

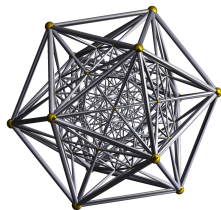
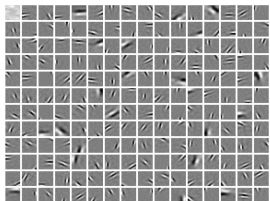
Computational Principles for High-dim Data Analysis

(Lecture Fifteen)

Yi Ma

EECS Department, UC Berkeley

October 21, 2021



Nonconvex Methods for Low-Dimensional Models

Dictionary Learning

- 1 Motivating Examples for Nonconvex Problems
- 2 Nonlinearity, Nonconvexity, and Symmetry
- 3 Rotational Symmetry (brief)
- 4 Discrete Symmetry: Dictionary Learning

“The mathematical sciences particularly exhibit order, symmetry, and limitations; and these are the greatest forms of the beautiful.”

– Aristotle, *Metaphysica*

Example: Magnetic Resonance Imaging

Simplified linear measurement model for MRI:

$$y = \mathcal{F}[I](\mathbf{u}) = \int_{\mathbf{v}} I(\mathbf{v}) \exp(-i 2\pi \mathbf{u}^* \mathbf{v}) d\mathbf{v} \in \mathbb{C}. \quad (1)$$

Real physical measurements as modulus:

$$y = |\mathcal{F}[I](\mathbf{u})| \in \mathbb{R}_+. \quad (2)$$

Fourier phase retrieval from multiple **nonlinear** real measurements:

$$\mathbf{y} = \left| \mathcal{F} \left(\begin{array}{c} \mathbf{x} \\ \text{unknown signal} \end{array} \right) \right| \in \mathbb{R}_+^m. \quad (3)$$



Example: Low-rank Matrix Completion

We observe:

$$\mathbf{Y} = \mathcal{P}_\Omega \left[\mathbf{X} \right].$$

Observed ratings Complete ratings

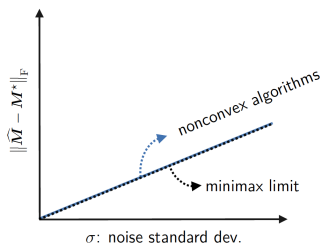
Matrix completion

via **bilinear** low-rank factorization¹:

$$\min_{\mathbf{U}, \mathbf{V}} f(\mathbf{U}, \mathbf{V}) = \sum_{(i,j) \in \Omega} [(UV^*)_{i,j} - \mathbf{Y}_{i,j}]^2 + \underbrace{\frac{\lambda}{2} \|\mathbf{U}\|_F^2 + \frac{\lambda}{2} \|\mathbf{V}\|_F^2}_{\text{reg}(\mathbf{U}, \mathbf{V})}.$$

$$\|\mathbf{M}\|_* = \min_{\mathbf{M} = \mathbf{UV}^*} \frac{\lambda}{2} \|\mathbf{U}\|_F^2 + \frac{\lambda}{2} \|\mathbf{V}\|_F^2$$

minimax limit	$\sigma\sqrt{n/p}$
nonconvex algorithms	$\sigma\sqrt{n/p}$ (optimal!)



¹figure courtesy from the lecture by Prof. Yuxin Chen of Princeton

Example: Dictionary for Image Representation

Image processing
(e.g. denoising or super-resolution)
against a known sparsifying dictionary:

$$I_{\text{noisy}} = \underset{\text{dictionary}}{\mathbf{A}} \times \underset{\text{sparse}}{\mathbf{x}} + \underset{\text{noise}}{\mathbf{z}}. \quad (4)$$



Dictionary learning: the motifs or atoms of the dictionary are **unknown**:

$$\underset{\text{data}}{\mathbf{Y}} = \underset{\text{dictionary}}{\mathbf{A}} \underset{\text{sparse}}{\mathbf{X}}. \quad (5)$$

- Band-limited signals: $\mathbf{A} = \mathbf{F}$, the Fourier transform;
- Piecewise smooth signals: $\mathbf{A} = \mathbf{W}$, the wavelet transforms;
- Natural images $\mathbf{A} = ?$ (How to **learn** \mathbf{A} from the data \mathbf{Y} ?)

