivation

## Linear Algebra Recap

Eigenvalues  
Singular values

## PCA

Overview  
Deflation  
Example

## Low-rank approximations

Problem  
Link with PCA  
Explained variance  
Factor models

## Extensions

Robust PCA  
**Sparse PCA**

- ▶ The Sparse PCA problem remains challenging due to the huge number of variables.
- ▶ SAFE technique does allow big reduction in problem size.
- ▶ Still area of active research. (Like SVD in the 70's-90's. . .)

# Sparse PCA

## Thresholded power iteration

*Efficient heuristic* to solve, for given  $n \times m$  matrix  $M$ :

$$\min_{p,q} \|M - pq^T\|_F : \mathbf{Card}(p) \leq k, \mathbf{Card}(q) \leq h.$$

Initialize  $p, q$  to be random and

$$p \rightarrow P(T_k(Mq)), \quad q \rightarrow P(T_h(M^T p)),$$

where

- ▶  $M$  is the (centered) data matrix,
- ▶  $P$  is the  $l_2$  normalization operator (for  $z \neq 0$ ,  $P(z) = z/\|z\|_2$ ),
- ▶ operator  $T_k$  removes all but the  $k$  largest-magnitude components of its input.
- ▶ Reduces to a standard method for PCA (power iteration) when  $k = n$ ,  $h = m$ .

For sparse data, can be used in very high dimensions.

### Motivation

### Linear Algebra Recap

Eigenvalues  
Singular values

### PCA

Overview  
Deflation  
Example

### Low-rank approximations

Problem  
Link with PCA  
Explained variance  
Factor models

### Extensions

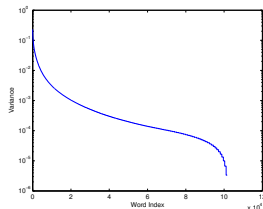
Robust PCA  
Sparse PCA

# Example

## Sparse PCA of New York Times headlines

*Data:* NYTimes text collection contains 300,000 articles and has a dictionary of 102,660 unique words.

The variance of the features (words) decreases very fast:



Sorted variances of 102,660 words in NYTimes data.

With a target number of words less than 10, SAFE allows to reduce the number of features from  $n \approx 100,000$  to  $n = 500$ .



# Example

## Sparse PCA of New York Times headlines

### Motivation

### Linear Algebra Recap

Eigenvalues  
Singular values

### PCA

Overview  
Deflation  
Example

### Low-rank approximations

Problem  
Link with PCA  
Explained variance  
Factor models

### Extensions

Robust PCA  
Sparse PCA

### Words associated with the top 5 sparse principal components in NYTimes

1st PC (6 words)	2nd PC (5 words)	3rd PC (5 words)	4th PC (4 words)	5th PC (4 words)
million	point	official	president	school
percent	play	government	campaign	program
business	team	united_states	bush	children
company	season	u.s	administration	student
market	game	attack		
companies				

**Note:** the algorithm found those terms without any information on the subject headings of the corresponding articles (unsupervised problem).

Thresholded PCA involves simply thresholding the principal components.

$k = 2$	$k = 3$	$k = 9$	$k = 14$
even	even	even	would
like	like	we	new
	states	like	even
		now	we
		this	like
		will	now
		united	this
		states	will
		if	united
			states
			world
			so
			some
			if

1st PC from Thresholded PCA for various cardinality  $k$ . The results contain a lot of non-informative words.

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA



G. Calafiore and L. El Ghaoui.

*Optimization Models.*

Cambridge University Press, 2014.



E.J. Candes, X. Li, Y. Ma, and J. Wright.

Robust principal component analysis.

*Arxiv preprint ArXiv:0912.3599*, 2009.



L. El Ghaoui.

Livebook: Optimization models and applications, 2016.

(Register to the livebook web site).



G.H. Golub and C.F. Van Loan.

*Matrix computations*, volume 3.

Johns Hopkins Univ Pr, 1996.



G. Strang.

*Introduction to linear algebra.*

Wellesley Cambridge Pr, 2003.



Madeleine Udell, Corinne Horn, Reza Zadeh, Stephen Boyd, et al.

Generalized low rank models.

*Foundations and Trends® in Machine Learning*, 9(1):1–118, 2016.

Motivation

Linear Algebra Recap

Eigenvalues

Singular values

PCA

Overview

Deflation

Example

Low-rank  
approximations

Problem

Link with PCA

Explained variance

Factor models

Extensions

Robust PCA

Sparse PCA