

EE4.66 Topics in Large Dimensional Data Processing

Wei Dai and Stefan Vlaski

Last Compiled: October 16, 2021

Contents

Contents	ii
1 Syllabus	1
1.1 What is this module about	1
1.2 The miscellaneous	1
1.3 Assessment	1
1.4 Syllabus	2
2 Introduction	5
2.1 Key Challenges	5
2.2 Data Points Are Isolated	5
2.3 False Structure	6
2.4 Overfitting	7
2.4.1 Univariate Polynomial Regression	7
2.4.2 Multivariate Polynomial Regression	9
3 Least Squares	10
3.1 Problem Statement	10
3.2 Matrix Decomposition	12
3.3 Projecton	15
3.4 The Underdetermined Case	17
3.5 Numerical Stability	18
4 Sparse Inverse Problems: Introduction	19
4.1 Problem Statement	19
4.2 Exhaustive Search	20
4.2.1 The Algorithm	20
4.2.2 Performance Guarantees	22
4.3 Applications	23
5 Sparse Inverse Problems: L0 Methods	26
5.1 Some Easy Problems	26
5.2 OMP	27
5.2.1 Algorithm	27
5.2.2 Performance Guarantees	28
5.3 SP and Similar Algorithms	30
5.3.1 Algorithms	30
5.3.2 Performance Guarantees	32

5.4	Majorization Minimization (MM)	37
5.5	Appendix	39
6	Convex Optimization: Introduction	41
6.1	Convexity	41
6.1.1	Definitions	41
6.1.2	Conditions of Convexity	42
6.2	Unconstrained Convex Optimization	44
6.3	Line Search Methods	45
7	Sparse Inverse Problems: ℓ_1-minimization	47
7.1	Convex Relaxation	47
7.2	LASSO Solvers	50
7.2.1	Soft Thresholding Function	50
7.2.2	Cyclic Coordinate Descent	51
7.2.3	Iterative Shrinkage Thresholding Algorithm (ISTA)	51
7.2.4	Majorization Minimization	52
7.2.5	Least Angle Regression (LAR)	52
8	Bilinear Inverse Problems	55
8.1	Bilinear Inverse Problems	55
8.1.1	Matrix Norms	57
8.2	Examples	58
8.2.1	Low-Rank Matrix Completion	58
8.2.2	Blind Deconvolution	60
8.2.3	Sparse PCA and Dictionary Learning	61
8.3	Methods	61
8.3.1	Alternating Optimization	61
8.3.2	Greedy Algorithms	62
8.3.3	Convex Relaxation	62
8.4	Appendix	63
8.4.1	Proof of Eckart-Young-Mirsky Theorem	63
9	Constrained Convex Optimization	66
9.1	General Form	66
9.1.1	Sublevel Sets	66
9.2	Duality	67
9.3	KKT Conditions	69
9.4	Examples	70
10	ADMM	73
10.1	Precursors	73
10.1.1	Dual Ascent Method	73

10.1.2	Method of Multipliers	74
10.2	ADMM	75
10.2.1	Optimality Conditions	77
11	Stochastic Optimization	78
11.1	Motivation	78
11.2	Stochastic Least Squares	78
11.2.1	Gradient Descent for LS	78
11.2.2	Stochastic Gradient Descent for LS	79

Syllabus

1

1.1 What is this module about

It is *not* about

- ▶ Programming
- ▶ Computer architecture

It is about

- ▶ Concepts, Principles, and Technical Tools
 - Main Topics
 - * Sparse linear inverse problems
 - * Sparse bilinear inverse problems
 - Technical Tools
 - * Linear algebra (matrix computations).
 - * Convex and nonconvex optimisation
 - Selected applications
 - * Google PageRank
 - * Denoising.
 - * Face recognition with block occlusion.
 - * Recommendation engine: Netflix problem.
 - * Dictionary learning
 - * Blind deconvolution.
 - * ...

1.1 What is this module about	1
1.2 The miscellaneous	1
1.3 Assessment	1
1.4 Syllabus	2

1.2 The miscellaneous

GTA

Miss Farwa Abbas, Mr Zhengang Guo, and Mr Yifan Ran

1.3 Assessment

Coursework 100%¹

- ▶ Project 1: ℓ_0 Methods for Sparse Recovery (16%)
- ▶ Project 2: ℓ_1 Methods for Sparse Recovery (20%)
- ▶ Project 3: Bilinear Inverse Problems (26%)
- ▶ Project 4: Read and Tell (30%)
- ▶ Peer Marking (8%)

1: More details on coursework will be announced.

1.4 Syllabus

The purpose of the module

The aim of the module is to introduce formulations and methods (convex and nonconvex optimization) arisen to process and analyse large dimensional data, and expose students to various application domains including image/video processing, machine learning, collaborative filtering, etc.

Prerequisites

Students are expected to have basic knowledge of linear algebra and algorithms.

Balance between theory and practice:

The lectures will cover ideas, formulations, and algorithms behind large dimensional data processing.

Some representative applications will be discussed in the lectures to show that difficult problems can be solved by clever formulation and appropriately designed methods.

Coursework will help students

- ▶ use modern and popular software tools for big data processing and machine learning, including VSCode, Jupyter Notebook, Markdown, LaTex, and Julia;
- ▶ design and implement algorithmic solutions to test taught concepts and solve real-world problems;
- ▶ explore recent research results, expand knowledge, get inspired by modern applications and novel solutions, and develop future career.

Learning outcomes

1. Knowledge:

Students will learn the theory underlying large dimensional data processing and analysis. We will

- a) present ideas, formulations, and algorithms for large dimensional data processing;
- b) demonstrate how basic tools in linear algebra and optimization can be applied to significantly improve the efficiency and accuracy of data processing;
- c) expose students to relevant real applications, including imaging processing, online recommendation, machine learning, etc.

2. Subject specific skills

a) Intellectual skills:

Students will learn essential paradigms and tools to handle large and high-dimensional data (e.g. in engineering, commerce, machine learning)

b) Practical Skills:

Students will learn to use modern tools popular for data science and machine learning, design and implement efficient algorithms for big data application, and get exposed to modern data processing methods and applications.

c) Transferable / Key Skills:

Students will learn to apply their knowledge to real-world situations. Through examples where real-world problems are addressed by different approaches, students will witness critical thinking and innovations.

Selected Applications

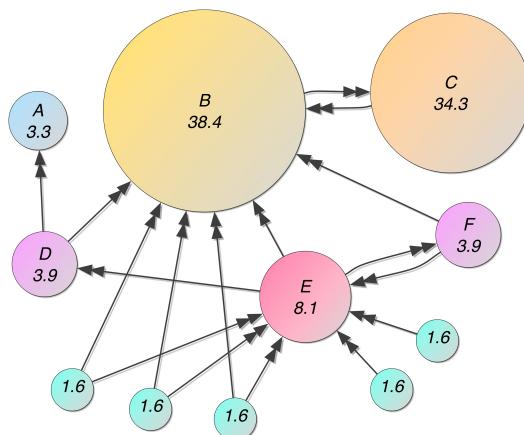


Figure 1.1: Google pagerank.

By 345Kai, Public Domain, <https://en.wikipedia.org/w/index.php?curid=14424767>



Figure 1.2: Image Denoising.

	Movie 1	Movie 2	Movie 3	Movie 4
User 1	★★★★★	★★★☆☆		
User 2	★★★★☆		★★★★★	
User 3		★★★★☆☆		★★★★☆
User 4	★★★★★			★★★★★

Figure 1.3: Netflix Problem.

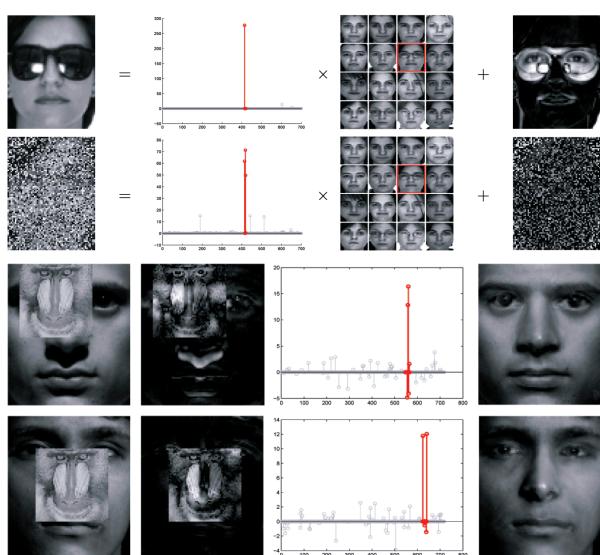


Figure 1.4: Face Recognition with Block Occlusion J. Wright et al., 2009.

2

Introduction

2.1 Key Challenges

What is data analysis?

Data analysis is to learn an unknown function f from training data x such that

$$y_i \approx f(x_i),$$

where $x_i \in \mathbb{R}^n$, $i = 1, 2, \dots, m$.

Curse of Dimensionality

Use some image database as an example.¹

In Classical data processing

- ▶ A small number n of parameters.
- ▶ A large number m of observations.

In many modern applications

- ▶ A large number n of parameters.
- ▶ A relatively small sample size m ($m \approx n$ or $m < n$).

Curse of dimensionality:

- ▶ The computational difficulty
- ▶ The intrinsic statistical difficulty
 - Data points are isolated.
 - False structures.
 - Overfitting (the inferred model describes the noise instead of the underlying relationship).

To address it: low dimensional structure.

2.2 Data Points Are Isolated

Example 2.2.1 (K-nearest neighbors algorithm) Suppose we have training data $\{x_1, y_1\}, \{x_2, y_2\}, \dots, \{x_m, y_m\}$ taking values in $\mathbb{R}^n \times \{0, 1\}$, where y is the class label of x for classification, or the value of the underlying (unknown) function for regression.

For a given test data x_t , reorder the training data such that $\|x_{(1)} - x_t\| \leq \dots \leq \|x_{(m)} - x_t\|$.

- ▶ For classification: majority vote using K-nearest neighbors.
- ▶ For regression: average value of the K-nearest neighbors.

2.1 Key Challenges	5
2.2 Data Points Are Isolated .	5
2.3 False Structure	6
2.4 Overfitting	7
2.4.1 Univariate Polynomial Regression	7
2.4.2 Multivariate Polynomial Regression	9

1:

- ▶ THE MNIST database of handwritten digits.
 - $28 \times 28 = 784$ pixels.
 - 10 categories, 60,000 training samples, 10,000 test samples.
- ▶ Caltech 256: Pictures of objects belonging to 256 categories.
 - Pictures of various sizes: normally $100 \times 100 = 10,000$ pixels.
 - 256 categories, 30,607 images in total.
- ▶ Modern database
 - Including Imagenet, Labelme, etc.
 - Pictures of various sizes, including HD ones.
 - Less training images per category in general.

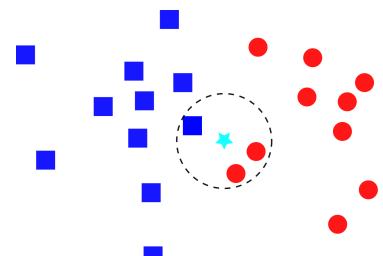


Figure 2.1: KNN Method

The performance of KNN is decided by how dense the training points are.

Question: In the n -dimensional space, how many training samples are needed so that for any given test sample, there exists at least one training sample that is *at most distance 1 away*?

Mathematical Question: how many unit balls are needed to cover the whole space the hypercube $[0, 1]^n$?

Cover the hypercube $[0, 1]^n$ by unit balls:

- The volume of $V_n(r)$ of a n -dimensional ball of radius $r > 0$:

$$V_n(r) = \frac{\pi^{n/2}}{\Gamma(n/2 + 1)} r^n \approx \left(\frac{2\pi e r^2}{n}\right)^{n/2} (n\pi)^{-1/2}.$$

- To cover the hypercube $[0, 1]^n$ by unit balls, one must have

$$[0, 1]^n \subset \bigcup_{i=1}^m B_n(x^{(i)}, 1).$$

- That is, $1 \leq m V_n(1)$, or

$$m \geq \frac{\Gamma(n/2 + 1)}{\pi^{n/2}} \approx \left(\frac{n}{2\pi e}\right)^{n/2} \sqrt{n\pi}.$$

Required number of data points for covering:

n	20	30	50	100	150
m	39	45630	5.7×10^{12}	42×10^{39}	1.28×10^{72}

2.3 False Structure

Task: Estimate empirical covariance matrix.

Given an n -dimensional random vector X , the covariance matrix is defined as

$$\Sigma := E[(X - E[X])(X - E[X])^T].$$

Suppose that the random variable X is normally distributed. Given a sample consisting of m independent observations x_1, \dots, x_m , the maximum likelihood estimator of the covariance matrix is given by

$$\begin{aligned} \hat{\Sigma} &:= \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})^T \\ &= \frac{1}{m} \tilde{X} \tilde{X}^T, \end{aligned}$$

where $\bar{x} = \frac{1}{m} \sum_i x_i$ is the sample average and $\tilde{X} = [x_1 - \bar{x}, \dots, x_m - \bar{x}]$ is the sample covariance matrix.

Case 1: If n is fixed and $m \rightarrow \infty$, by the law of large numbers²

$$\hat{\Sigma} \rightarrow I.$$

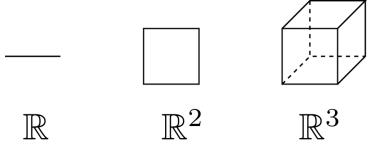


Figure 2.2: Unit tube in n -dimensional space

2: Interested readers see https://en.wikipedia.org/wiki/Law_of_large_numbers for details.

Case 2: What happens if m and n are at the same scale?

Example 2.3.1 Let $X \sim \mathcal{N}(\mathbf{0}, I)$. Compute the eigenvalues of the sample covariance matrix $\hat{\Sigma}$ with $n = 200$ and $m = 10^5$ ($m \gg n$). The histogram of the eigenvalues is shown in Figure 2.3.

Let $X \sim \mathcal{N}(\mathbf{0}, I)$. Compute the eigenvalues of $\hat{\Sigma}$ with $n = 2000$ and $m = 10^4$ ($m \gtrapprox n$). The histogram of the eigenvalues is shown in Figure 2.4.

Remark 2.3.1 The behaviour of eigenvalues of the empirical covariance can be predicted by random matrix theory.

If $n, m \rightarrow \infty$ proportionally with $n/m \rightarrow r \in \mathbb{R}^+$, then the empirical distribution of eigenvalues of empirical covariance matrix $\hat{\Sigma}$ converges to [Marchenko-Pastur distribution](#) (see https://en.wikipedia.org/wiki/Marchenko-Pastur_distribution for more details). See Figure 2.5 for simulations.

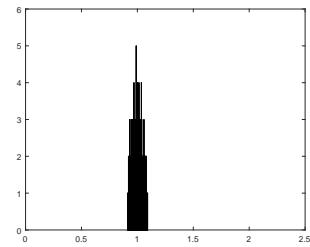


Figure 2.3: Eigenvalue histogram when number of samples is much more than the sample dimension.

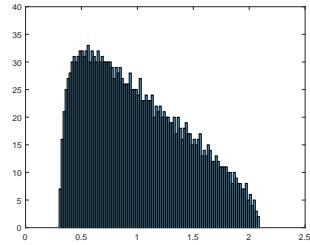


Figure 2.4: Eigenvalue histogram when number of samples is roughly equal to the sample dimension.

2.4 Overfitting

Definition 2.4.1 (Linear Regression) Given training samples (x_i, y_i) , $i = 1, \dots, m$, want to estimate a linear function represented by $\alpha \in \mathbb{R}^n$ s.t. $y_i \approx \langle x_i, \alpha \rangle$.

Solution: Let $\mathbf{y} = [y_1, \dots, y_m]^T$, $\mathbf{X} = [x_1, \dots, x_m]^T$, and $\mathbf{e} = [e_1, \dots, e_m]$. Write

$$\mathbf{y} = \mathbf{X}\alpha + \mathbf{e},$$

Then solve for α .³

Issue: when $m < n$, there are infinite many valid solutions to $\mathbf{y} = \mathbf{X}\alpha$.

2.4.1 Univariate Polynomial Regression

Task: Suppose that the input data $x \in \mathbb{R}$ (the univariate case). Assume that the unknown function f can be approximated by a degree S univariate polynomial:

$$f(x) = \sum_{s=0}^S \alpha_s x^s.$$

Based on the input-output pairs (x_i, y_i) , $i = 1, 2, \dots, m$, one would find the unknown polynomial so that $y_i \approx f(x_i)$.

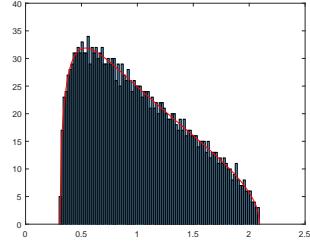


Figure 2.5: Empirical distribution of eigenvalues converges to Marchenko-Pastur distribution.

3: See https://en.wikipedia.org/wiki/Linear_regression for a more detailed description.

Fact

Given m distinct samples, there exists a polynomial of degree $m - 1$ to match the data perfectly.

Proof. Since

$$\sum_{\ell=0}^{m-1} \alpha_\ell x_i^\ell = y_i, \quad i = 1, 2, \dots, m,$$

It holds that

$$\underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}}_y = \underbrace{\begin{bmatrix} 1 & x_1 & \cdots & x_1^{m-1} \\ 1 & x_2 & \cdots & x_2^{m-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_m & \cdots & x_m^{m-1} \end{bmatrix}}_X \underbrace{\begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_{m-1} \end{bmatrix}}_\alpha.$$

The matrix X is a Vandermonde matrix and has full rank.⁴ Hence, the unknown vector α can be computed from $\alpha = X^{-1}y$. \square

4: See https://en.wikipedia.org/wiki/Vandermonde_matrix for the definition and property of Vandermonde matrix.

Issue: The polynomial obtained using the above method can be an overfitting of noise/numerical errors.

Example 2.4.1 Fit a univariate polynomial to the data points shown in Figure 2.6.

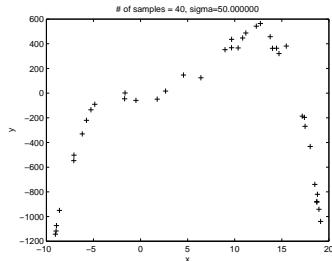


Figure 2.6: The data to be fit with.

Now we use the above method to fit a univariate polynomial to the data. It is clear that $f(x_i) = y_i$ from the results shown in Figure 2.7.

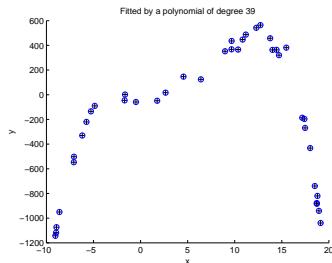


Figure 2.7: Fitted polynomial at data points.

However, the fitted polynomial is actually an overfitting. See Figure 2.8 for details.

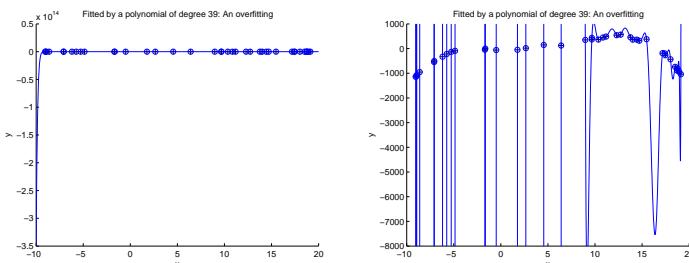


Figure 2.8: Fitted polynomial is an overfitting.

The reason of overfitting,⁵ in this case, is the poor condition number of the matrix X .⁶

$$y = X\alpha_0 + w \Rightarrow \hat{\alpha} = X^{-1}y = \alpha_0 + X^{-1}w.$$

Solution: A sparse solution of α .

α is **sparse** (only a few nonzeros): force $\alpha_5 = \alpha_6 = \dots = \alpha_{39} = 0$.

As presented in Figure 2.9, the sparse solution provides a much better fitting.

5: See <https://en.wikipedia.org/wiki/Overfitting> for overfitting.

6: See https://en.wikipedia.org/wiki/Condition_number. There is also an exercise in the first coursework.

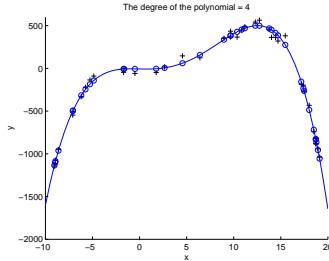


Figure 2.9: A sparse solution gives a much better fitting.

2.4.2 Multivariate Polynomial Regression

Task: Suppose that the input data $x \in \mathbb{R}^d$ (the multivariate case). Assume that the unknown function f can be approximated by a multivariate polynomial.

Here, we are particularly interested in a degree 2 multivariate polynomial

$$\begin{aligned} f(x) = & \alpha_0 + \alpha_1 x_1 + \dots + \alpha_d x_d \\ & + \alpha_{d+1} x_1^2 + \alpha_{d+2} x_1 x_2 + \dots + \alpha_{2d} x_1 x_d \\ & + \alpha_{2d+1} x_2^2 + \dots + \alpha_{n-2} x_{d-1} x_d \\ & + \alpha_{n-1} x_d^2. \end{aligned}$$

The regression task:

Given $(x(j), y(j)), j = 1, 2, \dots, m$, try to find $\alpha = [\alpha_0, \alpha_1, \dots, \alpha_{n-1}]$.

$$\underbrace{\begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ y(m) \end{bmatrix}}_y = \underbrace{\begin{bmatrix} 1 & x_1(1) & \dots & x_d(1) \\ 1 & x_1(2) & \dots & x_d^2(2) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1(m) & \dots & x_d^2(m) \end{bmatrix}}_X \underbrace{\begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_{n-1} \end{bmatrix}}_\alpha.$$

Two issues:

- No sufficient data:
 - Number of terms $n = O(d^2) \gg m$.
 - Even when data is abundant: need to avoid overfitting.

Sparsity plays a key role: Enforce most of the elements of α to be zero.

3

Least Squares

3.1 Problem Statement

Consider the model

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w},$$

where \mathbf{y} is the measurement vector, \mathbf{x} is the unknown signal/parameter vector, and \mathbf{w} is the noise vector.

The **least squares** estimation of \mathbf{x} is given by

$$\hat{\mathbf{x}} = \min_{\mathbf{x}} \underbrace{\frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2}_{f(\mathbf{x})}.$$

It can be verified that

$$\nabla f(\mathbf{x}) = \mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{y}).$$

The **optimal** solution (in least squares sense) is obtained by setting $\nabla f(\mathbf{x}) = \mathbf{0}$. One has

$$(\mathbf{A}^T\mathbf{A})\hat{\mathbf{x}} = \mathbf{A}^T\mathbf{y}.$$

We have two cases.

1. $\mathbf{A}^T\mathbf{A}$ is of full rank. In this case,

$$\hat{\mathbf{x}} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{y}.$$

2. $\mathbf{A}^T\mathbf{A}$ is rank deficient. In this case, the global minimizer \mathbf{x}^* is not unique. We will study one particular global minimizer defined as follows.¹ Consider the compact singular value decomposition (SVD) $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$. Set

$$\hat{\mathbf{x}} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T\mathbf{y}.$$

For both cases, we can write

$$\hat{\mathbf{x}} = \mathbf{A}^\dagger\mathbf{y} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T\mathbf{y},$$

where \mathbf{A}^\dagger is the **Moore-Penrose pseudoinverse**.

This chapter discusses

- ▶ Matrix decomposition: How to compute \mathbf{A}^\dagger especially when $\mathbf{A}^T\mathbf{A}$ is not of full rank.
 - Application: Google PageRank
- ▶ Projection: $\mathbf{A}\hat{\mathbf{x}} = \text{proj}(\mathbf{y}, \text{span}(\mathbf{A}))$.
- ▶ The underdetermined case: $\hat{\mathbf{x}} = \mathbf{A}^\dagger\mathbf{y}$ has the least length among all global minimizers.
- ▶ Numerical stability

3.1 Problem Statement	10
3.2 Matrix Decomposition	12
3.3 Projecton	15
3.4 The Underdetermined Case	17
3.5 Numerical Stability	18

1: The details will be covered in the rest of this chapter.

Pseudoinverse

Linear Independence and Dependence: Vectors $v_1, \dots, v_n \in \mathbb{R}^m$ are **linearly independent** if

$$\sum \lambda_i v_i = \mathbf{0} \Rightarrow \lambda_i = 0, \forall i.$$

In matrix format,

$$[v_1, \dots, v_n] \lambda = \mathbf{0} \in \mathbb{R}^m \Rightarrow \lambda = \mathbf{0} \in \mathbb{R}^n.$$

Vectors $v_1, \dots, v_n \in \mathbb{R}^m$ are **linearly dependent** if $\exists \lambda \neq \mathbf{0}$ such that

$$[v_1, \dots, v_n] \lambda = \mathbf{0}.$$

Matrix Rank: Let $A \in \mathbb{R}^{m \times n}$ be a given matrix.

Column rank: the maximum number of linearly independent columns.

Row rank: the maximum number of linearly independent rows.

Rank: For every matrix, column rank = row rank = rank.²

2: In Julia, use function `rank` (requiring `LinearAlgebra` package) to calculate the rank of a matrix.

Definition 3.1.1 (Matrix Inverse and Pseudoinverse)

- A matrix $A \in \mathbb{R}^{n \times n}$ is invertible if there exists a matrix $B \in \mathbb{R}^{n \times n}$ such that

$$AB = BA = I.$$

The matrix B is the inverse of A .

- For a matrix A , its Moore-Penrose pseudoinverse A^\dagger is defined to satisfy

- $AA^\dagger A = A$, and $A^\dagger AA^\dagger = A^\dagger$.
- $(AA^\dagger)^T = AA^\dagger$, and $(A^\dagger A)^T = A^\dagger A$.

3

3: In Julia, use functions `inv` and `LinearAlgebra > pinv` for inverse and pseudoinverse respectively.

- A is invertible $\Leftrightarrow A$ is of full rank.
- Moore-Penrose pseudoinverse is well defined no matter whether A has full rank or not.

Example 3.1.1

It is easy to compute the inverse and Moore-Penrose pseudoinverse of diagonal matrices.

