# Large-scale Robust Optimization and Applications Part II: Applications

Laurent El Ghaoui

EECS and IEOR Departments UC Berkeley

SESO 2015 Tutorial

June 22, 2015

Large-scale Robust Optimization Part II

Unsupervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction

Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● のへで

# Outline

#### Overview of Machine Learning Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Robust Optimization for Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation Resource allocation Likelihood uncertainty models Reduction to a 1D problem Numerical Experiments

References

#### Large-scale Robust Optimization Part II

Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Numerical Experiments

References

# Outline

#### Overview of Machine Learning Unsupervised learning Supervised learning

Sparse supervised learning

Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Robust Optimization for Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation Resource allocation Likelihood uncertainty mode Reduction to a 1D problem Numerical Experiments

References

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Dimensionali

Robust low-rank LP Low-rank LASSO

#### Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Numerical Experiments

References

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 善臣 - 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.

#### Large-scale Robust Optimization Part II

#### Overviev

#### Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D problem

Numerical Experiments

References

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

# What is supervised learning?

In supervised learning, the data points are associated with "side" information that can "guide" (supervise) the learning process.

- In linear regression, each data point x<sub>i</sub> is associated with a real number y<sub>i</sub> (the "response"); the goal of learning is to fit the response vector to (say, linear) function of the data points, *e.g.* y<sub>i</sub> ≈ w<sup>T</sup>x<sub>i</sub>.
- In classification, the side information is a Boolean "label" (typically y<sub>i</sub> = ±1); the goal is to find a set of coefficients such that the sign of a linear function w<sup>T</sup>x<sub>i</sub> matches the values y<sub>i</sub>.
- In structured output models, the side information is a more complex structure, such a tree.

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

- Motivation Example SAFE Relaxation Algorithms Examples Variants
- Dimensionality Reduction Robust low-rank LP Low-rank LASSO

#### Robust Resource Allocation

- Resource allocation Likelihood uncertainty models Reduction to a 1D problem
- References

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ のの⊙

### Popular loss functions

Squared loss: (for linear least-squares regression)

$$L(z, y) = ||z - y||_2^2.$$

Hinge loss: (for SVMs)

$$L(z, y) = \sum_{i=1}^{m} \max(0, 1 - y_i z_i)$$

Logistic loss: (for logistic regression)

$$L(z, y) = -\sum_{i=1}^{m} \log(1 + e^{-y_i z_i})$$

#### Large-scale Robust Optimization Part II

#### Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Variants

#### Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

#### Robust Resource Allocation

Resource allocation Likelihood uncertainty

Reduction to a 1D problem Numerical Experiments

References

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

# Outline

Overview of Machine Learning Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

**Sparse PCA** 

Motivation Example SAFE Relaxation Algorithms Examples

Robust Optimization for Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation Resource allocation Likelihood uncertainty model Reduction to a 1D problem Numerical Experiments

References

#### Large-scale Robust Optimization Part II

Overviev

Unsupervised learning Supervised learning

#### Sparse supervised learning

Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

## Generic sparse learning problem

Optimization problem with cardinality penalty:

$$\min_{w} L(X^{T}w) + \lambda \|w\|_{0}$$

- ▶ Data:  $X \in \mathbf{R}^{n \times m}$ .
- Loss function L is convex.
- ▶ Cardinality function  $||w||_0 := |\{j : w_j \neq 0\}|$  is non-convex.
- λ is a penalty parameter allowing to control sparsity.

- Arises in many applications, including (but not limited to) machine learning.
- Computationally intractable.

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning

#### Basics

Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation Resource allocation

models Reduction to a 1D problem

Numerical Experiments

Reference

# **Classical** approach

A now classical approach is to replace the cardinality function with an  $l_1$ -norm:

$$\min_{w} L(X^{T}w) + \lambda \|w\|_{1}.$$

### Pros:

- Problem becomes convex, tractable.
- Often works very well in practice.
- Many "recovery" results available.

Cons: may not work!

#### Large-scale Robust Optimization Part II

#### Overviev

Unsupervised learning Supervised learning

Sparse supervised learning

#### Basics

Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

#### Dimensionality Reduction

Low-rank LASSO

#### Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D proble

Numerical Experiments

References



Consider the sparse learning problem

$$\min_{x} \|w\|_0 : X^T w = y.$$

Assume optimal point is unique, let  $w^{(0)}$  be the optimal point.

Now solve *I*<sub>1</sub>-norm approximation

$$w^{(1)} := \arg\min_{x} \|w\|_{1} : X^{T}w = y.$$

Since  $w^{(1)}$  is feasible, we have  $X^{T}(w^{(1)} - w^{(0)}) = 0$ .

Facts: (see [?])

- Set of directions that decrease the norm from  $w^{(1)}$  form a cone.
- If the nullspace of  $X^T$  does not intersect the cone, then  $w^{(1)} = w^{(0)}$ .

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

```
Sparse supervised
learning
```

#### Recovery

Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Dimensionality Reduction Robust low-rank LP

Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D problem

.

