

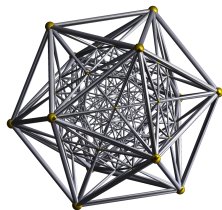
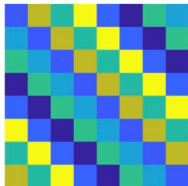
Computational Principles for High-dim Data Analysis

(Lecture Nineteen)

Yi Ma

EECS Department, UC Berkeley

November 4, 2021



Structured Nonlinear Low-Dimensional Models

Sparsity in Convolution and Deconvolution

- 1 Convolution for Image Modeling
- 2 Convolution and Circulant Matrix
- 3 The Blind Short-and-Sparse Deconvolution

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

– Aristotle, *Metaphysica*

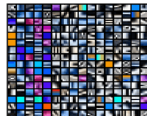
Importance of Mathematical Modeling

**If you formulate a problem correctly,
you have more than halfway solved it!**

Sparsity in Appearance of Image Patches

Patch-level image modeling
(e.g. denoising or super-resolution)
with a sparsifying dictionary:

$$I_{\text{patch}} = \underset{\text{dictionary}}{\mathbf{A}} \times \underset{\text{sparse}}{\mathbf{x}} + \underset{\text{noise}}{\mathbf{z}}. \quad (1)$$



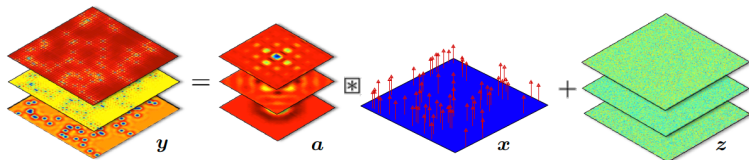
Dictionary learning: the motifs or atoms of the dictionary are **unknown**:

$$\underset{\text{data}}{\mathbf{Y}} = \underset{\text{dictionary}}{\mathbf{A}} \underset{\text{sparse}}{\mathbf{X}}. \quad (2)$$

- Band-limited signals: $\mathbf{A} = \mathbf{F}$, the Fourier transform (JPEG);
- Piecewise smooth: $\mathbf{A} = \mathbf{W}$, the wavelet transforms (JPEG2000);
- For natural images \mathbf{A} can be **learned** from patch samples \mathbf{Y} .

Sparsity in Occurrence of Patch Motif(s)

The same motif $\mathbf{a} \in \mathbf{A}$ occurs at a **sparse** number of locations $(i_1, j_1), \dots, (i_k, j_k)$ in space:



The overall observation y can be modeled as a superposition of translated versions of the motif a , one for each of locations (i_ℓ, j_ℓ) :

$$\underset{\text{data}}{y(i, j, e)} = \sum_{\ell=1}^k \underset{\text{translated motif}}{a(i - i_\ell, j - j_\ell, e)} + \underset{\text{noise}}{z(i, j, e)}. \quad (3)$$

One could generalize this to multiple motifs.

Modeling Translational Occurrence by Convolution

Define a two-dimensional sparse signal $\mathbf{x} \in \mathbb{R}^{w \times h}$, which takes on value 1 at locations (i_ℓ, j_ℓ) and zero elsewhere:

$$\mathbf{y}(\cdot, \cdot, e) = \mathbf{a}(\cdot, \cdot, e) * \mathbf{x} + \mathbf{z}(\cdot, \cdot, e). \quad (4)$$

Combining these equations for all energy levels e , the observed data \mathbf{y} is a convolution of the motif \mathbf{a} and a field \mathbf{x} of sparse spikes:

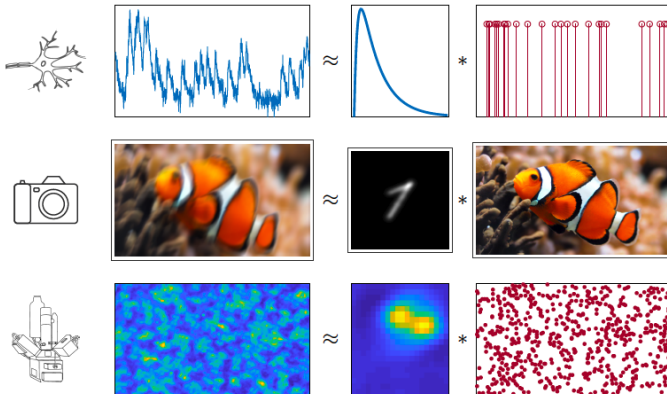
$$\underset{\text{data}}{\mathbf{y}} = \underset{\text{motif}}{\mathbf{a}} * \underset{\text{sparse spikes}}{\mathbf{x}} + \underset{\text{noise}}{\mathbf{z}}, \quad (5)$$

\mathbf{x} could also take different values other than 1 to model the intensity or weight of the motif at each location.