- $A = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \Rightarrow A^{-1} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1 \end{bmatrix} = A^\dagger.$
- $A = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix} \Rightarrow A^\dagger = \begin{bmatrix} 1/2 & 0 \\ 0 & 0 \end{bmatrix}.$

Its inverse is not well defined.

Matrix decomposition helps compute the inverse and Moore-Penrose pseudoinverse of a non-diagonal matrix.

3.2 Matrix Decomposition

Eigenvalue decomposition

Definition 3.2.1 (**Eigendecomposition**, spectral decomposition) A non-zero vector $v \in \mathbb{R}^n$ is an **eigenvector** of a **square** matrix $A \in \mathbb{R}^{n \times n}$ if there is a constant λ such that

$$Av = \lambda v,$$

where λ is called the **eigenvalue** corresponding to v .

4

4: In Julia, use function `LinearAlgebra > eigen` for eigenvalue decomposition.

Example 3.2.1

This example is linked to Google PageRank discussed in Section ??.

1. Let

$$A = \begin{bmatrix} 5 & 3 \\ 2 & 6 \end{bmatrix}.$$

Show that A has an eigenvalue 8 and the corresponding eigenvector is $v = [1, 1]^T$.

Use `Julia > LinearAlgebra > eigen` to numerically compute the eigenvalues and eigenvectors of A . Verify the above claim. Note that the two eigenvectors are linearly independent but not orthogonal.

2. Let $A \in \mathbb{R}^{n \times n}$ be a square matrix such that the sum of every row is the same, i.e.,

$$c = \sum_j A_{i,j} \quad \forall i.$$

Show that A has an eigenvalue c and the corresponding eigenvector is $v = [1, 1, \dots, 1]^T$.

3. Suppose that the eigenvalues and eigenvectors of A are given by $\lambda_1, \dots, \lambda_n$, and q_1, \dots, q_n , respectively.

Let $\Lambda = \text{diag}([\lambda_1, \dots, \lambda_n])$ and $Q = [q_1, \dots, q_n]$. It can be verified that

$$AQ = Q\Lambda \Rightarrow A = Q\Lambda Q^{-1}.$$

Now let $B = A^T$. Denote the eigenvalues and eigenvectors of B by $\lambda_{B,1}, \dots, \lambda_{B,n}$, and $q_{B,1}, \dots, q_{B,n}$, respectively. Let $\Lambda_B = \text{diag}([\lambda_{B,1}, \dots, \lambda_{B,n}])$ and $Q_B = [q_{B,1}, \dots, q_{B,n}]$.

What can you say about the relationship between Λ and Λ_B , and that between Q and Q_B ?

4. With the above, you can prove that the adjacent matrix in the Google PageRank example (Example 3.2.4) has an eigenvalue 1.

Singular Value Decomposition

Definition 3.2.2 (Singular Value Decomposition) For an arbitrary matrix $A \in \mathbb{R}^{m \times n}$, there exists a factorization of the form

$$A = U\Sigma V^T = \sum_i \sigma_i u_i v_i^T,$$

where $U \in \mathbb{R}^{m \times m}$ is unitary (contains m orthonormal columns), $\Sigma \in \mathbb{R}^{m \times n}$ is diagonal matrix with non-negative diagonal entries, and $V \in \mathbb{R}^{n \times n}$ is unitary.

The m columns of U and the n columns of V are called the **left-singular vectors** and **right-singular vectors** of M , respectively. The diagonal entries σ_i of Σ are known as the **singular values** of M .

5

5: In Julia, use function `LinearAlgebra.svd` for singular value decomposition.

Remark 3.2.1 1. Let $\sigma_i \geq 0$.

Note that if $A = \sum_i \sigma_i u_i v_i^T$ is a valid SVD, then so is $A = \sum_i -\sigma_i (-u_i) v_i^T$. WLOG, it is a convention to assume $\sigma_i \geq 0$.

2. List the singular values in descending order, i.e., $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min(m,n)} \geq 0$. In this case, the diagonal matrix Σ is uniquely determined.
3. The matrix U and V are not unique.

This is clear by observing that $A = \sum_i \sigma_i u_i v_i^T = \sum_i \sigma_i (-u_i) (-v_i)^T$.

Remark 3.2.2 (Compact SVD) Suppose that the matrix $A \in \mathbb{R}^{m \times n}$ has the rank $r < \min(m, n)$. Its compact SVD is given by

$$A = U_r \Sigma_r V_r^T,$$

where $U_r = U_{:,1:r} \in \mathbb{R}^{m \times r}$, $\Sigma_r = \Sigma_{1:r,1:r} \in \mathbb{R}^{r \times r}$, and $V_r = V_{:,1:r} \in \mathbb{R}^{n \times r}$.

Example 3.2.2 (Full and Compact SVDs)

$$\begin{aligned} & \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix} \\ &= \underbrace{\frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}}_U \underbrace{\begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}}_{\Sigma} \underbrace{\frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}}_{V^T} \quad (\text{Full SVD}) \\ &= \underbrace{\frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}}_U \underbrace{\begin{bmatrix} 2 \\ 0 \end{bmatrix}}_{\Sigma} \underbrace{\frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}}_{V^T} \quad (\text{Compact SVD}) \end{aligned}$$

Example 3.2.3 Find full SVD of the matrices

$$A = \frac{1}{2} \begin{bmatrix} 1 & 1 & -2 \\ 1 & 1 & 2 \\ 1 & -1 & 2 \\ 1 & -1 & -2 \end{bmatrix}$$

and

$$\mathbf{B} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 \\ -1 & 1 & 1 & -1 \\ 2 & -2 & -2 & 2 \end{bmatrix}.$$

6

6: A computing software is not needed.
Mind the pattern of columns of \mathbf{A} and \mathbf{B} .

EVD and SVD

Let

$$\mathbf{M} = \mathbf{U}\Sigma\mathbf{V}^T.$$

Then

$$\mathbf{M}\mathbf{M}^T = \mathbf{U}\Sigma^2\mathbf{U}^T \text{ and } \mathbf{M}^T\mathbf{M} = \mathbf{V}\Sigma^2\mathbf{V}^T.$$

- The left-singular vectors \mathbf{u}_i 's of \mathbf{M} are eigenvectors of $\mathbf{M}\mathbf{M}^T$.
- The right-singular vectors \mathbf{v}_i 's of \mathbf{M} are eigenvectors of $\mathbf{M}^T\mathbf{M}$.
- The singular values σ_i 's of \mathbf{M} are the square roots of the eigenvalues of both $\mathbf{M}\mathbf{M}^T$ and $\mathbf{M}^T\mathbf{M}$. That is, $\lambda_i = \sigma_i^2$.

Decomposition and Inverse

Eigenvalue decomposition and inverse: If matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ can be eigendecomposed and if none of its eigenvalues are zero, then \mathbf{A} is invertible (nonsingular) and its inverse is given by

$$\mathbf{A}^{-1} = \mathbf{Q}\Lambda^{-1}\mathbf{Q}^{-1},$$

where Λ is a diagonal matrix with

$$\Lambda_{i,i} = \lambda_i.$$

Singular value decomposition and pseudoinverse: Consider the **compact SVD** of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ given by

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T.$$

Then it holds that

$$\mathbf{A}^\dagger = \mathbf{V}\Sigma^\dagger\mathbf{U}^T. \quad (3.1)$$

It is straightforward to verify the following⁷

7: A nice exercise for you.

- An equivalent definition is

$$\mathbf{A}^\dagger = (\mathbf{A}^T\mathbf{A})^\dagger\mathbf{A}^T.$$

- The definition in (3.1) is equivalent to Definition 3.1.1.

•

$$\mathbf{A}\mathbf{A}^\dagger\mathbf{A} = \mathbf{A}, \mathbf{A}^\dagger\mathbf{A}\mathbf{A}^\dagger = \mathbf{A}^\dagger.$$

•

$$\mathbf{A}\mathbf{A}^\dagger = \mathbf{U}\mathbf{U}^T, \mathbf{A}^\dagger\mathbf{A} = \mathbf{V}\mathbf{V}^T. \quad (3.2)$$

Application

Google PageRank: Assign importance scores (ranks) to webpages.

Mathematical Model: Denote the rank of a webpage i by x_i . The goal is to find x_i such that

$$x_i = \sum_{j \in C_i} \frac{1}{N_j} x_j,$$

where

$$C_i = \{j : \text{There is a link from } j \text{ to } i\}$$

N_i = Number of links from j

Example 3.2.4 The graph in Figure 3.1 can be represented by the adjacent matrix

$$A = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & \frac{1}{4} \\ \frac{1}{3} & 0 & 0 & \frac{1}{2} & \frac{1}{2} & \frac{1}{4} \\ 0 & 1 & 0 & 0 & 0 & \frac{1}{4} \\ \frac{1}{3} & 0 & 0 & 0 & 0 & 0 \\ \frac{1}{3} & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{4} \\ 0 & 0 & 0 & 0 & \frac{1}{2} & 0 \end{bmatrix}.$$

It can be proved that A has an eigenvalue 1 and the corresponding eigenvector $\mathbf{x} = [x_1, x_2, \dots, x_n]^\top$ satisfies that

$$A\mathbf{x} = \mathbf{x} \Rightarrow x_i = \sum_{j \in C_i} \frac{1}{N_j} x_j.$$

Compute the eigenvalue decomposition of A using a scientific computation software. One gets

$$\mathbf{x} = [0.57, 0.49, 0.53, 0.19, 0.32, 0.16]^\top.$$

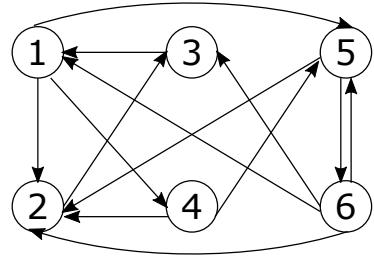


Figure 3.1: A simple directed graph as an example

3.3 Projecton

Definition 3.3.1 (Linear subspace) Let $\mathcal{B} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\} \subset \mathbb{R}^m$ containing linearly independent vectors.

The **linear span** of \mathcal{B} is defined as

$$\text{span}(\mathcal{B}) = \left\{ \sum_{i=1}^n \lambda_i \mathbf{v}_i : \lambda_i \in \mathbb{R} \right\}.$$

It is a **linear subspace** of \mathbb{R}^m .

- The set \mathcal{B} is a **basis** for the linear subspace $\mathcal{S} = \langle \mathcal{B} \rangle$.

- The basis \mathcal{B} **may not be unique**, but its dimension is.
- $\dim(\mathcal{S}) = n$: the # of vectors in a basis.

- \mathcal{B} is orthonormal if $\langle v_i, v_j \rangle = 0$ for $i \neq j$ and $\langle v_i, v_i \rangle = 1$.
- For convenience, we use $\text{span}(\mathcal{B})$ and $\text{span}(B)$ interchangeably where $B = [v_1, \dots, v_n]$.

Definition 3.3.2 (Projection) The *projection* of $v \in \mathbb{R}^m$ onto the subspace $\text{span}(A)$ is defined as

$$v_p = \text{proj}(v, A) = \arg \min_{u \in \text{span}(A)} \|v - u\|_2.$$

And the *projection residue* is given by

$$v_r = \text{resid}(v, A) = v - v_p.$$

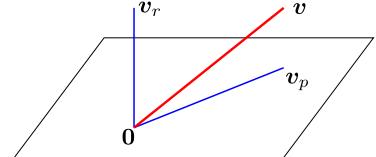


Figure 3.2: Projection

Example 3.3.1 Let

$$v = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}, B = \begin{bmatrix} \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & -\frac{2}{\sqrt{6}} \end{bmatrix}.$$

Find $\text{proj}(v, A)$ and $\text{proj}(v, B)$.

Note that the above example is easy to solve because both matrices A and B are orthonormal. The following proposition provides a way to calculate projection for general tall matrix A .

Proposition 3.3.1 Consider the compact SVD $A = U\Sigma V^T$. Then

1. $v_p = AA^\dagger v = UU^T v$.
2. $v_r \perp \text{span}(A)$.

Proof. 1. Any $u \in \text{span}(A)$ can be written as $u = Ax$ for some x .

Hence

$$\begin{aligned} v_p &= \arg \min_{u \in \text{span}(A)} \|v - u\|_2^2 \\ &= A \arg \min_x \|v - Ax\|_2^2 \end{aligned}$$

By setting $\nabla \|v - Ax\|_2^2 = \mathbf{0}$, it is clear that

$$\hat{x} = A^\dagger v \text{ and } v_p = AA^\dagger v.$$

That $v_p = UU^T v$ follows from (3.2).

2.

$$\begin{aligned} A^T v_r &= A^T (v - AA^\dagger v) \\ &= V\Sigma U^T (v - UU^T v) \\ &= V\Sigma U^T v - V\Sigma U^T v \\ &= \mathbf{0} \end{aligned}$$

□

Remark 3.3.1 Now it is clear that the least squares solution is nothing but a projection.

3.4 The Underdetermined Case

Consider the problem $\min_x \|y - Ax\|^2$ where $A^T A$ does not have full rank.

Let $r < n$ be the rank of A . Consider the compact SVD $A = U\Sigma V^T$, where $U \in \mathbb{R}^{m \times r}$, $\Sigma \in \mathbb{R}^{r \times r}$, and $V \in \mathbb{R}^{n \times r}$. We have the following observations.

\hat{x} is a global minimizer.

This can be verified by the facts that $f(x) = \|y - Ax\|^2$ is convex in x and that

$$\begin{aligned}\nabla f(\hat{x}) &= A^T(Ax - y) = -A^T \text{resid}(y, A) \\ &= \mathbf{0}.\end{aligned}$$

The global minimizer is not unique.

Define the set of global minimizer as

$$\begin{aligned}\mathcal{X} &:= \left\{ x : \|y - Ax\|^2 = \min_z \|y - Az\|^2 \right\} \\ &= \{ x : x = \hat{x} + V_\perp c_\perp, c_\perp \in \mathbb{R}^{n-m} \},\end{aligned}$$

where $V_\perp \in \mathbb{R}^{n \times (n-r)}$ is an orthogonal complement to V , i.e., $[V, V_\perp] \in \mathbb{R}^{n \times n}$ ⁸ is orthonormal.

A quick check: For all $x = \hat{x} + V_\perp c_\perp$, it holds that

$$\begin{aligned}\|y - Ax\|^2 &= \|y - A(\hat{x} + V_\perp c_\perp)\|^2 \\ &= \|y - A\hat{x} + \mathbf{0}\|^2 \\ &= \|y - A\hat{x}\|^2.\end{aligned}$$

The least squares solution $\hat{x} = A^\dagger y$ has the least length among all global minimizers.

The least squares solution $\hat{x} = A^\dagger y$ satisfies

$$\hat{x} = \arg \min \left\{ \|x\| : \|y - Ax\|^2 = \min_z \|y - Az\|^2 \right\}.$$

Proof. SVD provides a way to decompose any $x \in \mathbb{R}^n$ into two parts:

$$\begin{aligned}x &= Ix = [V, V_\perp] \begin{bmatrix} V^T \\ V_\perp^T \end{bmatrix} x \\ &= VV^T x + V_\perp V_\perp^T x.\end{aligned}$$

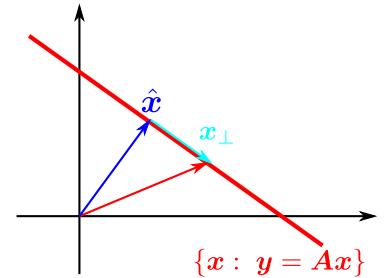


Figure 3.3: Solution set \mathcal{X}

8: $[V, V_\perp]$ can be calculated from a full SVD of A .

Consider an arbitrary $x \in \mathcal{X}$, it holds that

$$\begin{aligned}\|x\|^2 &= \|\hat{x} + V_{\perp}c_{\perp}\|^2 = \|Vc + V_{\perp}c_{\perp}\|^2 \\ &= \|Vc\|^2 + \langle Vc, V_{\perp}c_{\perp} \rangle + \|V_{\perp}c_{\perp}\|^2 \\ &= \|Vc\|^2 + \|V_{\perp}c_{\perp}\|^2 \\ &\geq \|Vc\|^2 = \|\hat{x}\|^2.\end{aligned}$$

□

3.5 Numerical Stability

Consider the model

$$y = Ax_0 + w,$$

where x_0 is the ground truth signal and $w \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$.

Assume that A is tall and has full rank. Consider the compact SVD $A = U\Sigma V^T$.⁹

9: What are the dimensions of U , Σ , and V ?

The least squares solution is given by

$$\begin{aligned}\hat{x} &= A^{\dagger}y = A^{\dagger}(Ax_0 + w) \\ &= (A^T A)^{-1} A^T (Ax_0 + w) \\ &= x_0 + V\Sigma^{-1}U^T w.\end{aligned}$$

Hence,

$$\|\hat{x} - x_0\|_2 = \|V\Sigma^{-1}U^T w\|.$$

The estimation error $\|\hat{x} - x_0\|$ can be then bounded by using the following steps:¹⁰

10: This is actually a coursework problem for you.

1. Show that $\|Vz\|^2 = \|z\|^2, \forall z \in \mathbb{R}^n$.
2. Show that $\forall z \in \mathbb{R}^n$,

$$\sigma_{\max}^{-1}\|z\| \leq \|\Sigma^{-1}z\| \leq \sigma_{\min}^{-1}\|z\|.$$

3. Show that¹¹

$$0 \leq \|U^T w\| \leq \|w\|,$$

11: What is the dimension of $U^T w$? Its distribution?

and

$$E[\|U^T w\|^2] = n\sigma^2.$$

4. Combine the results above, you should be able to derive bounds on $\|\hat{x} - x_0\|$ (in terms of $\|w\|$) and $E[\|\hat{x} - x_0\|^2]$ (in terms of σ^2).

In numerical analysis, the condition number of the matrix A ,¹² defined as

$$\kappa(A) := \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$$

12: See https://en.wikipedia.org/wiki/Condition_number for details.

is an important criterion to measure the stability of the least squares solution.

Sparse Inverse Problems: Introduction

4

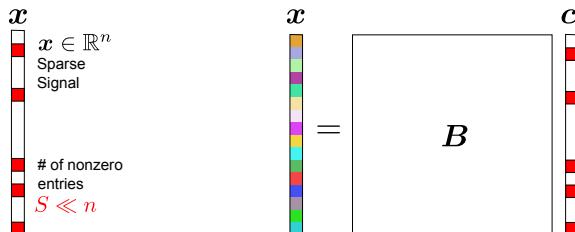
4.1 Problem Statement

Definition 4.1.1 (ℓ_0 -pseudonorm) *The ℓ_0 -pseudonorm counts the total number of nonzero elements of a vector.*

Definition 4.1.2 (Sparse signals) *A sparse signal is the one which contains only a small number of non-zero elements compared to its dimension.*

A signal $x \in \mathbb{R}^n$ is sparse if

- $\|x\|_0 = S \ll n$.
- $x = Bc$ and $\|c\|_0 = S \ll n$.



4.1 Problem Statement	19
4.2 Exhaustive Search	20
4.2.1 The Algorithm	20
4.2.2 Performance Guarantees .	22
4.3 Applications	23

In reality, sparse signals are rare but compressible signals are ubiquitous.

Compressible signals are those that can be approximated by sparse signals.

- Natural pictures are compressible under DCT/Wavelet transform.
- Communication signals are often compressible under Fourier transform.
- In function approximation, it is typically assumed that the unknown function can be well approximated by a few ‘kernel’ functions.

Figure 4.1: Sparse Signals

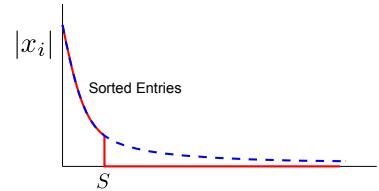


Figure 4.2: Compressible Signals

Sparse Inverse Problem

Given the observations y and the matrix A , try to find a sparse x such that $y \approx Ax$.

In mathematical terms:

$$\hat{x}_{\text{spa}} = \arg \min_{x: y \approx Ax} \|x\|_0 \quad (4.1)$$

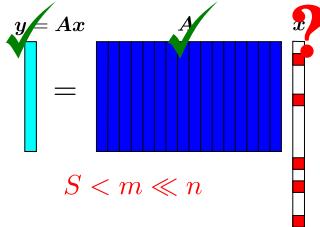


Figure 4.3: Sparse inverse problem

4.2 Exhaustive Search

Challenges

It is very challenging to solve sparse inverse problem (4.1).

A Look at Least Squares Solution:

- In big data applications, it is typically $m \ll n$. Hence there are infinitely many feasible solutions.

Define

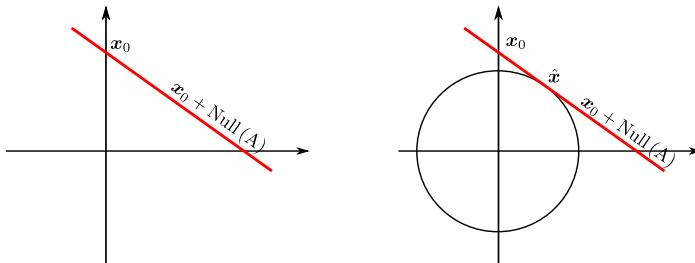
$$\mathcal{X} := \{x : y = Ax\}.$$

When $m < n$, $|\mathcal{X}| = \infty$.

- The least squares solution

$$\hat{x} = A^\dagger y = \arg \min_{x:y=Ax} \|x\|_2$$

is typically *not* sparse.

Figure 4.4: $|\mathcal{X}| = \infty$ and \hat{x} is not sparse.

A Look at ℓ_0 -pseudonorm

Normal optimization techniques cannot be applied to solve (4.1).

- ℓ_0 pseudo-norm is discontinuous¹ and nonconvex².
 - Discontinuity \Rightarrow Gradient based technique cannot be applied.
 - Nonconvexity \Rightarrow Global optimality is not guaranteed.

1: Draw contours of ℓ_0 -norm in 1-D and 2-D space. Convince yourself that ℓ_0 -norm is not continuous.

2: We will discuss convexity in detail later.

4.2.1 The Algorithm

Definition 4.2.1 Let $x \in \mathbb{R}^n$. Its support set is defined as

$$\text{supp}(x) = \{i : x_i \neq 0\}.$$

Example 4.2.1

$$\begin{aligned} \mathbf{x} &= [0.1, 0, 0, -1, 0]^T \\ \Rightarrow \mathcal{I} &= \{1, 4\}. \end{aligned}$$

Definition 4.2.2 Let $\mathcal{I} \subset \{1, 2, \dots, n\}$ be an index set. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{x} \in \mathbb{R}^n$.

- $\mathbf{A}_{\mathcal{I}}$: a matrix formed by columns of \mathbf{A} indexed by \mathcal{I} .
- $\mathbf{x}_{\mathcal{I}}$: a vector formed by entries of \mathbf{x} indexed by \mathcal{I} .

Example 4.2.2 Suppose that $\mathbf{y} = \mathbf{Ax}$ where \mathbf{x} is sparse. Let $\mathcal{I} = \text{supp}(\mathbf{x})$. Then $\mathbf{y} = \mathbf{Ax} = \mathbf{A}_{\mathcal{I}}\mathbf{x}_{\mathcal{I}}$.

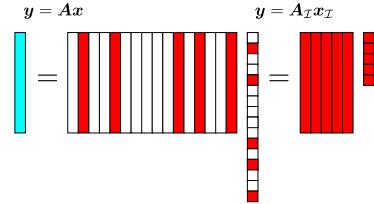


Figure 4.5: $y = Ax = A_{\mathcal{I}}x_{\mathcal{I}}$

Exhaustive Search Algorithm

1. For $s = 1, 2, \dots$
 - a) Try all $\mathcal{I} \subset [n] \triangleq \{1, 2, \dots, n\}$ s.t. $|\mathcal{I}| = s$.
 - b) If $\mathbf{y} = \mathbf{A}_{\mathcal{I}}\hat{\mathbf{x}}_{\mathcal{I}}$, then terminate the search.
Otherwise, continue.
2. Finalization: Set $\mathbf{x}_{\mathcal{I}} = \mathbf{A}_{\mathcal{I}}^{\dagger}\mathbf{y}$ and $\mathbf{x}_{\mathcal{I}^c} = 0$.

Exhaustive search is not practical.

- Exhaustive search terminates when $S^{\#} = \|\mathbf{x}\|_0$.
- The computational complexity is approximately

$$\sum_{s=1}^{S^{\#}} \binom{n}{s} \geq \binom{n}{S^{\#}} = \frac{n!}{S^{\#}!(n-S^{\#})!}.$$

- Use Stirling approximation: $n! \approx \sqrt{2\pi n} n^{n+1/2} e^{-n}$.
It holds that the complexity is

$$O(n^{S^{\#}}).$$

Exhaustive search is not practical for large n .

Example 4.2.3 (Finding Sparse Solutions) Given

1.

$$\begin{bmatrix} 2 \\ 2 \\ -2 \\ -2 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 & 2 & 0 & 0 & 0 \\ 1 & -1 & 1 & -1 & 0 & 2 & 0 & 0 \\ 1 & 1 & -1 & -1 & 0 & 0 & 2 & 0 \\ 1 & -1 & -1 & 1 & 0 & 0 & 0 & 2 \end{bmatrix} \mathbf{x}.$$

2.

$$\begin{bmatrix} 2 \\ 2 \\ 0 \\ 0 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 & 2 & 0 & 0 & 0 \\ 1 & -1 & 1 & -1 & 0 & 2 & 0 & 0 \\ 1 & 1 & -1 & -1 & 0 & 0 & 2 & 0 \\ 1 & -1 & -1 & 1 & 0 & 0 & 0 & 2 \end{bmatrix} \mathbf{x}.$$

3.

$$\begin{bmatrix} 2 \\ 3 \\ 0 \\ 0 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 & 2 & 0 & 0 & 0 \\ 1 & -1 & 1 & -1 & 0 & 2 & 0 & 0 \\ 1 & 1 & -1 & -1 & 0 & 0 & 2 & 0 \\ 1 & -1 & -1 & 1 & 0 & 0 & 0 & 2 \end{bmatrix} \mathbf{x}.$$

- What would the sparsest solution be?
- Is the sparsest solution unique?

