Short Course Robust Optimization and Machine Learning

Lecture 4: Optimization in Unsupervised Learning

Laurent El Ghaoui

EECS and IEOR Departments UC Berkeley

Spring seminar TRANSP-OR, Zinal, Jan. 16-19, 2012

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity Penalized

maximum-likelinoo

References

▲□▶▲□▶▲□▶▲□▶ ▲□ ● ●

Outline

Overview of Unsupervised Learning Unsupervised learning models Matrix facts

Principal Component Analysis

Motivations

- Variance maximization
- Deflation
- Factor models
- Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection Sparse graphical models Penalized maximum-likelihood Example

References

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PC/

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE

- Relaxation
- Algorithms

Examples

Variants

Sparse Covariance Selection

Sparsit

Penalized maximum-likelihood

Exampl

References

・ロト・日本・日本・日本・日本・日本

Outline

Overview of Unsupervised Learning Unsupervised learning models Matrix facts

Principal Component Analysis

Motivations

- Variance maximization
- Deflation
- Factor models
- Example

Sparse PCA

- Basics SAFE Relaxatior Algorithms
- Examples
- variants

Sparse Covariance Selection Sparse graphical models Penalized maximum-likelihood Example

References

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics

- Relaxatio
- Algorithm
- Examples
- Variants

Sparse Covariance Selection

- Sparsit
- Penalized maximum-likelihood

Exampl

References

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = 悪 = のへで

What is unsupervised learning?

In unsupervised learning, we are given a matrix of data points $X = [x_1, \ldots, x_m]$, with $x_i \in \mathbf{R}^n$; we wish to learn some condensed information from it.

Examples:

- Find one or several direction of maximal variance.
- Find a low-rank approximation or other structured approximation.
- Find correlations or some other statistical information (*e.g.*, graphical model).
- Find clusters of data points.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

> Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

The empirical covariance matrix Definition

Given $p \times n$ data matrix $A = [a_1, ..., a_m]$ (each row representing say a log-return time-series over *m* time periods), the *empirical covariance matrix* is defined as the $p \times p$ matrix

$$S = rac{1}{m} \sum_{i=1}^m (a_i - \hat{a}) (a_i - \hat{a})^T, \;\; \hat{a} := rac{1}{m} \sum_{i=1}^m a_i.$$

We can express S as

$$S=\frac{1}{m}A_{c}A_{c}^{T},$$

where A_c is the *centered data matrix*, with *p* columns $(a_i - \hat{a})$, i = 1, ..., m.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Dverview

Unsupervised learning

Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics

Delevetie

neiaxatioi

Algorithms

Voriente

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ のの⊙

The empirical covariance matrix

Link with directional variance

The (empirical) variance along direction x is

$$\operatorname{var}(x) = \frac{1}{m} \sum_{i=1}^{m} [x^{T}(a_{i} - \hat{a})]^{2} = x^{T} S x = \frac{1}{m} \|A_{c}x\|_{2}^{2}.$$

where A_c is the centered data matrix.

Hence, covariance matrix gives information about variance along any direction.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Jnsupervised learning

Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxatio Algorithm Examples

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

・ロト・西ト・山田・山田・山下

Eigenvalue decomposition for symmetric matrices

Theorem (EVD of symmetric matrices) We can decompose any symmetric $p \times p$ matrix Q as

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

where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$, with $\lambda_1 \ge \dots \ge \lambda_n$ 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 of Q. That is:

$$Qu_i = \lambda_i u_i, \quad i = 1, \ldots, p.$$

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning

Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Exampl

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

Singular Value Decomposition (SVD)

r

Theorem (SVD of general matrices) We can decompose any non-zero $p \times m$ matrix A as

$$\boldsymbol{A} = \sum_{i=1}^{T} \sigma_i \boldsymbol{u}_i \boldsymbol{v}_i^T = \boldsymbol{U} \boldsymbol{\Sigma} \boldsymbol{V}^T, \ \boldsymbol{\Sigma} = \textbf{diag}(\sigma_1, \dots, \sigma_r, \underbrace{0, \dots, 0}_{n-r \text{ times}})$$

where $\sigma_1 \ge ... \ge \sigma_r > 0$ are the singular values, and $U = [u_1, ..., u_m]$, $V = [v_1, ..., v_p]$ are square, orthogonal matrices $(U^T U = I_p, V^T V = I_m)$. The first r columns of U, V contains the left-and right singular vectors of A, respectively, that is:

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

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Jnsupervised learning

Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

Links between EVD and SVD

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 $U \wedge U^{T}$, with $\Lambda = \Sigma^{2}$.
 - The EVD of $A^T A$ is $V \wedge V^T$.