Challenges of Nonconvex Optimization – Pessimistic Views

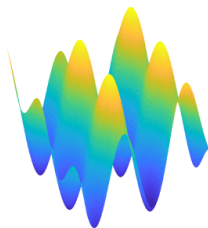
Consider the problem of minimizing a general nonlinear function:

$$\min_z \varphi(\mathbf{z}), \quad \mathbf{z} \in \mathcal{C}. \quad (6)$$

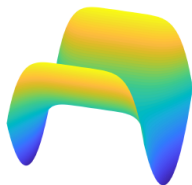
In **the worst case**, even finding a *local* minimizer can be NP-hard².

Hence typically people seek to work with relatively benign functions with benign guarantees (Chapter 9):

- ① convergence to some critical point $\bar{\mathbf{z}}$ such that $\nabla\varphi(\bar{\mathbf{z}}) = \mathbf{0}$;
- ② or convergence to some local minimizer $\nabla^2\varphi(\bar{\mathbf{z}}) \succeq \mathbf{0}$.



Spurious local minimizers

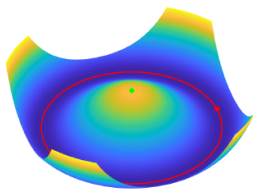


Flat saddle points

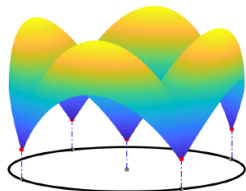
²Some NP-complete problems in quadratic and nonlinear programming, K.G Murty and S. N. Kabadi, 1987

Opportunities – Optimistic Views

However, nonconvex problems that arise from natural physical, geometrical, or statistical origins typically have **nice** structures, in terms of **symmetries!**



Rotational symmetry



Discrete symmetry

The function φ is **invariant** under certain group action:

- for phase recovery, invariant under a continuous rotation:

$$\varphi(e^{i\theta} \mathbf{x}) = \varphi(\mathbf{x}), \quad \forall \theta \in [0, 2\pi) = \mathbb{S}^1,$$

- for dictionary learning, invariant under signed permutations:

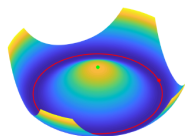
$$\varphi((\mathbf{A}, \mathbf{X})) = \varphi((\mathbf{A}\mathbf{\Pi}, \mathbf{\Pi}^* \mathbf{X})), \quad \forall \mathbf{\Pi} \in \text{SP}(n),$$

Optimization under Symmetry

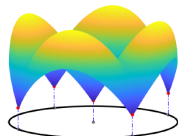
Definition (Symmetric Function)

Let \mathbb{G} be a group acting on \mathbb{R}^n . A function $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^{n'}$ is \mathbb{G} -symmetric if for all $\mathbf{z} \in \mathbb{R}^n$, $\mathbf{g} \in \mathbb{G}$, $\varphi(\mathbf{g} \circ \mathbf{z}) = \varphi(\mathbf{z})$.

Most symmetric objective functions that arise in structure signal recovery **do not** have spurious local minimizers or flat saddles.



Rotational symmetry



Discrete symmetry

Slogan 1: the (only!) local minimizers are symmetric versions of the ground truth.

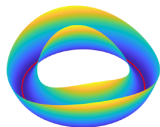
Slogan 2: any local critical point has negative curvature in directions that break symmetry.

Taxonomy of Symmetric Nonconvex Problems

Nonconvex Problems with Rotational Symmetries

Eigenspace Computation

Compute the principal subspace of a symmetric matrix.

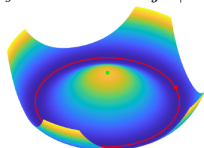


$$\min_{X \cdot X=I} -\frac{1}{2} \text{trace}[X^* A X].$$

Symmetry: $X \mapsto XR$
 $\mathbb{G} = O(r)$

Generalized Phase Retrieval

Recover a complex vector x_o from magnitude measurements $y = |Ax_o|$.