4.2.2 Performance Guarantees

The question:

- Let $\mathbf{y} = A\mathbf{x}_0$ where $\|\mathbf{x}_0\|_0 = S < m$.
- Let $\hat{\mathbf{x}}$ be a sparsest solution of $\mathbf{y} = A\mathbf{x}$.

Under which conditions $\hat{\mathbf{x}} = \mathbf{x}_0$?

Assumption:

Every m -column sub-matrix of A is of full rank.

- Let $A \in \mathbb{R}^{m \times n}$ ($m < n$) be generated from standard Gaussian distribution.
- The generated A satisfies the assumption with probability one.

Theorem 4.2.1 (Worst Case Performance Guarantees) *If $m > 2S$, then $\hat{\mathbf{x}} = \mathbf{x}_0$ for all S -sparse signals \mathbf{x}_0 .*

Proof. Suppose that there exists an \mathbf{x}_1 such that

$$\mathbf{x}_1 \neq \mathbf{x}_0, \|\mathbf{x}_1\|_0 \leq S, \mathbf{y} = A\mathbf{x}_1.$$

Then

$$\mathbf{0} = A(\mathbf{x}_0 - \mathbf{x}_1).$$

Note that

$$\|\mathbf{x}_0 - \mathbf{x}_1\|_0 \leq 2S < m$$

This contradicts that every m -columns of A are linearly independent. \square

Theorem 4.2.2 (Average Case Performance Guarantees) *If $m > S$, then $\hat{\mathbf{x}} = \mathbf{x}_0$ for nearly all S -sparse signals \mathbf{x}_0 .*

Proof. Suppose that there exists an \mathbf{x}_1 such that

$$\mathbf{x}_1 \neq \mathbf{x}_0, \|\mathbf{x}_1\|_0 \leq S, \mathbf{y} = A\mathbf{x}_1.$$

Let $\mathcal{T}_0 = \text{supp}(\mathbf{x}_0)$ and $\mathcal{T}_1 = \text{supp}(\mathbf{x}_1)$.

1. That $\mathbf{x}_1 \neq \mathbf{x}_0$ implies that $\mathcal{T}_0 \neq \mathcal{T}_1$ and $\text{span}(A_{\mathcal{T}_0}) \neq \text{span}(A_{\mathcal{T}_1})$. Otherwise, $\mathbf{x}_0 = A_{\mathcal{T}_0}^\dagger \mathbf{y} = A_{\mathcal{T}_1}^\dagger \mathbf{y} = \mathbf{x}_1$.

2. $y \in \text{span}(A_{\mathcal{T}_0}) \cap \text{span}(A_{\mathcal{T}_1}) =: \mathcal{A}_{\cap}$.

This follows from that

$$\begin{aligned} y &= Ax_0 = A_{\mathcal{T}_0}x_{0,\mathcal{T}_0} \\ &= Ax_1 = A_{\mathcal{T}_1}x_{1,\mathcal{T}_1}. \end{aligned}$$

3. $\dim(\mathcal{A}_{\cap}) < \dim(\text{span}(A_{\mathcal{T}_0})) = \dim(\text{span}(A_{\mathcal{T}_1}))$.

This follows from that $\text{span}(A_{\mathcal{T}_0}) \neq \text{span}(A_{\mathcal{T}_1})$.

4. For any x_0 generated from a continuous distribution,

$$\Pr(y \in \text{span}(A_{\mathcal{T}_0}) \cap \text{span}(A_{\mathcal{T}_1})) = 0.$$

□

4.3 Applications

Compressed Sensing

Challenges: In many applications, the cost (in physical space, time, spectrum, or operations) of sensing can be very high.

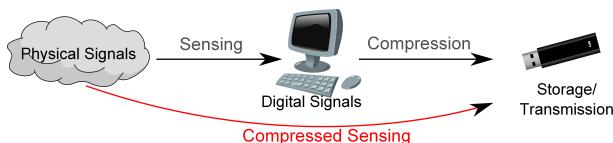
- Magnetic Resonance Imaging (MRI):



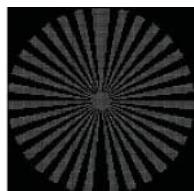
- Infrared sensing:

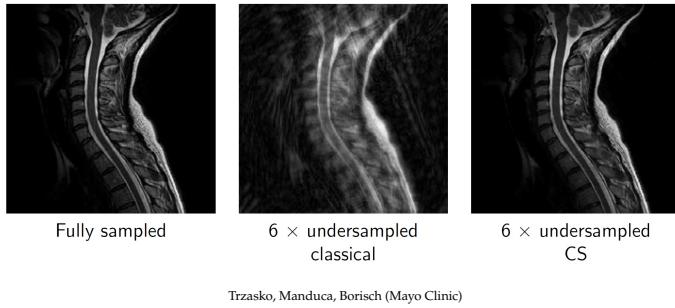


Compressed sensing: A paradigm shift to reduce the cost/time of sensing.



An example in MRI:

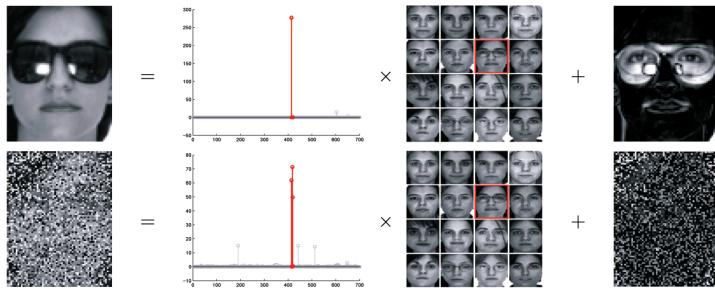




Trzasko, Manduca, Borisch (Mayo Clinic)

$$\min_{\mathbf{x}} \|\mathbf{x}\|_0 \text{ subject to } \mathbf{y} \approx \mathbf{F}_{\Omega,:} \mathbf{x}.$$

Robust Face Recognition



[Wright et al., 2009]

The setup:

- ▶ A set of training samples $\{\phi_i, l_i\}$
 - $\phi_i \in \mathbb{R}^m$ is the vector representation of the images.
 - $l_i \in \{1, 2, \dots, C\}$ label for the C subjects.
- ▶ Test sample \mathbf{y}
- ▶ Assumption:
 - For simplicity, assume a good face alignment.

Face recognition via sparse regression:

- ▶ Images of the same subject i form a low-dimensional linear subspace in \mathbb{R}^m .

$$\mathbf{y} \approx \sum_{\{j|l_j=i\}} \phi_j c_j =: \Phi_i \mathbf{c}_i,$$

and

$$\mathbf{y} \approx [\Phi_1, \Phi_2, \dots, \Phi_C] \mathbf{c} = \Phi \mathbf{c} \in \mathbb{R}^m$$

where $\mathbf{c} = [\dots, \mathbf{0}^T, \mathbf{c}_i^T, \mathbf{0}^T, \dots]^T$.

- ▶ Face recognition: the coefficient vector \mathbf{c} is group sparse.

$$\min_{\mathbf{c}} \|\mathbf{c}\|_{2,0} \text{ s.t. } \mathbf{y} \approx \Phi \mathbf{c}.$$

Robust face recognition:

- Model corruption and occlusion as

$$\mathbf{y} \approx \Phi \mathbf{c} + \mathbf{e},$$

where \mathbf{e} is an unknown error vector whose entries can be very large.

- Both coefficient vector \mathbf{c} and noise vector \mathbf{e} are sparse.

$$\min_{\mathbf{c}, \mathbf{e}} \|\mathbf{c}\|_{2,0} + \|\mathbf{e}\|_0 \text{ s.t. } \mathbf{y} \approx \Phi \mathbf{c} + \mathbf{e}.$$

Sparse Inverse Problems: L0 Methods

5

5.1 Some Easy Problems

L0 Indicator Function

Definition 5.1.1 Define indicator function

$$\delta(\cdot) = \begin{cases} 0 & \text{if the statement is true,} \\ \infty & \text{otherwise.} \end{cases}$$

Example 5.1.1 Find the value of the indicator function $\delta(\|x\|_0 \leq 1)$ for

- $x = [0, 0, 0]^\top$,
- $x = [1, 0, 0]^\top$, and
- $x = [1, 2, 0]^\top$.

We consider some easy ℓ_0 minimization problems, which help design efficient algorithms for general sparse inverse problems.

Example 5.1.2 Consider the optimization problem

$$\min_x \delta(\|x\|_0 \leq 1) + \frac{1}{2} \|x - z\|^2.$$

Find the optimal solution x^* and the corresponding optimal value of the objective function for the following choices of z .

- $z = [0, 0, 0]^\top$,
- $z = [1, 2, 3]^\top$, and
- $z = [1, 2, 2]^\top$.

Definition 5.1.2 (Moreau Envelope and Proximal Function) For a given function f and a given constant $\gamma > 0$, its Moreau envelope is defined as

$$\mathcal{M}_{\gamma f}(z) := \min_x f(x) + \frac{1}{2\gamma} \|x - z\|_2^2,$$

and the proximal operator is given by¹

$$x^* = \text{Prox}_{\gamma f}(z) := \arg \min_x f(x) + \frac{1}{2\gamma} \|x - z\|_2^2.$$

Example 5.1.3 Let $f(x) = \delta(\|x\|_0 \leq S)$ for some $S \in \mathbb{Z}_+$. Find the corresponding Moreau envelope $\mathcal{M}_{\gamma f}(z)$ and proximal operator $\text{Prox}_{\gamma f}(z)$ for some $z \in \mathbb{Z}^n$. Discuss the continuity of them. Discuss how the parameter γ affects the solutions.²

5.1	Some Easy Problems	26
5.2	OMP	27
5.2.1	Algorithm	27
5.2.2	Performance Guarantees	28
5.3	SP and Similar Algorithms	30
5.3.1	Algorithms	30
5.3.2	Performance Guarantees	32
5.4	Majorization Minimization (MM)	37
5.5	Appendix	39

1: Moreau envelope and proximal operator are very important in modern optimization theory and algorithms.

2: For any given vector $z \in \mathbb{R}^n$, let $|z_{i_1}| \geq |z_{i_2}| \geq \dots \geq |z_{i_n}|$ be the sequence of magnitude-sorted elements. The closed forms of Moreau envelope and proximal operator can be written using this sequence.

L0 Function

Example 5.1.4 Let $f(x) = \|x\|_0$. Find the corresponding Moreau envelope $\mathcal{M}_{\gamma f}(z)$ and proximal operator $\text{Prox}_{\gamma f}(z)$ for some $z \in \mathbb{R}^n$. Discuss the continuity of them. Discuss how the parameter γ affects the solutions.

If you have difficulties in solving this problem, the following simplified case may give you clues.

Consider the scalar case, i.e., $x \in \mathbb{R}$. One has $f(x) = 0$ if $x = 0$ and $f(x) = 1$ if $x \neq 0$. Let $\gamma = 1$. Then

$$\mathcal{M}_{\|\cdot\|_0}(z) = \min_x \|x\|_0 + \frac{1}{2}(x - z)^2.$$

1. Show that

$$\mathcal{M}_{\|\cdot\|_0}(z) \leq \|z\|_0 + \frac{1}{2}(z - z)^2 \leq 1.$$

2. For $z = 0, 1, 2$, find $\mathcal{M}_{\|\cdot\|_0}(z)$ and $x^* = \text{Prox}_{\|\cdot\|_0}(z)$ respectively.

Towards General Problems

Example 5.1.5 Let $A \in \mathbb{R}^{m \times n}$ be a matrix of which columns are normalized, i.e., $\|A_{:,n}\|_2 = 1$. Consider the optimization problem

$$\min_x \delta(\|x\|_0 \leq 1) + \frac{1}{2}\|y - Ax\|^2.$$

Find the optimal solution x^* and the corresponding optimal value of the objective function.³

3: The results in Example ?? are useful to justify the solutions.

While the above problem can be easily solved, the problem

$$\min_x \delta(\|x\|_0 \leq S) + \frac{1}{2}\|y - Ax\|^2$$

for $S > 1$ is much more difficult to solve.

5.2 OMP

5.2.1 Algorithm

The intuition: We solve the difficult problem

$$\min_x \delta(\|x\|_0 \leq S) + \frac{1}{2}\|y - Ax\|^2$$

by solving a sequence of easier problems

$$\min_x \delta(\|x\|_0 \leq 1) + \frac{1}{2}\|y - Ax\|^2.$$

Algorithm 1: Orthogonal Matching Pursuit (OMP)

Require: y, A, S

Ensure: \hat{x} such that $\|\hat{x}\|_0 \leq S$

- ```

1: Initialization: $\ell = 1$, $\hat{x} = 0$, $\mathcal{T}^\ell = \phi$, and $y_r = y$.
2: while $\ell \leq S$ do
3: $i_\ell = \arg \max_j |\langle a_j, y_r \rangle|$
4: $\mathcal{T}^\ell = \mathcal{T}^{\ell-1} \cup \{i_\ell\}$ {Add one index}
5: $\hat{x}_{\mathcal{T}^\ell} = A_{\mathcal{T}^\ell}^\dagger y$.
6: $y_r = y - A\hat{x}$.
7: $\ell = \ell + 1$.
8: end while

```

### 5.2.2 Performance Guarantees

**Definition 5.2.1** (Mutual coherence) The mutual coherence of a matrix  $A \in \mathbb{R}^{m \times n}$ , denoted by  $\mu(A)$ , is the maximal correlation (in magnitude) between two (normalized) columns.

$$\mu(A) = \max_{i \neq j} \frac{|\langle a_i, a_j \rangle|}{\|a_i\|_2 \|a_j\|_2}.$$

When  $\|\mathbf{a}_i\|_2 = 1$ ,  $\forall i \in [n]$ , then  $\mu(\mathbf{A}) = \max_{i \neq j} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|$ .

**Theorem 5.2.1** Suppose that  $A$  satisfies that

$$\mu < \frac{1}{2S}.$$

Then the OMP algorithm is guaranteed to exactly recover all  $S$ -sparse  $\mathbf{x}$ . That is,  $\hat{\mathbf{x}} = \mathbf{x}_0$ .

The key for the proof: To show  $\hat{x} = x_0$ :

- Want to show that  $\text{supp}(\hat{x}) = \text{supp}(x_0)$ .
  - Or show that at the  $\ell$ -th iteration of OMP, the chosen index  $i_\ell \in \mathcal{T}_0 := \text{supp}(x_0)$ .

The proof needs Cauchy-Schwartz Inequality:<sup>4</sup>

4: Also see Theorem 5.5.2.

$$\|x\|_1 = \sum_{i=1}^n |x_i| \leq \sqrt{n} \|x\|_2.$$

*Proof.* The proof is divided into two steps. The first step concerns the first iteration and proves that  $i_1 := \arg \max_i |\langle a_i, y \rangle| \in \mathcal{T}_0$ . In the second step, a very similar technique is used to show that  $i_\ell \in \mathcal{T}_0$ .

Step 1:  $i_1 \in \mathcal{T}_0$ .

It holds that  $\forall i$ ,

$$\begin{aligned} |\langle \mathbf{a}_i, \mathbf{y} \rangle| &= \left| \left\langle \mathbf{a}_i, \sum_{j \in \mathcal{T}_0} \mathbf{a}_j x_{0,j} \right\rangle \right| \\ &= \left| \sum_{j \in \mathcal{T}_0} x_{0,j} \langle \mathbf{a}_i, \mathbf{a}_j \rangle \right|. \end{aligned}$$

1. For all  $i \notin \mathcal{T}_0$ :

$$\begin{aligned} |\langle \mathbf{a}_i, \mathbf{y} \rangle| &= \left| \sum_{j \in \mathcal{T}_0} x_{0,j} \langle \mathbf{a}_i, \mathbf{a}_j \rangle \right| \leq \sum_{j \in \mathcal{T}_0} |x_{0,j}| |\langle \mathbf{a}_i, \mathbf{a}_j \rangle| \\ &\leq \mu \sum_{j \in \mathcal{T}_0} |x_{0,j}| \stackrel{(a)}{\leq} \mu \sqrt{S} \|\mathbf{x}\|_2 \end{aligned}$$

where (a) follows from Cauchy–Schwartz Inequality (Theorem 5.5.2). Hence,

$$\max_{i \notin \mathcal{T}_0} |\langle \mathbf{a}_i, \mathbf{y} \rangle| \leq \mu \sqrt{S} \|\mathbf{x}\|_2. \quad (5.1)$$

2. For all  $i \in \mathcal{T}_0$ :

$$\begin{aligned} |\langle \mathbf{a}_i, \mathbf{y} \rangle| &= \left| \sum_{j \in \mathcal{T}_0} x_{0,j} \langle \mathbf{a}_i, \mathbf{a}_j \rangle \right| \geq |\mathbf{x}_{0,i} \langle \mathbf{a}_i, \mathbf{a}_i \rangle| - \left| \sum_{j \neq i} x_{0,j} \langle \mathbf{a}_i, \mathbf{a}_j \rangle \right| \\ &\geq |x_{0,i}| - \mu \sum_{j \neq i} |x_{0,j}| \stackrel{(a)}{\geq} |x_{0,i}| - \mu \sqrt{S} \|\mathbf{x}\|_2, \end{aligned}$$

where (a) follows from Cauchy–Schwartz Inequality.

As a result,

$$\max_{i \in \mathcal{T}_0} |\langle \mathbf{a}_i, \mathbf{y} \rangle| \geq \frac{1}{\sqrt{S}} \|\mathbf{x}\|_2 - \mu \sqrt{S} \|\mathbf{x}\|_2, \quad (5.2)$$

where we have used the fact that

$$\frac{1}{\sqrt{S}} \|\mathbf{x}\|_2 = \frac{(\sum x_i^2)^{\frac{1}{2}}}{\sqrt{S}} \leq \frac{\left( \sum (\max_i |x_i|)^2 \right)^{\frac{1}{2}}}{\sqrt{S}} = \max_{i \in \mathcal{T}_0} |x_i|.$$

3. Now suppose that  $\mu < \frac{1}{2S}$  (the assumption in Theorem 5.2.1). Then

$$\frac{1}{\sqrt{S}} \|\mathbf{x}\|_2 > 2\mu \sqrt{S} \|\mathbf{x}\|_2,$$

Or equivalently,

$$\max_{i \in \mathcal{T}_0} |\langle \mathbf{a}_i, \mathbf{y} \rangle| \geq \frac{1}{\sqrt{S}} \|\mathbf{x}\|_2 - \mu \sqrt{S} \|\mathbf{x}\|_2 > \mu \sqrt{S} \|\mathbf{x}\|_2 \geq \max_{i \notin \mathcal{T}_0} |\langle \mathbf{a}_i, \mathbf{y} \rangle|.$$

One concludes that

$$i_1 \in \mathcal{T}_0.$$

**Step 2:**  $i_\ell \in \mathcal{T}_0$ . (Mathematical induction)

Let  $i_1, \dots, i_{\ell-1}$  be the indices chosen in the first  $\ell - 1$  iterations. Let  $\mathcal{T}^{\ell-1} = \{i_1, \dots, i_{\ell-1}\}$ . Assume that  $\mathcal{T}^{\ell-1} \subset \mathcal{T}_0$ .

Then

$$\mathbf{y}_r = \mathbf{y} - \mathbf{A}_{\mathcal{T}^{\ell-1}} \mathbf{A}_{\mathcal{T}^{\ell-1}}^\dagger \mathbf{y} = \mathbf{y} - \mathbf{A}_{\mathcal{T}^{\ell-1}} \hat{\mathbf{x}}_{\ell-1} \in \text{span}(\mathbf{A}_{\mathcal{T}_0}).$$

Or

$$\mathbf{y}_r = \mathbf{A}_{\mathcal{T}_0} \tilde{\mathbf{v}}_{\mathcal{T}_0}.$$

for some  $\tilde{\mathbf{v}}_{\mathcal{T}_0}$ .

Use the same arguments as before,  $i_\ell \in \mathcal{T}_0$ .

**Final Step:** At the same time,  $\mathbf{A}_{\mathcal{T}^{\ell-1}}^\top \mathbf{y}_r = \mathbf{0}$  and hence  $i_\ell \notin \mathcal{T}^{\ell-1}$ . So  $|\mathcal{T}^\ell| = \ell$ .

OMP algorithm needs  $S$  iterations to recover  $S$ -sparse signals.  $\square$

## 5.3 SP and Similar Algorithms

**The intuition:** If

- ▶ columns of  $\mathbf{A}$  are normalised,
- ▶ columns of  $\mathbf{A}$  are near orthogonal,

then  $\mathbf{A}^\top \mathbf{y} = \mathbf{A}^\top \mathbf{A} \mathbf{x} \approx \mathbf{I} \mathbf{x}$ . We then try to solve

$$\min_{\mathbf{x}} \delta(\|\mathbf{x}\|_0 \leq S) + \frac{1}{2} \|\mathbf{x} - \mathbf{A}^\top \mathbf{y}\|^2$$

### Key Differences Between OMP and Other Greedy Algorithms

- ▶ OMP
  - One index is added per iteration.
  - Analysis is based on near orthogonality of column pairs.
- ▶ SP, CoSaMP, IHT
  - Multiple indices are updated per iteration.
  - Analysis is based on near orthogonality of sub-matrices.

#### 5.3.1 Algorithms

Hard thresholding function  $H_S(\mathbf{a})$ :

$$H_S(\mathbf{a}) = \text{Prox}_{\delta(\|\cdot\|_0 \leq S)}.$$

Set all but the largest (in magnitude)  $S$  elements of  $\mathbf{a}$  to zero.

**Example 5.3.1**  $\mathbf{a} = [3, -4, 1]^T$ .

- ▶  $H_1(\mathbf{a}) = [0, -4, 0]^T$  and  $\text{supp}(H_1(\mathbf{a})) = \{2\}$ .

### Subspace pursuit (SP) algorithm

---

**Algorithm 2:** Subspace Pursuit (SP)

---

**Require:**  $y, A, S$

**Ensure:**  $\hat{x}$  such that  $\|\hat{x}\|_0 \leq S$

- 1: **Initialization:**  $\mathcal{T}^0 = \text{supp}(H_S(A^T y))$ , and  $y_r = \text{resid}(y, A_{\mathcal{T}^0})$ .
  - 2: **while**  $\ell = 1, 2, \dots$  until exit criteria are satisfied **do**
  - 3:    $\tilde{\mathcal{T}}^\ell = \mathcal{T}^{\ell-1} \cup \text{supp}(H_S(A^T y_r))$ . {Expand support}
  - 4:   Let  $b_{\tilde{\mathcal{T}}^\ell} = A_{\tilde{\mathcal{T}}^\ell}^T y$  and  $b_{(\tilde{\mathcal{T}}^\ell)^c} = 0$ . {Estimate a  $2S$ -sparse signal}
  - 5:   Set  $\mathcal{T}^\ell = \text{supp}(H_S(b))$ . {Shrink support}
  - 6:   Let  $x_{\mathcal{T}^\ell}^\ell = A_{\mathcal{T}^\ell}^T y$  and  $x_{(\mathcal{T}^\ell)^c}^\ell = 0$ . {Estimate  $S$ -sparse signal}
  - 7:   Let  $y_r = y - Ax^\ell$ . {Compute estimation error}
  - 8: **end while**
- 

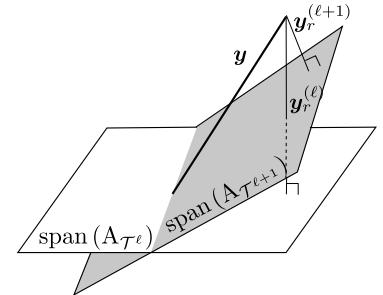


Figure 5.1: A geometric interpretation

### Compressive Sampling Matching Pursuit (CoSaMP) Algorithm

---

**Algorithm 3:** Compressive Sampling Matching Pursuit (CoSaMP)

---

**Require:**  $y, A, S$

**Ensure:**  $\hat{x}$  such that  $\|\hat{x}\|_0 \leq S$

- 1: **Initialization:**  $x^0 = 0$  and  $y_r = y$ .
  - 2: **while**  $\ell = 1, 2, \dots$  until exit criteria are satisfied **do**
  - 3:    $\tilde{\mathcal{T}}^\ell = \mathcal{T}^{\ell-1} \cup \text{supp}(H_{2S}(A^T y_r))$ . {Expand support}
  - 4:   Let  $b_{\tilde{\mathcal{T}}^\ell} = A_{\tilde{\mathcal{T}}^\ell}^T y$  and  $b_{(\tilde{\mathcal{T}}^\ell)^c} = 0$ . {Estimate a  $3S$ -sparse signal}
  - 5:    $x^\ell = H_S(b)$ . ( $\mathcal{T}^\ell = \text{supp}(H_S(b))$ ). {Shrink support}
  - 6:    $y_r = y - Ax^\ell$ . {Compute estimation error}
  - 7: **end while**
- 

### Iterative Hard Thresholding (IHT) Algorithm

---

**Algorithm 4:** Iterative Hard Thresholding (IHT)

---

**Require:**  $y, A, S$

**Ensure:**  $\hat{x}$  such that  $\|\hat{x}\|_0 \leq S$

- 1: **Initialization:**  $x^0 = 0$ .
- 2: **while**  $\ell = 1, 2, \dots$  until exit criteria are satisfied **do**
- 3:

$$x^\ell = H_S(x^{\ell-1} + A^T(y - Ax^{\ell-1}))$$

- 4: **end while**
- 

A more general form: for some step size  $\tau > 0$ ,

$$x^\ell = H_S(x^{\ell-1} + \tau A^T(y - Ax^{\ell-1})).$$

**Remark 5.3.1 (History)**

- ▶ MP: Friedman and Stuetzle, 1981; Mallat and Zhang, 1993; Qian and Chen, 1994.
- ▶ OMP: Chen, et al., 1989; Pati, et al., 1993; Davis, et al., 1994.

Analysed by Tropp, 2004.