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

Robust Optimization & Machine Learning 4. Unsupervised Learning

Verview

Unsupervised learning

Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

- Basics SAFE Relaxation
- Algorithm
- Examples
- Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

Variational characterizations

Largest and smallest eigenvalues and singular values

If Q is square, symmetric:

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

If A is a general rectangular matrix:

$$\sigma_{\max}(A) = \max_{x : \|x\|_2 = 1} \|Ax\|_2.$$

Similar formulae for minimum eigenvalues and singular values.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning

Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics

- SAFE
- Relaxation
- Algorithms

Examples

Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

・ロト・西ト・山下・山下・ 日・ もくの

Variational characterizations

Other eigenvalues and singular values

If Q is square, symmetric, the k-th largest eigenvalue satisfies

$$\lambda^k = \max_{x \in \mathcal{S}^k, : \, \|x\|_2 = 1} \, x^T Q x,$$

where S^k is the subspace spanned by $\{u_k, \ldots, u_p\}$.

A similar result holds for singular values.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Jnsupervised learning

Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics

Bolovatia

Telaxatio

Algorithms

Voriente

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ のの⊙

Low-rank approximation

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

$$A^{(k)} := \arg\min_{X} ||X - A||_{F} : \operatorname{Rank}(X) \leq k,$$

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 $A = U\Sigma V^{T}$ is an SVD of the matrix A.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Jnsupervised learning

Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE

Relaxation

Algorithms

Examples

Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ のの⊙

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

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

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.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Jnsupervised learning

Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = 悪 = のへで

Low-rank approximation

Interpretation: low-rank case

When A is rank k, that is,

$$A = UV^T$$
, $U \in \mathbf{R}^{p \times k}$, $V \in \mathbf{R}^{m \times k}$, $k \ll m, p$,

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

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

Thus, each time-series is the sum of scaled copies of *k* time-series represented by v_1, \ldots, v_k , with scaling factors given in u_1, \ldots, u_k . We can think of v_i 's as the few "factors" that drive all the time-series.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Jnsupervised learning

Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

Outline

Overview of Unsupervised Learning Unsupervised learning models Matrix facts

Principal Component Analysis

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection Sparse graphical models Penalized maximum-likelihood Example

References

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation

Examples

Vandins

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Exampl

References

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = 悪 = のへで

Motivation



Votes of US Senators, 2002-2004. The plot is impossible to read...

- Can we project data on a lower dimensional subspace?
- If so, how should we choose a projection?

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

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, medecine.
- Engineering design, signal compression and image processing.
- Search engines, data mining.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations

Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxatio

Algoritrim

Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

Solution principles

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

- PCs are computed via EVD of covariance matrix.
- Can be interpreted as a "factor model" of original data matrix.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations

Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxatio

Algorithm

Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

Variance maximization problem

Definition

Let us normalize the direction in a way that does not favor any direction.

Variance maximization problem:

$$\max_{x} \, \mathbf{var}(x) \, : \, \|x\|_2 = 1.$$

A non-convex problem!

Solution is easy to obtain via the eigenvalue decomposition (EVD) of S, or via the SVD of centered data matrix A_c .

Robust Optimization & Machine Learning 4. Unsupervised Learning

verview Jnsupervised learnir

PCA

Motivations Variance maximization

Factor models Example

Sparse PCA

Basics

Relaxatio

Algorithm

Examples

Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

・ロト・西ト・山田・山田・山下

Variance maximization problem

Variance maximization problem:

$$\max_{x} x^{T} S x : \|x\|_{2} = 1.$$

Assume the EVD of S is given:

$$\boldsymbol{S} = \sum_{i=1}^{p} \lambda_i \boldsymbol{u}_i \boldsymbol{u}_i^{T},$$

with $\lambda_1 \ge \ldots \lambda_p$, and $U = [u_1, \ldots, u_p]$ is orthogonal ($U^T U = I$). Then $\arg \max_{x : ||x||_2 = 1} x^T S x = u_1,$

where u_1 is any eigenvector of *S* that corresponds to the largest eigenvalue λ_1 of *S*.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview Unsupervised learning Matrix facts

PCA

Motivations Variance maximization

Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

・ロト・日本・日本・日本・日本・日本

Variance maximization problem

Example: US Senators voting data





Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview Unsupervised learnin Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example Sparse PCA Pasion

Basics SAFE Relaxation Algorithms Examples

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

Projection of US Senate voting data on random direction (left panel) and direction of maximal variance (right panel). The latter reveals party structure (party affiliations added after the fact). Note also the much higher range of values it provides.