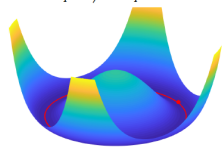


$$\min_x \frac{1}{2} \|y^2 - |Ax|^2\|_2^2.$$

Symmetry: $x \mapsto xe^{i\phi}$
 $\mathbb{G} = S^1 \cong O(2)$

Matrix Recovery

Recover a low-rank matrix $X = UV^*$ from incomplete/corrupted observations



$$\min_{U, V} \mathcal{L}(Y - A[UV^*]) + \rho(U, V).$$

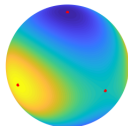
Symmetry: $(U, V) \mapsto (U\Gamma, V\Gamma^{-*})$
 $\mathbb{G} = GL(r)$ or $\mathbb{G} = O(r)$

Taxonomy of Symmetric Nonconvex Problems

Nonconvex Problems with Discrete Symmetries

Eigenvector Computation

Maximize a quadratic form over the sphere.

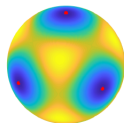


$$\max_{\mathbf{x} \in \mathbb{S}^{n-1}} \frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x}.$$

Symmetry: $\mathbf{x} \mapsto -\mathbf{x}$
 $\mathbb{G} = \{\pm 1\}$

Tensor Decomposition

Determine components \mathbf{a}_i of an orthogonal decomposable tensor $\mathbf{T} = \sum_i \mathbf{a}_i \otimes \mathbf{a}_i \otimes \mathbf{a}_i \otimes \mathbf{a}_i$

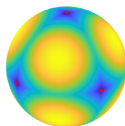


$$\max_{\mathbf{X} \in \mathcal{O}(n)} \sum_i \mathbf{T}(\mathbf{x}_i, \mathbf{x}_i, \mathbf{x}_i, \mathbf{x}_i).$$

Symmetry: $\mathbf{X} \mapsto \mathbf{X}\Gamma$
 $\mathbb{G} = \mathcal{P}(n)$

Dictionary Learning

Approximate a given matrix \mathbf{Y} as $\mathbf{Y} = \mathbf{A}\mathbf{X}$, with \mathbf{X} sparse

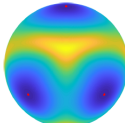


$$\min_{\mathbf{A} \in \mathcal{A}, \mathbf{X}} \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1.$$

Symmetry: $(\mathbf{A}, \mathbf{X}) \mapsto (\mathbf{A}\Gamma, \mathbf{X}\Gamma^*)$
 $\mathbb{G} = \text{SP}(n)$

Short-and-Sparse Deconvolution

Recover a short \mathbf{a} and a sparse \mathbf{x} from their convolution $\mathbf{y} = \mathbf{a} \otimes \mathbf{x}$.



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

Symmetry: $(\mathbf{a}, \mathbf{x}) \mapsto (\alpha s_\tau[\mathbf{a}], \alpha^{-1} s_{-\tau}[\mathbf{x}])$
 $\mathbb{G} = \mathbb{Z}_n \times \mathbb{R}_* \text{ or } \mathbb{G} = \mathbb{Z}_n \times \{\pm 1\}$

Dictionary Learning: the Minimal Case

Dictionary Learning with **one sparsity**:

$$\mathbf{Y} = \mathbf{A}_o \mathbf{X}_o. \quad (7)$$

data orthogonal dictionary 1-sparsity coefficients

Signed permutation symmetry:

$$\mathbf{Y} = \mathbf{A}_o \mathbf{X}_o = \mathbf{A}_o \mathbf{\Gamma} \mathbf{\Gamma}^* \mathbf{X}_o, \quad \forall \mathbf{\Gamma} \in \text{SP}(n).$$

Search for an orthogonal \mathbf{A} such that $\mathbf{A}^* \mathbf{Y}$ is *as sparse as possible*:

$$\min h(\mathbf{A}^* \mathbf{Y}) \quad \text{such that} \quad \mathbf{A} \in \text{O}(m), \quad (8)$$

where $h(\mathbf{X}) = \sum_{ij} h(\mathbf{X}_{ij})$ is a function that promotes sparsity.