- SP: Dai and Milenkovic, 2009. (Online available 06/03/2008)
- CoSaMP: Needell and Tropp, 2009. (Online available 17/03/2008)
- IHT: Blumensath and Davies, 2009. (Online available 05/05/2008)

### 5.3.2 Performance Guarantees

**Comparison:** Assume that the matrix  $A$  is generated as a random Gaussian matrix.

|                   | # of measurements             | # of iterations                     |
|-------------------|-------------------------------|-------------------------------------|
| Exhaustive Search | $2S + 1$                      | $\binom{n}{S} = O(n^S)$             |
| OMP               | $O(S^2 \log n)$               | $S$                                 |
| SP, CoSaMP, IHT   | $O(S \cdot \log \frac{n}{S})$ | Typically $O(\log S)$ , at most $S$ |

#### Restricted Isometry Property (RIP)

Intuitively, RIP means near orthogonality of submatrices.

**Definition 5.3.1** (Restricted isometry property (RIP) and restricted isometry constant (RIC)) A matrix  $A \in \mathbb{R}^{m \times n}$  is said to satisfy the **RIP** with parameters  $(K, \delta)$ , if for all  $\mathcal{T} \subset [n]$  such that  $|\mathcal{T}| \leq K$  and for all  $\mathbf{q} \in \mathbb{R}^{|\mathcal{T}|}$ , it holds that

$$(1 - \delta) \|\mathbf{q}\|_2^2 \leq \|A_{\mathcal{T}} \mathbf{q}\|_2^2 \leq (1 + \delta) \|\mathbf{q}\|_2^2.$$

The **RIC**  $\delta_K$  is defined as the smallest constant  $\delta$  for which the  $K$ -RIP holds, i.e.,

$$\delta_K = \inf \left\{ \delta : (1 - \delta) \|\mathbf{q}\|_2^2 \leq \|A_{\mathcal{T}} \mathbf{q}\|_2^2 \leq (1 + \delta) \|\mathbf{q}\|_2^2, \forall |\mathcal{T}| \leq K, \forall \mathbf{q} \in \mathbb{R}^{|\mathcal{T}|} \right\}.$$

**Remark 5.3.2** (RIP, Eigenvalues and Singular Values) Let  $B \in \mathbb{R}^{m \times K}$  be a tall matrix, i.e.  $m \geq K$ . Then the following statements are equivalent.

- For all  $\mathbf{q} \in \mathbb{R}^K$ ,

$$(1 - \delta_K) \|\mathbf{q}\|_2^2 \leq \|B \mathbf{q}\|_2^2 \leq (1 + \delta_K) \|\mathbf{q}\|_2^2.$$

- $1 - \delta_K \leq \lambda_{\min}(B^T B) \leq \lambda_{\max}(B^T B) \leq 1 + \delta_K.$
- $\sqrt{1 - \delta_K} \leq \sigma_{\min}(B) \leq \sigma_{\max}(B) \leq \sqrt{1 + \delta_K}.$

*Proof of the remark.*

$$\begin{aligned}\|Bq\|_2^2 &= \|U\Sigma V^T q\|_2^2 \\ &= q^T V \Sigma U^T U \Sigma V^T q \\ &= q^T V \Sigma^2 V^T q \\ &= q'^T \Sigma^2 q' \\ &= \sum_{i=1}^K \sigma_i^2 q_i'^2,\end{aligned}$$

where  $q' := V^T q$  satisfies  $\|q'\|_2^2 = q^T V V^T q = \|q\|_2^2$ .  $\square$

**Theorem 5.3.1** (Monotonicity of RIC) *It holds that  $\delta_1 \leq \delta_2 \leq \delta_3 \leq \dots (\delta_K \leq \delta_{K'} \text{ for all } K \leq K')$ .*

*Proof.* Let  $\mathcal{Q}_K = \{q \in \mathbb{R}^n : \|q\|_0 \leq K, \|q\|_2 \leq 1\}$ . It is clear that  $\mathcal{Q}_K \subset \mathcal{Q}_{K'}$  if  $K \leq K'$ .

Then it holds that

$$\delta_K := \sup_{q \in \mathcal{Q}_K} (\|Aq\|_2^2 - 1) \leq \sup_{q \in \mathcal{Q}_{K'}} (\|Aq\|_2^2 - 1) =: \delta_{K'}.$$

$\square$

**Theorem 5.3.2** (Near Orthogonality of Submatrices) *Let  $\mathcal{I}, \mathcal{J} \subset [n]$  be two disjoint sets, i.e.,  $\mathcal{I} \cap \mathcal{J} = \emptyset$ . For all  $a \in \mathbb{R}^{|\mathcal{I}|}$  and  $b \in \mathbb{R}^{|\mathcal{J}|}$ ,*

$$|\langle A_{\mathcal{I}} a, A_{\mathcal{J}} b \rangle| \leq \delta_{|\mathcal{I}|+|\mathcal{J}|} \|a\|_2 \|b\|_2, \quad (5.3)$$

and

$$\|A_{\mathcal{I}}^T A_{\mathcal{J}} b\|_2 \leq \delta_{|\mathcal{I}|+|\mathcal{J}|} \|b\|_2. \quad (5.4)$$

*Proof.* (5.3) obviously holds when either  $a$  or  $b$  is zero. Assume  $a \neq 0$  and  $b \neq 0$ . Define

$$\begin{aligned}a' &= a / \|a\|_2, & b' &= b / \|b\|_2, \\ x' &= A_{\mathcal{I}} a', & y' &= A_{\mathcal{J}} b'.\end{aligned}$$

Then RIP implies that

$$\begin{aligned}2(1 - \delta_{|\mathcal{I}|+|\mathcal{J}|}) &\leq \|x' + y'\|_2^2 = \left\| \begin{bmatrix} A_{\mathcal{I}} & A_{\mathcal{J}} \end{bmatrix} \begin{bmatrix} a' \\ b' \end{bmatrix} \right\|_2^2 \leq 2(1 + \delta_{|\mathcal{I}|+|\mathcal{J}|}), \\ 2(1 - \delta_{|\mathcal{I}|+|\mathcal{J}|}) &\leq \|x' - y'\|_2^2 = \left\| \begin{bmatrix} A_{\mathcal{I}} & A_{\mathcal{J}} \end{bmatrix} \begin{bmatrix} a' \\ -b' \end{bmatrix} \right\|_2^2 \leq 2(1 + \delta_{|\mathcal{I}|+|\mathcal{J}|}),\end{aligned}$$

Hence,

$$\begin{aligned}\langle \mathbf{x}', \mathbf{y}' \rangle &= \frac{\|\mathbf{x}' + \mathbf{y}'\|_2^2 - \|\mathbf{x}' - \mathbf{y}'\|_2^2}{4} \leq \delta_{|\mathcal{I}|+|\mathcal{J}|} \\ -\langle \mathbf{x}', \mathbf{y}' \rangle &= \frac{\|\mathbf{x}' - \mathbf{y}'\|_2^2 - \|\mathbf{x}' + \mathbf{y}'\|_2^2}{4} \leq \delta_{|\mathcal{I}|+|\mathcal{J}|}\end{aligned}$$

Therefore,

$$\frac{|\langle A_{\mathcal{I}} \mathbf{a}, A_{\mathcal{J}} \mathbf{b} \rangle|}{\|\mathbf{a}\|_2 \|\mathbf{b}\|_2} = |\langle \mathbf{x}', \mathbf{y}' \rangle| \leq \delta_{|\mathcal{I}|+|\mathcal{J}|}.$$

From (5.3) to (5.4):

$$\begin{aligned}\|A_{\mathcal{I}}^T A_{\mathcal{J}} \mathbf{b}\|_2 &= \max_{q: \|q\|_2=1} |\langle q, A_{\mathcal{I}}^T A_{\mathcal{J}} \mathbf{b} \rangle| \\ &= \max_{q: \|q\|_2=1} |q^T A_{\mathcal{I}}^T A_{\mathcal{J}} \mathbf{b}| \\ &\leq \max_{q: \|q\|_2=1} \delta_{|\mathcal{I}|+|\mathcal{J}|} \|q\|_2 \|\mathbf{b}\|_2 \\ &= \delta_{|\mathcal{I}|+|\mathcal{J}|} \|\mathbf{b}\|_2\end{aligned}$$

□

Why RIP:

- ▶ Near orthogonality between columns:
  - $|\langle \mathbf{a}_i, \mathbf{a}_j \rangle|$  is small.
- ▶ Near orthogonality between submatrices:
  - $\sigma_{\max}(A_{\mathcal{I}}^T A_{\mathcal{J}}) \leq \delta_{|\mathcal{I}|+|\mathcal{J}|}$  is small.
- ▶ Near-orthogonality of submatrices implies near-orthogonality of columns:
  - $\delta_2$  is small.
- ▶ Near-orthogonality of columns does not mean near-orthogonality of submatrices.
  - Suppose that

$$A_{\mathcal{I}}^T A_{\mathcal{J}} = \begin{bmatrix} \frac{1}{\ell} & \frac{1}{\ell} & \cdots & \frac{1}{\ell} \\ \frac{1}{\ell} & \frac{1}{\ell} & \cdots & \frac{1}{\ell} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\ell} & \frac{1}{\ell} & \cdots & \frac{1}{\ell} \end{bmatrix} \in \mathbb{R}^{\ell \times \ell}$$

Then  $\sigma(A_{\mathcal{I}}^T A_{\mathcal{J}}) = 1, 0, \dots, 0$ .

**Performance Guarantees: IHT**

**Theorem 5.3.3** Suppose that  $A$  satisfies the RIP with  $\delta_{3S} < 1/\sqrt{32}$ , then the  $k^{th}$  iteration of IHT obeys

$$\|\mathbf{x}_0 - \mathbf{x}^k\|_2 \leq 2^{-k} \|\mathbf{x}_0\|_2 + 5 \|\mathbf{w}\|_2.$$

That is, IHT estimates  $\mathbf{x}$  with accuracy

$$\|\mathbf{x}_0 - \mathbf{x}^k\|_2 \leq 6 \|\mathbf{w}\|_2, \text{ if } k > k^* = \left\lceil \log_2 \left( \frac{\|\mathbf{x}_0\|_2}{\|\mathbf{w}\|_2} \right) \right\rceil. \quad (5.5)$$

**Remark 5.3.3** (Optimality) No recovery method can perform fundamentally better.

Suppose that an oracle tells us the support  $\mathcal{T}_0$  of  $\mathbf{x}_0$ . Then

$$\begin{cases} \hat{\mathbf{x}}_{\mathcal{T}_0} = (\mathbf{A}_{\mathcal{T}_0}^\top \mathbf{A}_{\mathcal{T}_0})^{-1} \mathbf{A}_{\mathcal{T}_0}^\top \mathbf{y}, \\ \hat{\mathbf{x}}_{\mathcal{T}_0^c} = \mathbf{0}. \end{cases}$$

It is clear that  $\hat{\mathbf{x}}_{\mathcal{T}_0^c} - \mathbf{x}_{0, \mathcal{T}_0^c} = \mathbf{0}$ . At the same time,

$$\hat{\mathbf{x}}_{\mathcal{T}_0} - \mathbf{x}_{0, \mathcal{T}_0} = (\mathbf{A}_{\mathcal{T}_0}^\top \mathbf{A}_{\mathcal{T}_0})^{-1} \mathbf{A}_{\mathcal{T}_0}^\top \mathbf{w}$$

By the RIP property,

$$\frac{1}{\sqrt{1 + \delta_S}} \|\mathbf{w}\|_2 \leq \|\hat{\mathbf{x}} - \mathbf{x}_0\|_2 \leq \frac{1}{\sqrt{1 - \delta_S}} \|\mathbf{w}\|_2.$$

**Lemma 5.3.4** Let  $\mathbf{r}^k := \mathbf{x}_0 - \mathbf{x}^k$  ( $\mathbf{r}^0 = \mathbf{x}_0$ ). Then

$$\|\mathbf{r}^{k+1}\|_2 \leq \sqrt{8\delta_{3S}} \|\mathbf{r}^k\|_2 + 2\sqrt{1 + \delta_S} \|\mathbf{w}\|_2.$$

In particular, if  $\delta_{3S} < 1/\sqrt{32}$ ,

$$\|\mathbf{r}^{k+1}\|_2 \leq 0.5 \|\mathbf{r}^k\|_2 + 2.17 \|\mathbf{w}\|_2.$$

Back to the main result in Theorem 5.3.3:

$$\begin{aligned} \|\mathbf{r}^k\|_2 &\leq \frac{1}{2} \|\mathbf{r}^{k-1}\|_2 + 2.17 \|\mathbf{w}\|_2 \\ &\leq \frac{1}{4} \|\mathbf{r}^{k-2}\|_2 + 2.17 \left(1 + \frac{1}{2}\right) \|\mathbf{w}\|_2 \\ &\cdots < \frac{1}{2^k} \|\mathbf{r}^0\|_2 + 4.34 \|\mathbf{w}\|_2. \end{aligned}$$

*Proof of Lemma 5.3.4.* Recall that

$$\mathbf{x}^{k+1} = H_S \left( \mathbf{x}^k + \mathbf{A}^T \left( \mathbf{y} - \mathbf{A}\mathbf{x}^k \right) \right).$$

Define

$$\begin{aligned} \mathbf{a}^{k+1} &:= \mathbf{x}^k + \mathbf{A}^T \left( \mathbf{y} - \mathbf{A}\mathbf{x}^k \right) \\ &= \mathbf{x}_0 - \mathbf{x}_0 + \mathbf{x}^k + \mathbf{A}^T \left( \mathbf{A}\mathbf{x}_0 + \mathbf{w} - \mathbf{A}\mathbf{x}^k \right) \\ &= \mathbf{x}_0 + (\mathbf{A}^T \mathbf{A} - \mathbf{I}) (\mathbf{x}_0 - \mathbf{x}^k) + \mathbf{A}^T \mathbf{w} \\ &= \mathbf{x}_0 + (\mathbf{A}^T \mathbf{A} - \mathbf{I}) \mathbf{r}^k + \mathbf{A}^T \mathbf{w}. \end{aligned} \quad (5.6)$$

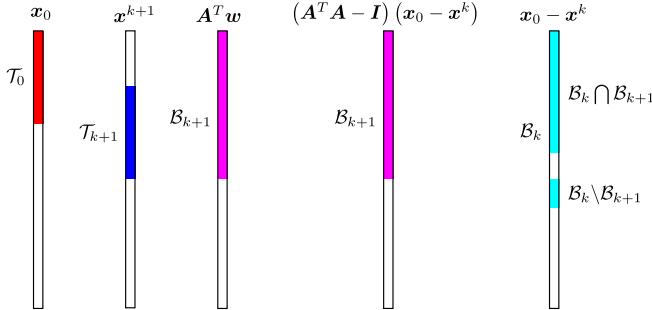
Then

$$\begin{aligned}\mathbf{x}^{k+1} &= H_S \left( \mathbf{x}_0 + (\mathbf{A}^T \mathbf{A} - \mathbf{I}) \mathbf{r}^k + \mathbf{A}^T \mathbf{w} \right) \\ &= \mathbf{x}_0 - \mathbf{r}^{k+1}\end{aligned}$$

We will show that  $\mathbf{r}^{k+1}$  is small.

Let  $\mathcal{T}_0 = \text{supp}(\mathbf{x}_0)$ ,  $\mathcal{T}^k = \text{supp}(\mathbf{x}^k)$ , and  $\mathcal{B}^k = \mathcal{T}_0 \cup \mathcal{T}^k$ . Note that

- $\mathbf{r}^k = \mathbf{x}_0 - \mathbf{x}^k$  is supported on  $\mathcal{B}^k$ .
- $\mathbf{r}^{k+1} = \mathbf{x}_0 - \mathbf{x}^{k+1}$  is supported on  $\mathcal{B}^{k+1}$ .



**Figure 5.2:**  $H_S \left( \mathbf{x}_0 + (\mathbf{A}^T \mathbf{A} - \mathbf{I}) \mathbf{r}^k + \mathbf{A}^T \mathbf{w} \right)$ .

We only need to focus on the set  $\mathcal{B}^{k+1}$ :

$$\begin{aligned}\|\mathbf{r}^{k+1}\|_2 &= \left\| \mathbf{x}_{0, \mathcal{B}^{k+1}} - \mathbf{x}_{\mathcal{B}^{k+1}}^{k+1} \right\|_2 \\ &= \left\| \mathbf{x}_{0, \mathcal{B}^{k+1}} - \mathbf{a}_{\mathcal{B}^{k+1}}^{k+1} + \mathbf{a}_{\mathcal{B}^{k+1}}^{k+1} - \mathbf{x}_{\mathcal{B}^{k+1}}^{k+1} \right\|_2 \\ &\stackrel{(a)}{\leq} \left\| \mathbf{x}_{0, \mathcal{B}^{k+1}} - \mathbf{a}_{\mathcal{B}^{k+1}}^{k+1} \right\|_2 + \left\| \mathbf{a}_{\mathcal{B}^{k+1}}^{k+1} - \mathbf{x}_{\mathcal{B}^{k+1}}^{k+1} \right\|_2 \\ &\stackrel{(b)}{\leq} 2 \left\| \mathbf{x}_{0, \mathcal{B}^{k+1}} - \mathbf{a}_{\mathcal{B}^{k+1}}^{k+1} \right\|_2,\end{aligned}\tag{5.7}$$

where (a) has used triangle inequality, and (b) follows from that  $\mathbf{x}_{\mathcal{B}^{k+1}}^{k+1}$  is the best  $s$ -term approximation to  $\mathbf{a}_{\mathcal{B}^{k+1}}^{k+1}$ .

We bound the noise term  $\mathbf{A}^T \mathbf{w}$ :

$$\|(A^T \mathbf{w})_{\mathcal{B}^{k+1}}\|_2 = \|A_{\mathcal{B}^{k+1}}^T \mathbf{w}\|_2 \leq \sqrt{1 + \delta_{2S}} \|\mathbf{w}\|_2.$$

The other term can be bounded as follows.

$$\begin{aligned}&\left( (\mathbf{I} - \mathbf{A}^T \mathbf{A}) \mathbf{r}^k \right)_{\mathcal{B}^{k+1}} \\ &= \mathbf{r}_{\mathcal{B}^{k+1}}^k - \mathbf{A}_{\mathcal{B}^{k+1}}^T \mathbf{A} \mathbf{r}^k \\ &= \mathbf{r}_{\mathcal{B}^{k+1}}^k - \mathbf{A}_{\mathcal{B}^{k+1}}^T \mathbf{A}_{\mathcal{B}^{k+1}} \mathbf{r}_{\mathcal{B}^{k+1}}^k - \mathbf{A}_{\mathcal{B}^{k+1}}^T \mathbf{A}_{\mathcal{B}^k \setminus \mathcal{B}^{k+1}} \mathbf{r}_{\mathcal{B}^k \setminus \mathcal{B}^{k+1}}^k \\ &= \left( \mathbf{I} - \mathbf{A}_{\mathcal{B}^{k+1}}^T \mathbf{A}_{\mathcal{B}^{k+1}} \right) \mathbf{r}_{\mathcal{B}^{k+1}}^k - \mathbf{A}_{\mathcal{B}^{k+1}}^T \mathbf{A}_{\mathcal{B}^k \setminus \mathcal{B}^{k+1}} \mathbf{r}_{\mathcal{B}^k \setminus \mathcal{B}^{k+1}}^k.\end{aligned}$$

Hence

$$\begin{aligned}&\left\| \left( (\mathbf{I} - \mathbf{A}^T \mathbf{A}) \mathbf{r}^k \right)_{\mathcal{B}^{k+1}} \right\|_2 \\ &\leq \delta_{2S} \left\| \mathbf{r}_{\mathcal{B}^{k+1}}^k \right\|_2 + \delta_{3S} \left\| \mathbf{r}_{\mathcal{B}^k \setminus \mathcal{B}^{k+1}}^k \right\|_2 \\ &\leq \sqrt{2} \delta_{3S} \left\| \mathbf{r}^k \right\|_2,\end{aligned}$$

where

- The 1st term follows from  $|\mathcal{B}^{k+1}| \leq 2S$  and RIP.
- The 2nd term follows from Theorem 5.3.2.
- The last term uses  $\delta_{2S} \leq \delta_{3S}$  (Theorem 5.3.1) and Cauchy-Schwartz Inequality

$$\begin{aligned} & \left\| \mathbf{r}_{\mathcal{B}^{k+1}}^k \right\|_2 + \left\| \mathbf{r}_{\mathcal{B}^k \setminus \mathcal{B}^{k+1}}^k \right\|_2 \\ & \leq \sqrt{2} \left( \left\| \mathbf{r}_{\mathcal{B}^{k+1}}^k \right\|_2^2 + \left\| \mathbf{r}_{\mathcal{B}^k \setminus \mathcal{B}^{k+1}}^k \right\|_2^2 \right)^{1/2} \\ & = \sqrt{2} \left\| \mathbf{r}_{\mathcal{B}^k \cup \mathcal{B}^{k+1}}^k \right\|_2 = \sqrt{2} \left\| \mathbf{r}^k \right\|_2. \end{aligned}$$

Finally,

$$\left\| \mathbf{r}^{k+1} \right\|_2 \leq 2 \left\| \mathbf{x}_{0, \mathcal{B}^{k+1}} - \mathbf{a}_{\mathcal{B}^{k+1}}^{k+1} \right\|_2 \leq \sqrt{8} \delta_{3S} \left\| \mathbf{r}^k \right\|_2 + \sqrt{1 + \delta_{3S}} \left\| \mathbf{w} \right\|_2.$$

□

## 5.4 Majorization Minimization (MM)

Consider an unconstrained optimization problem

$$\min_x f(x).$$

**Majorization minimization (MM) algorithm**<sup>5</sup> constructs a series of surrogate functions that are easier to solve and guarantee the decrease of the original objective function. The algorithm is described below.

5: [https://en.wikipedia.org/wiki/MM\\_algorithm](https://en.wikipedia.org/wiki/MM_algorithm)

---

### Algorithm 5: Majorization Minimization (MM)

---

- 1: **Initialization:** Choose an  $x^0$ .
- 2: **while**  $\ell = 1, 2, \dots$  until stopping criteria are satisfied **do**
- 3:   **Majorization step:** construct a surrogate function  $f(\cdot|x^\ell)$  s.t.

$$g(x|x^\ell) \geq f(x) + c^\ell,$$

where  $g(\cdot|x^\ell) - f(x)$  is minimized at  $x^\ell$ .

- 4:   **Minimization step:**

$$x^{\ell+1} \in \arg \min_x g(x|x^\ell).$$

- 5: **end while**
-

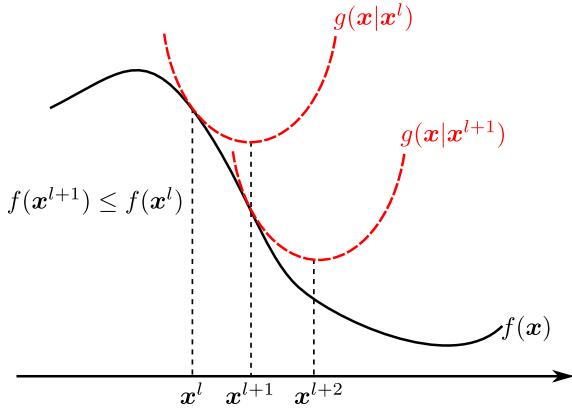


Figure 5.3: MM algorithm

**Claim:** the objective function decreases monotonically

$$f(\mathbf{x}^{l+1}) \stackrel{(a)}{\leq} g(\mathbf{x}^{l+1}|\mathbf{x}^l) - c^l \stackrel{(b)}{\leq} g(\mathbf{x}^l|\mathbf{x}^l) - c^l \stackrel{(c)}{\leq} f(\mathbf{x}^l),$$

where (a) and (c) follow from the majorization step ( $g(\cdot|\mathbf{x}^l) - f(\mathbf{x})$  is minimized at  $\mathbf{x}^l$ ), and (b) is due to the minimization step.

To understand MM algorithm, let us study the following example.

**Example 5.4.1** Consider the unconstrained optimization problem

$$\min_{\mathbf{x}} f(\mathbf{x}),$$

where

$$f(\mathbf{x}) = \delta(\|\mathbf{x}\|_0 \leq S) + \frac{\alpha}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2 \quad (5.8)$$

for some  $\alpha > 0$ .

As discussed in Section 5.1, this problem can be difficult to solve.

We construct a surrogate function as follows

$$\begin{aligned} g(\mathbf{x}|\mathbf{x}^l) &= f(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{x}^l\|_M^2 \\ &= \delta(\|\mathbf{x}\|_0 \leq S) + \frac{\alpha}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2 + \frac{1}{2} \|\mathbf{x} - \mathbf{x}^l\|_M^2, \end{aligned}$$

where

$$\frac{1}{2} \|\mathbf{x} - \mathbf{x}^l\|_M^2 := \frac{1}{2} (\mathbf{x} - \mathbf{x}^l)^T \mathbf{M} (\mathbf{x} - \mathbf{x}^l),$$

and  $\mathbf{M} \geq 0$  is a properly chosen matrix.

In particular, choose a  $\gamma$  s.t.  $0 < \gamma < 1/\sigma_{\max}^2(\mathbf{A})$  and define

$$\mathbf{M} := \frac{1}{\gamma} \mathbf{I} - \mathbf{A}^T \mathbf{A}.$$

It is straightforward to verify that



$$\mathbf{M} \succeq \sigma_{\max}^2(\mathbf{A}) \mathbf{I} - \mathbf{A}^T \mathbf{A} \succeq 0.$$

▶

$$g(\mathbf{x}|\mathbf{x}^l) = f(\mathbf{x}) + \frac{1}{2}\|\mathbf{x} - \mathbf{x}^l\|_M^2 \geq f(\mathbf{x}), \forall \mathbf{x}$$

▶

$$g(\mathbf{x}^l|\mathbf{x}^l) = f(\mathbf{x}^l)$$

and  $g(\mathbf{x}|\mathbf{x}^l) - f(\mathbf{x})$  is minimized at  $\mathbf{x}^l$ .