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.
- Fin 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 $k \ll p$).

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization

Deflation

Factor models Example

Sparse PCA

Basics SAFE Relaxatio Algorithm Examples

Sparse Covaria

Selection

Sparsity

Penalized maximum-likelihood

Example

References

・ロト・西ト・山田・山田・山下

Finding orthogonal directions Result

It turns out that the direction that solves

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

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, \ldots, u_k .

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization

Deflation

Factor models Example

Sparse PCA

Basics

- SAFE
- Relaxatio

Algorithm

Examples

Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

・ロト・西ト・山下・山下・ 日・ もくの

Factor models

PCA allows to build a low-rank approximation to the data matrix:

$$\boldsymbol{A} = \sum_{i=1}^{k} \sigma_i \boldsymbol{u}_i \boldsymbol{v}_i^{\mathsf{T}}$$

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

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation

Factor models

Example

Sparse PCA

Basics

- SAFE
- Relaxation

Algorithms

Examples

Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

・ロト・日本・日本・日本・日本・日本

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).

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models

Example

- Sparse PC/
- Basics SAFE Relaxation Algorithms Examples

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Exampl

References

・ロト・日本・日本・日本・日本

Outline

Overview of Unsupervised Learning Unsupervised learning models Matrix facts

Principal Component Analysis Motivations

Variance maximization

Deflation

Factor models

Example

Sparse PCA

Basics
SAFE
Relaxation
Algorithms
Examples
Variants

Sparse Covariance Selection Sparse graphical models Penalized maximum-likelihood Example

References

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Exampl

References

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = 悪 = のへで

Motivation

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.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics

SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

Sparse PCA Problem definition

Modify the variance maximization problem:

$$\max_{x} x^{T} S x - \lambda \operatorname{Card}(x) : ||x||_{2} = 1,$$

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

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

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ □臣 = のへで

Safe feature elimination

Express *S* as $S = R^T R$, with $R = [r_1, ..., r_p]$ (each r_i corresponds to one feature).

Theorem (Safe feature elimination [2]) *We have*

$$\max_{x: \|x\|_{2}=1} x^{T} S x - \lambda \operatorname{Card}(x) = \max_{z: \|z\|_{2}=1} \sum_{i=1}^{p} \max(0, (r_{i}^{T} z)^{2} - \lambda).$$

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models

Sparse PCA

Basics

SAFE

Relaxation

Algorithm

Examples

Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

▲□▶ ▲□▶ ▲目▶ ▲目▶ ▲□ ● のへで

SAFE

Corollary

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

- The presence of the penalty parameter allows to prune out dimensions in the problem.
- In practice, we want λ high as to allow better interpretability.
- Hence, interpretability requirement makes the problem easier in some sense!

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics

SAFE

Relaxation

Algorithm

Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

・ロト・日本・日本・日本・日本・日本

Relaxation for sparse PCA

Step 1: I1-norm bound

Sparse PCA problem:

$$\phi(\lambda) := \max_{x} x^{T} S x - \lambda \operatorname{Card}(x) : ||x||_{2} = 1,$$

First recall Cauchy-Schwartz inequality:

$$\|x\|_1 \leq \sqrt{\operatorname{Card}(x)} \|x\|_2,$$

hence we have the upper bound

$$\phi(\lambda) \leq \overline{\phi}(\lambda) := \max_{x} x^{T} S x - \lambda \|x\|_{1}^{2} : \|x\|_{2} = 1.$$

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basic

Relaxation

Algorithms Examples

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

Relaxation for sparse PCA

Step 2: lifting and rank relaxation

Next we rewrite problem in terms of (PSD, rank-one) $X := xx^T$:

 $\overline{\phi} = \max_{\mathbf{x}} \operatorname{Tr} SX - \lambda \|X\|_{1} : X \succeq 0, \quad \operatorname{Tr} X = 1, \quad \operatorname{Rank}(X) = 1.$

Drop the rank constraint, and get the upper bound

$$\overline{\lambda} \leq \psi(\lambda) := \max_{X} \operatorname{Tr} SX - \lambda \|X\|_{1} : X \succeq 0, \ \operatorname{Tr} X = 1.$$

- Upper bound is a semidefinite program (SDP).
- In practice, X is found to be (close to) rank-one at optimum.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Jnsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basic

Relaxation

Algorithms Examples

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ のの⊙

Sparse PCA Algorithms

- The Sparse PCA problem remains challenging due to the huge number of variables.
- Second-order methods become quickly impractical as a result.
- SAFE technique often allows huge reduction in problem size.
- Dual block-coordinate methods are efficient in this case [7].
- Still area of active research. (Like SVD in the 70's-90's...)

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PC/

Basics SAFE Relaxati

Algorithms

Example: Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ のの⊙

Example 1

Sparse PCA of New York Times headlines

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

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





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.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics

Deleveti

.....

Examples

Variants

Sparse Covariance Selection Sparsity Penalized

maximum-likelihood

Example

References

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 business	play team	government united_states	campaign bush	program children
company market	season game	u_s attack	administration	student
companies				

Words associated with the top 5 sparse principal components in NYTimes

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

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basic

SAFE

Relaxatio

Algorithm

Examples

Variants

Sparse Covariance Selection

Sparsity

Penalized

European la

References

▲□▶▲□▶▲□▶▲□▶ ▲□ ● ●

NYT Dataset

Comparison with thresholded PCA

Thresholded PCA involves simply thresholding the principal components.

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

Robust Optimization & Machine Learning 4. Unsupervised Learning

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

Robust 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 [1] 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?

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE

Relaxatio

Algorithm

Examples

Variants

Sparse Covariance Selection

Penalized

Example

References

・ロト・日本・日本・日本・日本・日本

Robust PCA model

In robust PCA, we solve the convex problem

 $\min_{N} \|\boldsymbol{A} - \boldsymbol{N}\|_{*} + \lambda \|\boldsymbol{N}\|_{1}$

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.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity Penalized maximum-likelihood Example

References

・ロト・西・・田・・田・・日・

CVX syntax

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 λ exist in the workspace:

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

Not the use of norm_nuc, which stands for the nuclear norm.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity Penalized maximum-likelihood

Example

References

・ロト・日本・日本・日本・日本・日本

Outline

Overview of Unsupervised Learning Unsupervised learning models Matrix facts

Principal Component Analysis

Motivations

- Variance maximization
- Deflation
- Factor models
- Example

Sparse PCA

Algorithms

Variants

Sparse Covariance Selection

Sparse graphical models Penalized maximum-likelihood Example

References

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PC/

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

- Basics SAFE Relaxation Algorithms Examples
- Variants

Sparse Covariance Selection

- Sparsit
 - Penalized maximum-likelihood
- Exampl

References

◆□▶ ◆□▶ ◆三▶ ◆三▶ ・三甲 のへぐ

Motivation

We'd like to draw a graph that describes the links between the features (*e.g.*, words).

- Edges in the graph should exist when some strong, natural metric of similarity exist between features.
- ► For better interpretability, a *sparse* graph is desirable.
- Various motivations: portfolio optimization (with sparse risk term), clustering, etc.

Here we focus on exploring *conditional independence* within features.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PC/

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood Example

References

・ロト・四ト・ヨト・ヨー もくら

Gaussian assumption

Let us assume that the data points are zero-mean, and follow a multi-variate Gaussian distribution: $x \simeq \mathcal{N}(0, \Sigma)$, with Σ a $p \times p$ covariance matrix. Assume Σ is positive definite.

Gaussian probability density function:

$$p(x) = \frac{1}{(2\pi \det \Sigma)^{p/2}} \exp((1/2)x^T \Sigma^{-1} x).$$

where $X := \Sigma^{-1}$ is the *precision* matrix.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE

Relaxatio

Algorithm

Examples

Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

References

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

Conditional independence

The pair of random variables x_i , x_j are *conditionally independent* if, for x_k fixed ($k \neq i, j$), the density can be factored:

 $p(x) = p_i(x_i)p_j(x_j)$

where p_i , p_j depend also on the other variables.