Find One Atom at a Time

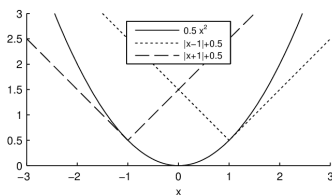
Take h to be **the Huber function**:

$$h_\lambda(x) = \begin{cases} \lambda|x| - \lambda^2/2 & |x| > \lambda, \\ x^2/2 & |x| \leq \lambda. \end{cases} \quad (9)$$

This can be viewed as a differentiable surrogate for the ℓ^1 norm.

For the dictionary $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_m]$, find the columns \mathbf{a}_i one at a time:

$$\min \varphi(\mathbf{a}) \doteq h_\lambda(\mathbf{a}^* \mathbf{Y}) \quad \text{such that} \quad \mathbf{a} \in \mathbb{S}^{m-1}. \quad (10)$$



Dictionary Learning: the Simplest Case

WLOG, assume $\mathbf{A}_o = \mathbf{I}$, and $\mathbf{X}_o = \mathbf{I}$ (uniformly random sampling).

$$\min \varphi(\mathbf{a}) \doteq h_\lambda(\mathbf{a}) \quad \text{such that} \quad \mathbf{a} \in \mathbb{S}^{m-1}. \quad (11)$$

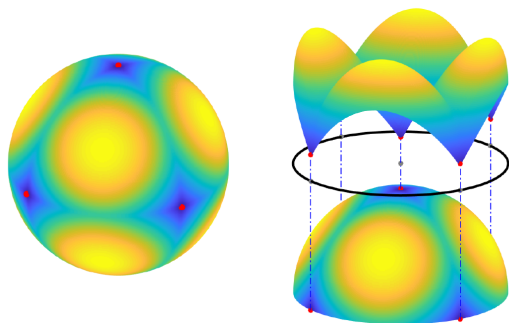


Figure: $h_\lambda(\mathbf{u})$ as a function on the sphere \mathbb{S}^2 .

First Order Characteristics of The Simplest Case

Critical Points of φ .

The gradient of φ :

$$\nabla\varphi(\mathbf{a}) = \lambda \text{sign}(\mathbf{a}) \odot \mathbb{1}_{|\mathbf{a}|>\lambda} + \mathbf{a} \odot \mathbb{1}_{|\mathbf{a}|\leq\lambda}, \quad (12)$$

where \odot denotes element-wise multiplication.

The Riemannian gradient is (tangent to the sphere \mathbb{S}^{m-1}):

$$\text{grad}[\varphi](\mathbf{a}) = \mathbf{P}_{\mathbf{a}^\perp} \nabla\varphi(\mathbf{a}). \quad (13)$$

The Riemannian gradient vanishes iff $\nabla\varphi(\mathbf{a}) \propto \mathbf{a}$, which occurs whenever

$$\mathbf{a} \propto \text{sign}(\mathbf{a}). \quad (14)$$

Second Order Characteristics of the Simplest Case

Hessian at Critical Points of φ .

The *Riemannian Hessian* is given by³

$$\begin{aligned} \text{Hess}[\varphi](\mathbf{a}) &= P_{\mathbf{a}^\perp} \left(\underbrace{\nabla^2 \varphi(\mathbf{a})}_{\text{curvature of } \varphi} - \underbrace{\langle \nabla \varphi(\mathbf{a}), \mathbf{a} \rangle \mathbf{I}}_{\text{curvature of the sphere}} \right) P_{\mathbf{a}^\perp} \\ &= P_{\mathbf{a}_{l,\sigma}^\perp} \left(P_{|\mathbf{a}_{l,\sigma}| \leq \lambda} - \lambda \|\mathbf{I}\| \right) P_{\mathbf{a}_{l,\sigma}^\perp}. \end{aligned}$$

At critical points $\mathbf{a}_{l,\sigma}$ the Hessian exhibits $(\|\cdot\| - 1)$ negative eigenvalues, and $m - \|\cdot\|$ positive eigenvalues.