The sparse occurrence/convolution model does generalize to other transformation groups, such as rotation etc.

Modeling Translational Occurrence by Convolution

Examples: Neuron, Camera, and Microscopy



The sparse occurrence/convolution model does generalize to other transformation groups, such as rotation etc.

Background: Convolution and Circulant Matrix

Given a vector $\mathbf{a} = [a_0, a_1, \dots, a_{n-1}]^* \in \mathbb{R}^n$, we may arrange all its circularly shifted versions in a circulant matrix form as

$$\mathbf{A} \doteq \text{circ}(\mathbf{a}) = \begin{bmatrix} a_0 & a_{n-1} & \dots & a_2 & a_1 \\ a_1 & a_0 & a_{n-1} & \dots & a_2 \\ \vdots & a_1 & a_0 & \ddots & \vdots \\ a_{n-2} & \vdots & \ddots & \ddots & a_{n-1} \\ a_{n-1} & a_{n-2} & \dots & a_1 & a_0 \end{bmatrix} \in \mathbb{R}^{n \times n}. \quad (6)$$

It is easy to see that the multiplication of such a circulant matrix \mathbf{A} with a vector \mathbf{x} gives a (circular) convolution $\mathbf{A}\mathbf{x} = \mathbf{a} \circledast \mathbf{x}$ with:

$$(\mathbf{a} \circledast \mathbf{x})_i = \sum_{j=0}^{n-1} x_j a_{i+n-j \bmod n}. \quad (7)$$

Fact: all circulant matrices share the same set of eigenvectors!

Background: Eigenvectors of Circulant Matrices

Let $i = \sqrt{-1}$ and $\omega_n := \exp(-\frac{2\pi i}{n})$ and we define the matrix:

$$\mathbf{F}_n \doteq \frac{1}{\sqrt{n}} \begin{bmatrix} \omega_n^0 & \omega_n^0 & \dots & \omega_n^0 & \omega_n^0 \\ \omega_n^0 & \omega_n^1 & \dots & \omega_n^{n-2} & \omega_n^{n-1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \omega_n^0 & \omega_n^{n-2} & \dots & \omega_n^{(n-2)^2} & \omega_n^{(n-2)(n-1)} \\ \omega_n^0 & \omega_n^{n-1} & \dots & \omega_n^{(n-2)(n-1)} & \omega_n^{(n-1)^2} \end{bmatrix} \in \mathbb{C}^{n \times n}. \quad (8)$$

\mathbf{F}_n is known as the *discrete Fourier transform* (DFT), with $\mathbf{F}_n \mathbf{F}_n^* = \mathbf{I}$.

Theorem (Eigenvectors of Circulant Matrix)

An $n \times n$ matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ is a circulant matrix if and only if it is diagonalizable by the unitary matrix \mathbf{F}_n :

$$\mathbf{F}_n^* \mathbf{A} \mathbf{F}_n = \mathbf{D}_a \quad \text{or} \quad \mathbf{A} = \mathbf{F}_n \mathbf{D}_a \mathbf{F}_n^*, \quad (9)$$

where \mathbf{D}_a is a diagonal matrix of (possibly complex) eigenvalues.

Probably the reason why our brain computes in spectral domain.

The Blind Deconvolution Problem

Problem: how to recover both the motif \mathbf{a} and sparse spikes \mathbf{x} from the observed data \mathbf{y} :

$$\underset{\text{data}}{\mathbf{y}} = \underset{\text{motif}}{\mathbf{a}} * \underset{\text{sparse spikes}}{\mathbf{x}} + \underset{\text{noise}}{\mathbf{z}}, \quad (10)$$

This problem is under-determined (Why?).