- Interpretation: if all the other variables are fixed then x_i, x_j are independent.
- Example: Gray hair and shoe size are independent, conditioned on age.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood Example

References

・ロト・四ト・日本・日本・日本・日本

Conditional independence

C.I. and the precision matrix

Theorem (C.I. for Gaussian RVs)

The variables x_i , x_j are conditionally independent if and only if the *i*, *j* element of the precision matrix is zero:

$$(\Sigma^{-1})_{ij}=0.$$

Proof.

The coefficient of $x_i x_j$ in log p(x) is $(\Sigma^{-1})_{ij}$.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

- Basics
- SAFE
- Relaxation
- Algorithms
- Examples
- Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

References

・ロト・日本・日本・日本・日本・日本

Sparse precision matrix estimation

Let us encourage sparsity of the precision matrix in the maximum-likelihood problem:

 $\max_{X} \log \det X - \operatorname{Tr} SX - \lambda \|X\|_{1},$

with $||X||_1 := \sum_{i,j} |X_{ij}|$, and $\lambda > 0$ a parameter.

- The above provides an invertible result, even if S is not positive-definite.
- The problem is convex, and can be solved in a large-scale setting by optimizing over column/rows alternatively.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

References

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ のの⊙

Dual

Sparse precision matrix estimation:

$$\max_{X} \log \det X - \operatorname{Tr} SX - \lambda \|X\|_{1}.$$

Dual:

$$\min_{U} - \log \det(S + U) : \|U\|_{\infty} \le \lambda.$$

Block-coordinate descent: Minimize over one column/row of *U* cyclically. Each step is a QP.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics

- Belavatia
- Tielaxatio

Algorithm

Examples

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

. . .

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Example Data: Interest rates





Using covariance matrix ($\lambda = 0$).

Using $\lambda = 0.1$.

The original precision matrix is dense, but the sparse version reveals the maturity structure.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxation Algorithms Examples Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

▲□▶ ▲□▶ ▲目▶ ▲目▶ ▲□ ● ● ●

Example Data: US Senate voting, 2002-2004



Again the sparse version reveals information, here political blocks within each party.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE Relaxatic

Algorithm

Examples

Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

・ロト・西ト・ヨト・ヨー もくの

Outline

Overview of Unsupervised Learning Unsupervised learning models Matrix facts

Principal Component Analysis

Motivations

Variance maximization

Deflation

Factor models

Example

Sparse PCA

Basics SAFE Relaxation

Examples

Variants

Sparse Covariance Selection Sparse graphical models Penalized maximum-likelihood Example

References

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PCA

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basics SAFE

Relaxation

Algorithm

Examples

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihood

Example

References

・ロト・日本・日本・日本・日本・日本

References



Emmanuel J. Candès, Xiaodong Li, Yi Ma, and John Wright.

Robust principal component analysis? 2009.



L. El Ghaoui.

On the quality of a semidefinite programming bound for sparse principal component analysis. arXiv:math/060144, February 2006.



Olivier Ledoit and Michael Wolf.

A well-conditioned estimator for large-dimensional covariance matrices. Journal of Multivariate Analysis, 88:365–411, February 2004.



O.Banerjee, L. El Ghaoui, and A. d'Aspremont.

Model selection through sparse maximum likelihood estimation for multivariate gaussian or binary data. Journal of Machine Learning Research, 9:485–516, March 2008.



Optimization for Machine Learning. MIT Press, 2011.



Y. Zhang, A. d'Aspremont, and L. El Ghaoui.

Sparse PCA: Convex relaxations, algorithms and applications.

In M. Anjos and J.B. Lasserre, editors, Handbook on Semidefinite, Cone and Polynomial Optimization: Theory, Algorithms, Software and Applications. Springer, 2011. To appear.



Large-scale sparse principal component analysis and application to text data. December 2011.

Robust Optimization & Machine Learning 4. Unsupervised Learning

Overview

Unsupervised learning Matrix facts

PC/

Motivations Variance maximization Deflation Factor models Example

Sparse PCA

Basic

Relavatio

Algorithm

Examples

Variants

Sparse Covariance Selection

Sparsity

Penalized maximum-likelihoo

Example

References

▲□▶ ▲圖▶ ▲国▶ ▲国▶ - 国 - のへで