

Mathematical Models and Methods For Image Processing

Rao

Politecnico di Milano

Originally written in: **2025-02-17**

Last updated at: **2025-07-20**

Contents

1	Introduction to Sparsity	4
1.1	Sparsity in Statistics	4
1.2	Sparsity in Signal Processing	4
2	Image Prior	5
2.1	Image Denoising Problem	5
2.2	Local Constancy Prior	6
2.3	Non-local Means Prior	7
2.4	Block-Matching 3D (BM3D) Algorithm	9
2.5	Sparsity-Based Image Prior	11
3	Discrete Cosine Transform (DCT)	12
3.1	1D DCT	12
3.2	2D DCT	12
3.3	DCT-Based Denoising Pipeline	13
3.4	Noise Standard Deviation Estimation	13
3.5	Sliding DCT Algorithm	14
3.6	Wiener Filter	16
3.7	The Sparsity Problem	17
3.8	Experimental Demonstration	17
4	Introduction to Dictionary Learning	18
4.1	Theoretical Properties of Overcomplete Systems	18
4.2	Regularization and Sparse Recovery	19
4.3	Towards Sparsity: ℓ_0 and ℓ_1 Regularization	20
5	Sparse Coding In ℓ_0 sense	21
5.1	The Sparse Coding Problem	21
5.2	Matching Pursuit Algorithm	23
5.3	Orthogonal Matching Pursuit	25
5.4	OMP-Based Image Denoising	26
5.5	Linearity Analysis of Sparse Coding Algorithms	27
5.6	Uniqueness Guarantees for Sparse Solutions	28
5.7	Application to Image Inpainting	29
6	Sparse Coding In ℓ_1 sense	30
6.1	The ℓ_1 Optimization Problem	30
6.2	Problem Components Analysis	30
6.3	Proximal Gradient Methods	31
6.4	FISTA	34
7	Structured Sparsity	36
7.1	Joint Sparsity and Mixed Norms	37
7.2	Joint Sparse Coding	38
7.3	Group Sparsity and Extensions	39
7.4	LASSO	40
7.5	Elastic Net	41
7.6	Algorithmic Summary	41
8	Dictionary Learning	42

8.1	Introduction to Dictionary Learning	42
8.2	Problem Formulation	42
8.3	Sparse Coding Phase	43
8.4	Dictionary Update Phase	43
9	Local Polynomial Approximation	47
9.1	Local Polynomial Model	47
9.2	Weighted Local Polynomial Approximation	49
9.3	QR Decomposition Approach	50
9.4	Convolution Implementation	51
9.5	Special Cases	51
9.6	Statistical Properties and Performance Analysis	52
9.7	Adaptive Neighborhood Selection	53
9.8	The Intersection of Confidence Intervals (ICI) Rule	53
10	Robust Fitting	55
10.1	Limitations of Vertical Distance Minimization	55
10.2	The Direct Linear Transformation (DLT) Algorithm	56
10.3	Robust Estimation Methods	59
10.4	RANSAC: Random Sample Consensus	60
10.5	MSAC and MLESAC Variants	61
11	Multiple Model Fitting	62
11.1	Outlier Complexity in Multi-Model Settings	62
11.2	Sequential Multi-Model Fitting	63
11.3	Simultaneous Multi-Model Fitting	65
12	Appendix: Fundamental Concepts in Linear Algebra	68
12.1	Vector Spaces and Linear Combinations	68
12.2	Linear Independence and Basis	68
12.3	Orthogonal Vectors	69
12.4	The ℓ_0 Norm	69
12.5	Matrix Spark	70
13	Appendix: Optimization Theory for Non-Differentiable Functions	71
13.1	Lipschitz Continuity	71
13.2	Majorization-Minimization	72

1 Introduction to Sparsity

The principle of sparsity or “parsimony” consists in representing some phenomenon with as few variable as possible. Stretch back to philosopher William Ockham in the 14th century, Wrinch and Jeffreys relate simplicity to parsimony:

The **existence of simple laws** is, then, apparently, to be regarded as **a quality of nature**; and accordingly we may infer that it is justifiable to **prefer a simple law to a more complex one that fits our observations slightly better**.

1.1 Sparsity in Statistics

Sparsity is used to **prevent overfitting and improve interpretability of learned models**. In model fitting, the number of parameters is typically used as a criterion to perform model selection. See Bayes Information Criterion (BIC), Akaike Information Criterion (AIC), ..., Lasso.

1.2 Sparsity in Signal Processing

Signal Processing: similar concepts but different terminology. **Vectors** corresponds to **signals** and data modeling is crucial for performing various operations such as **restoration, compression, solving inverse problems**.

Signals are approximated by sparse linear combinations of **prototypes**(basis elements / atoms of a dictionary), resulting in simpler and compact model.

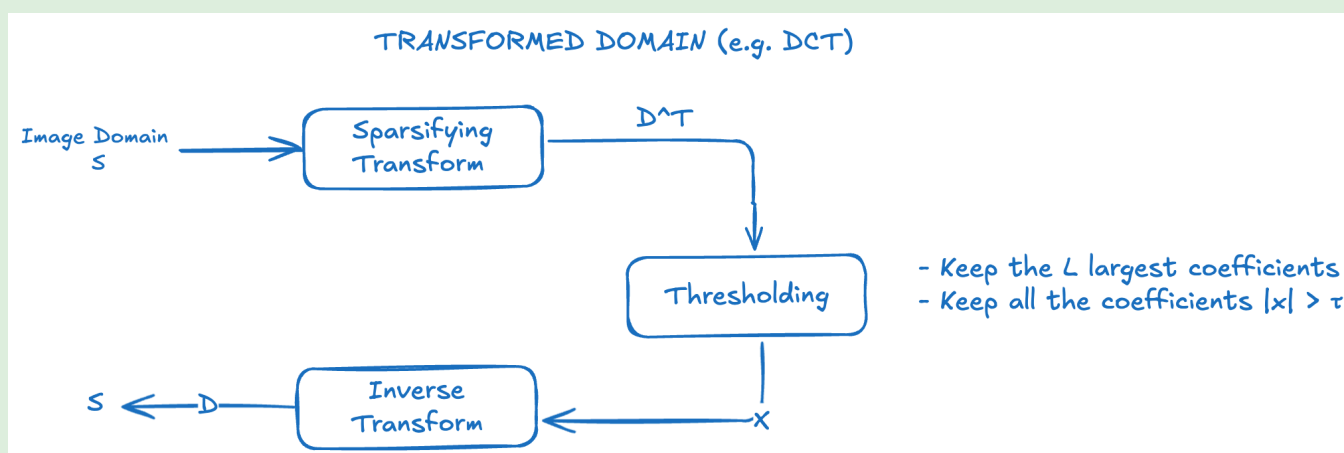


Figure 1.1: Enforce sparsity in signal processing

2 Image Prior

2.1 Image Denoising Problem

Image denoising provides a simple and clear problem formulation example. The **observation model** is:

$$z(x) = y(x) + \eta(x) \quad (2.1)$$

where:

- $z(x)$ is noisy observation at pixel coordinate x
- $y(x)$ is ideal (noise-free) image
- $\eta(x)$ is the noise component
- x is pixel coordinate

This formulation assumes *additive white Gaussian noise* (AWGN).

We assume that the noise is:

- **Additive Gaussian noise:** $\eta(x) \sim N(0, \sigma^2)$.
- **Independent and identically distributed (i.i.d.):** Noise realizations at different pixels are independent.

While real-world noise may exhibit correlations or non-Gaussian characteristics, the AWGN model remains prevalent for algorithm design. Practical systems often transform raw sensor data to approximate this model. The denoising objective is formalized as finding an estimator \hat{y} minimizing the the mean squared error (MSE):

$$\text{MSE}(\hat{y}) = \mathbb{E}[\|\hat{y} - y\|_2^2] \quad (2.2)$$

The observation model provides a **prior on noise** but we also need a **prior on images**.

Definition 2.1.1 (Noise Prior): A **noise prior** refers to the assumed statistical properties or characteristics of the noise itself that is corrupting the image. Knowing how the noise behaves helps in modeling and subsequently removing it. For AWGN, the true image y is likely to be a in a circular neighborhood around the observation $z(x)$.

Definition 2.1.2 (Image Prior): An **image prior** is a statistical model that captures the expected structure of natural images. It is used to regularize the denoising process, guiding the estimator towards plausible solutions.

In this chapter, we will explore various image priors and their applications in denoising. The goal is to leverage the statistical properties of natural images to improve the quality of denoised outputs.

2.2 Local Constancy Prior

Assumption: Images are locally constant within small patches. For a constant signal corrupted by Gaussian noise:

$$\hat{y} = \frac{1}{M} \sum_{i=1}^M z(x_i) \quad (2.3)$$

Properties of this estimator:

- **Unbiased:** $\mathbb{E}[\hat{y}] = y$ (true signal)
- **Reduced Variance:** $\text{Var}[\hat{y}] = \frac{\sigma^2}{M}$

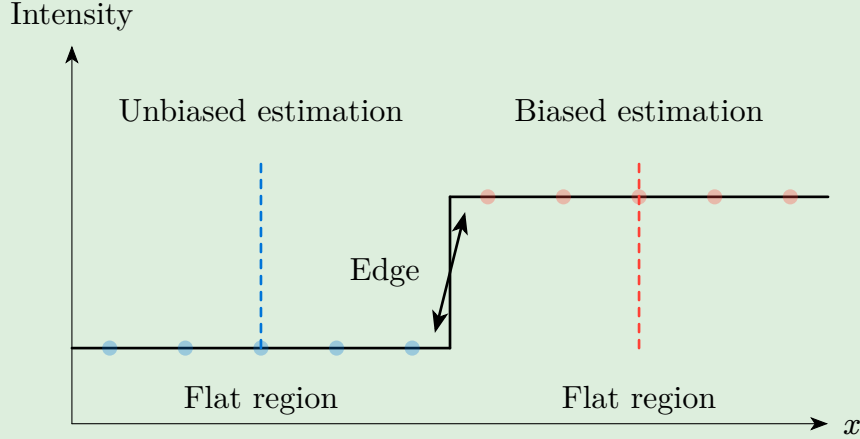


Figure 2.1: Bias-variance tradeoff in local averaging

Consider the classical approach of estimating pixel intensities through local averaging. For a pixel at location (i, j) , a naive estimate would be:

$$\hat{y}_{i,j} = \frac{1}{|U_{i,j}|} \sum_{(k,l) \in U_{i,j}} z_{k,l} \quad (2.4)$$

where $U_{i,j}$ denotes a neighborhood around pixel (i, j) , and $|U_{i,j}|$ is its cardinality.

This approach is equivalent to convolution with a normalized box kernel:

$$\hat{y} = z * h \quad (2.5)$$

where h is the box kernel with support $U_{i,j}$.

The primary limitation of local averaging becomes apparent near image discontinuities. Consider an ideal edge scenario where:

$$z_{i,j} = \begin{cases} a & \text{if } (i, j) \in \Omega_1 \\ b & \text{if } (i, j) \in \Omega_2 \end{cases} \quad (2.6)$$

where Ω_1 and Ω_2 are disjoint regions separated by an edge, and $a \neq b$.

For a pixel (i, j) near the edge boundary, the local average becomes:

$$\begin{aligned} \hat{y}_{i,j} &= \frac{1}{|U_{i,j}|} \left(\sum_{(k,l) \in U_{i,j} \cap \Omega_1} y_{k,l} + \sum_{(k,l) \in U_{i,j} \cap \Omega_2} y_{k,l} \right) \\ &\approx \frac{|\Omega_1 \cap U_{i,j}|}{|U_{i,j}|} a + \frac{|\Omega_2 \cap U_{i,j}|}{|U_{i,j}|} b \end{aligned} \quad (2.7)$$

This results in a blurred edge transition, motivating the need for adaptive filtering strategies.

2.3 Non-local Means Prior

The Non-Local Means (NLM) algorithm addresses the limitations of local averaging by exploiting the *self-similarity* property of natural images. **The key insight is that similar patches exist throughout the image, not just in local neighborhoods.**

Definition 2.3.1 (Self-Similarity Prior): For a natural image \mathbf{Z} , there exist multiple patches $\{P_k\}$ such that:

$$\|P_i - P_j\|_2 < \varepsilon \quad (2.8)$$

for some small threshold $\varepsilon > 0$, where P_k represents a patch extracted from the image.

The NLM estimate for pixel (i, j) is given by:

$$\hat{y}_{i,j} = \sum_{(k,l) \in S_{i,j}} w_{i,j}(k,l) \cdot z_{k,l} \quad (2.9)$$

where $S_{i,j}$ represents the search window around pixel (i, j) , and $w_{i,j}(k, l)$ are the similarity weights.

2.3.1 Weight Computation

The similarity weights are computed based on patch distances:

$$w_{i,j}(k, l) = \frac{1}{Z_{i,j}} \exp\left(-\frac{d^2(P_{i,j}, P_{k,l})}{h^2}\right) \quad (2.10)$$

where:

- $P_{i,j}$ and $P_{k,l}$ are patches centered at (i, j) and (k, l) respectively
- $d^2(P_{i,j}, P_{k,l})$ is the squared Euclidean distance between patches
- h is the filtering parameter controlling the decay rate
- $Z_{i,j}$ is the normalization constant ensuring $\sum w_{i,j}(k, l) = 1$

2.3.2 Patch Distance Computation

The patch distance is computed as:

$$d^2(P_{i,j}, P_{k,l}) = \frac{1}{|P|} \sum_{(u,v) \in P} |z_{i+u,j+v} - z_{k+u,l+v}|^2 \quad (2.11)$$

where P represents the patch domain and $|P|$ is the number of pixels in the patch.