We need to leverage low-dimensional structure in both \mathbf{a} and \mathbf{x} by assuming a **short-and-sparse** model (studied in the 90's):

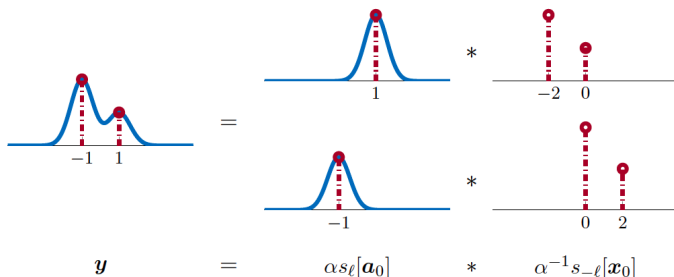
- ① \mathbf{a} is spatially localized, i.e., it is a *short* signal, whose spatial extent is small compared to that of \mathbf{y} ;
- ② \mathbf{x} is *sparse*, since it contains only one nonzero entry for each instance of the motif in \mathbf{y} . (Why not dense?)

Solution by Optimization

Simultaneously recover both \mathbf{a} and \mathbf{x} by the bilinear Lasso (BL):

$$\min_{\mathbf{a}, \mathbf{x}} \varphi_{\text{BL}}(\mathbf{a}, \mathbf{x}) \doteq \underbrace{\frac{1}{2} \|\mathbf{y} - \mathbf{a} * \mathbf{x}\|_F^2}_{\text{data fidelity}} + \underbrace{\lambda \|\mathbf{x}\|_1}_{\mathbf{x} \text{ sparse}} \quad \text{such that} \quad \underbrace{\mathbf{a} \in \mathcal{A}}_{\mathbf{a} \text{ short}} \quad (11)$$

Ambiguity due to a **scaling-shift symmetry**:

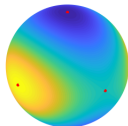


Taxonomy of Symmetric Nonconvex Problems

Nonconvex Problems with Discrete Symmetries

Eigenvector Computation

Maximize a quadratic form over the sphere.

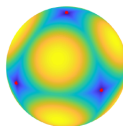


$$\max_{\mathbf{x} \in \mathbb{S}^{n-1}} \frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x}.$$

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

Dictionary Learning

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

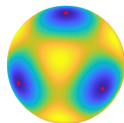


$$\min_{\mathbf{A} \in \mathcal{A}, \mathbf{X}} \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1.$$

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

Tensor Decomposition

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

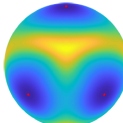


$$\max_{\mathbf{X} \in \mathcal{O}(n)} \sum_i \mathbf{T}(\mathbf{x}_i, \mathbf{x}_i, \mathbf{x}_i, \mathbf{x}_i).$$

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

Short-and-Sparse Deconvolution

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



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

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

Symmetry in Short-and-Sparse Deconvolution

Letting s_τ denote a shift by τ pixels, we have

$$\mathbf{y} = s_\tau[\mathbf{a}] * s_{-\tau}[\mathbf{x}] = \mathbf{a} * \mathbf{x}, \quad \text{with} \quad \underbrace{\|\mathbf{a}\|_F = 1}_{\text{normalization}}. \quad (12)$$

If \mathbf{a} is *shift incoherent*:

$$\underbrace{\mu_s = \max_{\tau \neq 0} |\langle \mathbf{a}, s_\tau[\mathbf{a}] \rangle| \approx 0}_{\text{incoherence}}, \quad \text{or} \quad \underbrace{\text{circ}(\mathbf{a}) \approx \mathbf{I}}_{\text{isometry}}.$$

the bilinear Lasso loss in (11) can be approximated as

$$\begin{aligned} \frac{1}{2} \|\mathbf{y} - \mathbf{a} * \mathbf{x}\|_F^2 &= \frac{1}{2} \|\mathbf{y}\|_F^2 + \frac{1}{2} \|\mathbf{a} * \mathbf{x}\|_F^2 - \langle \mathbf{y}, \mathbf{a} * \mathbf{x} \rangle \\ &\approx \frac{1}{2} \|\mathbf{y}\|_F^2 + \frac{1}{2} \|\mathbf{x}\|_F^2 - \langle \mathbf{y}, \mathbf{a} * \mathbf{x} \rangle. \end{aligned} \quad (13)$$