With the above particular choice of  $M \succeq 0$ , it holds that

$$\begin{aligned} g(\mathbf{x}|\mathbf{x}^l) &= \delta(\|\mathbf{x}\|_0 \leq S) + \frac{\alpha}{2}\|\mathbf{y} - A\mathbf{x}\|_2^2 + \frac{1}{2}\|\mathbf{x} - \mathbf{x}^l\|_M^2 \\ &= \delta(\|\mathbf{x}\|_0 \leq S) \\ &\quad + \frac{1}{2}\mathbf{x}^\top A^\top A\mathbf{x} - \mathbf{y}^\top A\mathbf{x} + \frac{1}{2}\mathbf{y}^\top \mathbf{y} \\ &\quad + \frac{1}{2}\mathbf{x}^\top M\mathbf{x} - \mathbf{x}^{l,\top} M\mathbf{x} + \frac{1}{2}\mathbf{x}^{l,\top} M\mathbf{x}^l \\ &= \delta(\|\mathbf{x}\|_0 \leq S) + \frac{1}{2\gamma} \left\| \mathbf{x} - \gamma \left( A^\top \mathbf{y} + M\mathbf{x}^l \right) \right\|_2^2 - c \\ &= \delta(\|\mathbf{x}\|_0 \leq S) + \frac{1}{2\gamma} \left\| \mathbf{x} - \left( \mathbf{x}^l + \gamma A^\top (\mathbf{y} - A\mathbf{x}^l) \right) \right\|_2^2 - c, \end{aligned}$$

where  $c$  is a constant independent of  $\mathbf{x}$  and hence does not affect the minimization step. Now it is clear that this choice of  $M$  admits a closed form solution of the minimization step (using the proximal operator in Example 5.1.3). MM algorithm can solve the nonconvex and non-smooth (unconstrained) optimization problem in (5.8).

**Example 5.4.2** Consider the unconstrained optimization problem

$$\min_{\mathbf{x}} f(\mathbf{x}),$$

where

$$f(\mathbf{x}) = \|\mathbf{x}\|_0 + \frac{\alpha}{2}\|\mathbf{y} - A\mathbf{x}\|_2^2$$

for some  $\alpha > 0$ .

## 5.5 Appendix

### Lp Norm

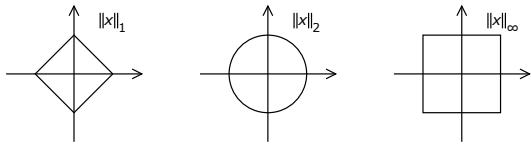
**Definition 5.5.1** ( $\ell_p$ -norm) For a real number  $p \geq 1$ , the  $\ell_p$ -norm of  $\mathbf{x} \in \mathbb{R}^n$  is given by

$$\|\mathbf{x}\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}.$$

**Example 5.5.1**

- ▶  $\ell_1$ -norm (Manhattan distance):  $\|\mathbf{x}\|_1 = \sum |x_i|$ .

- $\ell_2$ -norm (Euclidean norm):  $\|\mathbf{x}\| = \sqrt{\sum x_i^2}$ .
- $\ell_\infty$ -norm:  $\|\mathbf{x}\|_\infty = \max(|x_1|, \dots, |x_n|)$ .

Figure 5.4:  $\ell_p$ -norms

## The Cauchy-Schwartz Inequality

**Theorem 5.5.1** (The Hölder's inequality)

Let  $p, q \in [1, \infty]$  with  $1/p + 1/q = 1$ .

For all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , it holds that

$$\begin{aligned} \sum_{i=1}^n |x_i \cdot y_i| &\leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q \\ &= \left( \sum_{i=1}^n |x_i|^p \right)^{1/p} \cdot \left( \sum_{i=1}^n |y_i|^q \right)^{1/q}. \end{aligned}$$

The equality holds iff  $|x|^p$  and  $|y|^q$  are linear dependent, i.e.,  $\alpha |x_i|^p = \beta |y_i|^q$ ,  $\forall i$ .<sup>6</sup>

6: Hölder's inequality is a fundamental inequality for real analysis.

(Proof is omitted.)

**Theorem 5.5.2** (The Cauchy-Schwartz Inequality)

A special case of the Hölder's inequality is when  $p = q = 2$ .

$$\sum_{i=1}^n |x_i \cdot y_i| \leq \|\mathbf{x}\|_2 \cdot \|\mathbf{y}\|_2.$$

In particular, for all  $\mathbf{x} \in \mathbb{R}^n$ ,

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i| \leq \sqrt{n} \|\mathbf{x}\|_2,$$

where the equality holds iff  $|x_i| = |x_j|$ .

# Convex Optimization: Introduction

# 6

## 6.1 Convexity

### 6.1.1 Definitions

**Definition 6.1.1** A convex combination is a linear combination of points where all coefficients are non-negative and sum to 1.

More specifically, let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell \in \mathbb{R}^n$ . A convex combination of these points is of the form

$$\sum_{i=1}^{\ell} \lambda_i \mathbf{x}_i,$$

where the real coefficients  $\lambda_i$  satisfy  $\lambda_i \geq 0$  and  $\sum_{i=1}^n \lambda_i = 1$ .

**Definition 6.1.2** A set  $\mathcal{X}$  is a convex set if and only if the convex combination of any two points in the set belongs to the set.

That is,  $\mathcal{X} \subseteq \mathbb{R}^n$  is convex iff

$$\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}, \forall \lambda \in [0, 1], \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in \mathcal{X}.$$

### Example 6.1.1 (Convex Sets)

- A hyperplane:

$$\mathcal{H} = \{\mathbf{x} : \mathbf{a}^T \mathbf{x} = b\}$$

- A halfspace:

$$\mathcal{H}_+ = \{\mathbf{x} : \mathbf{a}^T \mathbf{x} \leq b\}$$

- A polyhedron:

$$\mathcal{P} = \{\mathbf{x} : \mathbf{a}_j^T \mathbf{x} \leq b_j, j = 1, \dots, m\}$$

- Intersections of convex sets are convex.

**Definition 6.1.3** The domain of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is defined as the set of the points where the function  $f$  is finite, i.e.,

$$\text{dom } f = \{\mathbf{x} \in \mathbb{R}^n : |f(\mathbf{x})| < \infty\}.$$

For example,  $\text{dom } \log x = \mathbb{R}^+$ .

**Definition 6.1.4 (Convex functions)** A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if for any  $\mathbf{x}_1, \mathbf{x}_2 \in \text{dom } f \subseteq \mathbb{R}^n$ ,  $\lambda \in [0, 1]$ , it holds

$$\lambda f(\mathbf{x}_1) + (1 - \lambda) f(\mathbf{x}_2) \geq f(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2).$$

|                                                 |    |
|-------------------------------------------------|----|
| 6.1 Convexity . . . . .                         | 41 |
| 6.1.1 Definitions . . . . .                     | 41 |
| 6.1.2 Conditions of Convexity . .               | 42 |
| 6.2 Unconstrained Convex Optimization . . . . . | 44 |
| 6.3 Line Search Methods . . . . .               | 45 |

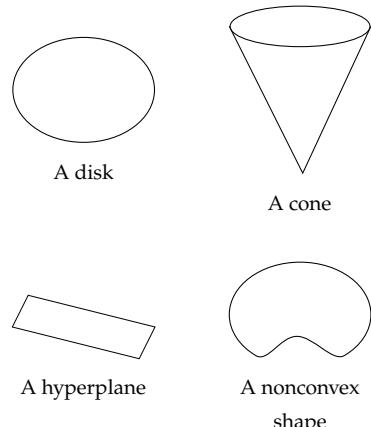


Figure 6.1: Convex and nonconvex sets

This definition implies that  $\text{dom } f$  is convex. However, for simplicity we assume that  $\text{dom } f = \mathbb{R}^n$  in this section.

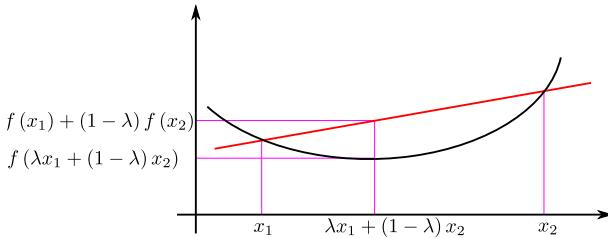


Figure 6.2: Convex functions

### 6.1.2 Conditions of Convexity

A useful mathematical tool to study conditions of convexity is Taylor's theorem.

**Theorem 6.1.1** (Taylor's Theorem) *If  $f$  is continuously differentiable, then*

$$f(\mathbf{x} + \mathbf{p}) = f(\mathbf{x}) + \nabla f(\mathbf{x} + t\mathbf{p})^\top \mathbf{p},$$

for some  $t \in (0, 1)$ .

Furthermore, if  $f$  is twice continuously differentiable, we have

$$\begin{aligned}\nabla f(\mathbf{x} + \mathbf{p}) &= \nabla f(\mathbf{x}) + \int_0^1 \nabla^2 f(\mathbf{x} + t\mathbf{p}) \mathbf{p} dt, \\ f(\mathbf{x} + \mathbf{p}) &= f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \mathbf{p} + \frac{1}{2} \mathbf{p}^\top \nabla^2 f(\mathbf{x} + t\mathbf{p}) \mathbf{p}\end{aligned}$$

for some  $t \in (0, 1)$ .

With Taylor's theorem, we are able to state and prove the first-order and second-order conditions of convexity.

**Theorem 6.1.2** (First-Order Condition of Convexity) *Suppose a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable. Then it is convex if and only if for all  $\mathbf{x}, \mathbf{y} \in \text{dom } f$ , it holds*

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}). \quad (6.1)$$

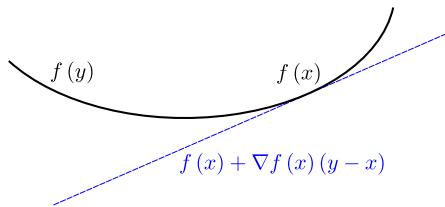


Figure 6.3: First-order condition of convexity

*Proof. Necessity:*

Assume first that  $f$  is convex and  $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$ . Since  $\text{dom}(f)$  is convex,  $\mathbf{x} + t(\mathbf{y} - \mathbf{x}) \in \text{dom}(f)$  for all  $0 < t \leq 1$ . By convexity of  $f$ ,

$$f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) \leq (1 - t)f(\mathbf{x}) + tf(\mathbf{y}).$$

Divide both sides by  $t$ . It holds

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \frac{f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{t}.$$

Take the limit as  $t \rightarrow 0$  yields (6.1).

#### Sufficiency:

To show the other direction (sufficiency), assume that (6.1) holds. Choose any  $\mathbf{x} \neq \mathbf{y}$  and  $\lambda \in [0, 1]$ . Let  $\mathbf{z} = \lambda\mathbf{x} + (1 - \lambda)\mathbf{y}$ . Applying (6.1) twice yields

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{z}) &\geq \nabla f(\mathbf{z})^T (\mathbf{x} - \mathbf{z}), \\ f(\mathbf{y}) - f(\mathbf{z}) &\geq \nabla f(\mathbf{z})^T (\mathbf{y} - \mathbf{z}). \end{aligned}$$

Multiply the first inequality by  $\lambda$  and the second by  $1 - \lambda$ , and then add them together. It holds

$$\begin{aligned} \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) - f(\mathbf{z}) \\ \geq \nabla f(\mathbf{z})^T (\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} - \mathbf{z}). \end{aligned}$$

By the definition of  $\mathbf{z}$ , the left side of the inequality is zero. Hence,

$$\lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \geq f(\mathbf{z}),$$

which proves that  $f$  is convex.  $\square$

**Theorem 6.1.3** (Second-Order Condition of Convexity) Suppose that  $f$  is twice differentiable. Then  $f$  is convex iff the Hessian of  $f$  is positive semidefinite (PSD), i.e.,  $\nabla^2 f(\mathbf{x}) \succeq 0$ .

*Proof.* The proof relies on Taylor series:

$$f(\mathbf{x} + \mathbf{v}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{v} + \frac{1}{2} \mathbf{v}^T \nabla^2 f(\mathbf{x} + \theta \mathbf{v}) \mathbf{v}$$

for some  $\theta \in [0, 1]$ .

- If  $\nabla^2 f \succeq 0$ , then  $f(\mathbf{x} + \mathbf{v}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{v}$ . By first-order condition for convexity,  $f$  is convex.
- If  $f$  is convex, then  $f(\mathbf{x} + \mathbf{v}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{v}$  for all  $\mathbf{v}$ . Hence,  $\mathbf{v}^T \nabla^2 f(\mathbf{x} + \theta \mathbf{v}) \mathbf{v} \geq 0$ . Let  $\|\mathbf{v}\|_2 \rightarrow 0$ . Then  $\nabla^2 f(\mathbf{x}) \succeq 0$ .

$\square$

#### Example 6.1.2 <sup>1</sup>

1.  $f(x) = \log(x)$  is concave on its domain  $x \in \mathbb{R}^+$ .
2.  $f(x) = x \log(x)$  is convex on its domain  $x \in \mathbb{R}^+$ .
3. Affine functions  $f(x) = \mathbf{A}\mathbf{x} + \mathbf{b}$  are both convex and concave.
4. Quadratic functions  $f(x) = \frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{x} + c$  are convex if and only if  $\mathbf{A} \succeq 0$ .
5.  $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$  is convex.

1: In this list, the first five examples can be verified using second-order condition of convexity.

6.  $f(\mathbf{x}) = \|\mathbf{x}\|_1 = \sum_i |x_i|$  is convex.
7. If  $f(\mathbf{x})$  is convex, then  $g(\mathbf{x}) := f(A\mathbf{x} + \mathbf{b})$  is convex.
  - $f(\mathbf{x}) = \|A\mathbf{x}\|_1$  is convex.
8. For large dimensional data processing, the frequently used function

$$\frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1,$$

where  $\lambda > 0$  is convex.

## 6.2 Unconstrained Convex Optimization

An unconstrained optimization problem is of the form

$$\min_{\mathbf{x}} f(\mathbf{x}).$$

**Definition 6.2.1** (global and local minimizer) A point  $\mathbf{x}^*$  is a *global minimizer* if  $f(\mathbf{x}^*) \leq f(\mathbf{x})$  for all  $\mathbf{x}$ .

A point  $\mathbf{x}^*$  is a *local minimizer* if there is a neighbourhood  $\mathcal{N}$  of  $\mathbf{x}^*$  such that  $f(\mathbf{x}^*) \leq f(\mathbf{x})$  for all  $\mathbf{x} \in \mathcal{N}$ .<sup>2</sup>

We are particularly interested in the case that  $f$  is convex due to some nice properties stated below.

2: A neighbourhood of  $\mathbf{x}^*$  is simply an open set that contains  $\mathbf{x}^*$ , for example,  $\mathcal{B}(\mathbf{x}^*, \epsilon) := \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\|_2 < \epsilon\}$ .

**Theorem 6.2.1** Suppose that a feasible point  $\mathbf{x}$  is locally optimal for a convex optimization problem. Then it is also globally optimal.

*Proof.* Suppose that  $\mathbf{x}$  is locally optimal but not globally optimal. Then there must exist a  $\mathbf{y} \neq \mathbf{x}$  such that  $f(\mathbf{y}) < f(\mathbf{x})$ .

Consider a point  $\mathbf{z}$  on the line segment between  $\mathbf{x}$  and  $\mathbf{y}$ , i.e.,

$$\mathbf{z} = (1 - \lambda)\mathbf{x} + \lambda\mathbf{y}, \quad \lambda \in (0, 1).$$

It is clear that

$$f(\mathbf{z}) \leq (1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y}) < f(\mathbf{x}),$$

where the first inequality follows from convexity of  $f$ , and the second (strict) inequality holds because  $f(\mathbf{y}) < f(\mathbf{x})$ . This contradicts with that  $\mathbf{x}$  is locally optimal and proves the global optimality of  $\mathbf{x}$ .  $\square$

**Theorem 6.2.2** (A Global Optimality Criterion) Suppose that the objective function  $f$  in a convex optimization problem is differentiable, i.e.,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f).$$

Then an  $\mathbf{x} \in \text{dom}(f)$  is optimal if and only if

$$\nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) \geq 0, \quad \forall \mathbf{y} \in \text{dom}(f).$$

Or equivalently,

$$\nabla f(\mathbf{x}) = \mathbf{0}$$

when  $f$  is closed.

*Proof.* **Sufficiency:** Suppose the inequality holds. Then for all  $\mathbf{y}$ ,

$$\begin{aligned} f(\mathbf{y}) &\geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) \\ &\geq f(\mathbf{x}). \end{aligned}$$

Hence,  $\mathbf{x}$  is globally optimal.

**Necessity:** suppose  $\mathbf{x}$  is optimal, but the inequality does not hold, i.e., for some  $\mathbf{y}$  we have

$$\nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) < 0.$$

Consider the point  $\mathbf{z}(t) = t\mathbf{y} + (1-t)\mathbf{x}$ ,  $t \in [0, 1]$ . Then

$$\begin{aligned} \frac{d}{dt} f(\mathbf{z}(t))|_{t=0} &= \nabla f(\mathbf{z}(0))^T \frac{d}{dt} \mathbf{z}(t)|_{t=0} \\ &= \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) < 0. \end{aligned}$$

This implies that for small  $t > 0$ , we have  $f(\mathbf{z}(t)) < f(\mathbf{x})$ , which contradicts the optimality of  $\mathbf{x}$ . The necessity is therefore proved.  $\square$

**Remark 6.2.1** The global optimal solution  $\mathbf{x}$  may not be unique.

### 6.3 Line Search Methods

In each iteration of linear search methods, the algorithm chooses a direction  $\mathbf{p}^l$  and searches along this direction for a lower function value. That is,

$$\mathbf{x}^{l+1} = \mathbf{x}^l + \alpha^l \mathbf{p}^l,$$

where the search direction  $\mathbf{p}^l$  and step size  $\alpha^l > 0$  are appropriately chosen so that  $f(\mathbf{x}^l + \alpha^l \mathbf{p}^l) < f(\mathbf{x}^l)$ .

#### Search Direction

Most line search methods require a search direction  $\mathbf{p}^l$  to be a descent direction:

$$\mathbf{p}^{l,\top} \nabla f(\mathbf{x}^l) < 0.$$

Based on Taylor's theorem, it is straightforward to show that the function  $f$  will be reduced along this direction as long as the step size  $\alpha_k$  is sufficiently small.

Typical choices of the search direction  $\mathbf{p}^l$  are of the form

$$\mathbf{p}^l = -\mathbf{B}^l \nabla f(\mathbf{x}^l),$$

where  $\mathbf{B}^l \succeq 0$  so that

$$\mathbf{p}^{l,\top} \nabla f(\mathbf{x}^l) = -\nabla f(\mathbf{x}^l)^\top \mathbf{B}^l \nabla f(\mathbf{x}^l) < 0.$$

Examples include

1. Steepest descent method:  $\mathbf{B}^l = \mathbf{I}$ .
2. Newton's method:  $\mathbf{B}^l = \nabla^2 f(\mathbf{x}^l)$ .
3. Quasi-Newton method:  $\mathbf{B}^l$  is an approximation to the exact Hessian  $\nabla^2 f(\mathbf{x}^l)$ .

## Step Size

### Exact Line Search

The step size  $\alpha^l$  is optimized at each iteration:

$$\alpha^l = \arg \min_{\alpha > 0} f(\mathbf{x}^l + \alpha \mathbf{p}^l).$$

In practice, it is often not possible to do this minimization exactly.

### Backtracking Line Search

A more practical strategy is to choose an  $\alpha^l > 0$  which gives *sufficient decrease* in the objective function  $f$ . For example, the popular *Armijo Condition* requires that

$$f(\mathbf{x}^l + \alpha^l \mathbf{p}^l) \leq f(\mathbf{x}^l) + c \alpha^l \mathbf{p}^{l,\top} \nabla f(\mathbf{x}^l),$$

for some constant  $c \in (0, 1)$ . In practice,  $c$  is chosen to be small, say  $c = 10^{-4}$ .

Armijo condition itself does not prevent tiny step sizes. To ensure reasonable progress of the algorithm, the *curvature condition* requires  $\alpha^l$  to satisfy

$$\mathbf{p}^{l,\top} \nabla f(\mathbf{x}^l + \alpha^l \mathbf{p}^l) \geq c' \mathbf{p}^{l,\top} \nabla f(\mathbf{x}^l)$$

for some constant  $c' \in (c, 1)$ .

The combination of Armijo and curvature conditions is called *Wolfe Conditions*:

$$\begin{aligned} f(\mathbf{x}^l + \alpha^l \mathbf{p}^l) &\leq f(\mathbf{x}^l) + c \alpha^l \mathbf{p}^{l,\top} \nabla f(\mathbf{x}^l), \\ \mathbf{p}^{l,\top} \nabla f(\mathbf{x}^l + \alpha^l \mathbf{p}^l) &\geq c' \mathbf{p}^{l,\top} \nabla f(\mathbf{x}^l) \end{aligned}$$

with  $0 < c < c' < 1$ .

# Sparse Inverse Problems: $\ell_1$ -minimization

7

## 7.1 Convex Relaxation

One popular approach to solve the sparse linear inverse problem is to replace the nonconvex objective functions mentioned in Chapter 5 with a convex one (convex relaxation).

### L<sub>p</sub>-norm

**Definition 7.1.1** Given a vector space  $\mathcal{V}$  over the field  $\mathbb{F}$  of complex (real) numbers, a norm on  $\mathcal{V}$  is a function  $p : \mathcal{V} \rightarrow \mathbb{R}$  with the following properties.

For all  $a \in \mathbb{F}$  and all  $\mathbf{u}, \mathbf{v} \in \mathcal{V}$ ,

1.  $p(a\mathbf{v}) = |a| p(\mathbf{v})$ , (absolute scalability)
2.  $p(\mathbf{u} + \mathbf{v}) \leq p(\mathbf{u}) + p(\mathbf{v})$ , (triangle inequality)
3. if  $p(\mathbf{v}) = 0$  then  $\mathbf{v}$  is the zero vector. (separates points)

Positivity follows from this definition.

*Proof.* By the first axiom,  $p(\mathbf{0}) = 0$  and  $p(-\mathbf{v}) = p(\mathbf{v})$ . Then by triangle inequality,

$$0 \leq p(\mathbf{v}) + p(-\mathbf{v}) = 2p(\mathbf{v}) \Rightarrow 0 \leq p(\mathbf{v}).$$

□

**Lemma 7.1.1** A norm is a convex function.

*Proof.* For any given  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$  and  $\lambda \in [0, 1]$ , it holds that

$$\begin{aligned} \|\lambda\mathbf{u} + (1 - \lambda)\mathbf{v}\| &\leq \|\lambda\mathbf{u}\| + \|(1 - \lambda)\mathbf{v}\| \\ &= \lambda \|\mathbf{u}\| + (1 - \lambda) \|\mathbf{v}\|, \end{aligned}$$

where we have used the triangle inequality and the absolute scalability. This establishes the convexity of the norm. □

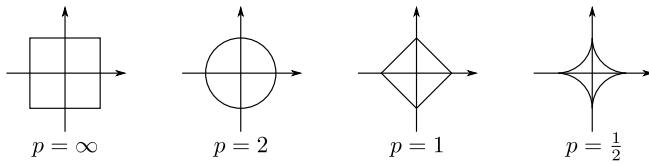
Recall Definition 5.5.1 of  $\ell_p$ -norm. For  $p \geq 1$ ,

$$\|\mathbf{x}\|_p = \left( \sum_i |x_i|^p \right)^{1/p}.$$

In this definition,  $\ell_p$ -norm is a proper norm iff  $p \geq 1$ .

The following picture draws the unit  $\ell_p$ -ball.

|                                                           |    |
|-----------------------------------------------------------|----|
| 7.1 Convex Relaxation . . . . .                           | 47 |
| 7.2 LASSO Solvers . . . . .                               | 50 |
| 7.2.1 Soft Thresholding Function                          | 50 |
| 7.2.2 Cyclic Coordinate Descent                           | 51 |
| 7.2.3 Iterative Shrinkage Thresholding Algorithm (ISTA) . | 51 |
| 7.2.4 Majorization Minimization                           | 52 |
| 7.2.5 Least Angle Regression (LAR) . . . . .              | 52 |

Figure 7.1: Unit  $\ell_p$  ball

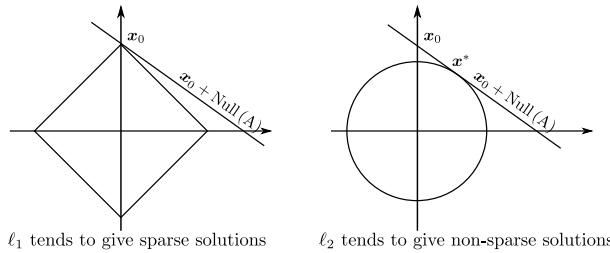
## LASSO Problem

The most well-known convex relaxation for sparse linear inverse problem is LASSO (least absolute shrinkage and selection operator):

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2 + \lambda \|\mathbf{x}\|_1. \quad (7.1)$$

The first term is for data fidelity and the second one is an  $\ell_1$  regularization.

Figure 7.2 gives the geometric intuition why  $\ell_1$ -norm promotes sparsity while  $\ell_2$ -norm does not. Note that the set of feasible solutions of  $\mathbf{y} = \mathbf{Ax}$  is given by  $\mathcal{X} = \mathbf{x}_0 + \text{Null}(A)$ .

Figure 7.2:  $\ell_1$ -norm promotes sparse solutions

## Difficulty in Solving LASSO

The main difficulty in solving LASSO is that the objective function is not differentiable! Gradient descent method cannot be directly applied.

One naive approach is to define the so called subgradient and apply a subgradient line search method.

**Definition 7.1.2** Let  $f : \mathcal{U} \rightarrow \mathbb{R}$  be a convex function defined on a convex open set  $\mathcal{U} \subset \mathbb{R}^n$ . A vector  $v \in \mathbb{R}^n$  is called a **subgradient** at a point  $\mathbf{x} \in \mathcal{U}$  if

$$f(\mathbf{y}) - f(\mathbf{x}) \geq v^\top (\mathbf{y} - \mathbf{x}), \forall \mathbf{y} \in \mathcal{U}.$$

The set of all subgradients at  $\mathbf{x}$  is called the **subdifferential** at  $\mathbf{x}$  and is denoted  $\partial f(\mathbf{x})$ .