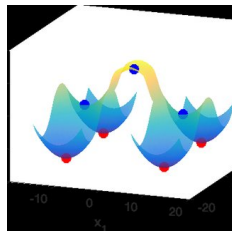
³can be derived by calculating $\left. \frac{d^2}{dt^2} \right|_{t=0} \varphi(\mathbf{a} \cos t + \boldsymbol{\delta} \sin t)$, with any direction $\boldsymbol{\delta} \in T_{\mathbf{a}} \mathbb{S}^{m-1}$ and $\|\boldsymbol{\delta}\| = 1$.

General Messages from the Simplest Case

Symmetric copies of the ground truth are minimizers. The objective function is strongly convex in the vicinity of local minimizers $\mathbf{a} = \pm \mathbf{e}_i$.

Negative curvature in symmetry breaking directions. Saddle points are balanced superpositions of target solutions: $\mathbf{a}_{l,\sigma} = \frac{1}{\sqrt{|l|}} \sum_{i \in l} \sigma_i \mathbf{e}_i$ with l and signs $\sigma_i \in \{\pm 1\}$. There is negative curvature in directions $\delta \in \text{span}(\{\mathbf{e}_i \mid i \in l\})$ that break the balance between target solutions.

Cascade of saddle points. Downstream negative curvature directions are the image of upstream negative curvature directions under gradient flow. Worst case, such as the “octopus function” shown in the Figure⁴, **never** occurs!



⁴Gradient Descent Can Take Exponential Time to Escape Saddle Points, S. Du et al, NeurIPS 2017.

Dictionary Learning: General Case

A Fundamental Problems in Data Analysis:

Given an n -dimensional signal: $\mathbf{y} \in \mathbb{R}^n$, find a transformation $\mathcal{T} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ or its “inverse” $\mathbf{D} : \mathbb{R}^m \rightarrow \mathbb{R}^n$, such that

$$\mathbf{x} = \mathcal{T}[\mathbf{y}], \quad \text{or} \quad \mathbf{y} = \mathbf{D}\mathbf{x}$$

where \mathbf{x} highly compressible or the sparsest possible.

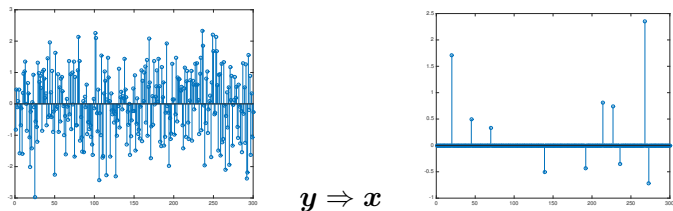


Figure: Sparse Representation Left: a *generic* vector $\mathbf{y} \in \mathbb{R}^n$, Right: a *sparse* representation $\mathbf{x} = \mathcal{T}[\mathbf{y}]$, after a proper transformation \mathcal{T} .

Introduction: History of Finding Good Transform



- **Fourier Transform** $D = F$
- **Wavelet Transform** $D = W$
- **Dictionary Learning**

Figure: Joseph Fourier, 1768 – 1830

Introduction: Fourier Transform

Assumption:

The signal y is **band-limited** and **sparse** in frequency domain: $y_k =$

$$\sum_{l=0}^{n-1} x_l \cdot e^{-\frac{i2\pi}{n}kl} \quad (y = Fx.)$$

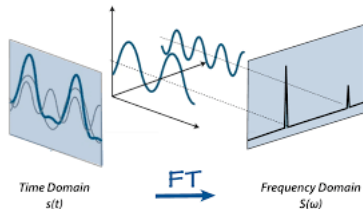


Figure: Fourier Transform



Figure: Lena Compression using Discrete Cosine Transform (JPEG) [pip18]

Introduction: History of Finding Good Transform



Figure: Alfred Haar, 1855 – 1933

- Fourier Transform $D = F$
- **Wavelet Transform** $D = W$
- Dictionary Learning

Introduction: Wavelet Transform

Assumption:

Signal y is piece-wise smooth, scale-invariant, etc: $y = Wx$, $W^*W = I$.