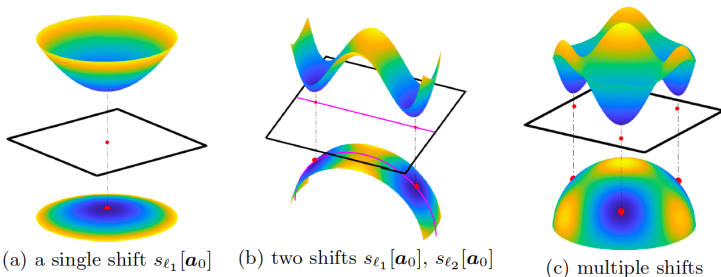
This gives:

$$\varphi_{\text{ABL}}(\mathbf{a}, \mathbf{x}) \doteq \frac{1}{2} \|\mathbf{y}\|_F^2 + \frac{1}{2} \|\mathbf{x}\|_F^2 - \langle \mathbf{y}, \mathbf{a} * \mathbf{x} \rangle + \lambda \|\mathbf{x}\|_1, \quad \|\mathbf{a}\|_F = 1. \quad (14)$$

Landscape of the Objective Function

Geometry of the approximate bilinear Lasso (ABL) objective:

$$\varphi_{\text{ABL}}(\mathbf{a}, \mathbf{x}) \doteq \frac{1}{2}\|\mathbf{y}\|_F^2 + \frac{1}{2}\|\mathbf{x}\|_F^2 - \langle \mathbf{y}, \mathbf{a} * \mathbf{x} \rangle + \lambda\|\mathbf{x}\|_1, \quad \mathbf{a} \in \mathcal{A}. \quad (15)$$



Notice: equivalent (symmetric) solutions are local minimizers, and there is negative curvature in symmetry breaking directions.

Sparsity and Shift-Coherence Tradeoff

Solving the sparse-and-short deconvolution (SaSD) from:

$$\min_{\mathbf{a}, \mathbf{x}} \varphi_{\text{BL}}(\mathbf{a}, \mathbf{x}) \doteq \frac{1}{2} \|\mathbf{y} - \mathbf{a} * \mathbf{x}\|_F^2 + \lambda \|\mathbf{x}\|_1 \quad \text{such that} \quad \mathbf{a} \in \mathcal{A}. \quad (16)$$

A sparsity-coherence tradeoff: Smaller $\mu_s(\mathbf{a}_0)$ allows higher $\theta(\mathbf{x})$.

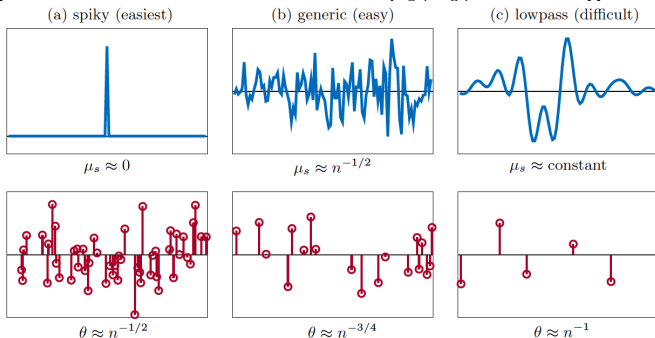


Figure: In order of increasing difficulty: (a) when \mathbf{a}_0 is a Dirac delta function, $\mu_s(\mathbf{a}_0) = 0$; (b) when \mathbf{a}_0 is uniform on the sphere \mathbb{S}^{n-1} , its shift-coherence is roughly $\mu_s(\mathbf{a}_0) \approx n^{-1/2}$; (c) when \mathbf{a}_0 is low-pass, $\mu_s(\mathbf{a}_0) \rightarrow \text{const.}$ as n grows.