### Mean width

Let  $S \subseteq \mathbf{R}^n$  be a convex set, with support function

$$S_C(d) = \sup_{x \in S} d^T x$$

Then  $S_C(d) + S_C(-d)$  measures "width along direction d".



*Mean width:* with  $S^{n-1}$  be the unit Euclidean ball in  $\mathbf{R}^n$ ,

$$\omega(C) := \mathbf{E}_u S_C(u) = \int_{u \in S^{n-1}} S_C(u) du.$$

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning

Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Variants

Dimensionality Reduction Robust low-rank LF

Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D problem

Numerical Experiments

References

▲□▶▲□▶★□▶★□▶ □ のへで

## Gordon's escape theorem

When does a random subspace  $\mathcal{A} \in \mathbf{R}^n$  intersect a convex cone C only at the origin?

Theorem: (Gordon, 1988) If

 $\operatorname{codim}(\mathcal{A}) \geq n \cdot \omega (\mathcal{C} \cap \mathcal{S}^{n-1})^2$ ,

then with high probability:  $\mathcal{A} \cap \mathcal{C} = \{0\}$ .

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning

Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE

......

- - - -

Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

# Bounding mean width

A duality approach

$$\begin{split} \omega(C \cap S^{n-1}) &= & \mathbf{E}_u \max_{x \in C, \|x\| = 1} u^T x \\ &\leq & \mathbf{E}_u \max_{x \in C, \|x\| \le 1} u^T x \\ &= & \mathbf{E}_u \min_{v \in C^*} \|u - v\|, \end{split}$$

where  $C^*$  is the polar cone:

$$\mathcal{C}^* := \left\{ v \; : \; v^T u \leq 0 \text{ for every } u \in \mathcal{C} 
ight\}.$$

Name of the game is to *choose* an appropriate *v*.

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics

Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation

Algorithms

Examples

Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty

Reduction to a 1D problem Numerical Experiments

References

▲□▶▲□▶▲□▶▲□▶ □ のQ@

### **Recovery rates**

*Fact:* ([?]) Assume that the solution to cardinality problem with n variables and m constraints:

$$w^{(0)} = \arg\min_{x} \|w\|_{0} : X^{T}w = y$$

is unique and has sparsity s. Using the  $l_1$ -norm approximation

$$w^{(1)} = \arg\min_{x} \|w\|_{1} : X^{T}w = y,$$

the condition

$$m \ge 2s \log \frac{n}{s} + \frac{5}{4}s$$

guarantees that with high probability,  $w^{(1)} = w^{(0)}$ .

Similar results hold for a variety of norms (not just  $l_1$ ).

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning

#### Recovery

Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms

Variants

#### Dimensionality Reduction Robust low-rank Lf

Low-rank LASSO

#### Robust Resource Allocation

Resource allocation Likelihood uncertainty

models Reduction to a 1D problem

Numerical Experiments

References

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

Basic idea

LASSO and its dual

"Square-root" LASSO:

$$\min_{w} \|\boldsymbol{X}^{T}\boldsymbol{w}-\boldsymbol{y}\|_{2}+\lambda\|\boldsymbol{w}\|_{1}.$$

with  $X^T = [a_1, \ldots, a_n] \in \mathbf{R}^{m \times n}$ ,  $y \in \mathbf{R}^m$ , and  $\lambda > 0$  are given. (Each  $a_i \in \mathbf{R}^m$  corresponds to a variable in *w*, *i.e.* a "feature".)

### Dual:

$$\max_{\theta} \theta^T y : \|\theta\|_2 \leq 1, \ |a_i^T \theta| \leq \lambda, \ i = 1, \dots, n.$$

From optimality conditions, if at optimum in the dual the *i*-constraint is not active:

 $|\boldsymbol{a}_{i}^{T}\boldsymbol{\theta}| < \lambda$ 

then  $w_i = 0$  at optimum in the primal.

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

```
Sparse supervised
learning
Basics
Recovery
Safe Feature Elimination
```

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

```
Dimensionality
Reduction
Robust low-rank LP
Low-rank LASSO
```

Robust Resource

Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ のの⊙

### Basic idea

Safe Feature Elimination (SAFE)

From optimality:

 $|\boldsymbol{a}_i^T\boldsymbol{\theta}| < \lambda \Longrightarrow \boldsymbol{w}_i = \boldsymbol{0}.$ 

Since the dual problem involves the constraint  $\|\theta\|_2 \leq 1$ , the condition

$$\forall \, \theta, \ \|\theta\|_2 \leq 1 \ : \ |\boldsymbol{a}_i^T \theta| < \lambda$$

ensures that  $w_i = 0$  at optimum.

SAFE condition:

$$\|\boldsymbol{a}_i\|_2 < \lambda \Longrightarrow \boldsymbol{w}_i = \boldsymbol{0}.$$

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE

Relaxatio

Algorithm

Examples

Variants

Dimensionality Reduction

Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

# Advanced SAFE tests

Test can be strenghtened:

- Exploit optimal solution to problem for a higher value of λ.
- ► Use idea within the loop of a coordinate-descent (CD) algorithm.
- Allows to eliminate variables on the go.