Noise-Aware Distance In the presence of noise, the distance can be corrected as:

$$d_{\text{corrected}}^2(P_{i,j}, P_{k,l}) = \max(d^2(P_{i,j}, P_{k,l}) - 2\sigma^2, 0) \quad (2.12)$$

This correction accounts for the noise contribution to the patch distance.

2.3.3 Normalization and Properties

The normalization constant is given by:

$$Z_{i,j} = \sum_{(k,l) \in S_{i,j}} \exp\left(-\frac{d^2(P_{i,j}, P_{k,l})}{h^2}\right) \quad (2.13)$$

Theorem 2.3.3.1 (NLM Consistency): For a noiseless image, the NLM estimate satisfies:

$$\lim_{\sigma \rightarrow 0} \hat{y}_{i,j} = y_{i,j} \quad (2.14)$$

Proof: As $\sigma \rightarrow 0$, the patch distances approach their true values, and the weight distribution becomes increasingly concentrated around patches identical to the reference patch.

2.3.4 Algorithm Framework

Algorithm 1: Non-Local Means Algorithm

```

1: procedure NLM(Noisy image  $y$ , patch size  $p$ , search window size  $s$ , filtering parameter  $h$ )
2:    $\triangleright$  Input: Noisy image  $y$ , patch size  $p$ , search window size  $s$ , filtering parameter  $h$ 
3:    $\triangleright$  Output: Denoised image  $\hat{x}$ 
4:
5:    $\triangleright$  For each pixel  $(i, j)$  in the image:
6:   for each pixel  $(i, j)$  in the image do
7:      $\triangleright$  Initialize weight matrix  $W = 0$ 
8:      $W \leftarrow 0$ 
9:
10:     $\triangleright$  Define search window  $S_{i,j}$  of size  $s \times s$ 
11:     $S_{i,j} \leftarrow$  search window of size  $s \times s$  around  $(i, j)$ 
12:
13:     $\triangleright$  For each pixel  $(k, l)$  in  $S_{i,j}$ :
14:    for each pixel  $(k, l)$  in  $S_{i,j}$  do
15:       $\triangleright$  Extract patches  $P_{i,j}$  and  $P_{k,l}$  of size  $p \times p$ 
16:       $P_{i,j}, P_{k,l} \leftarrow$  patches of size  $p \times p$ 
17:
18:       $\triangleright$  Compute distance  $d^2(P_{i,j}, P_{k,l})$  using patch distance formula
19:
20:       $d^2(P_{i,j}, P_{k,l}) \leftarrow \frac{1}{|P|} \sum_{(u,v) \in P} |y_{i+u,j+v} - y_{k+u,l+v}|^2$ 
21:
22:       $\triangleright$  Compute weight  $w_{i,j}(k, l)$  using NLM weights formula
23:
24:       $w_{i,j}(k, l) \leftarrow \exp\left(-\frac{d^2(P_{i,j}, P_{k,l})}{h^2}\right)$ 
25:
26:       $W(k, l) \leftarrow w_{i,j}(k, l)$ 
27:    end
28:
29:     $\triangleright$  Normalize weights:  $W = \frac{W}{\sum W}$ 
30:     $W \leftarrow \frac{W}{\sum W}$ 
31:
32:     $\triangleright$  Compute estimate:  $\hat{x}_{i,j} = \sum_{(k,l)} W(k, l) \cdot y_{k,l}$ 
33:     $\hat{x}_{i,j} \leftarrow \sum_{(k,l)} W(k, l) \cdot y_{k,l}$ 
34:  end
35:
36:  return  $\hat{x}$ 
37: end

```

Near image boundaries, patches may extend beyond the image domain. Common strategies include:

1. **Symmetric Padding:** Reflect image content across boundaries
2. **Periodic Padding:** Wrap image content periodically
3. **Zero Padding:** Extend with zeros (not recommended)

The symmetric padding approach is preferred as it maintains image statistics:

$$\tilde{x}_{i,j} = \begin{cases} x_{i,j} & \text{if } (i, j) \text{ is inside image} \\ x_{2N-i,j} & \text{if } i > N \text{ (bottom boundary)} \\ x_{i,2M-j} & \text{if } j > M \text{ (right boundary)} \end{cases} \quad (2.15)$$

2.3.5 Parameter Selection

Typical Parameter Values

- Patch size: $p = 7 \times 7$ (provides good balance between detail and computational cost)
- Search window: $s = 21 \times 21$ (sufficient for finding similar patches)
- Filtering parameter: $h = 10\sigma$ (empirically determined)

Adaptive Parameter Selection

The filtering parameter can be adapted based on local image characteristics:

$$h_{i,j} = \alpha \cdot \sigma \cdot \sqrt{\text{LocalVariance}(P_{i,j})} \quad (2.16)$$

where α is a scaling factor and LocalVariance measures local image activity.

2.4 Block-Matching 3D (BM3D) Algorithm

The Block-Matching 3D (BM3D) algorithm extends the self-similarity concept by combining it with sparsity-based denoising. The key innovations include:

1. **Grouping:** Collect similar patches into 3D arrays
2. **Collaborative Filtering:** Process groups jointly using 3D transforms
3. **Aggregation:** Combine processed patches back into the image

2.4.1 Patch Grouping

For a reference patch $P_{i,j}$, we define the group $G_{i,j}$ as:

$$G_{i,j} = \{P_{k,l} : d^2(P_{i,j}, P_{k,l}) < \tau\} \quad (2.17)$$

where τ is a similarity threshold.

2.4.2 3D Array Construction

The group is arranged as a 3D array $\mathcal{G} \in \mathbb{R}^{p \times p \times K}$, where $K = |G_{i,j}|$ is the number of patches in the group.

The BM3D algorithm applies a 3D transform to exploit both spatial and inter-patch correlations:

$$\mathcal{T} = \mathcal{W}_{3D} \cdot \mathcal{G} \quad (2.18)$$

where \mathcal{W}_{3D} represents the 3D transform operator.

2.4.3 Separable Transform

The 3D transform is typically implemented as a separable transformation:

$$\begin{aligned} \mathcal{T} &= \mathcal{W}_1 \cdot (\mathcal{W}_2 \cdot (\mathcal{W}_3 \cdot \mathcal{G})) \\ &= (\mathcal{W}_1 \times \mathcal{W}_2 \times \mathcal{W}_3) \cdot \mathcal{G} \end{aligned} \quad (2.19)$$

where \mathcal{W}_1 and \mathcal{W}_2 are 2D transforms (e.g., DCT) applied to each patch, and \mathcal{W}_3 is a 1D transform applied across the grouping dimension.

2.4.4 Hard Thresholding Stage

In the first stage, BM3D applies hard thresholding to the transform coefficients:

$$\hat{\mathcal{T}} = \mathcal{H}_\lambda(\mathcal{T}) \quad (2.20)$$

where the hard thresholding operator is defined as:

$$\mathcal{H}_\lambda(t) = \begin{cases} t & \text{if } |t| > \lambda \\ 0 & \text{if } |t| \leq \lambda \end{cases} \quad (2.21)$$

2.4.5 Threshold Selection

The threshold is typically chosen as:

$$\lambda = \beta \cdot \sigma \quad (2.22)$$

where β is a parameter controlling the aggressiveness of denoising (typically $\beta = 2.7$).

2.4.6 Collaborative Filtering Stage

The second stage performs collaborative filtering using both the noisy and first-stage estimates:

$$\hat{\mathcal{J}}^{(2)} = \frac{|\mathcal{J}^{(1)}|^2}{|\mathcal{J}^{(1)}|^2 + \sigma^2} \cdot \mathcal{J}^{(0)} \quad (2.23)$$

where:

- $\mathcal{J}^{(0)}$ represents the transform of the noisy group
- $\mathcal{J}^{(1)}$ represents the transform of the first-stage estimate
- $|\cdot|^2$ denotes element-wise squared magnitude

2.4.7 Aggregation and Weighting

After processing, patches must be aggregated back into the image. Due to overlapping patches, multiple estimates exist for each pixel:

$$\hat{x}_{i,j} = \frac{\sum_{G \ni (i,j)} w_G \cdot \hat{x}_{i,j}^{(G)}}{\sum_{G \ni (i,j)} w_G} \quad (2.24)$$

where the sum is over all groups G containing pixel (i, j) .

Sparsity-Aware Weighting

The weights are chosen based on the sparsity of the processed group:

$$w_G = \frac{1}{\|\mathcal{J}_G\|_0} \quad (2.25)$$

where $\|\cdot\|_0$ denotes the number of non-zero coefficients.

2.5 Sparsity-Based Image Prior

Natural images have **sparse representations** in certain **transform domains** (e.g., DCT), as evidenced by the success of JPEG compression. There are two types of sparsity models:

- **Transform-domain sparsity:** Images can be represented as sparse linear combinations of basis functions in a specific transform domain.
- **Synthesis Sparse Model:** The synthetic sparse model assumes that a signal (such as an image) can be synthesized through a linear combination of a small number of “atoms.” These “atoms” are drawn from a collection known as a **dictionary**.

2.5.1 Transform-domain Sparsity

Transform-domain sparsity is a fundamental concept in signal and image processing, asserting that many real-world signals and images, though complex in their original form, can be represented much more efficiently in a different mathematical domain. This efficiency is achieved when, after a specific mathematical transformation, the majority of the resulting coefficients are zero or close to zero, leaving only a few significant, non-zero values that capture the essential information of the original data. Some of the most common and powerful transforms include:

- **Discrete Cosine Transform (DCT):** Widely used in image compression (e.g., JPEG), it transforms spatial domain data into frequency domain, emphasizing low-frequency components.
- **Wavelet Transform:** Decomposes signals into different frequency components, allowing for multi-resolution analysis.
- **Fourier Transform (DFT):** Converts signals from time domain to frequency domain, revealing periodic structures.

2.5.2 Synthesis Sparse Model

The **synthesis sparse model** is a mathematical framework used to represent signals or images as linear combinations of a small number of basis functions, known as “atoms.” This model is particularly useful in applications like image compression, denoising, and reconstruction, where the goal is to efficiently represent data while preserving essential features. There are two main approaches to synthesis sparse modeling:

- **Sparse Coding:** In this approach, the signal is expressed as a linear combination of a dictionary of atoms. The goal is to find a sparse representation where only a few coefficients are non-zero, indicating that only a small number of atoms are used to reconstruct the signal. This is often achieved through optimization techniques that minimize the reconstruction error while enforcing sparsity constraints.
- **Dictionary Learning:** This approach involves learning a dictionary of atoms from a set of training signals. The dictionary is optimized to capture the underlying structure of the data, allowing for efficient representation. The learned dictionary can then be used to represent new signals or images in a sparse manner. Dictionary learning algorithms aim to find a set of basis functions that best represent the training data, often through iterative optimization techniques.

3 Discrete Cosine Transform (DCT)

3.1 1D DCT

Generate the DCT basis according to the following formula, the k -th atom of the DCT basis in dimension M is defined as:

$$\text{DCT}_{k(n)} = c_k \cos\left(k\pi \frac{2n+1}{2M}\right) \quad n, k = 0, 1, \dots, M-1 \quad (3.1)$$

where $c_0 = \sqrt{\frac{1}{M}}$ and $c_k = \sqrt{\frac{2}{M}}$ for $k \neq 0$.

For each $k = 0, \dots, M-1$, just sample each function

$$\text{DCT}_{k(n)} = \cos\left(k\pi \frac{2n+1}{2M}\right) \quad (3.2)$$

at $n = 0, \dots, M-1$, obtain a vector. Ignore the normalization coefficient. Divide each vector by its ℓ_2 norm.

Mathematically, suppose the image signal is $s \in \mathbb{R}^M$.

$$x = \text{dct2}(s) = D^T s \quad (3.3)$$

where D^T represents the **DCT basis matrix**. x contains the **DCT coefficients**, which are a **sparse representation** of s .

The **inverse DCT transformation** reconstructs s from x :

$$s = \text{idct2}(x) = Dx \quad (3.4)$$

3.2 2D DCT

2D Discrete Cosine Transform (DCT) can be used as a dictionary for representing image patches. A small patch of an image is extracted, represented as s , with dimension $p \times p$. This patch can be **flattened** into a vector of length $M = p^2$, meaning each patch is reshaped into a vector of length M . The **2D-DCT** is used to transform the patch s into DCT coefficients x .

Suppose the image signal is $S \in \mathbb{R}^{M \times N}$. The 2D DCT can be decomposed into two **1D DCT** operations:

1. **Column-wise DCT**: apply **1D DCT** to each column of the image patch: $Z = D^T S$
2. **Row-wise DCT**: apply **1D DCT** to each row of the image patch: $X^T = D^T Z^T \rightarrow X = D^T S D$

Example 3.2.1 (JPEG Compression): The image is divided into non-overlapping 8×8 blocks. Each block is treated separately during the compression process.

For each 8×8 block, the **DCT** is applied, transforming pixel values into frequency-domain coefficients. Each 8×8 block's coefficients are checked against a compression threshold τ , coefficients with absolute values below τ are **discarded**(set to zero). The larger the threshold τ , the more coefficients are discarded, leading to **higher compression**.

The compression ratio is defined as:

$$\text{Comp Ratio} = 1 - \frac{\# \text{Non-zero coefficients}}{\# \text{Pixels in the image}} \quad (3.5)$$

To measure how much the image quality is degraded after compression, **Peak Signal-to-Noise Ratio (PSNR)** is used:

$$\text{PSNR} = 10 \log_{10} \left(\frac{1}{\text{MSE}(Y, \hat{Y})} \right) \quad (3.6)$$

3.3 DCT-Based Denoising Pipeline

Step 1 : Analysis

$$X = D^T S \quad (3.7)$$

where:

- S is the vectorized image patch
- D is the DCT basis matrix
- X is the DCT coefficients vector

Step 2: Enforce Sparsity (Thresholding)

$$\hat{X}_i = \begin{cases} X_i & \text{if } |X_i| \geq \gamma \\ 0 & \text{if } |X_i| < \gamma \end{cases} \quad (3.8)$$

Important: Apply thresholding only to the coefficients $i \geq 1$ (preserve the DC component).