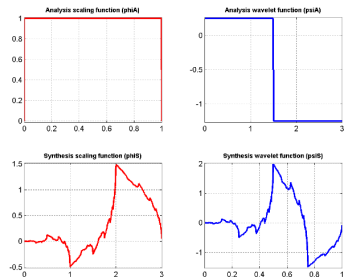


Figure: Haar & Daubechies Wavelets

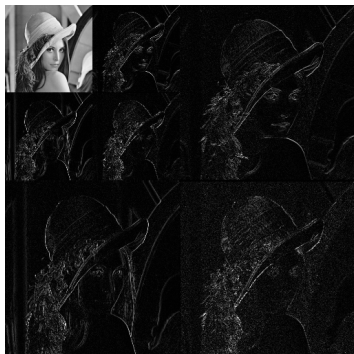


Figure: Lena Compression using Wavelet Transform (JPEG2000) [Jor06]

Why Dictionary Learning?

Limitations of Traditional “By Design” Methods

- A transform is not optimal for signals that do not satisfy the conditions under which the transform is designed (e.g. DCT not ideal for images).
- For different classes of signals, we need to design different transforms (e.g. all the x-lets), which may not even be possible if the properties are not clear.

Why Dictionary Learning?

Limitations of Traditional “By Design” Methods

- A transform is not optimal for signals that do not satisfy the conditions under which the transform is designed (e.g. DCT not ideal for images).
- For different classes of signals, we need to design different transforms (e.g. all the x-lets), which may not even be possible if the properties are not clear.

For a given class of signals, can we directly “learn” the corresponding optimal transform, from its samples?

Dictionary Learning: General Case

Given n -dimensional input data: $\{\mathbf{y}_1, \dots, \mathbf{y}_p\}$, $\forall i \in [p]$, $\mathbf{y}_i \in \mathbb{R}^n$, find a dictionary $\mathbf{D} \in \mathbb{R}^{n \times m}$ and its corresponding coefficients $\{\mathbf{x}_1, \dots, \mathbf{x}_p\}$, $\mathbf{x}_i \in \mathbb{R}^m$, such that

$$\mathbf{y}_i = \mathbf{D}\mathbf{x}_i, \quad \forall i \in [p], \quad (15)$$

and \mathbf{x}_i is sufficiently sparse. That is to factor the data matrix \mathbf{Y} into **two structured unknowns**: a matrix \mathbf{D} and a sparse matrix \mathbf{X} :

$$\mathbf{Y} = \underbrace{\begin{pmatrix} | & & | \\ \mathbf{y}_1 & \dots & \mathbf{y}_p \\ | & & | \end{pmatrix}}_{\text{Observations}} = \underbrace{\begin{pmatrix} d_{1,1} & \dots & d_{1,m} \\ \vdots & \ddots & \vdots \\ d_{n,1} & \dots & d_{n,m} \end{pmatrix}}_{\text{Dictionary } \mathbf{D}} \underbrace{\begin{pmatrix} | & & | \\ \mathbf{x}_1 & \dots & \mathbf{x}_p \\ | & & | \end{pmatrix}}_{\mathbf{X} \text{ is sparse, } \|\mathbf{x}_i\|_0 \ll m} = \mathbf{D}\mathbf{X}.$$

Dictionary Learning: General Case

Challenges

- **Computational Complexity**

Optimizing a nonconvex bilinear problem is NP-hard.

- **Sample Complexity**

Combinatorial possible outcomes for k -sparse x .

- **Signed Permutation Ambiguities**

$\forall P \in SP(m)$,⁵ (D_*P, P^*X_*) and (D_*, X_*) are equally sparse.

⁵ $SP(m)$ denote m dimensional signed permutation group, a group of orthogonal matrices whose entries contain only $0, \pm 1$.

Dictionary Learning: General Case

Challenges

- **Computational Complexity**
Optimizing a nonconvex bilinear problem is NP-hard.
- **Sample Complexity**
Combinatorial possible outcomes for k -sparse x .
- **Signed Permutation Ambiguities**
 $\forall P \in SP(m)$,⁵ (D_*P, P^*X_*) and (D_*, X_*) are equally sparse.