Test is cheap:

- SAFE test costs as much as one iteration of gradient or CD method.
- Typically involves matrix-vector multiply X<sup>T</sup>w, with w a sparse vector.

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Bobust low-rank L

Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models Beduction to a 1D proble

Numerical Experiments

References

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

# Experiment

*Data:* KDD 2010b, 30M features, 20M documents. Target cardinality is 50.



- Applying SAFE in the loop of a coordinate-descent algorithm.
- Graph shows number of features involved to attain a given sparsity level.

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PC Motivation Example SAFE Relaxation Algorithms Examples Variants

Reduction Robust low-rank LP

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

# Outline

#### Overview of Machine Learning Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Robust Optimization for Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation Resource allocation Likelihood uncertainty model Reduction to a 1D problem Numerical Experiments

References

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Numerical Experiments

References

# Principal Component Analysis



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?

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised earning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty

Reduction to a 1D problem Numerical Experiments

References

◆□▶ ◆□▶ ◆三▶ ◆三▶ ●□ ● ●

# 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.

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

#### Motivation

Example SAFE Relaxation Algorithm: Examples Variants

#### Dimensionality Reduction Robust low-rank Lf

Low-rank LASSO

#### Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D probler

Numerical Experiments

References

# 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.

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation

Example SAFE Relaxation Algorithms Examples

Variants

Dimensionality Reduction

Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ のの⊙

# Variance maximization problem

Definition

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

### Variance maximization problem:

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

A non-convex problem!

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

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation

Example SAFE Relaxation Algorithm: Examples

Variants

#### Dimensionality Reduction

Low-rank LASSO

#### Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D problem

Jumerical Experiments

Reference

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ = 臣 = のへで

# Variance maximization problem

Variance maximization problem:

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

Assume the EVD of S is given:

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

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  $\lambda_1$  of *S*.

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

#### Motivation Example SAFE Relaxation Algorithms Examples Variants

#### Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

#### Robust Resource Allocation

Resource allocation

models

Reduction to a 1D problem Numerical Experiments

References

◆□▶ ◆□▶ ◆三▶ ◆三▶ ・三 ・ つへぐ

## Variance maximization problem

Example: US Senators voting data





Large-scale Robust Optimization Part II



Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

> Dimensionality Reduction Robust low-rank LI

Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Numerical Experiments

Reference

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  $k \ll p$ ).

#### Large-scale Robust Optimization Part II

#### Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation

Algorithms Examples Variants

> Dimensionality Reduction

Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D problem

Numerical Experiments

References

### Finding orthogonal directions Result

It turns out that the direction that solves

$$\max_{x} \operatorname{var}(x) : x^{T} u_{1} = 0$$

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

After *k* steps of the deflation process, the directions returned are  $u_1, \ldots, u_k$ .

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation

Example SAFE Relaxation Algorithms Examples

Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty

Reduction to a 1D problem Numerical Experiments

References

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● のへで

### Factor models

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

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

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

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation

Example SAFE Relaxatio Algorithm

Examples

Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

▲□▶▲□▶▲□▶▲□▶ □ のQ@

### 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).

・ロト ・ ( 目 ト ・ 目 ト ・ 日 - )

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation

Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LF

Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

# Sparse PCA: 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.

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation

#### Example

SAFE Relaxation Algorithms Examples Variants

> Dimensionality Reduction Robust low-rank LF

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Numerical Experiments

References

▲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.

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation

Example

SAFE Relaxation Algorithms Examples Variants

> Dimensionality Reduction Robust low-rank LP

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Numerical Experiments

References

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● 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 [?]) *We have*

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

#### Large-scale Robust Optimization Part II

#### Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation

#### SAFE

Relaxation Algorithms

Examples

. Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● 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 \u03c6 high as to allow better interpretability.
- Hence, interpretability requirement makes the problem easier in some sense!

#### Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example

SAFE

Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models Beduction to a 1D proble

Numerical Experiments

References

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

# 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

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

#### Large-scale Robust Optimization Part II

#### Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivatio

#### Relaxation

Algorithms Examples Variants

#### Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

#### Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

## Relaxation for sparse PCA

Step 2: lifting and rank relaxation

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

 $\overline{\phi} = \max \operatorname{Tr} SX - \lambda \|X\|_1 : X \succeq 0, \quad \operatorname{Tr} X = 1, \quad \operatorname{Rank}(X) = 1.$ 

Drop the rank constraint, and get the upper bound

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

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

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivatio Example

SAFE

Relaxation

Algorithms Examples Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ 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 [?].
- Still area of active research. (Like SVD in the 70's-90's...)

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivatio Example SAFE Belaxatio

#### Algorithms

Examples Variants

Dimensionality Reduction Robust low-rank LF

Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ 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:



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.

# Large-scale Robust Optimization Part II

## Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation

Examples Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

▲□▶ ▲圖▶ ▲国▶ ▲国▶ - 国 - のへで

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

# 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).

# Large-scale Robust Optimization Part II

## Overviev

Unsupervised learning Supervised learning

Sparse supervised earning Basics Recovery Safe Feature Elimination

## Sparse PCA

Motivation Example SAFE Relaxatio

Examples

Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D probler

Numerical Experiments

Reference

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ = 臣 = のへで

# NYT Dataset

Comparison with thresholded PCA

Thresholded PCA involves simply thresholding the principal components.