Step 3: Synthesis

$$\hat{S} = D \hat{X} \quad (3.9)$$

Theorem 3.3.1 (Universal Thresholding): For AWGN with variance σ^2 , the optimal threshold for denoising is given by:

$$\gamma = \sigma \sqrt{2 \log(n)} \quad (3.10)$$

where:

- σ is the noise standard deviation
- n is the dimension of coefficients vector

For 8×8 patches: $\gamma \approx 3\sigma$

3.4 Noise Standard Deviation Estimation

To use the universal thresholding, we need to estimate σ from the noisy image itself.

Theorem 3.4.1 (Robust Estimation of Noise Standard Deviation): Given a noisy image Z , the noise standard deviation can be estimated using the following robust method:

$$\hat{\sigma} = \frac{\text{MAD}}{0.6745 \times \sqrt{2}} \quad (3.11)$$

where MAD = Median Absolute Deviation, defined as:

$$\text{MAD}(D) = \text{median}(|D - \text{median}(D)|) \quad (3.12)$$

where D denotes the horizontal differences of the image.

3.5 Sliding DCT Algorithm

3.5.1 Non-Overlapping Tiles (No Aggregation)

The simplest approach processes the image in non-overlapping 8×8 blocks. Let Z be the noisy image of size $M \times N$, partitioned into non-overlapping blocks $B_{i,j}$ of size 8×8 . For each block:

$$X_{i,j} = \text{DCT}_2(B_{i,j}) \quad (3.13)$$

Apply hard thresholding with universal threshold $\gamma = 3\sigma$:

$$\hat{X}_{i,j}(u, v) = \begin{cases} X_{i,j}(u, v) & \text{if } |X_{i,j}(u, v)| \geq \gamma \\ 0 & \text{otherwise} \end{cases} \quad (3.14)$$

Reconstruct each block:

$$\hat{S}_{i,j} = \text{IDCT}_2(\hat{X}_{i,j}) \quad (3.15)$$

The final denoised image is the union of all processed blocks:

$$\hat{Y} = \bigcup_{i,j} \hat{S}_{i,j} \quad (3.16)$$

Properties:

- Complexity: $O(N)$ operations
- Blocking artifacts due to independent processing
- Fast but lower quality

3.5.2 Sliding Window with Uniform Weights

For overlapping patches, each pixel receives multiple estimates that must be aggregated.

For each patch position (i, j) with step size $\text{STEP} = 1$, extract $p \times p$ patch $S_{i,j}$ from the noisy image and do the same DCT processing as before, getting the reconstructed patch $\hat{S}_{i,j}$.

The denoised image \hat{I} is obtained by weighted aggregation:

$$\hat{I}(m, n) = \frac{\sum_{i,j \in \Omega_{m,n}} w \cdot \hat{S}_{i,j}(m - i, n - j)}{\sum_{i,j \in \Omega_{m,n}} w + \varepsilon} \quad (3.17)$$

where $w = 1.0$ is the uniform weight, $\varepsilon = 10^{-8}$ prevents division by zero, and $\Omega_{m,n}$ denotes the set of patch positions (i, j) that contain pixel (m, n) .

For an image of size $M \times N$, each pixel (m, n) can be covered by at most $p \times p$ overlapping patches when using unit step size, but the actual number depends on the pixel's position relative to image boundaries.

3.5.3 Sparsity-Adaptive Weight Aggregation

For overlapping patches, each pixel receives multiple estimates that must be aggregated with weights adapted to the sparsity of the thresholded DCT coefficients.

For each patch position (i, j) with step size $\text{STEP} = 1$, extract $p \times p$ patch $S_{i,j}$ from the noisy image, getting the reconstructed patch $\hat{S}_{i,j}$.

The sparsity-adaptive weight for each patch is computed based on the number of **non-zero coefficients** after thresholding:

$$w_{i,j} = \text{nnz}(\hat{X}_{i,j}) \quad (3.18)$$

where $\text{nnz}(\hat{X}_{i,j})$ counts the number of non-zero elements in the thresholded DCT coefficient matrix.

The denoised image \hat{I} is obtained by sparsity-weighted aggregation:

$$\hat{I}(m, n) = \frac{\sum_{i,j \in \Omega_{m,n}} w_{i,j} \cdot \hat{S}_{i,j}(m - i, n - j)}{\sum_{i,j \in \Omega_{m,n}} w_{i,j} + \varepsilon} \quad (3.19)$$

where $\varepsilon = 10^{-8}$ prevents division by zero, and $\Omega_{m,n}$ denotes the set of patch positions (i, j) that contain pixel (m, n) .

3.6 Wiener Filter

The **Wiener filter** is a powerful tool for image denoising, particularly when the noise characteristics are known. It operates in the frequency domain, leveraging the DCT coefficients to perform adaptive filtering based on local statistics.

3.6.1 Empirical Wiener Filter

Let $\hat{\mathbf{y}}^{\text{HT}}$ be the **hard threshold estimate**, with DCT coefficients:

$$\hat{\mathbf{x}}^{\text{HT}} = D^T \hat{\mathbf{y}}^{\text{HT}} \quad (3.20)$$

The empirical Wiener filter attenuates the DCT coefficients as:

$$\hat{x}_i^{\text{Wie}} = \frac{(\hat{x}_i^{\text{HT}})^2}{(\hat{x}_i^{\text{HT}})^2 + \sigma^2} x_i \quad (3.21)$$

The empirical Wiener estimate is thus:

$$\hat{\mathbf{y}}^{\text{HT}} = D \hat{\mathbf{x}}^{\text{Wie}} \quad (3.22)$$

3.6.2 Transform Domain Patch Processing

Given an image \mathbf{Y} of size $M \times M$, we extract overlapping patches $\mathbf{P}_{i,j}$ of size $p \times p$ centered at pixel (i, j) :

$$\mathbf{P}_{i,j} = \mathbf{Y} \left[i - \left\lfloor \frac{p}{2} \right\rfloor : i + \left\lfloor \frac{p}{2} \right\rfloor, j - \left\lfloor \frac{p}{2} \right\rfloor : j + \left\lfloor \frac{p}{2} \right\rfloor \right] \quad (3.23)$$

For each patch $\mathbf{P}_{i,j}$, we apply the following procedure:

1. **Vectorization:** Convert patch to vector $\mathbf{p}_{i,j} \in \mathbb{R}^{p^2}$
2. **Transformation:** Apply orthogonal transform $\tilde{\mathbf{p}}_{i,j} = \mathbf{T} \mathbf{p}_{i,j}$
3. **Preliminary Estimation:** Obtain initial estimate $\hat{\mathbf{p}}_{i,j}^{(0)}$ using a simple denoising method
4. **Wiener Filtering:** Apply empirical Wiener filter coefficient-wise
5. **Inverse Transform:** Reconstruct patch $\hat{\mathbf{p}}_{i,j} = \mathbf{T}^{-1} \hat{\mathbf{p}}_{i,j}$

3.7 The Sparsity Problem

While orthonormal bases provide computational convenience and guarantee unique representations, they suffer from a fundamental limitation: *no single orthonormal basis can provide sparse representations for all signals of interest.*

Example 3.7.1 (DCT Basis Limitation): Consider a signal $\mathbf{s}_0 \in \mathbb{R}^n$ that admits a sparse representation with respect to the Discrete Cosine Transform (DCT) basis \mathbf{D}_{DCT} :

$$\mathbf{s}_0 = \mathbf{D}_{\text{DCT}} \mathbf{x}_0 \quad (3.24)$$

where \mathbf{x}_0 is sparse (most entries are zero).

Now consider the modified signal:

$$\mathbf{s} = \mathbf{s}_0 + \lambda \mathbf{e}_j \quad (3.25)$$

where \mathbf{e}_j is the j -th canonical basis vector and $\lambda \in \mathbb{R}$ is a scaling factor.

The DCT representation of \mathbf{s} becomes:

$$\mathbf{x} = \mathbf{D}_{\text{DCT}}^T \mathbf{s} = \mathbf{D}_{\text{DCT}}^T \mathbf{s}_0 + \lambda \mathbf{D}_{\text{DCT}}^T \mathbf{e}_j = \mathbf{x}_0 + \lambda \mathbf{D}_{\text{DCT}}^T \mathbf{e}_j \quad (3.26)$$

Since $\mathbf{D}_{\text{DCT}}^T \mathbf{e}_j$ is typically dense (all entries are non-zero), the addition of a single spike destroys the sparsity of the representation.

3.8 Experimental Demonstration

The experimental verification of this limitation involves:

1. Generate a sparse signal \mathbf{s}_0 with respect to DCT basis
2. Add a single spike: $\mathbf{s} = \mathbf{s}_0 + \lambda \mathbf{e}_j$
3. Compute DCT coefficients of both signals
4. Observe the loss of sparsity in the modified signal

The results consistently show that the addition of a single spike causes all DCT coefficients to become significant, effectively destroying the sparse structure that denoising algorithms rely upon.

4 Introduction to Dictionary Learning

The solution to the sparsity limitation lies in abandoning the constraint of orthonormality and embracing redundancy. Instead of using a single $n \times n$ orthonormal basis, we construct an $n \times m$ dictionary matrix \mathbf{D} where $m > n$.

Definition 4.1 (Overcomplete Dictionary): An *overcomplete dictionary* is a matrix $\mathbf{D} \in \mathbb{R}^{n \times m}$ with $m > n$ such that:

$$\text{span}\{\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_m\} = \mathbb{R}^n \quad (4.1)$$

where \mathbf{d}_i are the columns of \mathbf{D} .

For the DCT-spike example, we construct the overcomplete dictionary by concatenating the DCT basis with the **canonical basis**:

$$\mathbf{D} = (\mathbf{D}_{\text{DCT}} \ \mathbf{I}) \in \mathbb{R}^{n \times 2n} \quad (4.2)$$

This construction ensures that:

- Signals sparse in DCT domain remain sparse
- Signals sparse in canonical domain remain sparse
- Mixed signals (DCT-sparse + spikes) admit sparse representations

Example 4.1 (Sparse Representation with Overcomplete Dictionary): Consider the signal $\mathbf{s} = \mathbf{s}_0 + \lambda \mathbf{e}_j$ where $\mathbf{s}_0 = \mathbf{D}_{\text{DCT}} \mathbf{x}_0$ with sparse \mathbf{x}_0 .

The representation with respect to the overcomplete dictionary is:

$$\mathbf{s} = \mathbf{D} \begin{pmatrix} \mathbf{x}_0 \\ \lambda \mathbf{e}_j \end{pmatrix} \quad (4.3)$$

The coefficient vector $\begin{pmatrix} \mathbf{x}_0 \\ \lambda \mathbf{e}_j \end{pmatrix} \in \mathbb{R}^{2n}$ is sparse, containing only the non-zero entries of \mathbf{x}_0 plus the single entry λ at position j in the second block.

4.1 Theoretical Properties of Overcomplete Systems

Theorem 4.1.1 (Rouché-Capelli Theorem): Consider the linear system $\mathbf{D}\mathbf{x} = \mathbf{s}$ where $\mathbf{D} \in \mathbb{R}^{n \times m}$ and $\mathbf{s} \in \mathbb{R}^n$. The system admits a solution if and only if:

$$\text{rank}(\mathbf{D}) = \text{rank}((\mathbf{D} \ \mathbf{s})) \quad (4.4)$$

When $m > n$ and $\text{rank}(\mathbf{D}) = n$, the system has infinitely many solutions forming an affine subspace of dimension $m - n$.

Theorem 4.1.2 (Solution Space Dimension): If $\mathbf{D} \in \mathbb{R}^{n \times m}$ with $m > n$ and $\text{rank}(\mathbf{D}) = n$, then for any $\mathbf{s} \in \mathbb{R}^n$, the solution set of $\mathbf{D}\mathbf{x} = \mathbf{s}$ forms an affine subspace of dimension $m - n$.

The abundance of solutions in overcomplete systems necessitates additional criteria for solution selection. This is where **regularization theory** becomes essential.

4.2 Regularization and Sparse Recovery

Definition 4.2.1 (Regularization): Given an **ill-posed** problem $D\mathbf{x} = \mathbf{s}$ with multiple solutions, *regularization* involves solving:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} J(\mathbf{x}) \quad \text{subject to} \quad D\mathbf{x} = \mathbf{s} \quad (4.5)$$

where $J : \mathbb{R}^m \rightarrow \mathbb{R}_+$ is a regularization functional encoding our prior knowledge about the desired solution.

4.2.1 ℓ_2 Regularization: Ridge Regression

The most mathematically tractable regularization is the ℓ_2 norm:

$$J(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2 = \frac{1}{2} \sum_{i=1}^m x_i^2 \quad (4.6)$$

This leads to the constrained optimization problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x}\|_2^2 \quad \text{subject to} \quad D\mathbf{x} = \mathbf{s} \quad (4.7)$$

Alternatively, we can formulate the unconstrained version:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2} \|D\mathbf{x} - \mathbf{s}\|_2^2 + \frac{\lambda}{2} \|\mathbf{x}\|_2^2 \quad (4.8)$$

Theorem 4.2.1.1 (Ridge Regression Solution): The solution to the ridge regression problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2} \|D\mathbf{x} - \mathbf{s}\|_2^2 + \frac{\lambda}{2} \|\mathbf{x}\|_2^2 \quad (4.9)$$

is given by:

$$\hat{\mathbf{x}} = (D^T D + \lambda I)^{-1} D^T \mathbf{s} \quad (4.10)$$

where $\lambda > 0$ ensures the matrix $(D^T D + \lambda I)$ is invertible.