Some heuristic algorithms

- K-SVD [AEB⁺06]
- Alternative Direction Methods [SQW17]

⁵ $SP(m)$ denote m dimensional signed permutation group, a group of orthogonal matrices whose entries contain only $0, \pm 1$.

Dictionary Learning: General Case

Challenges

- **Computational Complexity**
Optimizing a nonconvex bilinear problem is NP-hard.
- **Sample Complexity**
Combinatorial possible outcomes for k -sparse x .
- **Signed Permutation Ambiguities**
 $\forall P \in SP(m)$,⁵ (D_*P, P^*X_*) and (D_*, X_*) are equally sparse.

Some heuristic algorithms

- K-SVD [AEB⁺06]
- Alternative Direction Methods [SQW17]

Learn the dictionary with tractable algorithms and sample size?

⁵ $SP(m)$ denote m dimensional signed permutation group, a group of orthogonal matrices whose entries contain only $0, \pm 1$.

Complete Dictionary Learning – Prior Arts

A Random Model:

For complete dictionary learning, [SWW12] assumes data \mathbf{Y} is generated by a **complete**⁶ dictionary \mathbf{D}_o and sparse coefficients \mathbf{X}_o :

$$\mathbf{Y} = \mathbf{D}_o \mathbf{X}_o,$$

where \mathbf{X}_o follows a Bernoulli Gaussian model:

$$\mathbf{X}_o = \mathbf{\Omega} \circ \mathbf{G}^7, \quad \Omega_{i,j} \sim_{iid} \text{Ber}(\theta), G_{i,j} \sim_{iid} \mathcal{N}(0, 1).$$

⁶square and invertible

⁷ \circ denote element-wise product: $\forall \mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times m}, \{\mathbf{A} \circ \mathbf{B}\}_{i,j} = a_{i,j} b_{i,j}$

Complete Dictionary Learning – Prior Arts

A Random Model:

For complete dictionary learning, [SWW12] assumes data \mathbf{Y} is generated by a **complete**⁶ dictionary \mathbf{D}_o and sparse coefficients \mathbf{X}_o :

$$\mathbf{Y} = \mathbf{D}_o \mathbf{X}_o,$$

where \mathbf{X}_o follows a Bernoulli Gaussian model:

$$\mathbf{X}_o = \mathbf{\Omega} \circ \mathbf{G}^7, \quad \Omega_{i,j} \sim_{iid} \text{Ber}(\theta), G_{i,j} \sim_{iid} \mathcal{N}(0, 1).$$

Preconditioning:

[SQW17] shows that learning a complete dictionary is equivalent with learning an **orthogonal** one through preconditioning

$$\bar{\mathbf{Y}} \leftarrow \left(\frac{1}{p\theta} \mathbf{Y} \mathbf{Y}^* \right)^{-\frac{1}{2}} \mathbf{Y} = \mathbf{D}_o \mathbf{X}_o, \quad \text{with } \mathbf{D}_o \in \mathbf{O}(n).$$

⁶square and invertible

⁷ \circ denote element-wise product: $\forall \mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times m}, \{\mathbf{A} \circ \mathbf{B}\}_{i,j} = a_{i,j} b_{i,j}$

Complete Dictionary Learning – Prior Arts

Complete dictionary learning can be reduced to find the sparsest direction in a subspace:

- 1 D_o is complete $\implies \boxed{\text{row}(\mathbf{Y}) = \text{row}(\mathbf{X}_o)}$
- 2 Rows of \mathbf{X}_o form a *sparse basis* of $\text{row}(\mathbf{Y})$.
- 3 Find \mathbf{x}_1 , the *sparsest vector* in the subspace $\text{row}(\mathbf{Y})$.
- 4 Find \mathbf{x}_i , the *sparsest vector* in $\text{row}(\mathbf{Y}) \setminus \{\mathbf{x}_1, \dots, \mathbf{x}_{i-1}\}$.
- 5 Recover D_o by: $D_o = \mathbf{Y} \mathbf{X}_o^* (\mathbf{X}_o \mathbf{X}_o^*)^{-1}$.