**Remark 7.1.1** If  $f$  is convex and its subdifferential at  $\mathbf{x}$  contains exactly one subgradient, then  $f$  is differentiable at  $\mathbf{x}$ .

**Example 7.1.1** Let  $f(x) = |x|$ . Then the subdifferential of  $f$  is given by

$$\partial f(x) = \partial|x| = \begin{cases} 1 & \text{if } x > 0, \\ [-1, 1] & \text{if } x = 0, \\ -1 & \text{if } x < 0. \end{cases}$$

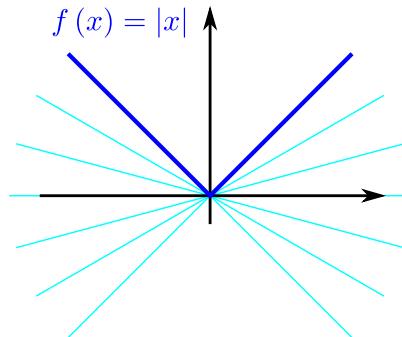


Figure 7.3: Subdifferential of  $|x|$

With the definition of subgradient, one can implement a subgradient line search method. However,

- Subgradient line search does not necessarily lead to a decrease of the objective function.
- Subgradient line search does not necessarily converge to a global minimizer.

**Example 7.1.2** Let  $x \in \mathbb{R}^2$  and

$$f(x) = 3|x_1| + |x_2|.$$

Consider the point  $x^0 = [0, 1]^T$ . Show that

- $g^1 = [0, 1]^T \in \partial f(x)$ .
- $g^2 = [3, 1]^T \in \partial f(x)$ .

Calculate

$$\lim_{\tau > 0, \tau \rightarrow 0} \frac{f(x^0 - \tau g) - f(x^0)}{\tau}$$

with  $g = g^1$  and  $g = g^2$  respectively. Comment on whether  $-g$  is always a descent direction or not.

**Example 7.1.3 (Wolfe's Example)** Let  $x \in \mathbb{R}^2$  and

$$f(x) = \begin{cases} 5(9x_1^2 + 16x_2^2)^{1/2} & \text{if } x_1 > |x_2|, \\ 9x_1 + 16|x_2| & \text{if } x_1 \leq |x_2|. \end{cases}$$

Show that

$$x^{\text{nd}} = [0, 0]^T$$

is a non-differentiable point.

Suppose that  $x^0 = [16/9, 1]^\top$ . Consider exact line search where

$$x^{l+1} = x^l - t^l \nabla f(x^l),$$

where

$$t^l = \arg \min_t f(x^l - t \nabla f(x^l)).$$

Show that

1.  $x^l$  are all differentiable points;
2.  $x^l \rightarrow [0, 0]^\top$  as  $l \rightarrow \infty$ , and the point  $[0, 0]^\top$  is not a global minimizer.

## 7.2 LASSO Solvers

### 7.2.1 Soft Thresholding Function

Proximal operator of  $\ell_1$ -norm plays an essential role in developing practical algorithms to solve LASSO problems.

We first study the simple case — the scalar case.

Let

$$x^* = \text{Prox}_{\gamma|\cdot|}(y) = \arg \min_x |x| + \frac{1}{2\gamma}(x - y)^2.$$

It is clear that the objective function

$$f_\gamma(x) := |x| + \frac{1}{2\gamma}(x - y)^2$$

is convex. Set its subgradient at the global minimizer  $x^*$  to zero. One has

$$\partial_x f_\gamma(x^*) = \partial_x |x^*| + \frac{1}{\gamma}(x^* - y).$$

Hence,

$$\begin{aligned} x^* &= y - \gamma \partial_x |x^*| \\ &= \begin{cases} y - \gamma & \text{if } y \geq \gamma, \\ 0 & \text{if } -\gamma < y < \gamma, \\ y + \gamma & \text{if } y \leq -\gamma, \end{cases} \\ &=: \eta(y; \gamma) = \eta_\gamma(y) \\ &= \text{sign}(y)(|y| - \gamma)_+, \end{aligned} \tag{7.2}$$

where  $(z)_+ = \max(z, 0)$ .

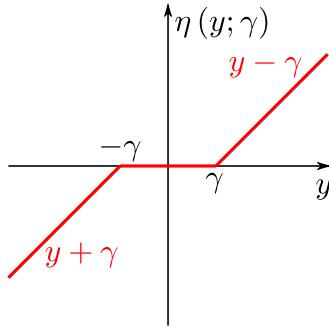


Figure 7.4: Soft thresholding function

The extension to vector case is straightforward. Let  $x^* = \text{Prox}_{\gamma\|\cdot\|}(\mathbf{y})$ . Then  $x_i^* = \eta(y_i; \gamma)$ .

This closed form solution of the  $\ell_1$ -norm proximal operator is called **soft thresholding function**.

### 7.2.2 Cyclic Coordinate Descent

Consider the LASSO problem

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2 + \lambda \|\mathbf{x}\|_1.$$

At each step of cyclic coordinate descent, we optimize the objective function with respect to  $x_i$  for some  $i$  by fixing all  $x_j, j \neq i$ . The objective function can be then written as

$$\begin{aligned} & \frac{1}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2 + \lambda \|\mathbf{x}\|_1 \\ &= \frac{1}{2} \left\| \mathbf{y} - \mathbf{a}_i x_i - \sum_{j \neq i} \mathbf{a}_j x_j \right\|_2^2 + \lambda |x_i| + \lambda \sum_{j \neq i} |x_j| \\ &= \frac{1}{2} \|\mathbf{r}_i - \mathbf{a}_i x_i\|_2^2 + \lambda |x_i| + c. \end{aligned}$$

The global minimizer  $x_i^*$  is given by

$$x_i^* = \eta(\mathbf{a}_i^\top \mathbf{r}_i; \lambda),$$

where

$$\begin{aligned} \mathbf{r}_i &:= \mathbf{y} - \sum_{j \neq i} \mathbf{a}_j \hat{x}_j \\ &= \mathbf{y} - \sum_j \mathbf{a}_j \hat{x}_j + \mathbf{a}_i \hat{x}_i. \end{aligned}$$

In practice, cyclic coordinate descent is slow as one needs to cyclically update all  $\hat{x}_i, i = 1, 2, \dots, N$ , many times.

### 7.2.3 Iterative Shrinkage Thresholding Algorithm (ISTA)

The key idea of ISTA is a Taylor approximation to allow the direct application of proximal operator.

Let  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2$ . Approximate the objective function locally by

$$\begin{aligned} & f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1 \\ & \approx f(\mathbf{x}^l) + \langle \mathbf{x} - \mathbf{x}^l, \nabla f(\mathbf{x}^l) \rangle + \frac{1}{2\tau^l} \|\mathbf{x} - \mathbf{x}^l\|_2^2 + \lambda \|\mathbf{x}\|_1 \\ & = \frac{1}{2\tau^l} \left\| \mathbf{x} - \left( \mathbf{x}^l - \tau^l \nabla f(\mathbf{x}^l) \right) \right\|_2^2 + \lambda \|\mathbf{x}\|_1 + c. \end{aligned}$$

Note that

$$\nabla f(\mathbf{x}^l) = -A^\top(\mathbf{y} - A\mathbf{x}^l).$$

The ISTA iterations are given as

$$\mathbf{x}^{l+1} = \eta \left( \mathbf{x}^l + \tau^l A^\top(\mathbf{y} - A\mathbf{x}^l); \lambda \tau^l \right).$$

In practice, the step size  $\tau^l$  is fixed as  $\tau^l = \tau \in (0, 2/\sigma_{\max}(A^\top A))$ .

#### 7.2.4 Majorization Minimization

Apply majorization minimization (Section 5.4) to the LASSO problem. In each iteration, one has objective function

$$\begin{aligned} & \lambda \|\mathbf{x}\|_1 + \frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \frac{1}{2} \|\mathbf{x} - \mathbf{x}^l\|_M^2 \\ & = \lambda \|\mathbf{x}\|_1 + \frac{1}{2\gamma} \left\| \mathbf{x} - \left( \mathbf{x}^l + \gamma A^\top(\mathbf{y} - A\mathbf{x}^l) \right) \right\|_2^2 + c, \end{aligned}$$

where  $M = \frac{1}{\gamma} I - A^\top A \succeq 0$  and  $\gamma > 1/\sigma_{\max}(A^\top A)$ . Hence in each iteration of majorization minimization, proximal operator can be directly used for the update:

$$\begin{aligned} \mathbf{x}^{l+1} &= \text{Prox}_{\lambda\gamma\|\cdot\|_1} \left( \mathbf{x}^l + \gamma A^\top(\mathbf{y} - A\mathbf{x}^l) \right) \\ &= \eta \left( \mathbf{x}^l + \gamma A^\top(\mathbf{y} - A\mathbf{x}^l); \lambda \gamma \right). \end{aligned}$$

#### 7.2.5 Least Angle Regression (LAR)

**Observation:** Different  $\lambda > 0$  results in different sparsity in the optimal solution

$$\mathbf{x}_\lambda := \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1.$$

►  $\lambda = 0$ :

$$\begin{aligned} \mathbf{x}_\lambda &= \arg \min_{\mathbf{x}} \|\mathbf{y} - A\mathbf{x}\|_2^2 \\ &= A^\dagger \mathbf{y} \end{aligned}$$

is not sparse.

►  $\lambda \rightarrow \infty$ :

$$\begin{aligned}\mathbf{x}_\lambda &= \arg \min_{\mathbf{x}} \|\mathbf{x}\|_1 \\ &= \mathbf{0}.\end{aligned}$$

**Least angle regression (LAR)**: Find a  $\lambda$  such that the optimal  $\mathbf{x}$  has a specific sparsity.

The derivation of LAR is based on the optimal condition:

$$\mathbf{A}^\top(\mathbf{y} - \mathbf{A}\mathbf{x}_\lambda) = \lambda \text{sign}(\mathbf{x}_\lambda).$$

Let  $\mathcal{I} := \text{supp}(\mathbf{x}_\lambda)$  and  $\mathcal{I}^c := \{1, \dots, N\} \setminus \text{supp}(\mathbf{x}_\lambda)$ . It holds that

$$\begin{cases} \mathbf{x}_{\lambda, \mathcal{I}} = (\mathbf{A}_{\mathcal{I}}^\top \mathbf{A}_{\mathcal{I}})^{-1} (\mathbf{A}_{\mathcal{I}}^\top \mathbf{y} - \lambda \text{sign}(\mathbf{x}_{\lambda, \mathcal{I}})), \\ \mathbf{x}_{\lambda, \mathcal{I}^c} = \mathbf{0}. \end{cases}$$

That is, when the sign pattern of  $\mathbf{x}_\lambda$  is fixed (so is  $\mathcal{I}$ ),  $\mathbf{x}_\lambda$  is linear in  $\lambda$ .

LAR: find the knots  $\lambda$  where  $\text{sign}(\mathbf{x}_\lambda)$  changes.

$$\lambda \text{sign}(\mathbf{x}_\lambda) = \mathbf{A}^\top(\mathbf{y} - \mathbf{A}\mathbf{x}_\lambda) = \mathbf{A}^\top(\mathbf{y} - \mathbf{A}_{\mathcal{I}}\mathbf{x}_{\lambda, \mathcal{I}}) = \mathbf{A}^\top \mathbf{r}_\lambda$$

- Keep track  $\mathbf{x}_\lambda$  and  $|\langle \mathbf{a}_j, \mathbf{r}_\lambda \rangle|$  until
  - Either  $|\mathcal{I}|$  increases

$$\max_{j \notin \mathcal{I}} |\langle \mathbf{a}_j, \mathbf{r}_\lambda \rangle| = \lambda.$$

Define  $i = \arg \max_{j \notin \mathcal{I}} |\langle \mathbf{a}_j, \mathbf{r}_\lambda \rangle|$  and set  $\mathcal{I} = \mathcal{I} \cup \{i\}$ .

- Or  $|\mathcal{I}|$  decreases  
For some  $i \in \mathcal{I}$ ,

$$|\langle \mathbf{a}_i, \mathbf{r}_\lambda \rangle| = \lambda \Rightarrow (\mathbf{x}_\lambda)_i = 0.$$

Set  $\mathcal{I} = \mathcal{I} \setminus \{i\}$ .

**Example**: Find  $\lambda_0$  such that

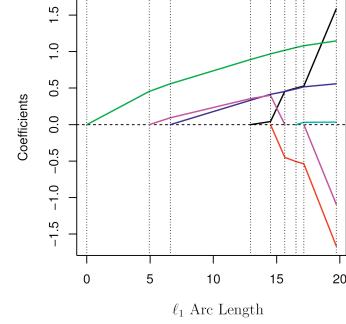
$$\begin{cases} \|\mathbf{x}_\lambda\|_0 = 0, & \forall \lambda > \lambda_0, \\ \|\mathbf{x}_\lambda\|_0 = 1, & \forall \lambda \in (\lambda_0 - \epsilon, \lambda_0), \end{cases}$$

for some  $\epsilon > 0$ .

Let  $\text{supp}(\mathbf{x}_\lambda) = \{n\}$  for  $\lambda \in (\lambda_0 - \epsilon, \lambda_0)$ .

The optimality of  $\mathbf{x}_\lambda$  suggests that

$$\begin{aligned}\mathbf{A}^\top(\mathbf{y} - \mathbf{A}\mathbf{x}_\lambda) &= \lambda \text{sign}(\mathbf{x}_\lambda) \\ \Rightarrow \mathbf{a}_n^\top(\mathbf{y} - \mathbf{a}_n \mathbf{x}_{\lambda, n}) &= \lambda \text{sign}(\mathbf{x}_{\lambda, n}) \\ \Rightarrow \mathbf{x}_{\lambda, n} &= \mathbf{a}_n^\top \mathbf{y} - \lambda \text{sign}(\mathbf{x}_{\lambda, n})\end{aligned}$$



**Figure 7.5:** Typical behaviour of  $\mathbf{x}_\lambda$ . (Source: T. Hastie, et al., Statistical learning with sparsity: the lasso and generalizations. Chapman and Hall/CRC, 2015: page 120.)

Hence,

$$x_i^* = (\|a_i^T y\| - \lambda) \operatorname{sign}(a_i^T y), \text{ and}$$
$$\lambda_0 = \max_i \|a_i^T y\|.$$

# 8

## Bilinear Inverse Problems

### 8.1 Bilinear Inverse Problems

A [linear map](#)  $\mathbb{L}$  satisfies

$$\mathbb{L}(\lambda_1 \mathbf{x}_1 + \lambda_2 \mathbf{x}_2) = \lambda_1 \mathbb{L}(\mathbf{x}_1) + \lambda_2 \mathbb{L}(\mathbf{x}_2)$$

for all  $\lambda_1, \lambda_2 \in \mathbb{R}$  and  $\mathbf{x}_1, \mathbf{x}_2$  in the domain of  $\mathbb{L}$ .

The general form of linear functions (domain of  $\mathbb{L}$  is finite dimensional) is given by

$$\mathbb{L}(\mathbf{x}) = A\mathbf{x}.$$

A noisy linear system is typically written as

$$\mathbf{y} = A\mathbf{x} + \mathbf{w}.$$

A linear inverse problem is given by

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2.$$

A [bilinear map](#)  $\mathbb{B}$  is a map such that

$$\begin{aligned} \mathbb{B}(\lambda_1 \mathbf{p}_1 + \lambda_2 \mathbf{p}_2, \mathbf{q}) &= \lambda_1 \mathbb{B}(\mathbf{p}_1, \mathbf{q}) + \lambda_2 \mathbb{B}(\mathbf{p}_2, \mathbf{q}) \\ \mathbb{B}(\mathbf{p}, \lambda_1 \mathbf{q}_1 + \lambda_2 \mathbf{q}_2) &= \lambda_1 \mathbb{B}(\mathbf{p}, \mathbf{q}_1) + \lambda_2 \mathbb{B}(\mathbf{p}, \mathbf{q}_2) \end{aligned}$$

for all  $\lambda_1, \lambda_2 \in \mathbb{R}$  and  $\mathbf{p}, \mathbf{p}_1, \mathbf{p}_2, \mathbf{q}, \mathbf{q}_1, \mathbf{q}_2$ , in the domain of  $\mathbb{B}$ .

The general form of bilinear map (domain of  $\mathbb{B}$  is finite dimensional) is given by

$$\begin{aligned} \mathbb{B}(\mathbf{p}, \mathbf{q}) &= [f_1(\mathbf{p}, \mathbf{q}), \dots, f_m(\mathbf{p}, \mathbf{q})]^\top \\ &= [\mathbf{p}^T \mathbf{A}_1 \mathbf{q}, \dots, \mathbf{p}^T \mathbf{A}_m \mathbf{q}]^\top, \end{aligned}$$

where  $f_i(\mathbf{p}, \mathbf{q}) = \mathbf{p}^T \mathbf{A}_i \mathbf{q}$  are bilinear functions.

The above bilinear form can be written in a more convenient form. For this, we need the following definitions.

**Definition 8.1.1 (Trace)** *The trace of an  $n \times n$  square matrix  $\mathbf{A}$  is defined as the sum of diagonal entries, i.e.,*

$$tr(\mathbf{A}) = \sum_{i=1}^n a_{i,i}.$$

**Definition 8.1.2 (Vectorization)** *The vectorization of a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is a linear map which converts the matrix into a column vector by stacking*

the columns of the matrix  $\mathbf{A}$  together:

$$\begin{aligned} \text{vec}(\mathbf{A}) &= [\mathbf{a}_{:,1}^\top, \dots, \mathbf{a}_{:,n}^\top]^\top \\ &= [a_{1,1}, \dots, a_{m,1}, \dots, a_{1,n}, \dots, a_{m,n}]^\top. \end{aligned}$$

**Lemma 8.1.1** (Properties of Trace)

1.

$$\text{tr}(\mathbf{A}) = \text{tr}(\mathbf{A}^\top).$$

2.

$$\text{tr}(\mathbf{A}^\top \mathbf{B}) = \text{tr}(\mathbf{B} \mathbf{A}^\top) = \sum_{i,j} a_{i,j} b_{j,i}.$$

3.

$$\text{tr}(\mathbf{A}^\top \mathbf{B}) = \text{tr}(\mathbf{A} \mathbf{B}^\top) = \text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B}).$$

**Definition 8.1.3** (Frobenius Inner Product) Given two matrices of the same size  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times m}$ , the Frobenius inner product is defined as

$$\begin{aligned} \langle \mathbf{A}, \mathbf{B} \rangle &:= \sum_{i,j} \overline{a_{i,j}} b_{i,j} \\ &= \text{tr}(\mathbf{A}^\text{H} \mathbf{B}), \end{aligned}$$

where the overline denotes the complex conjugate and the superscript  $\text{H}$  denotes Hermitian conjugate.

For two real matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ , the Frobenius inner product becomes

$$\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^\top \mathbf{B}).$$

We write a bilinear map using trace.

$$\begin{aligned} \mathbb{B}(\mathbf{p}, \mathbf{q}) &= [\mathbf{p}^\top \mathbf{A}_1 \mathbf{q}, \dots, \mathbf{p}^\top \mathbf{A}_m \mathbf{q}]^\top \\ &= [\text{tr}(\mathbf{p}^\top \mathbf{A}_1 \mathbf{q}), \dots, \text{tr}(\mathbf{p}^\top \mathbf{A}_m \mathbf{q})]^\top \\ &\stackrel{(a)}{=} [\text{tr}(\mathbf{A}_1 \mathbf{q} \mathbf{p}^\top), \dots, \text{tr}(\mathbf{A}_m \mathbf{q} \mathbf{p}^\top)]^\top \\ &\stackrel{(b)}{=} [\langle \mathbf{A}_1^\top, \mathbf{q} \mathbf{p}^\top \rangle, \dots, \langle \mathbf{A}_m^\top, \mathbf{q} \mathbf{p}^\top \rangle]^\top \\ &= \begin{bmatrix} \text{vect}(\mathbf{A}_1)^\top \\ \vdots \\ \text{vect}(\mathbf{A}_m)^\top \end{bmatrix} \text{vect}(\mathbf{p} \mathbf{q}^\top) \\ &= \mathbb{A}(\text{vect}(\mathbf{p} \mathbf{q}^\top)), \end{aligned}$$

where (a) follows from Lemma 8.1.1-2 and (b) comes from Definition 8.1.3.

Hence, a noisy bilinear system can be expressed as

$$\begin{aligned} \mathbf{y} &= \mathbb{A}(\mathbf{p} \mathbf{q}^\top) + \mathbf{w} \\ &= A \text{vect}(\mathbf{p} \mathbf{q}^\top) + \mathbf{w}. \end{aligned}$$

A bilinear inverse problem can be written as

$$\min_{p,q} \frac{1}{2} \|y - A \text{vect}(pq^T)\|_2^2$$

### 8.1.1 Matrix Norms

#### Frobenius Norm

$$\|A\|_F = \left( \sum_{i,j} |a_{i,j}|^2 \right)^{1/2} = \left( \sum_k \sigma_k^2 \right)^{1/2}.$$

#### Operator Norm

Operator norms of matrices are *induced* by vector norms.

##### Definition 8.1.4 (Operator Norm)

$$\|A\| = \sup\{\|Ax\| : \|x\| = 1\}.$$

##### Example 8.1.1

►  $p = 1$ :

$$\begin{aligned} \|A\|_1 &= \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{i,j}| \\ &= \max_{1 \leq j \leq n} \|\mathbf{a}_{:,j}\|_1. \end{aligned}$$

►  $p = \infty$ :

$$\begin{aligned} \|A\|_\infty &= \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{i,j}| \\ &= \max_{1 \leq i \leq m} \|\mathbf{a}_{i,:}\|_1. \end{aligned}$$

►  $p = 2$ :

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^\top A)} = \sigma_{\max}(A).$$

**Remark 8.1.1** Note the difference between  $\ell_2$ -norm and Frobenius norm of matrices.

$$\begin{aligned} \|A\|_2 &= \sigma_{\max}(A) \\ &\leq \|A\|_F = \left( \sum_{i,j} |a_{i,j}|^2 \right)^{1/2} = \left( \sum_k \sigma_k^2 \right)^{1/2}, \end{aligned}$$

where the equality holds if and only if  $A$  has a rank at most one.

**$L_{2,1}$ -norm**

$$\|A\|_{2,1} := \sum_{j=1}^n \|a_{:,j}\|_2 = \sum_{j=1}^n \left( \sum_{i=1}^m |a_{i,j}|^2 \right)^{1/2}$$

This norm has been widely used to promote sparsity in columns of a matrix.

**Schatten Norms**

The Schatten  $p$ -norm applies the  $p$ -norm to the singular value vector:

$$\|A\|_p = \left( \sum_i \sigma_i^p(A) \right)^{1/p}.$$

The most popular cases include

- $p = 1$ : nuclear norm.

$$\|A\|_* = \sum_i \sigma_i(A). \quad (8.1)$$

Nuclear norm is often used to prompt a low-rank solution.

- $p = 2$ : Frobenius norm.

## 8.2 Examples

The first example is designed to show that bilinear inverse problems are inherently nonconvex.

**Example 8.2.1** (A simple bilinear inverse problem) Consider the following bilinear inverse problem

$$f(u, v) = (1 - uv)^2.$$

- For a fixed  $u$ ,  $f(u, v)$  is convex in  $v$ .
- For a fixed  $v$ ,  $f(u, v)$  is convex in  $u$ .
- $f(u, v)$  is not convex in  $(u, v)$ . To see this, consider two points  $(u_1, v_1) = (0.1, 10)$  and  $(u_2, v_2) = (10, 0.1)$ .

See below for other examples of bilinear inverse problems.

### 8.2.1 Low-Rank Matrix Completion

**Example 8.2.2** (Netflix Problem)

|          | Movie 1 | Movie 2 | Movie 3 | Movie 4 |
|----------|---------|---------|---------|---------|
| Person 1 | ★★★★★   | ★★★★☆   | ★★★★☆   |         |
| Person 2 | ★★★★☆   |         | ★★★★★   |         |
| Person 3 |         | ★★★★☆   |         | ★★★★★   |
| Person 4 | ★★★★★   |         |         | ★★★★★   |

Figure 8.1: Netflix problem [https://en.wikipedia.org/wiki/Netflix\\_Prize](https://en.wikipedia.org/wiki/Netflix_Prize)

A low rank model: the full matrix is of low-rank.

The inverse problem is then given by

$$\begin{aligned} \min_{P \in \mathbb{R}^{m,r}, Q \in \mathbb{R}^{n,r}} \frac{1}{2} \|Y_\Omega - (PQ^T)_\Omega\|_F^2, \\ \equiv \\ \min_{X: \text{rank}(X) \leq r} \frac{1}{2} \|Y_\Omega - X_\Omega\|_F^2 \end{aligned}$$

where  $\Omega$  is the index set of revealed entries.

### Example 8.2.3 (Simple Examples of Low Rank Matrix Completion)

1. Find a rank-1 matrix given the observations

$$\begin{bmatrix} 1 & ? & 3 \\ 2 & -2 & ? \\ ? & -1 & ? \end{bmatrix}$$

2. Find a rank-2 matrix given the observations

$$\begin{bmatrix} 1 & 2 & ? & ? \\ 1 & 2 & 0 & 4 \\ 1 & 0 & 2 & 2 \\ ? & ? & 2 & 2 \end{bmatrix}$$

### Eckart-Young-Mirsky Theorem

Though bilinear inverse problem is nonconvex in general, we present a special case where a global minimizer can be easily found (via SVD).

Let  $Y \in \mathbb{R}^{m \times n}$  be a given matrix. Consider the following nonconvex optimization problem

$$\min_{X: \text{rank}(X) \leq r} \|Y - X\|_F^2 \equiv \min_{P \in \mathbb{R}^{m \times r}, Q \in \mathbb{R}^{n \times r}} \|Y - PQ^T\|_F^2$$

**Theorem 8.2.1** (The Eckart-Young-Mirsky Theorem) Consider SVD of the matrix  $Y = \sum_i \sigma_i u_i v_i^\top$ . The **best rank  $r$  approximation** of  $Y$  is given by the truncated SVD

$$X_r = \sum_{i=1}^r \sigma_i u_i v_i^\top.$$

**Remark 8.2.1** (Weighted Low-Rank Approximation Problems) Consider the weighted low-rank approximation problems of the form

$$\min_{X: \text{rank}(X) \leq r} \frac{1}{2} \text{vec}(Y - X)^\top W \text{vec}(Y - X).$$

The general weighted low-rank approximation problem does not admit an analytic solution. Global optimality of the solution is not guaranteed.