Proof: Define the objective function:

$$f(\mathbf{x}) = \frac{1}{2} \|D\mathbf{x} - \mathbf{s}\|_2^2 + \frac{\lambda}{2} \|\mathbf{x}\|_2^2 \quad (4.11)$$

Expanding the squared norms:

$$\begin{aligned} f(\mathbf{x}) &= \frac{1}{2} (D\mathbf{x} - \mathbf{s})^T (D\mathbf{x} - \mathbf{s}) + \frac{\lambda}{2} \mathbf{x}^T \mathbf{x} \\ &= \frac{1}{2} \mathbf{x}^T D^T D \mathbf{x} - D\mathbf{x}^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \mathbf{s} + \frac{\lambda}{2} \mathbf{x}^T \mathbf{x} \end{aligned} \quad (4.12)$$

Taking the gradient with respect to \mathbf{x} :

$$\nabla f(\mathbf{x}) = D^T D \mathbf{x} - D^T \mathbf{s} + \lambda \mathbf{x} \quad (4.13)$$

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

$$(D^T D + \lambda I) \mathbf{x} = D^T \mathbf{s} \quad (4.14)$$

Since $\lambda > 0$, the matrix $(D^T D + \lambda I)$ is positive definite and therefore invertible, yielding the stated solution.

4.2.2 Limitations of ℓ_2 Regularization

While ℓ_2 regularization provides a computationally efficient solution, it does not promote sparsity.

The solution $\hat{\mathbf{x}}$ typically has all non-zero entries, which contradicts our goal of sparse representation.

⚠ Sparsity vs. ℓ_2 Regularization

attention 4.2.1

The ℓ_2 norm penalizes large coefficients but does not drive small coefficients to zero. For sparse recovery, we need regularization functionals that promote sparsity, such as the ℓ_1 norm or ℓ_0 pseudo-norm.

4.3 Towards Sparsity: ℓ_0 and ℓ_1 Regularization

4.3.1 The ℓ_0 “Norm” and True Sparsity

The most natural regularization for sparse recovery is the ℓ_0 “norm” (technically a pseudo-norm):

$$\|\mathbf{x}\|_0 = |\{i : x_i \neq 0\}| \quad (4.15)$$

This counts the number of non-zero entries in \mathbf{x} . The corresponding optimization problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_0 \quad \text{subject to} \quad \mathbf{D}\mathbf{x} = \mathbf{s} \quad (4.16)$$

directly seeks the sparsest representation.

4.3.2 Computational Challenges

The ℓ_0 minimization problem is NP-hard in general, making it computationally intractable for large-scale problems. This has led to the development of convex relaxations and approximation algorithms.

5 Sparse Coding In ℓ_0 sense

5.1 The Sparse Coding Problem

Given the desire for sparse representations, a natural formulation is to seek the sparsest possible α such that:

$$\mathbf{x} = D\alpha \quad (5.1)$$

This leads to the following optimization problem:

$$\min_{\alpha \in \mathbb{R}^n} \|\alpha\|_0 \quad \text{subject to} \quad \mathbf{x} = D\alpha \quad (5.2)$$

This is often referred to as the **P_0 problem**.

Goal: Among all solutions α such that $D\alpha = \mathbf{x}$, select the sparsest one.

Challenge: The ℓ_0 norm is non-convex, discontinuous, and leads to combinatorial complexity. Solving Equation (5.2) exactly is NP-hard in general.

5.1.1 Union of Subspaces Interpretation

Let us assume $D \in \mathbb{R}^{m \times n}$ is a dictionary with $n > m$, i.e., an overcomplete dictionary.

Suppose we restrict α to have at most s non-zero entries. Then the image $D\alpha$ lies in a subspace spanned by s columns of D .

Definition 5.1.1.1 (Sparsity-Induced Subspace): If α has $\|\alpha\|_0 \leq s$, then $\mathbf{x} = D\alpha$ lies in a subspace $S_\omega := \text{span}(D_\omega)$, where ω is the support of α and D_ω denotes the sub-matrix of D restricted to columns indexed by ω .

Therefore, the space of all s -sparse representations is the union of all such s -dimensional subspaces:

$$\mathcal{M}_s := \bigcup_{\omega \subset \{1, \dots, n\}, |\omega| \leq s} \text{span}(D_\omega) \quad (5.3)$$

Interpretation: Sparse modeling corresponds to finding the best subspace (among exponentially many) in which to approximate the signal \mathbf{x} .

Key Point: Unlike PCA (which projects onto a single global subspace), sparse coding projects onto a *union of low-dimensional subspaces*, selected adaptively based on the input \mathbf{x} .

5.1.2 A 2D Illustration of Sparsity

Suppose we are working in \mathbb{R}^2 and D has 3 atoms:

$$D = (\mathbf{d}_1 \ \mathbf{d}_2 \ \mathbf{d}_3), \quad D \in \mathbb{R}^{2 \times 3} \quad (5.4)$$

Each \mathbf{d}_i is a column vector in \mathbb{R}^2 .

Let $\mathbf{x} \in \mathbb{R}^2$ be a signal we wish to approximate. If we restrict $\|\alpha\|_0 = 1$, then the approximant $D\alpha$ must lie along one of the directions \mathbf{d}_1 , \mathbf{d}_2 , or \mathbf{d}_3 .

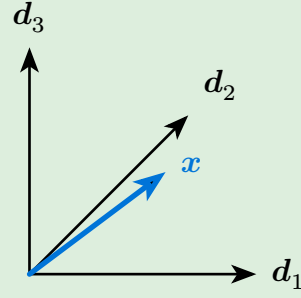


Figure 5.1: Approximation of \mathbf{x} via projections onto sparse atoms

Interpretation:

- We seek the atom \mathbf{d}_i such that the projection of \mathbf{x} onto $\text{span}(\mathbf{d}_i)$ minimizes the residual.
- This is the best 1-sparse approximation.

In general: If $\|\alpha\|_0 \leq s$, the approximation lives in a union of $\binom{n}{s}$ subspaces.

5.1.3 Combinatorial Intractability

Why is P_0 hard? To find the optimal s -sparse representation of \mathbf{x} , one must:

1. Enumerate all subsets $\omega \subset \{1, \dots, n\}$ of size s
2. Solve the least squares problem:

$$\alpha_\omega = \arg \min_{\mathbf{z} \in \mathbb{R}^s} \|D_\omega \mathbf{z} - \mathbf{x}\|_2^2 \quad (5.5)$$

3. Select the best ω minimizing the residual.

The number of subsets grows exponentially:

$$\#\text{subspaces} = \binom{n}{s} \sim \left(n \frac{e}{s}\right)^s \quad (5.6)$$

Illustrative Example: Let $n = 1000$, $s = 20$. Then:

$$\binom{1000}{20} \approx 10^{34} \quad (5.7)$$

Assuming a billion operations per second, exhaustive search would take more time than the age of the universe.

Conclusion: The ℓ_0 sparse coding problem is combinatorially explosive. Efficient approximations are necessary.

5.2 Matching Pursuit Algorithm

The **Matching Pursuit (MP)** algorithm embodies the greedy principle for sparse coding:

Input: Signal \mathbf{y} , dictionary \mathbf{D} (normalized), stopping criterion

Output: Sparse representation \mathbf{x}

1. **Initialize:**

$$\begin{aligned} \mathbf{x}^{(0)} &= \mathbf{0} && \text{(coefficient vector)} \\ \mathbf{r}^{(0)} &= \mathbf{y} && \text{(residual)} \\ \Omega^{(0)} &= \emptyset && \text{(active set)} \\ k &= 0 && \text{(iteration counter)} \end{aligned} \tag{5.8}$$

2. **Sweep Stage:** For each atom $j = 1, \dots, n$, compute the approximation error:

$$E_j^{(k)} = \|\mathbf{r}^{(k)} - \langle \mathbf{d}_j, \mathbf{r}^k \rangle \mathbf{d}_j\|_2^2 \tag{5.9}$$

3. **Atom Selection:** Choose the atom with minimum error:

$$j^* = \arg \min_{j=1, \dots, n} E_j^{(k)} \tag{5.10}$$

Equivalently (by maximizing correlation):

$$j^* = \arg \max_{j=1, \dots, n} \langle \mathbf{d}_j, \mathbf{r}^k \rangle \tag{5.11}$$

4. **Coefficient Update:** Compute the projection coefficient:

$$z_{j^*}^{(k)} = \langle \mathbf{d}_{j^*}, \mathbf{r}^k \rangle \tag{5.12}$$

5. **Solution Update:**

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + z_{j^*}^{(k)} \mathbf{e}_{j^*} \tag{5.13}$$

where \mathbf{e}_{j^*} is the j^* -th standard basis vector.

6. **Residual Update:**

$$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - z_{j^*}^{(k)} \mathbf{d}_{j^*} \tag{5.14}$$

7. **Active Set Update:**

$$\Omega^{(k+1)} = \Omega^{(k)} \cup \{j^*\} \tag{5.15}$$

8. **Stopping Criteria:** Terminate if:

- $|\Omega^{(k+1)}| \geq k_{\max}$ (maximum sparsity reached)
- $\|\mathbf{r}^{(k+1)}\|_2 \leq \varepsilon$ (residual threshold met)

Otherwise, set $k \leftarrow k + 1$ and return to step 2.

5.2.1 Residual Monotonicity

The Matching Pursuit algorithm produces a monotonically decreasing sequence of residual norms:

$$\|\mathbf{r}^{(k+1)}\|_2 \leq \|\mathbf{r}^{(k)}\|_2 \quad (5.16)$$

with strict inequality unless $\mathbf{r}^{(k)}$ is orthogonal to all dictionary atoms.

5.2.2 Atom Reselection

Unlike orthogonal methods, Matching Pursuit may select the same atom multiple times in successive iterations. This occurs because:

1. The algorithm does not enforce orthogonality of residuals to previously selected atoms
2. Residual components may align with previously selected atoms after updates
3. This can lead to slower convergence compared to orthogonal variants

5.2.3 Convergence Analysis

For any finite dictionary \mathbf{D} and signal \mathbf{y} , the Matching Pursuit algorithm converges in the sense that:

$$\lim_{k \rightarrow \infty} \|\mathbf{r}^{(k)}\|_2 = \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2 \quad (5.17)$$

Furthermore, if $\mathbf{y} \in \text{span}(\mathbf{D})$, then the algorithm achieves exact recovery in finite steps.

5.2.4 Approximation Quality

While Matching Pursuit provides computational tractability, it may not achieve the globally optimal sparse solution. The quality of approximation depends on the coherence structure of the dictionary.

The **coherence** of a dictionary \mathbf{D} with normalized columns is:

$$\mu(\mathbf{D}) = \max_{i \neq j} |\mathbf{d}_i^T \mathbf{d}_j| \quad (5.18)$$

Under certain conditions on dictionary coherence and signal sparsity, Matching Pursuit provides approximation guarantees. Specifically, if the true sparse representation has sparsity k and the dictionary satisfies appropriate coherence conditions, then MP recovers a solution with controlled approximation error.

5.2.5 Computational Complexity

Each iteration of Matching Pursuit requires:

- $\mathcal{O}(mn)$ operations for the sweep stage (computing all correlations)
- $\mathcal{O}(m)$ operations for residual update
- Total per-iteration complexity: $\mathcal{O}(mn)$

For k iterations, the total complexity is $\mathcal{O}(kmn)$, which is polynomial and practically feasible.

5.3 Orthogonal Matching Pursuit

Orthogonal Matching Pursuit (OMP) is a variant of Matching Pursuit that enforces orthogonality of the residuals to previously selected atoms. This leads to improved convergence properties and better approximation quality.

The OMP algorithm embodies the greedy principle for sparse coding with orthogonal projections:

Input: Signal \mathbf{y} , dictionary \mathbf{D} (normalized), stopping criterion

Output: Sparse representation \mathbf{x}

1. **Initialize:**

$$\begin{aligned} \mathbf{x}^{(0)} &= \mathbf{0} && \text{(coefficient vector)} \\ \mathbf{r}^{(0)} &= \mathbf{y} && \text{(residual)} \\ \Omega^{(0)} &= \emptyset && \text{(active set)} \\ k &= 0 && \text{(iteration counter)} \end{aligned} \tag{5.19}$$

2. **Atom Selection:** Choose the atom with maximum correlation:

$$j^* = \arg \max_{j=1, \dots, n} |\langle \mathbf{d}_j, \mathbf{r}^{(k)} \rangle| \tag{5.20}$$

3. **Active Set Update:**

$$\Omega^{(k+1)} = \Omega^{(k)} \cup \{j^*\} \tag{5.21}$$

4. **Orthogonal Projection:** Solve the least squares problem over the active set:

$$\boldsymbol{\alpha}_{\Omega}^{(k+1)} = \arg \min_{\boldsymbol{\alpha}} \|\mathbf{D}_{\Omega} \boldsymbol{\alpha} - \mathbf{y}\|_2^2 \tag{5.22}$$

where \mathbf{D}_{Ω} is the sub-matrix of \mathbf{D} with columns indexed by $\Omega^{(k+1)}$.

5. **Solution Update:**

$$\mathbf{x}_j^{(k+1)} = \begin{cases} \alpha_j^{(k+1)} & \text{if } j \in \Omega^{(k+1)} \\ 0 & \text{otherwise} \end{cases} \tag{5.23}$$

6. **Residual Update:** Compute the orthogonal residual:

$$\mathbf{r}^{(k+1)} = \mathbf{y} - \mathbf{D}_{\Omega} \boldsymbol{\alpha}_{\Omega}^{(k+1)} \tag{5.24}$$

7. **Stopping Criteria:** Terminate if:

- $|\Omega^{(k+1)}| \geq k_{\max}$ (maximum sparsity reached)
- $\|\mathbf{r}^{(k+1)}\|_2 \leq \varepsilon$ (residual threshold met)

Otherwise, set $k \leftarrow k + 1$ and return to step 2.