<i>k</i> = 2	<i>k</i> = 3	<i>k</i> = 9	<i>k</i> = 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.

# Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA Motivation

Example SAFE Relaxatio

Examples

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへで

# **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* [?] 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?

## Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation

Examples

Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty

Reduction to a 1D problem Numerical Experiments

References

# 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  $I_1$ -norm of the vector of singular values, and  $I_1$ -norm minimization encourages sparsity of its argument.

## Large-scale Robust Optimization Part II

## Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

## Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

## Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Numerical Experiments

## References

# 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  $\lambda$  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.

# Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms

#### Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

## Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

# Outline

# Overview of Machine Learning Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

# Robust Optimization for Dimensionality Reduction Robust low-rank LP Low-rank LASSO

# Robust Resource Allocation Resource allocation Likelihood uncertainty mode Reduction to a 1D problem Numerical Experiments

References

# Large-scale Robust Optimization Part II

#### Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

# Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

## Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

# Low-rank LP

Consider a linear programming problem in *n* variables with *m* constraints:

$$\min_{x} c^{\mathsf{T}} x : A x \leq b,$$

with  $A \in \mathbf{R}^{m \times n}$ ,  $b \in \mathbf{R}^m$ , and such that

- Many different problem instances involving the same matrix A have to be solved.
- The matrix A is close to low-rank.

- Clearly, we can approximate A with a low-rank matrix A<sub>ir</sub> once, and exploit the low-rank structure to solve many instances of the LP fast.
- In doing so, we cannot guarantee that the solutions to the approximated LP are even feasible for the original problem.

## Large-scale Robust Optimization Part II

# Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D problem

Numerical Experiments

Reference

# Approach: robust low-rank LP

For the LP

$$\min_{x} c^{\mathsf{T}} x : A x \leq b,$$

with many instances of b, c:

- Invest in finding a low-rank approximation A<sub>lr</sub> to the data matrix A, and estimate ∈ := ||A − A<sub>lr</sub>||.
- Solve the robust counterpart

$$\min_{x} c^{\mathsf{T}} x : (A_{\mathrm{lr}} + \Delta) x \leq b \ \forall \Delta, \ \|\Delta\| \leq \epsilon.$$

Robust counterpart can be written as SOCP

$$\min_{x,t} c^{T}x : A_{lr}x + t\mathbf{1} \le b, \ t \ge ||x||_{2}.$$

► We can exploit the low-rank structure of A<sub>ir</sub> and solve the above problem in time linear in m + n, for fixed rank.

## Large-scale Robust Optimization Part II

# Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

# Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource

Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

# A motivation: topic imaging

*Task:* find a short list of words that summarizes a topic in a large corpus. (StatNews project; see Miratrix et al, 2014)



Image of topic "Climate change" over time. Each square encodes the size of regression coefficient in LASSO. *Source:* People's Daily, 2000-2011.

# Interactive plot at

http://statnews.eecs.berkeley.edu/showcase/staircase\_economy/stair.html

# Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LP

Robust Resource

Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

Reference

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

# Low-rank LASSO

In many learning problems, we need to solve many instances of the LASSO problem

$$\min_{w} \|\boldsymbol{X}^{\mathsf{T}}\boldsymbol{w} - \boldsymbol{y}\|_2 + \lambda \|\boldsymbol{w}\|_1.$$

where

- For all the instances, the matrix X is a rank-one modification of the same matrix X.
- Matrix  $\tilde{X}$  is close to low-rank (hence, X is).

In the topic imaging problem:

- $\tilde{X}$  is a term-by-document matrix that represents the whole corpus.
- ▶ y is one row of X that encodes presence or absence of the topic in documents.
- X contains all remaining rows.

## Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LP

Robust Resource Allocation Resource allocation Likelihood uncertainty models Reduction to a 1D problen Numerical Experiments

Reference

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ のの⊙

# Robust low-rank LASSO

The robust low-rank LASSO

$$\min_{w} \max_{\|\Delta\| \leq \epsilon} \|(X_{\mathrm{lr}} + \Delta)^{\mathsf{T}} w - y\|_2 + \lambda \|w\|_1$$

is expressed as a variant of "elastic net":

$$\min_{\boldsymbol{w}} \|\boldsymbol{X}_{\mathrm{lr}}^{\mathsf{T}}\boldsymbol{w} - \boldsymbol{y}\|_{2} + \lambda \|\boldsymbol{w}\|_{1} + \epsilon \|\boldsymbol{w}\|_{2}.$$

- Solution can be found in time linear in m + n, for fixed rank.
- Solution has much better properties than low-rank LASSO, e.g. we can control the amount of sparsity.

# Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

## Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LP

# Robust Resource Allocation Resource allocation

likelihood uncertainty nodels

Numerical Experiments

References

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

# Example



Rank-1 LASSO (left) and Robust Rank-1 LASSO (right) with random data. The plot shows the elements of the solution as a function of the  $l_1$ -norm penalty parameter.