**Remark 8.2.2** (Uniqueness of the Solution)

$$\min_{X: \text{rank}(X) \leq r} \|Y - X\|_F^2 \equiv \min_{P \in \mathbb{R}^{m \times r}, Q \in \mathbb{R}^{n \times r}} \|Y - PQ^T\|_F^2$$

- Optimal solutions  $(P^*, Q^*)$  are not unique: scaling/rotation invariant.
- If  $(P^*, Q^*)$  is the optimal solution, so is  $(P^*S, Q^*S^{-T})$  for any invertible  $S \in \mathbb{R}^{r \times r}$ . This is because

$$P^*Q^{*\top} = P^*(SS^{-1})Q^{*\top} = (P^*S)(Q^*S^{-\top})^T$$

A solution  $(P^*, Q^*)$  is optimal/unique up to scaling / rotation / permutation.

- $X^* = P^*Q^{*\top}$  is unique.

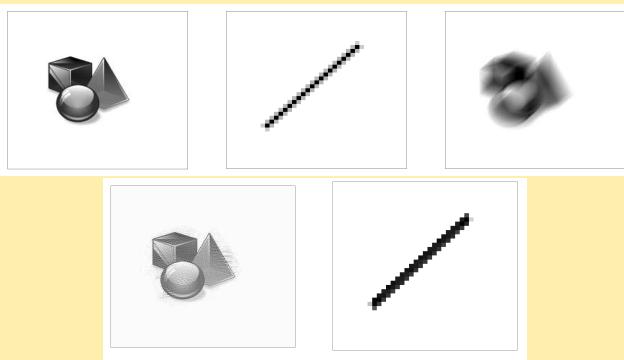
## 8.2.2 Blind Deconvolution

**Example 8.2.4** (Blind Deconvolution) Assume that the observed data is a convolution of an unknown source  $s$  and an unknown channel  $h$ , i.e.,

$$y = s \otimes h : y[n] = \sum_{\ell=0}^L s[n-\ell] h[\ell].$$

It is straightforward to verify that blind deconvolution problem is a bilinear inverse problem.

Figure 8.2 shows one simulation result from the seminal paper Ahmed, Recht, and Romberg, 2013.



**Figure 8.2:** Blind deconvolution (Source: Ahmed, Recht, and Romberg, 2013). The top row from left to right: source, channel, and convolution of source and channel. The bottom row from left to right: estimated source and channel via solving the corresponding bilinear inverse problem.

Blind deconvolution can be formulated as a low-rank matrix recovery problem as follows. Define

$$X := hs^T = \begin{bmatrix} h[0]s[0] & h[0]s[1] & h[0]s[2] & h[0]s[3] & \cdots \\ h[1]s[0] & h[1]s[1] & h[1]s[2] & h[1]s[3] & \cdots \\ h[2]s[0] & h[2]s[1] & h[2]s[2] & h[2]s[3] & \cdots \end{bmatrix}.$$

It is clear that the matrix  $X$  has rank one. The observations  $y$  can be

written as

$$\begin{aligned} \mathbf{y} &= \begin{bmatrix} y[0] \\ y[1] \\ y[2] \\ y[3] \\ \vdots \end{bmatrix} = \mathbb{A}(\mathbf{X}) \\ &= \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & \cdots \\ 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & \cdots \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & \cdots \\ \vdots & \ddots \end{bmatrix} \text{vect}(\mathbf{X}). \end{aligned}$$

The inverse problem can be written as

$$\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{y} - \mathbb{A}(\mathbf{X})\|_F^2.$$

### 8.2.3 Sparse PCA and Dictionary Learning

**Example 8.2.5** (Sparse PCA and Dictionary Learning)

$$\underbrace{\mathbf{Y}_{m \times n}}_{m \times n} = \mathbf{D}\mathbf{S}$$

$$= \underbrace{\left[ \begin{array}{cccc} \mathbf{d}_1 & \cdots & \mathbf{d}_k \end{array} \right]}_{m \times k} \underbrace{\mathbf{S}_{k \times n}}_{k \times n}$$

Typical assumption:  $\mathbf{S}$  is sparse.

- ▶ Sparse PCA
- ▶ Dictionary learning

The same formulation is applicable to blind source separation.<sup>1</sup>

1: [https://en.wikipedia.org/wiki/Signal\\_separation](https://en.wikipedia.org/wiki/Signal_separation)

## 8.3 Methods

### 8.3.1 Alternating Optimization

$$\min_{\mathbf{P} \in \mathbb{R}^{m \times r}, \mathbf{Q} \in \mathbb{R}^{n \times r}} \|\mathbf{y} - \mathbb{A}(\mathbf{P}\mathbf{Q}^T)\|_2^2$$

Alternating projection: in each iteration

1. Fix  $\mathbf{P}$ , optimise  $\mathbf{Q}$ :

$$\begin{aligned} \mathbf{Q}^{l+1} &= \arg \min_{\mathbf{Q}} \|\mathbf{y} - \mathbb{A}(\mathbf{P}^l \mathbf{Q}^T)\|_2^2 \\ &= \arg \min_{\mathbf{Q}} \|\mathbf{y} - \mathbb{A}_{\mathbf{P}^l}(\mathbf{Q}^T)\|_2^2 \end{aligned}$$

2. Fix  $\mathbf{Q}$ , optimise  $\mathbf{P}$ :

$$\begin{aligned}\mathbf{P}^{l+1} &= \arg \min_{\mathbf{P}} \|\mathbf{y} - \mathbb{A}(\mathbf{P}\mathbf{Q}^{l+1,T})\|_2^2 \\ &= \arg \min_{\mathbf{P}} \|\mathbf{y} - \mathbb{A}_{\mathbf{Q}^{l+1}}(\mathbf{P})\|_2^2\end{aligned}$$

**Remark 8.3.1**

1. Each step is to solve a least squares problem of which the global optimal solution can be computed.
2. The objective function is monotonically decreasing. Convergence is guaranteed.

### 8.3.2 Greedy Algorithms

OMP, SP, CoSaMP and IHT can be adapted to solve low-rank matrix recovery problems. In the following, we use IHT as an example.

For greedy algorithms, the optimization problem can be formulated as

$$\min_{X: \text{rank}(X) \leq r} \|\mathbf{Y} - \mathbb{A}(X)\|_F^2.$$

The proximal operator of rank indicator function is given as

$$\begin{aligned}\mathbf{Z}^* &= \text{Prox}_{\gamma\delta(\text{rank}(\cdot) \leq r)}(\mathbf{Z}) \\ &= \arg \min_{\mathbf{X}} \delta(\text{rank}(\mathbf{X}) \leq r) + \frac{1}{2\gamma} \|\mathbf{X} - \mathbf{Z}\|_F^2 \\ &= \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top \\ &=: H_{\text{rank}-r}(\mathbf{Z}),\end{aligned}$$

where

$$\mathbf{Z} = \sum_i \sigma_i \mathbf{u}_i \mathbf{v}_i^\top.$$

The IHT algorithm for low-rank matrix recovery is then given by

$$\mathbf{X}^{l+1} = H_{\text{rank}-r} \left( \mathbf{X}^l + \tau^l \mathbb{A}^* \left( \mathbf{y} - \mathbb{A}(\mathbf{X}^l) \right) \right),$$

where  $\mathbb{A}^*$  is the adjoint operator of the linear map  $\mathbb{A}$ : if we represent the linear map  $\mathbb{A}$  by a matrix  $A$ , then  $\mathbb{A}^* = A^\top$  if  $A$  is a real matrix and  $\mathbb{A}^* = A^H$  if  $A$  is a complex matrix.

### 8.3.3 Convex Relaxation

Nuclear norm

$$\|A\|_* = \sum_i \sigma_i(A).$$

can be viewed as  $\ell_1$ -norm of the singular value vector, promoting sparsity in singular value vector, i.e., low-rank solutions.

A Lasso type formulation to promote low-rank solutions:

$$\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{y} - \mathbb{A}(\mathbf{X})\|_2^2 + \lambda \|\mathbf{X}\|_*.$$

An ISTA algorithm can be developed as follows. Note that

$$\partial \|\mathbf{X}\|_* = \sum_{i=1}^{\min(m,n)} \text{sign}(\sigma_i) \mathbf{u}_i \mathbf{v}_i^T$$

The proximal operator of nuclear norm  $\|\cdot\|_*$  is given by

$$\begin{aligned} \mathbf{X}^* &= \text{Prox}_{\gamma \|\cdot\|_*}(\mathbf{Z}) \\ &= \arg \min_{\mathbf{X}} \|\mathbf{X}\|_* + \frac{1}{2\gamma} \|\mathbf{X} - \mathbf{Z}\|_F^2 \\ &= \sum_{i=1}^{\min(m,n)} \eta(\sigma_i; \gamma) \mathbf{u}_i \mathbf{v}_i^T, \end{aligned}$$

where

- $\mathbf{Z} = \sum_{i=1}^{\min(m,n)} \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ , and
- $\eta(\sigma_i; \gamma) = \max(0, \sigma_i - \gamma)$  is the soft thresholding function.<sup>2</sup>

Let  $f(\mathbf{X}) = \frac{1}{2} \|\mathbf{y} - \mathbb{A}(\mathbf{X})\|_2^2$ . Note that

$$\begin{aligned} \nabla f(\mathbf{X}) &= -\mathbb{A}^*(\mathbf{y} - \mathbb{A}(\mathbf{X})) \\ &= -\mathbf{A}^T (\mathbf{y} - \text{Avec}(\mathbf{X})), \end{aligned}$$

where the linear map  $\mathbb{A}$  is represented by a matrix  $\mathbf{A}$ . Approximate the objective function locally by

$$\begin{aligned} f(\mathbf{X}) + \lambda \|\mathbf{X}\|_* &\approx f(\mathbf{X}^l) + \langle \mathbf{X} - \mathbf{X}^l, \nabla f(\mathbf{X}^l) \rangle + \frac{1}{2\tau^l} \|\mathbf{X} - \mathbf{X}^l\|_2^2 + \lambda \|\mathbf{X}\|_* \\ &= \frac{1}{2\tau^l} \left\| \mathbf{X} - \left( \mathbf{X}^l - \tau^l \nabla f(\mathbf{X}^l) \right) \right\|_2^2 + \lambda \|\mathbf{X}\|_* + c. \end{aligned}$$

2: See (7.2). Here, as by default singular values are non-negative, the form of soft-thresholding function can be simplified.

Hence, the ISTA iterations are given as

$$\mathbf{X}^{l+1} = \eta_\sigma \left( \mathbf{X}^l + \tau^l \mathbb{A}^* \left( \mathbf{y} - \mathbb{A}(\mathbf{X}^l) \right); \lambda \tau^l \right),$$

where the notation  $\eta_\sigma$  emphasizes that the soft-thresholding function acts on the singular values of the matrix (rather than the elements of the matrix). In practice, the step size  $\tau^l$  is fixed as  $\tau^l = \tau \in (0, 2/\sigma_{\max}(\mathbb{A}^* \mathbb{A})) = (0, 2/\sigma_{\max}(\mathbf{A}^T \mathbf{A}))$ .

## 8.4 Appendix

### 8.4.1 Proof of Eckart-Young-Mirsky Theorem

We first prove Eckart-Young-Mirsky theorem for spectral norm and then for Frobenius norm.

Consider singular value decomposition

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top = \sum_i \sigma_i \mathbf{u}_i \mathbf{v}_i^\top.$$

Consider the rank- $k$  approximation generated by SVD given by

$$\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^\top.$$

Consider any given (arbitrary) rank- $k$  approximation

$$\mathbf{B}_k = \mathbf{P}\mathbf{Q}^\top,$$

where both  $\mathbf{P} \in \mathbb{R}^{m \times k}$  and  $\mathbf{Q} \in \mathbb{R}^{n \times k}$  have  $k$  columns.

To prove Eckart–Young–Mirsky theorem for spectral norm, it suffices to show that

$$\|\mathbf{A} - \mathbf{A}_k\|_2 = \sigma_{k+1} \leq \|\mathbf{A} - \mathbf{B}_k\|_2.$$

Towards this end, note that  $\mathbf{Q}$  has only  $k$  columns. There must exist a vector  $\mathbf{w}$  such that

1.  $\mathbf{Q}^\top \mathbf{w} = \mathbf{0}$ ;
2.  $\mathbf{w} = \gamma_1 \mathbf{v}_1 + \cdots + \gamma_{k+1} \mathbf{v}_{k+1}$ ;
3.  $\|\mathbf{w}\|_2 = 1$ ; This implies that  $\|\gamma\|_2^2 = \sum_{i=1}^{k+1} \gamma_i^2 = 1$ .

Then

$$\begin{aligned} \|\mathbf{A} - \mathbf{B}_k\|_2^2 &\geq \|(\mathbf{A} - \mathbf{B}_k)\mathbf{w}\|_2^2 \\ &= \|\mathbf{A}\mathbf{w}\|_2^2 = \gamma_1^2 \sigma_1^2 + \cdots + \gamma_{k+1}^2 \sigma_{k+1}^2 \\ &\geq \sigma_{k+1}^2, \end{aligned}$$

which concludes the proof for the spectral norm.

For Frobenius norm, the goal is to establish

$$\|\mathbf{A} - \mathbf{A}_k\|_F^2 = \sum_{i=k+1}^n \sigma_i^2 \leq \|\mathbf{A} - \mathbf{B}_k\|_F^2.$$

Let  $\mathbf{A}' = \mathbf{A} - \mathbf{B}_k$  and  $\mathbf{A}'' = \mathbf{B}_k$ . Let  $\mathbf{A}'_i$  and  $\mathbf{A}''_i$  be the rank- $i$  approximation of  $\mathbf{A}'$  and  $\mathbf{A}''$  by using SVD, respectively. Then for any  $i, j \geq 1$ ,

$$\begin{aligned} \sigma_i(\mathbf{A}') + \sigma_j(\mathbf{A}'') &= \sigma_1(\mathbf{A}' - \mathbf{A}'_{i-1}) + \sigma_1(\mathbf{A}'' - \mathbf{A}''_{j-1}) \\ &\stackrel{(a)}{\geq} \sigma_1 \left( \mathbf{A} - (\mathbf{A}'_{i-1} + \mathbf{A}''_{j-1}) \right) \\ &\stackrel{(b)}{\geq} \sigma_1 (\mathbf{A} - \mathbf{A}_{i+j-2}) \\ &= \sigma_{i+j-1}(\mathbf{A}), \end{aligned}$$

where (a) comes from triangle inequality with the spectral norm:

$$\begin{aligned} & \sigma_1 \left( A - (A'_{i-1} + A''_{j-1}) \right) \\ &= \sigma_1 \left( A' - A'_{i-1} + A'' - A''_{j-1} \right) \\ &\leq \sigma_1 (A' - A'_{i-1}) + \sigma_1 (A'' - A''_{j-1}), \end{aligned}$$

and (b) holds by combining the fact that

$$\text{rank} (A'_{i-1} + A''_{j-1}) \leq \text{rank} (A_{i+j-2}),$$

and Eckart–Young–Mirsky theorem for spectral norm. Since  $\sigma_{k+1}(B_k) = 0$ , one has

$$\sigma_i(A - B_k) \geq \sigma_{k+i}(A).$$

This leads to

$$\begin{aligned} \|A - B_k\|_F^2 &= \sum_{i=1}^n \sigma_i(A - B_k)^2 \\ &\geq \sum_{i=k+1}^n \sigma_i(A)^2 \\ &= \|A - A_k\|_F^2, \end{aligned}$$

which concludes the proof.

# Constrained Convex Optimization

9

## 9.1 General Form

A constrained optimization problem of the form

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad (9.1)$$

subject to  $h_i(\mathbf{x}) \leq 0, i = 1, \dots, m,$   
 $\ell_i(\mathbf{x}) = 0, i = 1, \dots, r,$

where

- ▶ the objective function  $f_0$  is convex,
- ▶  $h_i$ 's are convex, and
- ▶  $\ell_i$ 's are affine, i.e., in the form of  $\mathbf{a}_i^T \mathbf{x} + b_i = 0$ .

In the following, we show that convex  $h_i$  and affine  $\ell_i$  lead to a *convex* feasible set.

### 9.1.1 Sublevel Sets

**Definition 9.1.1** (Sublevel Sets, a.k.a. Lower Contour Sets) *The  $\alpha$ -sublevel set of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is defined as*

$$\mathcal{C}_\alpha = \{\mathbf{x} \in \text{dom}(f) : f(\mathbf{x}) \leq \alpha\}.$$

**Lemma 9.1.1** *Sublevel sets of a convex function  $f$  are convex.*

*Proof.* We shall show that for all  $\mathbf{x}, \mathbf{y} \in \mathcal{C}_\alpha$ , it holds  $\lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in \mathcal{C}_\alpha$  for all  $\lambda \in [0, 1]$ .

By the definition of  $\mathcal{C}_\alpha$ ,  $f(\mathbf{x}) \leq \alpha$  and  $f(\mathbf{y}) \leq \alpha$ . By the convexity of  $f$ ,

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \leq \alpha,$$

which proves this lemma.  $\square$

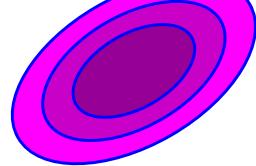


Figure 9.1: Sublevel sets of a convex function are convex

The converse of Lemma 9.1.1 is not true. That sublevel sets of a function  $f$  are convex does not imply that  $f$  is convex.

Consider the general form of constrained convex optimization (9.1). That  $h_i$  is convex ensures that

$$\{\mathbf{x} : h_i(\mathbf{x}) \leq 0\}$$

is convex by Lemma 9.1.1. Further note that

$$\{\mathbf{x} : \ell_i(\mathbf{x}) = 0\}$$

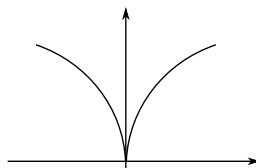


Figure 9.2: Converse of Lemma 9.1.1 is not true.

is convex. It can be concluded that

$$\left( \bigcap_i \{x : h_i(x) \leq 0\} \right) \cap \left( \bigcap_j \{x : \ell_j(x) = 0\} \right)$$

is convex.

## 9.2 Duality

Consider a general optimization problem

$$\begin{aligned} & \min_x f(x) \\ \text{s.t. } & h_i(x) \leq 0, \quad i = 1, \dots, m, \\ & \ell_j(x) = 0, \quad j = 1, \dots, r. \end{aligned}$$

where the objective function  $f$  needs not to be convex. Of course we pay special attention to the convex case. Its Lagrangian is defined as follows.

**Definition 9.2.1** (Lagrangian)

$$L(x, u, v) := f(x) + \sum_{i=1}^m u_i h_i(x) + \sum_{j=1}^r v_j \ell_j(x).$$

Here  $u \in \mathbb{R}^m$ ,  $v \in \mathbb{R}^r$ , and  $u \geq 0$ .

The corresponding Lagrange dual function is defined as

**Definition 9.2.2** (Lagrange Dual Function)

$$\begin{aligned} g(u, v) &:= \min_{x \in \mathbb{R}^n} L(x, u, v) \\ &= \min_{x \in \mathbb{R}^n} f(x) + \sum_{i=1}^m u_i h_i(x) + \sum_{j=1}^r v_j \ell_j(x). \end{aligned}$$

The following summarizes some important properties of Lagrangian and Lagrange dual function.

- For every feasible  $x$ ,  $L(x, u, v) \leq f(x)$ :

$$\begin{aligned} L(x, u, v) &= f(x) + \underbrace{\sum_{i=1}^m u_i h_i(x)}_{\leq 0} + \underbrace{\sum_{j=1}^r v_j \ell_j(x)}_{=0} \\ &\leq f(x). \end{aligned}$$

- Let  $\mathcal{X}$  denote the primal feasible set.

$$g(u, v) = \min_{x \in \mathbb{R}^n} L(x, u, v) \leq \min_{x \in \mathcal{X}} L(x, u, v) \leq f(x). \quad (9.2)$$

- $g(u, v)$  is concave in  $(u, v)$ . (See the lemma below.)

**Lemma 9.2.1** *The Lagrange dual function*

$$\begin{aligned} g(\mathbf{u}, \mathbf{v}) &= \min_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \mathbf{u}, \mathbf{v}) \\ &= \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + \sum_{i=1}^m u_i h_i(\mathbf{x}) + \sum_{j=1}^r v_j \ell_j(\mathbf{x}) \end{aligned}$$

is concave in  $(\mathbf{u}, \mathbf{v})$ .

We prove Lemma 9.2.1 by proving the following lemma.

**Lemma 9.2.2**

- Let  $f_\alpha(\mathbf{x})$  be concave. Then  $g(\mathbf{x}) = \inf_\alpha f_\alpha(\mathbf{x})$  is concave.
- Let  $f_\alpha(\mathbf{x})$  be convex. Then  $g(\mathbf{x}) = \sup_\alpha f_\alpha(\mathbf{x})$  is convex.

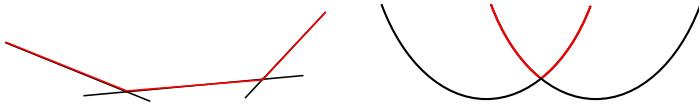


Figure 9.3: Pointwise Maximum

*Proof.* For any  $\lambda \in [0, 1]$ ,

$$\begin{aligned} g(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) &= \inf_\alpha f_\alpha(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \\ &\geq \inf_\alpha \lambda f_\alpha(\mathbf{x}) + (1 - \lambda) f_\alpha(\mathbf{y}) \\ &\geq \lambda \inf_\alpha f_\alpha(\mathbf{x}) + (1 - \lambda) \inf_\alpha f_\alpha(\mathbf{y}). \end{aligned}$$

where the first inequality comes from concavity of  $f_\alpha$ .  $\square$

*Proof of Lemma 9.2.1.* For any given  $\mathbf{x}$ ,  $L(\mathbf{x}, \mathbf{u}, \mathbf{v})$  is linear in  $(\mathbf{u}, \mathbf{v})$ , and hence concave in  $(\mathbf{u}, \mathbf{v})$ . Lemma 9.2.1 is then a direct consequence of Lemma 9.2.2.  $\square$

**Definition 9.2.3** (Lagrange Dual Problem) *Given the primal problem*

$$\begin{aligned} &\min_{\mathbf{x}} f(\mathbf{x}) \\ &\text{s.t. } h_i(\mathbf{x}) \leq 0, i = 1, \dots, m, \\ &\quad \ell_j(\mathbf{x}) = 0, j = 1, \dots, r, \end{aligned}$$

its Lagrange dual problem is

$$\begin{aligned} &\max_{\mathbf{u}, \mathbf{v}} g(\mathbf{u}, \mathbf{v}) \\ &\text{s.t. } \mathbf{u} \geq \mathbf{0}. \end{aligned}$$

- **Weak duality:** the dual optimal value  $g^*$  satisfies

$$f^* \geq g^*.$$

This is a direct consequence of (9.2).

- **Strong duality** is referred to as the case that

$$f^* = g^*.$$

- **Slater's condition:** if the primal is a convex problem (i.e.,  $f$  and  $g_i$ 's are convex and  $\ell_j$ 's are affine), and there exists at least one strictly feasible  $x \in \mathbb{R}^n$  satisfying

$$h_i(x) < 0, \forall i \in [m], \text{ and } \ell_j(x) = 0, \forall j \in [r],$$

then strong duality holds. (Proof is omitted.)

### 9.3 KKT Conditions

**Definition 9.3.1** (Karush-Kuhn-Tucker Conditions) *Given the optimization problem (9.1), the Karush-Kuhn-Tucker (KKT) conditions for optimal  $x^*, u^*, v^*$  are:*

- ▶  $0 \in \partial f(x^*) + \sum_{i=1}^m u_i^* \partial h_i(x^*) + \sum_{j=1}^r v_j^* \partial \ell_j(x^*)$ . (stationarity)
- ▶  $u_i^* h_i(x^*) = 0, \forall i$ . (complementary slackness)
- ▶  $h_i(x^*) \leq 0, \ell_j(x^*) = 0, \forall i, \forall j$ . (primal feasibility)
- ▶  $u_i^* \geq 0, \forall i$ . (dual feasibility)

KKT conditions are

- ▶ Always sufficient.
- ▶ Necessary under strong duality.

*Proof.* We prove sufficiency first and then necessity.

1. (Sufficiency) If  $x^*, u^*, v^*$  satisfy the KKT conditions, then

$$\begin{aligned} g(u^*, v^*) &= \min_x L(x, u^*, v^*) \\ &= f(x^*) + \sum_{i=1}^m u_i^* h_i(x^*) + \sum_{j=1}^r v_j^* \ell_j(x^*) \\ &= f(x^*), \end{aligned}$$

where the 2nd equality follows from stationarity, and the 3rd one follows from complementary slackness. Recall (9.2). This equality suggests the duality gap is zero and  $x^*, u^*$  and  $v^*$  are primal and dual optimal.

2. (Necessity) Suppose that the strong duality holds and that  $x^*$  and  $u^*, v^*$  are primal and dual solutions. One has

$$\begin{aligned} f(x^*) &= g(u^*, v^*) \\ &= \min_{x \in \mathbb{R}^n} f(x) + \sum_{i=1}^m u_i^* h_i(x) + \sum_{j=1}^r v_j^* \ell_j(x) \\ &\leq f(x^*) + \sum_{i=1}^m u_i^* h_i(x^*) + \sum_{j=1}^r v_j^* \ell_j(x^*) \\ &\leq f(x^*). \end{aligned}$$

In other words, all the inequalities are actually equalities. KKT conditions hold.