5.4 OMP-Based Image Denoising

The integration of OMP into image denoising frameworks requires careful consideration of patch processing, dictionary design, and aggregation strategies.

Natural images exhibit strong local correlations but varying global statistics. The patch-based approach decomposes the image into overlapping patches, each processed independently:

Algorithm 2: OMP-Based Image Denoising

```

1: procedure OMP-IMAGE-DENOISING( $\mathbf{Y}$ ,  $\mathbf{D}$ ,  $s$ )
2:    $\triangleright$  Input: Noisy image  $\mathbf{Y} \in \mathbb{R}^{N \times M}$ , dictionary  $\mathbf{D} \in \mathbb{R}^{n \times m}$ , sparsity level  $s$ 
3:    $\triangleright$  Output: Denoised image  $\hat{\mathbf{Y}}$ 
4:
5:    $\triangleright$  Patch Extraction:
6:    $\triangleright$  For image  $\mathbf{Y} \in \mathbb{R}^{N \times M}$ , extract patches  $(\mathbf{y}_i, i = 1, \dots, P)$ ,  $P$  is the number of patches
7:    $\triangleright$  where each  $\mathbf{y}_i \in \mathbb{R}^n$  represents a vectorized  $\sqrt{n} \times \sqrt{n}$  patch
8:   for  $i = 1, \dots, P$  do
9:      $\mathbf{y}_i \leftarrow$  extract  $\sqrt{n} \times \sqrt{n}$  patch and vectorize
10:  end
11:
12:   $\triangleright$  Mean Computation and Centering:
13:  for  $i = 1, \dots, P$  do
14:     $\triangleright$  Compute mean:  $\mu_i = \frac{1}{n} \sum_{j=1}^n y_{i,j}$ 
15:     $\mu_i \leftarrow \frac{1}{n} \sum_{j=1}^n y_{i,j}$ 
16:
17:     $\triangleright$  Center patch:  $\tilde{\mathbf{y}}_i = \mathbf{y}_i - \mu_i \mathbf{1}$ 
18:     $\tilde{\mathbf{y}}_i \leftarrow \mathbf{y}_i - \mu_i \mathbf{1}$ 
19:     $\triangleright$  where  $\mathbf{1} \in \mathbb{R}^n$  is the vector of all ones
20:  end
21:
22:   $\triangleright$  Sparse Coding:
23:  for  $i = 1, \dots, P$  do
24:     $\triangleright$  Apply OMP to each mean-centered patch
25:     $\hat{\boldsymbol{\alpha}}_i \leftarrow \text{OMP}(\mathbf{D}, \tilde{\mathbf{y}}_i, s)$ 
26:  end
27:
28:   $\triangleright$  Reconstruction:
29:  for  $i = 1, \dots, P$  do
30:     $\triangleright$  Compute denoised patches by adding back the mean
31:     $\hat{\mathbf{x}}_i \leftarrow \mathbf{D} \hat{\boldsymbol{\alpha}}_i + \mu_i \mathbf{1}$ 
32:  end
33:
34:   $\triangleright$  Aggregation:
35:   $\triangleright$  Reconstruct the full image by averaging overlapping reconstructions
36:   $\triangleright$  at each pixel location
37:   $\hat{\mathbf{Y}} \leftarrow \text{Aggregate}(\{\hat{\mathbf{x}}_i\}_{i=1}^P)$ 
38:
39:  return  $\hat{\mathbf{Y}}$ 
40: end

```

5.5 Linearity Analysis of Sparse Coding Algorithms

A fundamental question in sparse coding concerns the linearity properties of the resulting algorithms. An algorithm \mathcal{A} is considered linear if and only if it satisfies:

Definition 5.5.1 (Linear Algorithm): An algorithm $\mathcal{A} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is linear if and only if for all $\alpha, \beta \in \mathbb{R}$ and $\mathbf{y}_1, \mathbf{y}_2 \in \mathbb{R}^m$:

$$\mathcal{A}(\alpha \mathbf{y}_1 + \beta \mathbf{y}_2) = \alpha \mathcal{A}(\mathbf{y}_1) + \beta \mathcal{A}(\mathbf{y}_2) \quad (5.25)$$

5.5.1 Nonlinearity of OMP

Despite the final solution taking the form of a linear projection:

$$\hat{\mathbf{x}}_{\mathcal{S}} = (\mathbf{D}_{\mathcal{S}}^T \mathbf{D}_{\mathcal{S}})^{-1} \mathbf{D}_{\mathcal{S}}^T \mathbf{y} \quad (5.26)$$

the OMP algorithm is fundamentally nonlinear due to the adaptive selection of the support set \mathcal{S} .

Proposition 5.5.1.1 (Nonlinearity of OMP): The OMP algorithm is nonlinear because the support selection depends on the input signal: $\mathcal{S}(\mathbf{y})$ is a function of \mathbf{y} .

Consider two signals \mathbf{y}_1 and \mathbf{y}_2 that would result in different support sets under OMP. The support of $\mathbf{y}_1 + \mathbf{y}_2$ may differ from the union of individual supports, violating the linearity condition.

Comparison with Linear Denoising Methods

Fixed linear methods like convolution-based filtering or PCA projection onto the first k principal components are linear but less adaptive. The nonlinearity of OMP enables signal-dependent subspace selection, providing superior performance for sparse signals.

5.6 Uniqueness Guarantees for Sparse Solutions

The following theorem provides conditions under which the solution to the ℓ_0 constraint problem is unique.

Theorem 5.6.1 (Uniqueness of ℓ_0 Solutions): Consider the system $D\mathbf{x} = \mathbf{y}$ where $D \in \mathbb{R}^{\{m \times n\}}$. If there exists a solution $\hat{\mathbf{x}}$ such that:

$$\|\hat{\mathbf{x}}\|_0 < \frac{1}{2} \text{spark}(D) \quad (5.27)$$

then $\hat{\mathbf{x}}$ is the unique solution to the ℓ_0 constraint problem.

Proof: Suppose, for the sake of contradiction, that there exists another solution $\tilde{\mathbf{x}} \neq \hat{\mathbf{x}}$ such that $D\tilde{\mathbf{x}} = \mathbf{y}$.

Since both $\hat{\mathbf{x}}$ and $\tilde{\mathbf{x}}$ satisfy the linear system:

$$\begin{aligned} D\hat{\mathbf{x}} &= \mathbf{y} \\ D\tilde{\mathbf{x}} &= \mathbf{y} \end{aligned} \quad (5.28)$$

Subtracting these equations yields:

$$D(\hat{\mathbf{x}} - \tilde{\mathbf{x}}) = \mathbf{0} \quad (5.29)$$

This shows that $\hat{\mathbf{x}} - \tilde{\mathbf{x}}$ is a non-zero solution to the homogeneous system. By [Lemma 12.5.1](#):

$$\text{spark}(D) \leq \|\hat{\mathbf{x}} - \tilde{\mathbf{x}}\|_0 \quad (5.30)$$

The definition of spark see Section 12.5.

Using the triangle inequality for the ℓ_0 pseudo-norm:

$$\|\hat{\mathbf{x}} - \tilde{\mathbf{x}}\|_0 \leq \|\hat{\mathbf{x}}\|_0 + \|\tilde{\mathbf{x}}\|_0 \quad (5.31)$$

Combining these inequalities:

$$\text{spark}(D) \leq \|\hat{\mathbf{x}}\|_0 + \|\tilde{\mathbf{x}}\|_0 \quad (5.32)$$

Since $\tilde{\mathbf{x}}$ is also assumed to be a solution to the ℓ_0 problem, and $\hat{\mathbf{x}}$ is the optimal solution:

$$\|\tilde{\mathbf{x}}\|_0 \geq \|\hat{\mathbf{x}}\|_0 \quad (5.33)$$

Therefore:

$$\text{spark}(D) \leq 2\|\hat{\mathbf{x}}\|_0 \quad (5.34)$$

This contradicts our assumption that $\|\hat{\mathbf{x}}\|_0 < \frac{1}{2} \text{spark}(D)$. Hence, $\hat{\mathbf{x}}$ is unique.

▲ Practical Limitations

attention 5.6.1

While theoretically elegant, the uniqueness conditions are often too restrictive in practice:

- Computing $\text{spark}(D)$ is computationally intractable for large matrices
- The bound $\|\hat{\mathbf{x}}\|_0 < \frac{1}{2} \text{spark}(D)$ is often very conservative
- Real-world signals may not satisfy the sparsity requirements

5.7 Application to Image Inpainting

Image inpainting addresses the reconstruction of missing or corrupted pixels in digital images. Using sparse coding theory, we can formulate inpainting as a sparse reconstruction problem.

Consider an image patch $\mathbf{s}_0 \in \mathbb{R}^n$ (vectorized) and its corrupted version $\mathbf{s} \in \mathbb{R}^n$ where some pixels are missing or corrupted. The relationship between them can be expressed as:

$$\mathbf{s} = \mathbf{\Omega} \mathbf{s}_0 \quad (5.35)$$

where $\mathbf{\Omega} \in \mathbb{R}^{n \times n}$ is a diagonal matrix with:

$$\Omega_{ii} = \begin{cases} 1 & \text{if pixel } i \text{ is known} \\ 0 & \text{if pixel } i \text{ is missing} \end{cases} \quad (5.36)$$

Assuming the original patch admits a sparse representation:

$$\mathbf{s}_0 = \mathbf{D} \mathbf{x}_0 \quad (5.37)$$

where $\mathbf{D} \in \mathbb{R}^{n \times m}$ is an overcomplete dictionary and \mathbf{x}_0 is sparse, the corrupted patch becomes:

$$\mathbf{s} = \mathbf{\Omega} \mathbf{D} \mathbf{x}_0 = \mathbf{D}_\Omega \mathbf{x}_0 \quad (5.38)$$

where $\mathbf{D}_\Omega = \mathbf{\Omega} \mathbf{D}$ represents the ‘‘inpainted dictionary.’’

The key insight is that the spark of the inpainted dictionary relates to the original dictionary:

Proposition 5.7.1 (Spark of Inpainted Dictionary): For the inpainted dictionary $\mathbf{D}_\Omega = \mathbf{\Omega} \mathbf{D}$:

$$\text{spark}(\mathbf{D}_\Omega) \geq \text{spark}(\mathbf{D}) \quad (5.39)$$

Proof: Removing rows (zeroing out pixels) from a matrix cannot decrease the spark, as linear dependencies between columns are preserved or potentially eliminated.

Let \mathbf{s}_0 be an image patch with sparse representation $\mathbf{s}_0 = \mathbf{D} \mathbf{x}_0$ where $\|\mathbf{x}_0\|_0 < \frac{1}{2} \text{spark}(\mathbf{D}_\Omega)$. Then:

1. The sparse coding problem $\min_{\mathbf{x}} \|\mathbf{x}\|_0$ subject to $\mathbf{D}_\Omega \mathbf{x} = \mathbf{s}$ has a unique solution \mathbf{x}_0
2. The reconstruction $\hat{\mathbf{s}}_0 = \mathbf{D} \mathbf{x}_0$ perfectly recovers the original patch

Proof: The proof follows directly from [Theorem 5.6.1](#) applied to the inpainted dictionary \mathbf{D}_Ω .

Sparse Coding Inpainting Algorithm

Input: Corrupted image patch \mathbf{s} , dictionary \mathbf{D} , mask $\mathbf{\Omega}$

1. **Construct inpainted dictionary:** $\mathbf{D}_\Omega = \mathbf{\Omega} \mathbf{D}$
2. **Solve sparse coding:** $\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_0$ subject to $\mathbf{D}_\Omega \mathbf{x} = \mathbf{s}$
3. **Reconstruct patch:** $\hat{\mathbf{s}}_0 = \mathbf{D} \hat{\mathbf{x}}$

Output: Inpainted patch $\hat{\mathbf{s}}_0$

In practice, step 2 is solved using OMP with the inpainted dictionary:

$$\hat{\mathbf{x}} = \text{OMP}(\mathbf{D}_\Omega, \mathbf{s}, K) \quad (5.40)$$

where K is a predetermined sparsity level. The key insight is that the synthesis step uses the original dictionary \mathbf{D} , not the inpainted dictionary \mathbf{D}_Ω .

6 Sparse Coding In ℓ_1 sense

6.1 The ℓ_1 Optimization Problem

The ℓ_0 norm, defined as the number of non-zero components in a vector, provides the most intuitive measure of sparsity. However, optimization problems involving the ℓ_0 norm are NP-hard due to their combinatorial nature. This computational intractability necessitates the exploration of alternative sparsity-promoting norms that maintain favorable optimization properties.

The ℓ_1 optimization problem for sparse coding can be formulated in two equivalent ways:

Definition 6.1.1 (Constrained ℓ_1 Problem):

$$\min_{\alpha} \|\alpha\|_1 \quad \text{subject to} \quad \|D\alpha - x\|_2 \leq \varepsilon \quad (6.1)$$

Definition 6.1.2 (Regularized ℓ_1 Problem):

$$\min_{\alpha} \frac{1}{2} \|D\alpha - x\|_2^2 + \lambda \|\alpha\|_1 \quad (6.2)$$

6.2 Problem Components Analysis

6.2.1 Data Fidelity Term

The term $\frac{1}{2} \|D\alpha - x\|_2^2$ serves as the data fidelity term, ensuring that the solution α produces a reconstruction $D\alpha$ that is close to the observed signal x .