Complete Dictionary Learning – Prior Arts

Finding the sparsest vector in $\text{row}(\mathbf{Y})$ can be naively formulated as

$$\min_{\mathbf{q}} \|\mathbf{q}^* \mathbf{Y}\|_0, \quad \text{such that } \mathbf{q} \neq \mathbf{0},$$

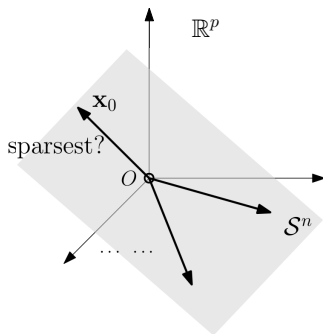


Figure: The sparsest direction in a subspace. Credit: Prof. Qing Qu.

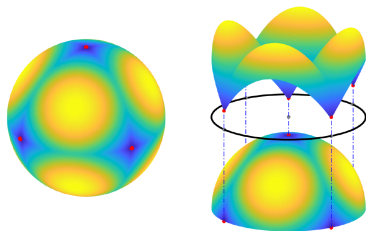
Related Works in Finding the Sparsest Direction

- Linear Programming [SWW12]:

$$\min_{\mathbf{q}} \|\mathbf{q}^* \mathbf{Y}\|_1, \quad \text{such that} \quad \|\mathbf{q}^* \mathbf{Y}\|_\infty = 1.$$

- Nonconvex Optimization on a Sphere [SQW17, BJS18]:

$$\min_{\mathbf{q}} \|\mathbf{q}^* \mathbf{Y}\|_1, \quad \text{such that} \quad \|\mathbf{q}\|_2 = 1.$$



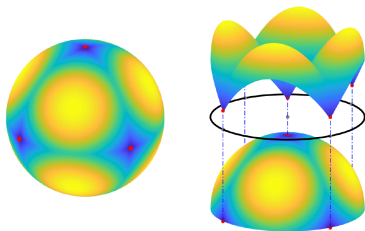
Related Works in Finding the Sparsest Direction

- Linear Programming [SWW12]:

$$\min_{\mathbf{q}} \|\mathbf{q}^* \mathbf{Y}\|_1, \quad \text{such that} \quad \|\mathbf{q}^* \mathbf{Y}\|_\infty = 1.$$

- Nonconvex Optimization on a Sphere [SQW17, BJS18]:

$$\min_{\mathbf{q}} \|\mathbf{q}^* \mathbf{Y}\|_1, \quad \text{such that} \quad \|\mathbf{q}\|_2 = 1.$$



Solving the same optimization n times (high computational cost)!

Assignments

- Reading: Section 7.1 - 7.3 of Chapter 7.
- Programming Homework #3.

References I



Michal Aharon, Michael Elad, Alfred Bruckstein, et al.

K-svd: An algorithm for designing overcomplete dictionaries for sparse representation.
IEEE Transactions on signal processing, 54(11):4311, 2006.



Yu Bai, Qijia Jiang, and Ju Sun.

Subgradient descent learns orthogonal dictionaries.
arXiv preprint arXiv:1810.10702, 2018.



Palle Jorgensen.

<http://homepage.divms.uiowa.edu/~jorgen/Haar.html>, 2006.



[pipo1995_2.](https://www.taringa.net/+info/como-una-foto-de-una-playboy-se-convirtio-en-el-formato-jpg_1ejzk6)

https://www.taringa.net/+info/como-una-foto-de-una-playboy-se-convirtio-en-el-formato-jpg_1ejzk6, 2018.



Ju Sun, Qing Qu, and John Wright.

Complete dictionary recovery over the sphere i: Overview and the geometric picture.
IEEE Transactions on Information Theory, 63(2):853–884, 2017.



Daniel A Spielman, Huan Wang, and John Wright.

Exact recovery of sparsely-used dictionaries.
In Conference on Learning Theory, pages 37–1, 2012.