□

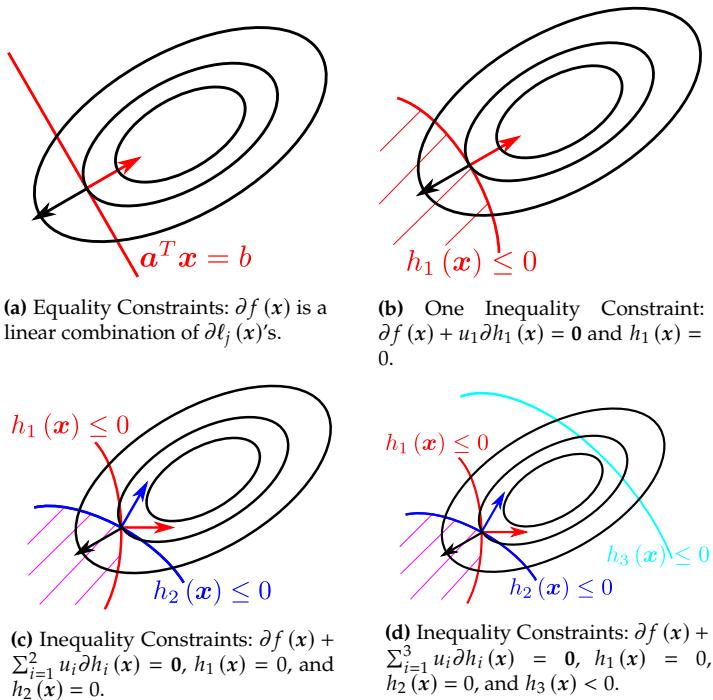


Figure 9.4: Geometric Interpretation of KKT Conditions

## 9.4 Examples

### Quadratic Programming with Equality Constraints

Consider a quadratic programming problem given by

$$\begin{aligned} \min_{\mathbf{x}} \quad & \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} \\ \text{s.t. } & \mathbf{C} \mathbf{x} = \mathbf{d}, \end{aligned}$$

where  $\mathbf{A} \succeq 0$ .

The corresponding Lagrangian is given by

$$L(\mathbf{x}, \mathbf{v}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + \langle \mathbf{v}, \mathbf{C} \mathbf{x} - \mathbf{d} \rangle.$$

By KKT conditions,  $\mathbf{x}$  is the minimizer if and only if

$$\left[ \begin{array}{cc} \mathbf{A} & \mathbf{C}^T \\ \mathbf{C} & \mathbf{0} \end{array} \right] \left[ \begin{array}{c} \mathbf{x} \\ \mathbf{v} \end{array} \right] = \left[ \begin{array}{c} -\mathbf{b} \\ \mathbf{d} \end{array} \right],$$

where the first set of linear equations come from the stationarity and the second set follows from the primal feasibility. The optimal  $\mathbf{x}^*$  can be obtained by solving the system of linear equations.

## Water Filling Power Allocation for Wireless Communications

Suppose independent parallel channels in a wireless communication system. The corresponding channel capacity is given by

$$\sum_i \log \left( 1 + \frac{x_i}{\alpha_i} \right),$$

where  $x_i \geq 0$  is the power allocated to the  $i$ -th channel and  $\alpha_i > 0$  is the given noise power for the  $i$ -th channel. The task is to find the optimal power allocation to maximize the channel capacity under the total power constraint  $\sum_i x_i = 1$ .

Towards this goal, the following constrained convex optimization problem can be formulated.

$$\begin{aligned} \min_{\mathbf{x}} \quad & - \sum_{i=1}^n \log(\alpha_i + x_i) \\ \text{s.t. } \mathbf{x} \geq \mathbf{0}, \quad & \mathbf{1}^T \mathbf{x} = 1. \end{aligned}$$

The corresponding Lagrangian is given by

$$L(\mathbf{x}, \mathbf{u}, v) = - \sum_{i=1}^n \log(\alpha_i + x_i) - \sum_i u_i x_i + v(\mathbf{1}^T \mathbf{x} - 1).$$

By KKT conditions, the optimal solution satisfies

- $-1/(\alpha_i + x_i^\star) - u_i^\star + v^\star = 0, \forall i.$
- $u_i^\star x_i^\star = 0, \forall i.$
- $\mathbf{x}^\star \geq \mathbf{0}, \mathbf{1}^T \mathbf{x}^\star = 1, \mathbf{u}^\star \geq \mathbf{0}.$

From the first two conditions, one has

$$\begin{aligned} 1/(\alpha_i + x_i^\star) &\leq v^\star, \text{ and} \\ x_i^\star (v^\star - 1/(\alpha_i + x_i^\star)) &= 0. \end{aligned}$$

Therefore, the optimal power allocation is given by

$$x_i^\star = \max(0, 1/v^\star - \alpha_i)$$

where  $v^\star$  is chosen such that

$$\sum_{i=1}^n \max(0, 1/v^\star - \alpha_i) = 1.$$

## Constrained $\ell_1$ -minimization

Consider the optimization problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & \|\mathbf{x}\|_1 \\ \text{s.t. } \mathbf{A}\mathbf{x} &= \mathbf{y}. \end{aligned}$$

The corresponding Lagrangian is given by

$$L(\mathbf{x}, \mathbf{v}) = \|\mathbf{x}\|_1 + \langle \mathbf{v}, A\mathbf{x} - \mathbf{y} \rangle.$$

The corresponding Lagrange dual function is then given by

$$\begin{aligned} g(\mathbf{v}) &= \min_{\mathbf{x}} \|\mathbf{x}\|_1 + \langle \mathbf{v}, A\mathbf{x} - \mathbf{y} \rangle \\ &= -\mathbf{v}^\top \mathbf{y} + \min_{\mathbf{x}} \|\mathbf{x}\|_1 + \langle A^\top \mathbf{v}, \mathbf{x} \rangle. \end{aligned}$$

Note that

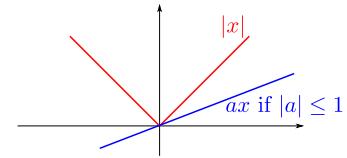
$$\min_{x_i} |x_i| + (A^\top \mathbf{v})_i x_i = \begin{cases} 0 & \text{if } |(A^\top \mathbf{v})_i| \leq 1 \\ -\infty & \text{if } |(A^\top \mathbf{v})_i| > 1. \end{cases}$$

One has

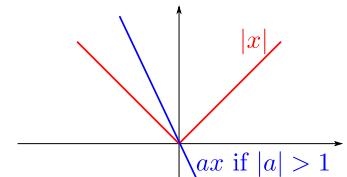
$$g(\mathbf{v}) = \begin{cases} -\mathbf{v}^\top \mathbf{y} & \text{if } \|A^\top \mathbf{v}\|_\infty \leq 1 \\ -\infty & \text{if } \|A^\top \mathbf{v}\|_\infty > 1 \end{cases}$$

The dual problem is given by

$$\begin{aligned} \max_{\mathbf{v}} & -\mathbf{v}^\top \mathbf{y} \\ \text{s.t. } & \|A^\top \mathbf{v}\|_\infty \leq 1. \end{aligned}$$



(a) Case 1:  $\min_{\mathbf{x}} |\mathbf{x}| + a\mathbf{x} = 0$ .



(b) Case 1:  $\min_{\mathbf{x}} |\mathbf{x}| + a\mathbf{x} = -\infty$ .

**Figure 9.5:** An illustration of the dual function for constrained  $\ell_1$ -minimization

**ADMM:** Alternating Direction Method of Multipliers.

**Key reference:** Boyd, Stephen, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. “Distributed optimization and statistical learning via the alternating direction method of multipliers.” *Foundations and Trends in Machine learning* 3, no. 1 (2011): 1-122.

|                                        |    |
|----------------------------------------|----|
| <b>10.1 Precursors . . . . .</b>       | 73 |
| 10.1.1 Dual Ascent Method . . . . .    | 73 |
| 10.1.2 Method of Multipliers . . . . . | 74 |
| <b>10.2 ADMM . . . . .</b>             | 75 |
| 10.2.1 Optimality Conditions . . . . . | 77 |

## 10.1 Precursors

### 10.1.1 Dual Ascent Method

Consider the convex optimization problem

$$\begin{aligned} & \min f(x) \\ & \text{s.t. } Ax = b \end{aligned} \tag{10.1}$$

Its Lagrangian is

$$L(x, v) = f(x) + v^\top (Ax - b).$$

The dual function is given by

$$g(v) = \inf_x L(x, v).$$

The dual problem is

$$\max g(v).$$

Apply dual ascent method for the dual problem. One updates the dual variable via

$$v^{k+1} = v^k + \alpha^k \nabla g(v^k).$$

From the definition of the dual function, one has

$$\begin{aligned} x^{k+1} &= \arg \min_x L(x, v^k) \\ v^{k+1} &= v^k + \alpha^k \nabla_v L(x^{k+1}, v^k) \\ &= v^k + \alpha^k (Ax^{k+1} - b). \end{aligned}$$

With appropriate chosen  $\alpha^k$ ,  $g(v^{k+1}) > g(v^k)$ . Dual ascent method converges under certain assumptions.

However, the required assumptions do not hold in many applications. For example, when  $f(x) = \|x\|_1$ , the  $x$ -update step fails as Lagrangian  $L$  is unbounded below for some choices of  $v$ .

### 10.1.2 Method of Multipliers

Instead of Lagrangian, we consider the augmented Lagrangian for (10.1) given by

$$L_\rho(\mathbf{x}, \mathbf{v}) = f(\mathbf{x}) + \mathbf{v}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) + \frac{\rho}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2,$$

where  $\rho > 0$  is called penalty parameter.

Augmented Lagrangian can be viewed as Lagrangian of the problem

$$\begin{aligned} \min_{\mathbf{x}} & f(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 \\ \text{s.t. } & \mathbf{A}\mathbf{x} = \mathbf{b}, \end{aligned}$$

which is equivalent to (10.1).

The associated dual function is given by

$$g_\rho(\mathbf{v}) = \inf_{\mathbf{x}} L_\rho(\mathbf{x}, \mathbf{v}).$$

Apply dual ascent method. The *method of multipliers* algorithm is given by

$$\begin{aligned} \mathbf{x}^{k+1} &= \arg \min_{\mathbf{x}} L_\rho(\mathbf{x}, \mathbf{v}^k) \\ \mathbf{v}^{k+1} &= \mathbf{v}^k + \rho (\mathbf{A}\mathbf{x}^{k+1} - \mathbf{b}). \end{aligned}$$

Note the fixed step size  $\rho$ .

The choice of the step size  $\rho$  can be motivated as follows. For simplicity, assume  $f$  is differentiable. The optimality conditions for (10.1) are primal and dual feasibility

$$\mathbf{A}\mathbf{x}^* - \mathbf{b} = \mathbf{0}, \quad \nabla f(\mathbf{x}^*) + \mathbf{A}^T \mathbf{v}^* = \mathbf{0}.$$

As  $\mathbf{x}^{k+1}$  minimizes  $L_\rho(\mathbf{x}, \mathbf{v}^k)$ , it holds that

$$\begin{aligned} \mathbf{0} &= \nabla_{\mathbf{x}} L_\rho(\mathbf{x}^{k+1}, \mathbf{v}^k) \\ &= \nabla_{\mathbf{x}} f(\mathbf{x}^{k+1}) + \mathbf{A}^T (\mathbf{v}^k + \rho (\mathbf{A}\mathbf{x}^{k+1} - \mathbf{b})) \\ &= \nabla_{\mathbf{x}} f(\mathbf{x}^{k+1}) + \mathbf{A}^T \mathbf{v}^{k+1}. \end{aligned}$$

Using  $\rho$  as step size,  $(\mathbf{x}^{k+1}, \mathbf{v}^{k+1})$  is dual feasible.

Compared with dual ascent method, the method of multipliers converges under much more general conditions. However, the  $\mathbf{x}$ -update step may be difficult to compute.

## 10.2 ADMM

ADMM solves problems in the form

$$\begin{aligned} & \min f(x) + g(z) \\ & \text{s.t. } Ax + Bz = c \end{aligned}$$

The corresponding augmented Lagrangian is given by

$$\begin{aligned} L_\rho(x, z, v) &= f(x) + g(z) + v^\top (Ax + Bz - c) \\ &+ \frac{\rho}{2} \|Ax + Bz - c\|_2^2 \\ &= f(x) + g(z) + \frac{\rho}{2} \|Ax + Bz - c + v/\rho\|_2^2 \\ &- \|v\|_2^2/(2\rho) \end{aligned}$$

for some  $\rho > 0$ . ADMM is an iterative algorithm with iterations:

$$\begin{aligned} x^{k+1} &= \arg \min_x L_\rho(x, z^k, v^k), \\ z^{k+1} &= \arg \min_z L_\rho(x^{k+1}, z, v^k), \\ v^{k+1} &= v^k + \rho (Ax^{k+1} + Bz^{k+1} - c). \end{aligned}$$

**Example 10.2.1** (Solving Lasso by ADMM) Consider the Lasso problem:

$$\min \frac{1}{2} \|y - Ax\|_2^2 + \lambda \|x\|_1.$$

An ADMM reformulation is given by

$$\begin{aligned} & \min \frac{1}{2} \|y - Ax\|_2^2 + \lambda \|z\|_1 \\ & \text{s.t. } x = z. \end{aligned}$$

The corresponding ADMM iterations are

$$\begin{aligned} x^{k+1} &= \arg \min_x \frac{1}{2} \|y - Ax\|_2^2 + \frac{\rho}{2} \left\| x - z^k + \frac{v^k}{\rho} \right\|_2^2 \\ z^{k+1} &= \arg \min_z \lambda \|z\|_1 + \frac{\rho}{2} \left\| x^{k+1} - z + \frac{v^k}{\rho} \right\|_2^2 \\ v^{k+1} &= v^k + \rho (x^{k+1} - z^{k+1}). \end{aligned}$$

### Remark 10.2.1

1. Each step of the above ADMM has a closed form solution.
2. ADMM is different from the method of multipliers. In the method

of multipliers, the variable update is

$$\begin{aligned} & (\mathbf{x}^{k+1}, \mathbf{z}^{k+1}) \\ &= \arg \min_{\mathbf{x}, \mathbf{z}} L_\rho(\mathbf{x}, \mathbf{z}, \mathbf{v}^k) \\ &= \arg \min_{\mathbf{x}, \mathbf{z}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{z}\|_1 + \frac{\rho}{2} \left\| \mathbf{x} - \mathbf{z} + \frac{\mathbf{v}^k}{\rho} \right\|_2^2. \end{aligned}$$

As  $\mathbf{x}$  and  $\mathbf{z}$  are not separate, there is no closed form for the variable update step.

3. The above ADMM converges for an arbitrary penalty constant  $\rho > 0$ . In ISTA (Section 7.2.3), the step size  $\tau$  needs to be in the range  $(0, 2/\sigma_{\max}(\mathbf{A}^\top \mathbf{A}))$  to guarantee a convergence.

**Example 10.2.2** (Constrained Lasso Solved by ADMM) Consider the problem

$$\begin{aligned} & \min \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1 \\ & \text{s.t. } \mathbf{B}\mathbf{x} \leq \mathbf{c}. \end{aligned}$$

An ADMM reformulation is given by

$$\begin{aligned} & \min \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{z}_1\|_1 + \delta_{\geq 0}(\mathbf{z}_2) \\ & \text{s.t. } \begin{bmatrix} \mathbf{I} \\ \mathbf{B} \end{bmatrix} \mathbf{x} + \begin{bmatrix} -\mathbf{I} \\ \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{c} \end{bmatrix}, \end{aligned}$$

where

$$\delta_{\geq 0}(z) = \begin{cases} 0 & \text{if } z \geq 0, \\ \infty & \text{if } z < 0. \end{cases}$$

Write the above constraint as  $\mathbf{B}'\mathbf{x} + \mathbf{D}'\mathbf{z} = \mathbf{c}'$ . The ADMM iterations are

$$\begin{aligned} \mathbf{x}^{k+1} &= \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \frac{\rho}{2} \left\| \mathbf{B}'\mathbf{x} + \mathbf{D}'\mathbf{z}^k - \mathbf{c}' + \frac{\mathbf{v}^k}{\rho} \right\|_2^2 \\ \mathbf{z}^{k+1} &= \arg \min_{\mathbf{z}_1, \mathbf{z}_2} \lambda \|\mathbf{z}_1\|_1 + \frac{\rho}{2} \left\| \mathbf{x}^{k+1} - \mathbf{z}_1 + \frac{\mathbf{v}_1^k}{\rho} \right\|_2^2 \\ &\quad + \delta_{\geq 0}(\mathbf{z}_2) + \frac{\rho}{2} \left\| \mathbf{B}\mathbf{x}^{k+1} + \mathbf{z}_2 - \mathbf{c} + \frac{\mathbf{v}_2^k}{\rho} \right\|_2^2 \\ \mathbf{v}^{k+1} &= \mathbf{v}^k + \rho \left( \mathbf{B}'\mathbf{x}^{k+1} + \mathbf{D}'\mathbf{z}^{k+1} - \mathbf{c}' \right). \end{aligned}$$

Each step admits a closed form.

### 10.2.1 Optimality Conditions

The necessary and sufficient conditions for optimality are primal feasibility

$$\mathbf{A}\mathbf{x}^* + \mathbf{B}\mathbf{z}^* - \mathbf{c} = \mathbf{0}. \quad (10.2)$$

and dual feasibility

$$\mathbf{0} \in \partial f(\mathbf{x}^*) + \mathbf{A}^T \mathbf{v}^* \quad (10.3)$$

$$\mathbf{0} \in \partial g(\mathbf{z}^*) + \mathbf{B}^T \mathbf{v}^*. \quad (10.4)$$

1. Consider primal feasibility. Define *primal residual* as

$$\mathbf{r}^{k+1} = \mathbf{A}\mathbf{x}^{k+1} + \mathbf{B}\mathbf{z}^{k+1} - \mathbf{c}.$$

We expect  $\mathbf{r}^k$  converges to zero.

2. Consider (10.4) in dual feasibility. It turns out that  $\mathbf{z}^{k+1}$  and  $\mathbf{v}^{k+1}$  always satisfy (10.4):

$$\begin{aligned} \mathbf{0} &\in \partial g(\mathbf{z}^{k+1}) + \mathbf{B}^T \mathbf{v}^k + \rho \mathbf{B}^T (\mathbf{A}\mathbf{x}^{k+1} + \mathbf{B}\mathbf{z}^{k+1} - \mathbf{c}) \\ &= \partial g(\mathbf{z}^{k+1}) + \mathbf{B}^T \mathbf{v}^{k+1}. \end{aligned}$$

3. Consider (10.3) in dual feasibility. The situation about  $\mathbf{x}^{k+1}$  is different. By definition  $\mathbf{x}^{k+1}$  minimizes  $L_\rho(\mathbf{x}, \mathbf{z}^k, \mathbf{v}^k)$ . It holds that

$$\begin{aligned} \mathbf{0} &\in \partial f(\mathbf{x}^{k+1}) + \mathbf{A}^T \mathbf{v}^k + \rho \mathbf{A}^T (\mathbf{A}\mathbf{x}^{k+1} + \mathbf{B}\mathbf{z}^k - \mathbf{c}) \\ &= \partial f(\mathbf{x}^{k+1}) + \mathbf{A}^T (\mathbf{v}^k + \rho \mathbf{r}^{k+1} + \rho \mathbf{B}(\mathbf{z}^k - \mathbf{z}^{k+1})) \\ &= \partial f(\mathbf{x}^{k+1}) + \mathbf{A}^T \mathbf{v}^{k+1} + \rho \mathbf{A}^T \mathbf{B}(\mathbf{z}^k - \mathbf{z}^{k+1}). \end{aligned}$$

Or equivalently

$$\rho \mathbf{A}^T \mathbf{B}(\mathbf{z}^{k+1} - \mathbf{z}^k) \in \partial f(\mathbf{x}^{k+1}) + \mathbf{A}^T \mathbf{v}^{k+1}.$$

Define *dual residual* as

$$\mathbf{s}^{k+1} = \rho \mathbf{A}^T \mathbf{B}(\mathbf{z}^{k+1} - \mathbf{z}^k).$$

Under mild conditions, ADMM converges:

- Primal residual convergence:  $\mathbf{r}^k \rightarrow \mathbf{0}$  as  $k \rightarrow \infty$ .
- Objective convergence:  $f(\mathbf{x}^k) + g(\mathbf{z}^k) \rightarrow p^*$  as  $k \rightarrow \infty$ .
- Dual variable convergence:  $\mathbf{v}^k \rightarrow \mathbf{v}^*$  as  $k \rightarrow \infty$ .

In practice, a reasonable criterion of terminating ADMM iterations is that the primal and dual residuals are small, i.e.,

$$\|\mathbf{r}^k\|_2 \leq \epsilon^{\text{pri}}, \quad \|\mathbf{s}^k\|_2 \leq \epsilon^{\text{dual}}.$$

## 11.1 Motivation

Throughout these notes, we have encountered a wide range of inference and recovery problems, all of which could be formulated as optimization problems. In the case of least squares estimation, treated in Chapter 3, the optimal solution to the optimization problem could be found in terms of a closed-form solution. In later chapters, as the optimization problems became more elaborate, such as for example in case of Lasso (treated in Chapter 7), the optimal solution was no longer available in closed-form. Instead, we pursued optimal solutions through iterative means, such as the Iterative Shrinkage Thresholding Algorithm (ISTA). We repeat the problems and solutions here for reference:

$$\hat{x} = \arg \min_x \|y - Ax\|_2^2 \implies \hat{x} = (A^T A)^{-1} A^T y \quad (11.1)$$

$$\hat{x} = \arg \min_x \|y - Ax\|_2^2 + \lambda \|x\|_1 \implies x^{l+1} = \text{Prox}_{\gamma \lambda \|\cdot\|_1} \left( x^l + \gamma A^T (y - Ax^l) \right) \quad (11.2)$$

In both cases, the expressions were *batch* solutions, meaning that the full batch of data, contained in the vector  $y$  and matrix  $A$  are leveraged either to compute the closed-form solution (11.1), or at every iteration (11.2), when evaluating the ISTA update. This is feasible when the amount of data, either in terms of problem dimension (i.e., the number of columns of  $A$ ), or in terms of available samples (number of rows of  $A$ ) is moderate when compared to the computational and storage resources. In many large-scale settings, the available data grows faster than computational resources, eventually reaching the point where we are no longer able to process the entire batch of data at once.

Motivated by these considerations, in this chapter, we resort to *stochastic* optimization techniques, which process smaller *mini-batches* of data at a time. We begin by deriving a stochastic version of least optimization, and will subsequently generalize to allow for arbitrary loss functions.

## 11.2 Stochastic Least Squares

### 11.2.1 Gradient Descent for LS

Let us revisit the least squares problem (11.1) and derive an alternative algorithm for pursuing its minimizer. We recall that the squared Euclidean norm of a vector is given by the sum of its elements squared, yielding the alternative expression:

$$\hat{x} = \arg \min_x \|y - Ax\|_2^2 = \arg \min_x \sum_{m=1}^M (y_m - a_m^T x)^2 \quad (11.3)$$

|                                                     |    |
|-----------------------------------------------------|----|
| 11.1 Motivation . . . . .                           | 78 |
| 11.2 Stochastic Least Squares . . . . .             | 78 |
| 11.2.1 Gradient Descent for LS . . . . .            | 78 |
| 11.2.2 Stochastic Gradient Descent for LS . . . . . | 79 |

where we denote by  $a_m$  the *column* vector corresponding to the  $m$ -th row of the data matrix  $A$ . The minimizing argument in (11.3) is invariant to scaling of the cost, and hence we can normalize equivalently:

$$\hat{x} = \arg \min_x \frac{1}{M} \sum_{m=1}^M (y_m - a_m^T x)^2 \quad (11.4)$$

While a direct closed-form expression for  $\hat{x}$  is available via (11.1), it will be instructive to derive a gradient based iterative algorithm and compare the computational cost. Using a constant step-size  $\alpha$  (i.e., no line search), we find:

$$x^{l+1} = x^l - \alpha \nabla \left( \frac{1}{M} \sum_{m=1}^M (y_m - a_m^T x)^2 \right) = x^l + \frac{\alpha}{M} \sum_{m=1}^M a_m (y_m - a_m^T x) \quad (11.5)$$

**Example 11.2.1** (Per-iteration complexity of Least-Squares Gradient Descent) We note that computing the gradient update in (11.5) involves the averaging of  $M$  smaller operations of the form  $a_m(y_m - a_m^T x)$ . Disregarding additions, evaluating  $y_m - a_m^T x$  requires  $N$  multiplications, where  $N$  denotes the dimension of the feature space (and hence also of the parameter vector  $x$ ). To find  $a_m(y_m - a_m^T x)$ , we require  $N$  more multiplications. All this is repeated  $M$  times to find the gradient  $\frac{1}{M} \sum_{m=1}^M a_m (y_m - a_m^T x)$ . In summary, we find:

$$\text{Per-iteration complexity of LS GD} \sim MN \quad (11.6)$$

We observe that the complexity of evaluating a single step of gradient descent, applied to the least mean squares cost (11.4) scales linearly in both the dimension  $N$ , and the size of the data set  $M$ . Of course this operation needs to be repeated a total of  $L$  times, where  $L$  denotes the total number of iterations. Hence:

$$\text{Total complexity of LS GD} \sim MNL \quad (11.7)$$

A natural question given the above consideration is then the following: Will the complexity (both per iteration, and in total) of a learning algorithm always scale at least linearly with the number of available samples? In either case, the answer would have far-reaching implications for large-scale data processing. Perhaps surprisingly, the answer is no. We will explore this in the context of least squares estimation first, and subsequently generalize.

### 11.2.2 Stochastic Gradient Descent for LS

A common strategy in engineering and statistics for dealing with intractable quantities of data is to sample subsets of the data and compute stochastic estimates of the quantity of interest. In the context of least

squares gradient descent, the quantity of interest is the gradient:

$$\nabla f(x) = \frac{1}{M} \sum_{m=1}^M a_m(y_m - a_m^T x) \quad (11.8)$$

Suppose now, at every iteration instead of evaluating the full gradient  $\nabla f(x)$  by averaging  $M$  gradients of the form  $a_m(y_m - a_m^T x)$ , we subsample the total data set  $\{x_m, y_m\}_{m=1}^M$ . In particular, we will sample  $B$  times from  $\{x_m, y_m\}_{m=1}^M$ , obtaining the *mini-batch*  $\{x_b^t, y_b^t\}_{b=1}^B$ .