Proposition 6.2.1.1 (Properties of Data Fidelity Term): The data fidelity term $g(x) = \frac{1}{2} \|D\alpha - x\|_2^2$ is:

1. Convex (as a composition of convex functions)
2. Differentiable with gradient $\nabla g(\alpha) = D^T(D\alpha - x)$
3. Strongly convex if D has full column rank

6.2.2 Regularization Term

The term $\lambda \|\alpha\|_1$ acts as a regularization term, promoting sparsity in the solution.

Proposition 6.2.2.1 (Properties of ℓ_1 Regularization): The regularization term $h(\alpha) = \|\alpha\|_1$ is:

1. Convex
2. Non-differentiable at $\alpha_i = 0$ for any component i
3. Promotes sparsity through its geometric properties

6.3 Proximal Gradient Methods

6.3.1 Proximal Operator

For the composite optimization problem $\min_{\alpha} f(\alpha) + g(\alpha)$, where f is differentiable and g is convex but not necessarily differentiable, we can use the proximal operator to handle the non-differentiable part (See Section 13.2 for details).

Definition 6.3.1.1 (Proximal Operator): The proximal operator of a convex function g is defined as:

$$\text{prox}_{\lambda g}(\mathbf{v}) = \arg \min_{\mathbf{x}} \left\{ \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{v}\|_2^2 + g(\mathbf{x}) \right\} \quad (6.3)$$

The aim of the proximal operator is to find a point that minimizes $g(\mathbf{x})$, while being close to \mathbf{v} .

6.3.2 Proximal Gradient Descent

The proximal gradient descent algorithm iteratively updates the solution by combining the gradient step of the differentiable part and the proximal operator of the non-differentiable part.

For the problem $\min_{\alpha} f(\alpha) + g(\alpha)$, the update rule is: For $k = 0, 1, 2, \dots$ until convergence:

$$\alpha^{k+1} = \text{prox}_{t g}(\alpha^k - t \nabla f(\alpha^k)) \quad (6.4)$$

where t is a step size parameter. Usually, t is chosen based on the Lipschitz constant of ∇f . $t = \frac{1}{L}$, where L is the Lipschitz constant equal to the largest eigenvalue of $\mathbf{D}^T \mathbf{D}$ (See Section 13.1 for details).

6.3.3 Iterative Soft Thresholding Algorithm (ISTA)

Theorem 6.3.3.1 (Soft Thresholding): The proximal operator of the ℓ_1 norm is given by the soft thresholding operator:

$$[\text{prox}_{\lambda \|\cdot\|_1}(\mathbf{v})]_i = \text{sign}(v_i) \max\{|v_i| - \lambda, 0\} \quad (6.5)$$

where $\text{sign}(v_i)$ is the sign function.

Proof: The proximal operator problem becomes:

$$\min_{\mathbf{x}} \left\{ \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{v}\|_2^2 + \|\mathbf{x}\|_1 \right\} \quad (6.6)$$

For each component i , we solve:

$$\min_{x_i} g_{i(x_i)} = \frac{1}{2\lambda} (x_i - v_i)^2 + |x_i| \quad (6.7)$$

The function $g_{i(x_i)}$ is convex but non-differentiable at $x_i = 0$. We use the subdifferential calculus:

- For $x_i > 0$: $\partial g_{i(x_i)} = \frac{1}{\lambda}(x_i - v_i) + 1$
- For $x_i < 0$: $\partial g_{i(x_i)} = \frac{1}{\lambda}(x_i - v_i) - 1$
- For $x_i = 0$: $\partial g_{i(0)} = -\frac{v_i}{\lambda} + [-1, 1]$

Setting the subdifferential to zero:

Case 1: If $x_i^* > 0$, then $\frac{1}{\lambda}(x_i^* - v_i) + 1 = 0$, giving $x_i^* = v_i - \lambda$. This is valid only if $v_i - \lambda > 0$, i.e., $v_i > \lambda$.

Case 2: If $x_i^* < 0$, then $\frac{1}{\lambda}(x_i^* - v_i) - 1 = 0$, giving $x_i^* = v_i + \lambda$. This is valid only if $v_i + \lambda < 0$, i.e., $v_i < -\lambda$.

Case 3: If $x_i^* = 0$, then $0 \in -\frac{v_i}{\lambda} + [-1, 1]$, which means $-1 \leq -\frac{v_i}{\lambda} \leq 1$, or equivalently $-\lambda \leq v_i \leq \lambda$.

Combining all cases:

$$x_i^* = \begin{cases} v_i - \lambda & \text{if } v_i > \lambda \\ 0 & \text{if } -\lambda \leq v_i \leq \lambda \\ v_i + \lambda & \text{if } v_i < -\lambda \end{cases} \quad (6.8)$$

This can be written compactly as:

$$x_i^* = \text{sign}(v_i) \max\{|v_i| - \lambda, 0\} \quad (6.9)$$

The solution is the soft thresholding operator due to the subdifferential analysis of the absolute value function.

Algorithm 3: Iterative Soft Thresholding Algorithm (ISTA)

```

1: procedure ISTA( $D, x, \lambda, \varepsilon$ )
2:    $\triangleright$  Input:
3:    $\triangleright D \in \mathbb{R}^{m \times n}$ : dictionary/measurement matrix
4:    $\triangleright x \in \mathbb{R}^m$ : observation vector
5:    $\triangleright \lambda > 0$ : regularization parameter
6:    $\triangleright \varepsilon > 0$ : convergence tolerance
7:    $\triangleright$  Output:
8:    $\triangleright \alpha^* \in \mathbb{R}^n$ : sparse coefficient vector
9:
10:   $\triangleright$  Initialize
11:   $\alpha^0 \leftarrow \mathbf{0}$ 
12:   $L \leftarrow$  largest eigenvalue of  $D^T D$ 
13:   $t \leftarrow \frac{1}{L}$ 
14:   $k \leftarrow 0$ 
15:
16:  while not converged do
17:     $\triangleright$  Gradient step
18:     $u^k \leftarrow \alpha^k - tD^T(D\alpha^k - x)$ 
19:
20:     $\triangleright$  Proximal step (soft thresholding)
21:    for  $i = 1 < n$  do
22:       $\alpha_i^{k+1} \leftarrow \text{sign}(u_i^k) \cdot \max\{|u_i^k| - \lambda t, 0\}$ 
23:    end
24:
25:     $\triangleright$  Check convergence
26:    if  $\|\alpha^{k+1} - \alpha^k\|_2 < \varepsilon$  then
27:      return  $\alpha^{k+1}$ 
28:    end
29:     $k \leftarrow k + 1$ 
30:  end
31:  return  $\alpha^k$ 
32: end

```

6.3.4 Backtracking Line Search

When computing eigenvalues is impractical, adaptive step size selection via backtracking provides a robust alternative:

Algorithm 4: Backtracking Line Search for ISTA

```

1: procedure BACKTRACKING-LINE-SEARCH-ISTA( $\mathbf{D}$ ,  $\mathbf{x}$ ,  $\lambda$ ,  $\varepsilon$ )
2:    $\triangleright$  Input:
3:    $\triangleright \mathbf{D} \in \mathbb{R}^{m \times n}$ : dictionary/measurement matrix
4:    $\triangleright \mathbf{x} \in \mathbb{R}^m$ : observation vector
5:    $\triangleright \lambda > 0$ : regularization parameter
6:    $\triangleright \varepsilon > 0$ : convergence tolerance
7:    $\triangleright$  Parameters:  $\beta \in (0, 1)$  (typically  $\beta = 0.5$ ),  $\eta \in (0, 1)$  (typically  $\eta = 0.9$ )
8:
9:    $\triangleright$  Output:
10:   $\triangleright \boldsymbol{\alpha}^* \in \mathbb{R}^n$ : sparse coefficient vector
11:
12:   $\triangleright$  Initialize
13:   $t \leftarrow t_0$ 
14:   $\triangleright$  initial step size, e.g.,  $t_0 = 1$ 
15:   $\boldsymbol{\alpha}^0 \leftarrow \mathbf{0}$ 
16:   $k \leftarrow 0$ 
17:
18:  while not converged do
19:     $\triangleright$  Compute gradient step
20:     $\mathbf{g}^k \leftarrow \mathbf{D}^T(\mathbf{D}\boldsymbol{\alpha}^k - \mathbf{x})$ 
21:
22:     $\triangleright$  Compute proximal operator (soft thresholding)
23:     $\boldsymbol{\alpha}^+ \leftarrow \text{prox}_{t\lambda \|\cdot\|_1}(\boldsymbol{\alpha}^k - t\mathbf{g}^k)$ 
24:
25:     $\triangleright$  Define objective function:  $F(\boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{D}\boldsymbol{\alpha} - \mathbf{x}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1$ 
26:     $\triangleright$  Backtracking line search
27:    while  $F(\boldsymbol{\alpha}^+) > F(\boldsymbol{\alpha}^k) - \frac{\eta}{t} \|\boldsymbol{\alpha}^+ - \boldsymbol{\alpha}^k\|_2^2$  do
28:       $\triangleright$  Decrease step size
29:       $t \leftarrow \beta t$ 
30:       $\triangleright$  Recompute proximal step with new step size
31:       $\boldsymbol{\alpha}^+ \leftarrow \text{prox}_{t\lambda \|\cdot\|_1}(\boldsymbol{\alpha}^k - t\mathbf{g}^k)$ 
32:    end
33:
34:     $\triangleright$  Check convergence
35:    if  $\|\boldsymbol{\alpha}^+ - \boldsymbol{\alpha}^k\|_2 < \varepsilon$  then
36:      return  $\boldsymbol{\alpha}^+$ 
37:    end
38:     $\triangleright$  Update for next iteration
39:     $\boldsymbol{\alpha}^{k+1} \leftarrow \boldsymbol{\alpha}^+$ 
40:     $k \leftarrow k + 1$ 
41:  end
42:  return  $\boldsymbol{\alpha}^k$ 
43: end

```

6.4 FISTA

While ISTA achieves $O(\frac{1}{k})$ convergence, Nesterov's acceleration technique can improve this to $O(\frac{1}{k^2})$ without additional computational cost per iteration. This acceleration is achieved through a momentum-like mechanism that exploits the history of iterates.

Algorithm 5: Fast Iterative Shrinkage-Thresholding Algorithm (FISTA)

```

1: procedure FISTA( $D, x, \lambda, \varepsilon$ )
2:    $\triangleright$  Input:
3:    $\triangleright D \in \mathbb{R}^{m \times n}$ : dictionary/measurement matrix
4:    $\triangleright x \in \mathbb{R}^m$ : observation vector
5:    $\triangleright \lambda > 0$ : regularization parameter
6:    $\triangleright \varepsilon > 0$ : convergence tolerance
7:
8:    $\triangleright$  Output:
9:    $\triangleright \alpha^* \in \mathbb{R}^n$ : sparse coefficient vector
10:
11:    $\triangleright$  Initialize
12:    $\alpha^0 \leftarrow \mathbf{0}$ 
13:    $y^1 \leftarrow \alpha^0$ 
14:    $t_1 \leftarrow 1$ 
15:    $L \leftarrow$  largest eigenvalue of  $D^T D$ 
16:    $\alpha \leftarrow \frac{1}{L}$ 
17:    $k \leftarrow 1$ 
18:
19:   while not converged do
20:      $\triangleright$  (a) Proximal gradient step:
21:      $\alpha^k \leftarrow \text{prox}_{\alpha\lambda\|\cdot\|_1}(y^k - \alpha D^T(Dy^k - x))$ 
22:
23:      $\triangleright$  (b) Update momentum parameter:
24:
25:      $t_{k+1} \leftarrow \frac{1 + \sqrt{1 + 4t_k^2}}{2}$ 
26:
27:      $\triangleright$  (c) Compute extrapolated point:
28:
29:      $y^{k+1} \leftarrow \alpha^k + \frac{t_k - 1}{t_{k+1}}(\alpha^k - \alpha^{k-1})$ 
30:
31:      $\triangleright$  Check convergence
32:     if  $\|\alpha^k - \alpha^{k-1}\|_2 < \varepsilon$  then
33:       return  $\alpha^k$ 
34:     end
35:      $k \leftarrow k + 1$ 
36:   end
37:   return  $\alpha^k$ 
38: end

```

6.4.1 The Momentum Sequence

The sequence $\{t_k\}$ satisfies the recurrence relation:

$$t_{k+1}^2 - t_{k+1} - t_k^2 = 0 \quad (6.10)$$

This yields the closed-form expression:

$$t_k = \frac{k+1}{2} + O(1) \approx \frac{k}{2} \text{ for large } k \quad (6.11)$$

6.4.2 The Extrapolation Step

The extrapolation coefficient:

$$\beta_k = \frac{t_k - 1}{t_{k+1}} \approx \frac{k-2}{k+1} \rightarrow 1 \text{ as } k \rightarrow \infty \quad (6.12)$$

This creates an “overshoot” effect that accelerates convergence by anticipating the trajectory of the iterates.

6.4.3 Convergence Rate

Theorem 6.4.3.1 (FISTA Convergence Rate): For FISTA with step size $\alpha = \frac{1}{L}$, the following bound holds:

$$F(\mathbf{x}^k) - F(\mathbf{x}^*) \leq \frac{2L \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2}{(k+1)^2} \quad (6.13)$$

The $O(\frac{1}{k^2})$ rate is optimal for first-order methods on the class of convex functions with Lipschitz continuous gradients.

7 Structured Sparsity

The classical sparse coding framework seeks representations with minimal non-zero coefficients, treating each coefficient independently. However, in many practical applications, the *locations* of non-zero coefficients exhibit inherent structure that standard sparsity models fail to exploit. Structured sparsity extends the sparse coding paradigm by incorporating prior knowledge about coefficient patterns, leading to more robust and interpretable representations.