- Without robustness (ϵ = 0), the cardinality is 1 for 0 < λ < λ<sub>max</sub>, where λ<sub>max</sub> is a function of data. For λ ≥ λ<sub>max</sub>, w = 0 at optimum. Hence the l<sub>1</sub>-norm fails to control the solution.
- With robustness (ε = 0.01), increasing λ allows to gracefully control the number of non-zeros in the solution.

# Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation likelihood uncertainty nodels Reduction to a 1D problem Jumerical Experiments

Reference

# Numerical experiments: low-rank approximation

RCV1V2 NYTIMES Dataset TMC2007 28,596 23.149300,000 8,200,000 n 46.236 d 49.060 102,660 141.043 Time (s)  $\sigma_{k+1}/\sigma_1$ Time (s)  $\sigma_{k+1}/\sigma_1$ Time (s)  $\sigma_{k+1}/\sigma_1$ Time (s)  $\sigma_{k+1}/\sigma_1$ k = 50.1539 0.26090.4095187 0.4072k = 101 0.1196 1 0.210050 0.3075451 0.3494k = 151 0.1010 1 0.1907 59 0.2709520 0.30410.2793 k = 202 0.0958 2 0.176973 0.2432589 k = 253 0.0909 3 0.1662 87 0.2312687 0.2680k = 30794 0.25804 0.0880 4 0.1615 93 0.2180k = 354 0.0858 4 0.15550.2098932 0.2477114 k = 405 0.0836 5 0.15070.20120.2354130 1150 6 5 0.2255k = 450.0826 0.1475142 0.19321208 0.2209 k = 500.0811 0.1430 158 0.18501862

Are real-world datasets approximately low-rank?

Runtimes<sup>1</sup> for computing a rank-k approximation to the whole data matrix.

# Large-scale Robust Optimization Part II

Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Dimensionality Reduction

Robust low-rank LP

Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D proble

Numerical Experiments

References

<sup>1</sup> Experiments are conducted on a personal work station: 16GB RAM, 2.6GHz quadeore Intel: 🕨 < 🚊 🔹 🖓 🔍

# Multi-label classification

In multi-label classification, the task involves the same data matrix X, but many different response vectors y.

- Treat each label as a single classification subproblem (one-vs-all).
- Evaluation metric: Macro-F1 measure.
- Datasets:
  - RCV1-V2: 23,149 training documents; 781,265 test documents; 46,236 features; 101 labels.
  - TMC2007: 28,596 aviation safety reports; 49,060 features; 22 labels.

# Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Example SAFE Relaxation Algorithms Examples Variants Dimensionali

Reduction Robust low-rank LE

Low-rank LASSO

Robust Resource Allocation Resource allocation

Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

# Multi-label classification

Plot performance vs. training times for various values of rank  $k = 5, 10, \ldots, 50$ .



# RCV1V2 data set



In both cases, the low-rank robust counterpart allows to recover the performance obtained with full-rank LASSO (red dot), for a fraction of computing time.

# Large-scale Robust Optimization Part II

## Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

## Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

# Robust Resource Allocation

Resource allocation Likelihood uncertainty nodels Reduction to a 1D problem Numerical Experiments

## Reference

# Topic imaging

- Labels are columns of whole data matrix X.
- Compute low-rank approximation of X
   When a column is removed.
- Evaluation: report predictive word lists for 10 queries.
- Datasets:
  - NYTimes: 300,000 documents; 102,660 features, file size is 1GB. Queries: 10 industry sectors.
  - PUBMED: 8,200,000 documents; 141,043 features, file size is 7.8GB. Queries: 10 diseases.
- ► In both cases we have pre-computed a rank k (k = 20) approximation using power iteration.

# Large-scale Robust Optimization Part II

## Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Bobust low-rank LE

# Low-rank LASSO

Robust Resource Allocation

Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

# Topic imaging

automotive	agriculture	technology	tourism	aerospace	defence	financial	healthcare	petroleum	gaming
car	government	company	tourist	boeing	afghanistan	company	health	oil	game
vehicle	farm	computer	hotel	aircraft	attack	million	care	prices	gambling
auto	farmer	system	business	space	forces	stock	cost	gas	casino
sales	food	web	visitor	program	military	market	patient	fuel	player
model	water	information	economy	jet	gulf	money	corp	company	online
driver	trade	internet	travel	plane	troop	business	al_gore	barrel	computer
ford	land	american	tour	nasa	aircraft	firm	doctor	gasoline	tribe
driving	crop	job	local	flight	terrorist	fund	drug	bush	money
engine	economic	product	room	airbus	president	investment	medical	energy	playstation
consumer	country	software	plan	military	war	economy	insurance	opec	video

The New York Times data: Top 10 predictive words for different queries corresponding to industry sectors.