Alternating Descent Algorithm for SaSD

Solving the sparse-and-short deconvolution (SaSD) from:

$$\min_{\mathbf{a}, \mathbf{x}} \varphi_{\text{BL}}(\mathbf{a}, \mathbf{x}) \doteq \frac{1}{2} \|\mathbf{y} - \mathbf{a} * \mathbf{x}\|_F^2 + \lambda \|\mathbf{x}\|_1 \quad \text{such that} \quad \mathbf{a} \in \mathcal{A}. \quad (17)$$

Fix \mathbf{a} and take a proximal gradient step on \mathbf{x} .

$$\text{Gradient w.r.t. } \mathbf{x} : \quad \nabla_{\mathbf{x}} \psi(\mathbf{a}, \mathbf{x}) = \iota_{\mathbf{x}}^* \check{\mathbf{a}} * (\mathbf{a} * \mathbf{x} - \mathbf{y}). \quad (18)$$

$$\text{Proximal gradient:} \quad \mathbf{x}_{k+1} = \text{prox}_{t\lambda g} [\mathbf{x}_k - t \nabla_{\mathbf{x}} \psi(\mathbf{a}_k, \mathbf{x}_k)]. \quad (19)$$

Fix \mathbf{x} and take a projected gradient step on $\mathbf{a} \in \mathcal{A}$ and $\|\mathbf{a}\|_2 = 1$.

$$\text{Gradient w.r.t. } \mathbf{a} : \quad \nabla_{\mathbf{a}} \psi(\mathbf{a}, \mathbf{x}) = \iota_{\mathbf{a}}^* \check{\mathbf{x}} * (\mathbf{a} * \mathbf{x} - \mathbf{y}). \quad (20)$$

$$\text{Proximal gradient:} \quad \mathbf{a}_{k+1} = \mathcal{P}_{\mathcal{A}} [\mathbf{a}_k - \tau_k \nabla_{\mathbf{a}} \psi(\mathbf{a}_k, \mathbf{x}_{k+1})]. \quad (21)$$

Additional Heuristics

In practice, the kernel \mathbf{a} might **not** be so shift incoherent.

Better Optimization Algorithm: Momentum Acceleration

$$\mathbf{w}_k = \mathbf{x}_k + \underbrace{\beta \cdot (\mathbf{x}_k - \mathbf{x}_{k-1})}_{\text{inertial term}}, \quad (22)$$

$$\mathbf{x}_{k+1} = \text{prox}_{t_k g} [\mathbf{w}_k - t_k \nabla_{\mathbf{x}} \psi(\mathbf{a}_k, \mathbf{w}_k)]. \quad (23)$$

Better Optimization Strategy: Homotopy Continuation

Gradually decreasing λ_n to produce the *solution path* $\{(\hat{\mathbf{a}}_n, \hat{\mathbf{x}}_n; \lambda_n)\}$. By ensuring that \mathbf{x} remains sparse along the solution path.

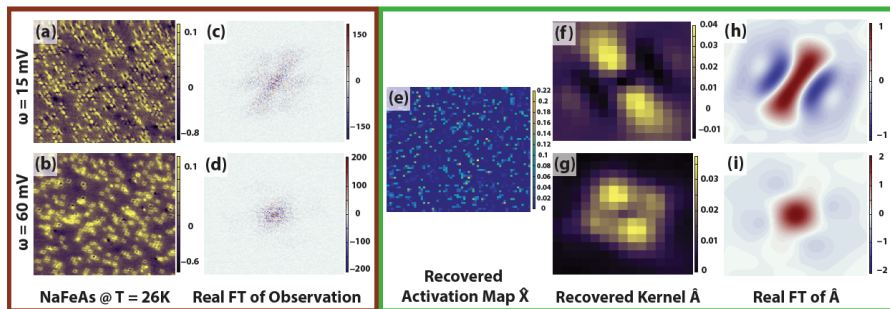
Better Initialization: from the Data

Small pieces of \mathbf{y} are superpositions of a few shifted copies of \mathbf{a}_0 . One could select a small window of \mathbf{y} and then normalizes it to initialize \mathbf{a} .

An Example of Scanning Tunneling Microscopy

Short and Sparse Deconvolution on Real NaFeAs Data¹

This dataset y consists of measurements across a $100 \times 100 nm^2$ area at $E = 41$ different bias voltages.



¹Dictionary learning in Fourier-transform scanning tunneling spectroscopy, Sky Cheung et. al., Nature Communications, 2020.

Assignments

- Reading: Section 7.3.3 and Chapter 12.
- Written Homework #4.