Consider the fundamental **sparse coding problem**:

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_0 \quad (7.1)$$

where $\mathbf{y} \in \mathbb{R}^m$ represents the observed signal, $\mathbf{D} \in \mathbb{R}^{m \times n}$ denotes an overcomplete dictionary with $n > m$, and $\mathbf{x} \in \mathbb{R}^n$ contains the sparse coefficients.

This formulation promotes sparsity uniformly across all coefficients. However, structured sparsity recognizes that coefficients often exhibit dependencies or groupings that reflect the underlying data generation process.

1. **Multi-channel Signal Processing:** When processing multiple signals acquired from the same source (e.g., multi-electrode recordings, hyperspectral imaging), the active atoms in the dictionary tend to be consistent across channels.
2. **Texture Analysis:** Patches extracted from textured images share common structural elements, suggesting that their sparse representations should utilize similar dictionary atoms.
3. **Statistical Variable Selection:** In high-dimensional regression problems, predictors often form natural groups (e.g., dummy variables for categorical features, measurements from the same instrument).

7.1 Joint Sparsity and Mixed Norms

7.1.1 Problem Formulation

Let $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_L] \in \mathbb{R}^{m \times L}$ represent a collection of L signals sharing common structural properties. The joint sparse coding problem seeks a coefficient matrix $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L] \in \mathbb{R}^{n \times L}$ such that:

$$\mathbf{Y} \approx \mathbf{D}\mathbf{X} \quad (7.2)$$

The key insight is that columns of \mathbf{X} should not only be individually sparse but should also share common support patterns.

7.1.2 Mixed Norms for Matrices

To enforce joint sparsity, we introduce the (p, q) -mixed norm for matrices.

Definition 7.1.2.1 ($\ell_{p,q}$ Mixed Norm): For a matrix $\mathbf{X} \in \mathbb{R}^{n \times L}$, the (p, q) -mixed norm is defined as:

$$\|\mathbf{X}\|_{p,q} = \left(\sum_{i=1}^n \left(\sum_{j=1}^L |X_{ij}|^p \right)^{\frac{q}{p}} \right)^{\frac{1}{q}} \quad (7.3)$$

This can be interpreted as:

$$\begin{aligned} \|\mathbf{X}\|_{p,q} &= \left(\|\mathbf{x}^{(1)}\|_p, \|\mathbf{x}^{(2)}\|_p, \dots, \|\mathbf{x}^{(n)}\|_p \right)^T \Big|_q \\ &= \|\mathbf{v}\|_q \end{aligned} \quad (7.4)$$

where $\mathbf{x}^{(i)}$ denotes the i -th row of \mathbf{X} , and $v_i = \|\mathbf{x}^{(i)}\|_p$.

Frobenius Norm: When $p = q = 2$:

$$\|\mathbf{X}\|_{2,2} = \sqrt{\sum_{i,j} |X_{ij}|^2} = \|\mathbf{X}\|_F \quad (7.5)$$

Entry-wise Norms: When $p = q$, the mixed norm reduces to the vectorized ℓ_p norm:

$$\|\mathbf{X}\|_{p,p} = \|\text{vec}(\mathbf{X})\|_p \quad (7.6)$$

Joint Sparsity Norm: The $\ell_{2,1}$ norm:

$$\|\mathbf{X}\|_{2,1} = \sum_{i=1}^n \sqrt{\sum_{j=1}^L |X_{ij}|^2} = \sum_{i=1}^n \|\mathbf{x}^{(i)}\|_2 \quad (7.7)$$

Proposition 7.1.2.1 (Joint Sparsity Property): The $\ell_{2,1}$ norm promotes row-sparsity in \mathbf{X} , meaning entire rows become zero. This corresponds to selecting the same dictionary atoms across all signals.

7.2 Joint Sparse Coding

The joint sparse coding problem with the $\ell_{2,1}$ norm regularization is formulated as:

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times L}} \frac{1}{2} \|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_{2,1} \quad (7.8)$$

This optimization problem combines a smooth, convex data fidelity term with a non-smooth but convex regularizer, making it amenable to proximal gradient methods.

The iterative solution via proximal gradient descent follows:

$$\mathbf{X}^{(k+1)} = \text{prox}_{\gamma\lambda \|\cdot\|_{2,1}}(\mathbf{X}^{(k)} - \gamma \nabla f(\mathbf{X}^{(k)})) \quad (7.9)$$

where $f(\mathbf{X}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_F^2$ is the smooth part, and $\gamma > 0$ is the step size.

The key computational challenge lies in evaluating the proximal mapping:

$$\text{prox}_{\tau \|\cdot\|_{2,1}}(\mathbf{Z}) = \arg \min_{\mathbf{X}} \left\{ \tau \|\mathbf{X}\|_{2,1} + \frac{1}{2} \|\mathbf{X} - \mathbf{Z}\|_F^2 \right\} \quad (7.10)$$

Theorem 7.2.1 (Row-wise Separability of $\ell_{2,1}$ Proximal Mapping): The proximal mapping of the $\ell_{2,1}$ norm can be computed row-wise:

$$\left[\text{prox}_{\tau \|\cdot\|_{2,1}}(\mathbf{Z}) \right]_i = \text{shrink}_{\tau}^{(2)}(\mathbf{z}^{(i)}) \quad (7.11)$$

where $\mathbf{z}^{(i)}$ is the i -th row of \mathbf{Z} , and $\text{shrink}_{\tau}^{(2)}$ is the multivariate soft-thresholding operator.

Proof: The objective function in Equation (7.10) can be rewritten as:

$$\begin{aligned} \tau \|\mathbf{X}\|_{2,1} + \frac{1}{2} \|\mathbf{X} - \mathbf{Z}\|_F^2 &= \tau \sum_{i=1}^n \|\mathbf{x}^{(i)}\|_2 + \frac{1}{2} \sum_{i=1}^n \|\mathbf{x}^{(i)} - \mathbf{z}^{(i)}\|_2^2 \\ &= \sum_{i=1}^n \left[\tau \|\mathbf{x}^{(i)}\|_2 + \frac{1}{2} \|\mathbf{x}^{(i)} - \mathbf{z}^{(i)}\|_2^2 \right] \end{aligned} \quad (7.12)$$

Since the objective decomposes into independent row-wise problems, the minimization can be performed separately for each row.

Definition 7.2.1 (Multivariate Soft-Thresholding): For a vector $\mathbf{v} \in \mathbb{R}^L$ and threshold $\tau > 0$:

$$\text{shrink}_{\tau}^{(2)}(\mathbf{v}) = \begin{cases} \frac{\tau}{\|\mathbf{v}\|_2} \cdot \max(0, \|\mathbf{v}\|_2 - \tau) & \text{if } \mathbf{v} \neq \mathbf{0} \\ \mathbf{0} & \text{if } \mathbf{v} = \mathbf{0} \end{cases} \quad (7.13)$$

This operator exhibits two key behaviors:

1. **Nullification:** If $\|\mathbf{v}\|_2 \leq \tau$, the entire vector is set to zero.
2. **Shrinkage:** If $\|\mathbf{v}\|_2 > \tau$, the vector is scaled down while preserving its direction.

7.3 Group Sparsity and Extensions

Consider a dictionary \mathbf{D} partitioned into G groups:

$$\mathbf{D} = [\mathbf{D}_1 \mid \mathbf{D}_2 \mid \dots \mid \mathbf{D}_G] \quad (7.14)$$

where $\mathbf{D}_g \in \mathbb{R}^{m \times n_g}$ contains atoms corresponding to group g , and $\sum_{g=1}^G n_g = n$.

The coefficient vector \mathbf{x} is correspondingly partitioned:

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_{[1]} \\ \mathbf{x}_{[2]} \\ \vdots \\ \mathbf{x}_{[G]} \end{pmatrix} \quad (7.15)$$

where $\mathbf{x}_{[g]} \in \mathbb{R}^{n_g}$ contains coefficients for group g .

The group sparse coding problem seeks representations that activate entire groups rather than individual atoms:

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2^2 + \lambda \sum_{g=1}^G w_g \|\mathbf{x}_{[g]}\|_2 \quad (7.16)$$

where $w_g > 0$ are **group-specific weights**, typically set as $w_g = \sqrt{n_g}$ to account for group size differences.

▲ Group Selection Property

attention 7.3.1

The group LASSO penalty induces sparsity at the group level: either all coefficients within a group are zero, or the group is active with potentially multiple non-zero coefficients.

The proximal mapping for the group LASSO penalty decomposes into group-wise operations:

$$\left[\text{prox}_{\tau \sum_g w_g \|\cdot\|_2}(\mathbf{z}) \right]_{[g]} = \text{shrink}_{\tau w_g}^{(2)}(\mathbf{z}_{[g]}) \quad (7.17)$$

This allows efficient computation by applying multivariate soft-thresholding to each group independently.

7.4 LASSO

In the statistical setting, we observe m samples with response variable y_i and n predictors $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{in}]^T$. The linear model assumes:

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, 2, \dots, m \quad (7.18)$$

In matrix notation:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (7.19)$$

where $\mathbf{y} \in \mathbb{R}^m$, $\mathbf{X} \in \mathbb{R}^{m \times n}$ is the design matrix, $\boldsymbol{\beta} \in \mathbb{R}^n$ contains regression coefficients, and $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$.

The classical least squares estimator:

$$\hat{\boldsymbol{\beta}}_{\text{LS}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (7.20)$$

This estimator is unbiased ($\mathbb{E}[\hat{\boldsymbol{\beta}}_{\text{LS}}] = \boldsymbol{\beta}$) but may have high variance, especially when predictors are correlated or n is large relative to m .

The Least Absolute Shrinkage and Selection Operator (LASSO) adds ℓ_1 regularization:

$$\hat{\boldsymbol{\beta}}_{\text{LASSO}} = \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \right\} \quad (7.21)$$

Theorem 7.4.1 (Variable Selection Property): For sufficiently large λ , the LASSO estimator $\hat{\boldsymbol{\beta}}_{\text{LASSO}}$ contains exact zeros, performing automatic variable selection.

The LASSO introduces bias to reduce variance:

$$\begin{aligned} \text{MSE}(\hat{\boldsymbol{\beta}}) &= \mathbb{E}[\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2^2] \\ &= \|\mathbb{E}[\hat{\boldsymbol{\beta}}] - \boldsymbol{\beta}\|_2^2 + \text{trace}(\text{Var}(\hat{\boldsymbol{\beta}})) \\ &= \text{Bias}^2 + \text{Variance} \end{aligned} \quad (7.22)$$

While OLS minimizes bias, LASSO accepts some bias in exchange for substantially reduced variance through sparsity.

Theorem 7.4.2 (LASSO in High Dimensions): When $n > m$ (more predictors than observations), OLS is not unique. However, LASSO provides a unique sparse solution for appropriate $\lambda > 0$.

This property makes LASSO particularly valuable in modern applications like genomics, where the number of features vastly exceeds the sample size.

7.5 Elastic Net

The LASSO has two notable limitations:

1. In the $n > m$ setting, it selects at most m variables
2. When predictors are highly correlated, LASSO tends to arbitrarily select one from each group

The Elastic Net addresses these issues by combining ℓ_1 and ℓ_2 penalties:

$$\hat{\beta}_{\text{EN}} = \arg \min_{\beta} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2 \right\} \quad (7.23)$$

7.5.1 Geometric Interpretation

The constraint region for Elastic Net is:

$$\mathcal{C}_{\text{EN}} = \{\beta : \alpha \|\beta\|_1 + (1 - \alpha) \|\beta\|_2^2 \leq t\} \quad (7.24)$$

This creates a compromise between the diamond-shaped ℓ_1 ball and the spherical ℓ_2 ball, maintaining sparsity-inducing corners while allowing smoother boundaries.

7.5.2 Proximal Gradient Solution

The Elastic Net optimization can be solved efficiently using proximal gradient methods. The key insight is that the ℓ_2 penalty can be absorbed into the smooth part:

$$f(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_2 \|\beta\|_2^2 \quad (7.25)$$

with gradient:

$$\nabla f(\beta) = \mathbf{X}^T(\mathbf{X}\beta - \mathbf{y}) + 2\lambda_2\beta = (\mathbf{X}^T\mathbf{X} + 2\lambda_2\mathbf{I})\beta - \mathbf{X}^T\mathbf{y} \quad (7.26)$$

The proximal gradient update becomes:

$$\beta^{(k+1)} = \text{shrink}_{\gamma\lambda_1}(\beta^{(k)} - \gamma\nabla f(\beta^{(k)})) \quad (7.27)$$

where shrink_{τ} is the element-wise soft-thresholding operator.

7.6 Algorithmic Summary

Problem	Regularizer	Proximal Operator
Standard Sparsity	$\ \mathbf{x}\ _1$	$\text{shrink}_{\tau}(x_i) = \text{sign}(x_i) \max(0, x_i - \tau)$
Joint Sparsity	$\ \mathbf{X}\ _{2,1}$	Row-wise multivariate soft-thresholding
Group Sparsity	$\sum_g w_g \ \mathbf{x}_{[g]}\ _2$	Group-wise multivariate soft-thresholding
Elastic Net	$\lambda_1 \ \mathbf{x}\ _1 + \lambda_2 \ \mathbf{x}\ _2^2$	Modified soft-thresholding with ℓ_2 in gradient

Table 1: Summary of sparsity-inducing regularizers and their proximal operators

8 Dictionary Learning

8.1 Introduction to Dictionary Learning

Definition 8.1.1 (Dictionary Learning Problem): Given a set of training signals $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N \in \mathbb{R}^n$, the dictionary learning problem seeks to find:

1. A dictionary matrix $\mathbf{D} \in \mathbb{R}^{n \times m}$ with $m > n$ (redundant dictionary)
2. Sparse coefficient vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \in \mathbb{R}^m$

such that $\mathbf{y}_i \approx \mathbf{D}\mathbf{x}_i$ for all $i = 1, 2, \dots, N$, where each \mathbf{x}_i has at most T_0 non-zero entries.