arthritis	asthma	cancer	depression	diabetes	gastritis	hiv	leukemia	migraines	parkinson
joint	bronchial	tumor	effect	diabetic	gastric	aid	cell	headache	treatment
synovial	asthmatic	treatment	treatment	insulin	h_pylori	infection	acute	headaches	effect
infection	children	carcinoma	disorder	level	chronic	cell	bone_marrow	pain	nerve
chronic	respiratory	cell	depressed	glucose	ulcer	hiv-1	leukemic	disorder	syndrome
pain	symptom	chemotherapy	pressure	control	acid	infected	tumor	women	disorder
treatment	allergic	survival	anxiety	plasma	stomach	antibodies	remission	chronic	neuron
fluid	infant	risk	symptom	diet	atrophic	risk	t_cell	duration	receptor
knee	inhalation	dna	drug	liver	antral	positive	antigen	symptom	alzheimer
acute	airway	malignant	neuron	renal	reflux	transmission	chemotherapy	gene	response
therapy	fev1	diagnosis	response	normal	treatment	drug	expression	therapy	brain

*PubMed* data: Top 10 predictive words for different queries corresponding to diseases.

# Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised earning Basics Recovery Safe Feature Elimination

## Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LP

Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D problem

Reference

# Outline

# Overview of Machine Learning Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery

Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Robust Optimization for Dimensionality Reduction Robust low-rank LP Low-rank LASSO

# **Robust Resource Allocation**

Resource allocation Likelihood uncertainty models Reduction to a 1D problem Numerical Experiments

References

## Large-scale Robust Optimization Part II

Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

# Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D problem

References

# **Resource allocation**

We consider resource allocation problems, of the form

 $\max_{w\in\mathcal{W}} U(w)$ 

where

$$\mathcal{W} := \left\{ \boldsymbol{w} \in \mathbf{R}^n : \boldsymbol{w} \ge \mathbf{0}, \ \boldsymbol{w}^T \mathbf{1} = \mathbf{1} \right\},\$$

and U is a concave utility function.

The vector w may represent

- A fraction of budget allocated across n different items;
- A proportion of time spent displaying an ad.

## Large-scale Robust Optimization Part II

#### Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

# Resource allocation

Likelihood uncertainty models Reduction to a 1D problem Numerical Experiments

References

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

# Robust resource allocation problem

Many resource allocation problems are of the form

$$\phi := \max_{w \in \mathcal{W}} \min_{r \in \mathcal{R}} r^T w,$$

where the "return vector" r is assumed to be unknown-but-bounded via a given "uncertainty set"  $\mathcal{R}$ .

The corresponding utility function

$$U(w) := \min_{r \in \mathcal{R}} r^T w$$

is concave, and positively homogeneous.

# Large-scale Robust Optimization Part II

(1)

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank Lf

Low-rank LASSO

Robust Resource Allocation

# Resource allocation

Likelihood uncertainty models Reduction to a 1D problem Numerical Experiments

References

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

# Challenges

Practical challenges:

- ▶ How to choose the uncertainty set *R*?
- Can we connect this choice to some probabilistic model of the return?
- Can we solve the problem fast, e.g., in linear time?

# Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Dimensionality Reduction

Low-rank LASSO

Robust Resource Allocation

# Resource allocation

Likelihood uncertainty models Reduction to a 1D problem Numerical Experiments

References

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● のへで

# Example: portfolio optimization

In finance, we consider *r* to be a "return" vector, and *w* represents a portfolio, with return  $r^T w$ . In practice, *r* is never fully known.

In our model, the return vector is assumed to be uncertain, and only known to be contained in the given set  $\mathcal{R}$ .

For example, we may assume that the set  $\mathcal{R}$  is an ellipsoid:

 $\mathcal{R} = \{\hat{r} + Ru : \|u\|_2 \leq \kappa\},\$ 

with  $\hat{r} \in \mathbf{R}^n$ , *R* a matrix, and  $\kappa$  a measure of the size of the ellipsoid.

#### Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LF

Robust Resource

#### Resource allocation

Likelihood uncertainty models Reduction to a 1D problem Numerical Experiments

References

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ○ □ ● のへで

# Connection with Gaussian models

In practice the ellipsoid  ${\mathcal R}$  can be derived from a Gaussian assumption on the return.

Specifically: if we assume that the returns are Gaussian, with mean  $\hat{r}$  and covariance matrix  $\Sigma$ . Factor  $\Sigma$  as  $\Sigma = RR^{T}$ , with R a matrix. Then the set  $\mathcal{R}$  is a set of confidence for the returns, based on the normal likelihood function.

The robust portfolio optimization problem reads

$$\max_{\boldsymbol{w}\in\mathcal{W}} \hat{\boldsymbol{r}}^T \boldsymbol{w} - \kappa \|\boldsymbol{R}^T \boldsymbol{w}\|_2$$

This is closely connected to the (more standard) mean-variance model (shown here with "risk aversion parameter"  $\sigma$ ):

$$\max_{\boldsymbol{w}\in\mathcal{W}} \hat{\boldsymbol{r}}^T \boldsymbol{w} - \sigma \|\boldsymbol{R}^T \boldsymbol{w}\|_2^2.$$

## Large-scale Robust Optimization Part II

## Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

## Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

> Dimensionality Reduction Robust low-rank LF

Low-rank LASSO

## Robust Resource Allocation

# Resource allocation

Likelihood uncertainty models Reduction to a 1D problem Numerical Experiments

## References

# Challenges

In practice, estimating  $\boldsymbol{\Sigma}$  in high dimensions is hard. Further, solving the problem

$$\max_{\boldsymbol{w}\in\mathcal{W}} \hat{\boldsymbol{r}}^{\mathsf{T}}\boldsymbol{w} - \kappa \|\boldsymbol{R}^{\mathsf{T}}\boldsymbol{w}\|_2,$$

or its more standard mean-variance version, requires  $O(n^3)$ , which may be prohibitive.

# Large-scale Robust Optimization Part II

## Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

# Resource allocation

Likelihood uncertainty models Reduction to a 1D problem Numerical Experiments

References

# ▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● のへで

# **Motivation**

We seek to derive the uncertainty set  $\ensuremath{\mathcal{R}}$  from a probabilistic model of the returns.

To this end, we assume that the set  $\mathcal{R}$  has the form

$$\mathcal{R} := \{r : H(r) \le \kappa\},\$$

with *H* the negative log-likelihood, and  $\kappa \ge 0$  is a measure of uncertainty.

The above uncertainty model is very natural as it corresponds to returns that are likely under the assumed probabilistic model.

# Large-scale Robust Optimization Part II

#### Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation

# Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

◆□▶ ◆□▶ ◆三▶ ◆三▶ ・三甲 のへぐ

# Decomposable uncertainty

We assume that the function H is convex, differentiable, and decomposable:

$$\forall r \in \mathbf{dom} \ h : \ H(r) = \sum_{i=1}^n h_i(r_i),$$

with  $h_i$ 's convex and differentiable. We make a few additional technical assumptions on H, seen next.

When H is a negative log-likelihood, the decomposability corresponds to assuming that the different components of the return vector r are independent.

## Large-scale Robust Optimization Part II

## Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

# Dimensionality Reduction Robust low-rank LF

Robust Resource

Resource allocation

Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

# **Technical assumptions**

- 1. The functions *h<sub>i</sub>* and their gradient can be easily computed anywhere on their respective domain.
- 2. The quantities

$$\tau_i^u := \arg\min_{\tau} h_i(\tau), \quad \kappa_i := h_i(\tau_i^u) = \min_{\tau} h_i(\tau)$$

are finite, and available.

3. The following condition holds:

$$\kappa > \kappa_{\min} := \min_{r} H(r) = \sum_{i=1}^{n} \kappa_i,$$

so that the equivalent problem

$$\phi = \min_{r \in \mathcal{R}(\kappa)} \max_{1 \le i \le n} r$$

is strictly feasible.

4. A lower bound on  $\phi$ ,  $\phi_{\min}$ , is available.

## Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

## Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Dimensionality Reduction

Low-rank LASSO

Robust Resource Allocation

Resource allocation

# Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References



The expressions

$$h_i(r_i)=\frac{1}{2\sigma_i^2}(r_i-\hat{r}_i)^2,$$

naturally arise when the returns are assumed to be Gaussian, with a diagonal covariance matrix. Here,  $\hat{r}_i \in \mathbf{R}$ ,  $\sigma_i \in \mathbf{R}_{++}$ , i = 1, ..., n are given.

- The diagonal covariance matrix corresponds to an independence assumptions.
- The constraint  $H(r) \le \kappa$  naturally "couples" the returns.
- Compare this with an "interval model"  $r_i \in [\hat{r}_i \kappa \sigma_i, \hat{r}_i + \kappa \sigma_i]$ , which would allow returns that are **jointly** very unlikely.

# Large-scale Robust Optimization Part II

# Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

# Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation

# Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

# References

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ○ □ ○ ○ ○ ○

# Comments

- The model couples the different components of r, even though the random variable r has uncorrelated components. This captures the fact that jointly observing large values for independent Gaussian scalars is a rare event.
- The model puts a very low burden on statistical estimation task, as only individual variances need be estimated, and does not require the knowledge of the full covariance matrix.

# Large-scale Robust Optimization Part II

# Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank Ll

Low-rank LASSO

Robust Resource Allocation

Resource allocation

Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

# ◆□▶ ◆□▶ ◆目▶ ◆目▶ ●目 ● のへで

# Example <sup>β</sup> distributions

The  $\beta$ -likelihood models arise with functions  $h_i$  with domain [0, 1], of the form

$$h_i(r_i) = -\alpha_i \log(r_i) - \beta_i \log(1 - r_i), \ r_i \in [0, 1]$$

and  $+\infty$  otherwise. This corresponds to a log-likelihood function for  $\beta$ -distributions, with  $\alpha_i \ge 1$ ,  $\beta_i \ge 1$  corresponding to event counts.

In this case,

$$\tau_i^u = \frac{\alpha_i}{\alpha_i + \beta_i}.$$

Such models are useful in the context of sparse data, since they allow to gracefully enforce non-negativity of returns.

# Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation

Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

# Main result

# Theorem

With the assumptions in place, the robust allocation problem can be solved as a one-dimensional one:

$$\phi = \min_{t} t : \sum_{i=1}^{n} h_i(\min(t,\tau_i^u)) \le \kappa.$$
(2)

Once the above problem is solved, the optimal weights are given as follows. Set  $\tau_i^* = \min(t^*, \tau_i^u), \eta_i^* = (-h'_i(\tau_i^*))_+, i = 1, ..., n$ . Then,  $\eta^* \neq 0$ , and

$$w_i^* = \frac{\eta_i^*}{\sum\limits_{j=1}^n \eta_j^*}, \ i = 1, \dots, n.$$
 (3)

▲ロト ▲周 ト ▲ ヨ ト ▲ ヨ ト つのの

# Large-scale Robust Optimization Part II

## Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation

Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

References

# **Bisection algorithm**

We can solve the problem with a simple bisection algorithm, provided we know upper and lower bounds on t,  $t^{u}$ ,  $t^{l}$ :

Input data:  $\kappa$ ,  $h_i(\cdot)$ , where  $i = 1, \ldots, n$ ; and  $\epsilon$ .

1. Compute  $\tau^{u}$ ,  $t^{l}$ ,  $t^{u}$  as detailed next.

2. Set 
$$t = (t^u + t')/2$$
.

• If 
$$\sum_{i=1}^{n} h_i(\min(t, \tau_i^u)) \le \kappa$$
, set  $t^u = t$ ;

• Otherwise, set t' = t.

3. If  $t^u - t' \le \epsilon$ , exit.

# Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation

models

Reduction to a 1D problem Numerical Experiments

References

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ○ □ ○ ○ ○ ○

# Initialization

For an upper bound, we note that the vector  $\tau^{u}$  is feasible:  $H(\tau^{u}) \leq \kappa$ , we have then  $\phi = t^{*} \leq t^{u} := \max_{1 \leq i \leq n} \tau_{i}^{u}$ .

For the lower bound, we have  $t^* \ge t^l := \max_i t^i$ , where  $t^i = \min_{r \in \mathcal{R}(\kappa)} r_i$ . The constraint translates as

$$h_i(r_i) \leq \eta_i := \kappa - \sum_{i=1}^n h_i(\tau_i^u).$$

We then have to solve the problems

$$t^i = \min_{\xi} \xi : h_i(\xi) \leq \eta_i.$$

Usually these can be solved in closed-form in specific instances. If the set  $\mathcal{R}(\kappa)$  is contained in the non-negative orthant, we simply set  $t^{l} = 0$ . In case the above problem is not easily solved, we can simply set  $t^{l} = \phi_{\min}$ , where  $\phi_{\min}$  is any lower bound on  $\phi$  (which we assumed is known).

Large-scale Robust Optimization Part II

#### Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples Variants

Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Reduction to a 1D problem Numerical Experiments

Reference

# Numerical experiment: robust bandit problem

- We have applied the decision model to a bandit problem with Bernoulli return rates uniformly sampled from the interval [0.18, 0.2].
- We compared different approaches (UCB and Thomson sampling) to ours.
- ▶ We have used a simple uncorrelated Gaussian model.
- ► The simulations run for  $T = 10^6$  rounds and the policies are only updated every 1000 rounds.
- We measure performance in terms of cumulative regret.

#### Large-scale Robust Optimization Part II

## Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithm Examples Variants

Dimensionality Reduction Robust low-rank LF

Robust Resource Allocation

Resource allocation Likelihood uncertainty

Reduction to a 1D problem

Numerical Experiments

References

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ のの⊙

# Results

# Large-scale Robust Optimization Part II



Mean regret for UCB, Thompson Sampling ('Thompson') and Robust policy with confidence levels 0.999 ('Robust 0.999'), 0.9 ('Robust 0.9') and 0.5 ('Robust 0.5'). The mean of the regret is computed with 20 repetitions. イロト イポト イヨト

Numerical Experiments

э
# Outline

### Overview of Machine Learning Unsupervised learning Supervised learning

Sparse supervised learning Basics Becovery

Safe Feature Elimination

# Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Robust Optimization for Dimensionality Reduction Robust low-rank LP Low-rank LASSO

Robust Resource Allocation Resource allocation Likelihood uncertainty model Reduction to a 1D problem Numerical Experiments

# References

Large-scale Robust Optimization Part II

Overview

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

Sparse PCA

Motivation Example SAFE Relaxation Algorithms Examples

Variants

Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

Robust Resource Allocation

Resource allocation Likelihood uncertainty models Reduction to a 1D problem

Numerical Experiments

References

# References I

### Large-scale Robust Optimization Part II

#### Overviev

Unsupervised learning Supervised learning

Sparse supervised learning Basics Recovery Safe Feature Elimination

#### Sparse PCA

Motivation Example SAFE Relaxation Algorithms

Variants

### Dimensionality Reduction

Robust low-rank LP Low-rank LASSO

### Robust Resource Allocation

Resource allocation Likelihood uncertainty models

Numerical Experiments

### References

# ◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ のへぐ