Dictionary learning employs a block coordinate descent strategy, alternating between two phases:

1. **Sparse Coding Phase:** Fix \mathbf{D} and solve for \mathbf{X}
2. **Dictionary Update Phase:** Fix \mathbf{X} and update \mathbf{D}

8.2 Problem Formulation

Let $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N] \in \mathbb{R}^{n \times N}$ denote the training matrix, where each column represents a training signal. Similarly, let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{m \times N}$ represent the sparse coefficient matrix.

The dictionary learning problem can be formulated as the following optimization:

$$\min_{\mathbf{D}, \mathbf{X}} \|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_F^2 \quad \text{subject to} \quad \|\mathbf{x}_i\|_0 \leq T_0, \quad \forall i = 1, 2, \dots, N \quad (8.1)$$

where $\|\cdot\|_F$ denotes the Frobenius norm and $\|\cdot\|_0$ is the ℓ_0 pseudo-norm counting non-zero entries.

(Normalization Constraint:) To resolve scaling ambiguities, we impose the constraint that each column of \mathbf{D} has unit ℓ_2 norm:

$$\|\mathbf{d}_j\|_2 = 1, \quad \forall j = 1, 2, \dots, m \quad (8.2)$$

The optimization problem Equation (8.1) presents several fundamental challenges:

1. **Non-convexity:** The objective function is non-convex in the joint variables \mathbf{D}, \mathbf{X} , even though it is convex in each variable individually when the other is fixed.
2. **Sparsity Constraint:** The ℓ_0 pseudo-norm constraint is non-convex and combinatorial, making direct optimization intractable.
3. **Solution ambiguity:** Multiple equivalent solutions exist due to:
 - Column permutations of \mathbf{D} with corresponding row permutations of \mathbf{X}
 - Sign ambiguities: $\mathbf{d}_j, \mathbf{x}_j \equiv (-\mathbf{d}_j, -\mathbf{x}_j)$

8.3 Sparse Coding Phase

The sparse coding phase solves the following problem for each training signal:

$$\mathbf{x}_i^{(k+1)} = \arg \min_{\mathbf{x}} \|\mathbf{y}_i - \mathbf{D}^{(k)} \mathbf{x}\|_2^2 \quad \text{subject to} \quad \|\mathbf{x}\|_0 \leq T_0 \quad (8.3)$$

This is precisely the sparse coding problem discussed in previous sections, which can be solved using greedy algorithms such as:

- **Orthogonal Matching Pursuit (OMP)**: Iteratively selects atoms that best correlate with the current residual
- **Matching Pursuit (MP)**: Similar to OMP but without orthogonalization
- **Basis Pursuit**: Convex relaxation using ℓ_1 norm

The sparse coding phase is computationally the most expensive part of dictionary learning, as it requires solving N sparse coding problems (one for each training signal) at each iteration.

8.4 Dictionary Update Phase

The dictionary update phase constitutes the core innovation of K-SVD. Rather than updating the entire dictionary simultaneously, K-SVD updates one column at a time while simultaneously updating the corresponding sparse coefficients.

8.4.1 Matrix Factorization Perspective

Consider the error matrix:

$$\mathbf{E} = \mathbf{Y} - \mathbf{D}\mathbf{X} \quad (8.4)$$

Using the fundamental matrix identity, we can decompose the product $\mathbf{D}\mathbf{X}$ as:

$$\mathbf{D}\mathbf{X} = \sum_{j=1}^m \mathbf{d}_j \mathbf{x}_j^T \quad (8.5)$$

where \mathbf{x}_j^T denotes the j -th row of \mathbf{X} .

8.4.2 Isolated Column Update

To update the j_0 -th column of \mathbf{D} , we rewrite Equation (8.5) as:

$$\mathbf{D}\mathbf{X} = \sum_{j \neq j_0} \mathbf{d}_j \mathbf{x}_j^T + \mathbf{d}_{j_0} \mathbf{x}_{j_0}^T \quad (8.6)$$

Define the error matrix excluding the j_0 -th atom:

$$\mathbf{E}_{j_0} = \mathbf{Y} - \sum_{j \neq j_0} \mathbf{d}_j \mathbf{x}_j^T \quad (8.7)$$

The update problem becomes:

$$\min_{\mathbf{d}_{j_0}, \mathbf{x}_{j_0}^T} \|\mathbf{E}_{j_0} - \mathbf{d}_{j_0} \mathbf{x}_{j_0}^T\|_F^2 \quad \text{subject to} \quad \|\mathbf{d}_{j_0}\|_2 = 1 \quad (8.8)$$

This is a rank-one matrix approximation problem, optimally solved using the Singular Value Decomposition (SVD).

8.4.3 SVD solution

Theorem 8.4.3.1 (Rank-One Matrix Approximation): Let $\mathbf{A} \in \mathbb{R}^{n \times N}$ be given. The solution to

$$\min_{\mathbf{u}, \mathbf{v}} \|\mathbf{A} - \mathbf{u}\mathbf{v}^T\|_F^2 \quad \text{subject to} \quad \|\mathbf{u}\|_2 = 1 \quad (8.9)$$

is given by $\mathbf{u} = \mathbf{u}_1$ and $\mathbf{v}^T = \sigma_1 \mathbf{v}_1^T$, where $\mathbf{A} = \sum_{i=1}^{\min(n, N)} \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ is the SVD of \mathbf{A} .

Proof: The Frobenius norm can be expressed as:

$$\begin{aligned}\|\mathbf{A} - \mathbf{u}\mathbf{v}^T\|_F^2 &= \|\mathbf{A}\|_F^2 - 2 \operatorname{trace}(\mathbf{A}^T \mathbf{u}\mathbf{v}^T) + \|\mathbf{u}\mathbf{v}^T\|_F^2 \\ &= \|\mathbf{A}\|_F^2 - 2\mathbf{v}^T \mathbf{A}^T \mathbf{u} + \|\mathbf{v}\|_2^2\end{aligned}\quad (8.10)$$

Since $\|\mathbf{u}\|_2 = 1$, maximizing $\mathbf{v}^T \mathbf{A}^T \mathbf{u}$ is equivalent to finding the leading singular vectors of \mathbf{A} .

8.4.4 Sparsity Preservation

A critical challenge in the dictionary update is preserving the sparsity structure of \mathbf{X} . The naive application of [Theorem 8.4.3.1](#) would yield a dense row vector $\mathbf{x}_{j_0}^T$, violating the sparse coding constraint.

Solution - Restricted SVD: K-SVD addresses this by restricting the update to only those training signals that actually use the j_0 -th atom:

$$\Omega_{j_0} = \{i : x_{j_0,i} \neq 0\} \quad (8.11)$$

Define the restricted error matrix:

$$\mathbf{E}_{j_0}^R = \mathbf{E}_{j_0}(:, \Omega_{j_0}) \quad (8.12)$$

The restricted update problem becomes:

$$\min_{\mathbf{d}_{j_0}, \mathbf{x}_{j_0}^R} \left\| \mathbf{E}_{j_0}^R - \mathbf{d}_{j_0} (\mathbf{x}_{j_0}^R)^T \right\|_F^2 \quad \text{subject to} \quad \|\mathbf{d}_{j_0}\|_2 = 1 \quad (8.13)$$

where $\mathbf{x}_{j_0}^R$ contains only the non-zero elements of $\mathbf{x}_{j_0}^T$.

Dictionary Update Algorithm for Column j_0 :

1. Compute error matrix: $\mathbf{E}_{j_0} = \mathbf{Y} - \sum_{j \neq j_0} \mathbf{d}_j \mathbf{x}_j^T$
2. Identify support: $\Omega_{j_0} = \{i : x_{j_0,i} \neq 0\}$
3. Extract restricted matrix: $\mathbf{E}_{j_0}^R = \mathbf{E}_{j_0}(:, \Omega_{j_0})$
4. Compute SVD: $\mathbf{E}_{j_0}^R = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$
5. Update dictionary: $\mathbf{d}_{j_0} = \mathbf{u}_1$
6. Update coefficients: $\mathbf{x}_{j_0}^R = \sigma_1 \mathbf{v}_1$
7. Restore to full representation: $\mathbf{x}_{j_0}^T(\Omega_{j_0}) = \mathbf{x}_{j_0}^R$

Algorithm 6: K-SVD Dictionary Learning Algorithm

```

1: procedure K-SVD( $\mathbf{Y}, m, T_0, K$ )
2:    $\triangleright$  Input: Training signals  $\mathbf{Y} \in \mathbb{R}^{n \times N}$ , dictionary size  $m$ , sparsity level  $T_0$ , max iterations  $K$ 
3:    $\triangleright$  Output: Dictionary  $\mathbf{D} \in \mathbb{R}^{n \times m}$  and sparse coefficients  $\mathbf{X} \in \mathbb{R}^{m \times N}$ 
4:
5:    $\triangleright$  Initialize dictionary:
6:    $\mathbf{D}^{(0)} \leftarrow$  random matrix with unit-norm columns
7:    $k \leftarrow 0$ 
8:
9:    $\triangleright$  Main K-SVD loop:
10:  while  $k < K \wedge$  not converged do
11:     $\triangleright$  SPARSE CODING PHASE:
12:     $\triangleright$  For each training signal  $\mathbf{y}_i$ , solve sparse coding problem
13:    for  $i = 1, \dots, N$  do
14:       $\mathbf{x}_i^{(k+1)} \leftarrow \arg \min_{\mathbf{x}} \|\mathbf{y}_i - \mathbf{D}^{(k)} \mathbf{x}\|_2^2$  subject to  $\|\mathbf{x}\|_0 \leq T_0$ 
15:       $\triangleright$  Use OMP, MP, or other sparse coding algorithm
16:    end
17:
18:     $\triangleright$  DICTIONARY UPDATE PHASE:
19:     $\triangleright$  Update each dictionary column sequentially
20:    for  $j = 1, \dots, m$  do
21:       $\triangleright$  (a) Compute error matrix excluding  $j$ -th atom:
22:       $\mathbf{E}_j \leftarrow \mathbf{Y} - \sum_{\ell \neq j} \mathbf{d}_\ell \mathbf{x}_\ell^T$ 
23:
24:       $\triangleright$  (b) Identify support set:
25:       $\Omega_j \leftarrow \{i : x_{j,i} \neq 0\}$ 
26:
27:       $\triangleright$  (c) Check if atom is used:
28:      if  $\Omega_j = \emptyset$  then
29:         $\triangleright$  Replace unused atom with random unit vector
30:         $\mathbf{d}_j \leftarrow$  random unit vector
31:      end
32:      else
33:         $\triangleright$  (d) Extract restricted error matrix:
34:         $\mathbf{E}_j^R \leftarrow \mathbf{E}_j(:, \Omega_j)$ 
35:
36:         $\triangleright$  (e) Compute SVD of restricted matrix:
37:         $\mathbf{U}, \Sigma, \mathbf{V} \leftarrow \text{SVD}(\mathbf{E}_j^R)$ 
38:
39:         $\triangleright$  (f) Update dictionary column:
40:         $\mathbf{d}_j \leftarrow \mathbf{u}_1$ 
41:         $\triangleright$  First left singular vector
42:
43:         $\triangleright$  (g) Update sparse coefficients:
44:         $\mathbf{x}_j^R \leftarrow \sigma_1 \mathbf{v}_1$ 
45:         $\triangleright$  Scaled first right singular vector
46:
47:         $\triangleright$  (h) Restore to full representation:
48:         $\mathbf{x}_j^T(\Omega_j) \leftarrow \mathbf{x}_j^R$ 
49:         $\mathbf{x}_j^T(\text{complement}(\Omega_j)) \leftarrow \mathbf{0}$ 
50:      end
51:    end
52:  end

```

```
53:   ▷ Check convergence:
54:   if  $\|Y - D^{(k+1)}X^{(k+1)}\|_F < \varepsilon$  then
55:     ▷ Convergence achieved
56:     BREAK
57:   end
58:
59:    $k \leftarrow k + 1$ 
60: end
61:
62: return  $D^{(k)}, X^{(k)}$ 
63: end
```

9 Local Polynomial Approximation

Local Polynomial Approximation (LPA) represents a fundamental shift in signal processing methodology, moving away from global sparsity constraints toward localized polynomial modeling. This approach recognizes that many real-world signals exhibit piecewise smooth behavior that can be effectively captured through local polynomial representations.

Traditional signal processing approaches often rely on global assumptions about signal structure, such as sparsity in transformed domains (e.g., DCT, wavelet transforms). While these methods have proven successful in many applications, they impose uniform constraints across the entire signal domain. In contrast, LPA methods adapt to local signal characteristics, providing more flexible and often more accurate approximations.

The transition from sparsity-based to polynomial-based modeling represents a significant conceptual advancement:

- **Sparsity-based approach:** Assumes signal can be represented as a linear combination of few basis functions from a fixed dictionary
- **Polynomial-based approach:** Assumes signal can be locally approximated by polynomials of appropriate degree

This shift enables adaptive processing where each spatial location can have its own optimal approximation parameters, leading to superior performance in regions with varying signal characteristics.

Definition 9.1 (Local Polynomial Approximation): Given a signal $f : \mathbb{R} \rightarrow \mathbb{R}$ and a point $x_0 \in \mathbb{R}$, a local polynomial approximation of degree L is a polynomial $P_L(x)$ of the form:

$$P_L(x) = \sum_{j=0}^L \beta_j (x - x_0)^j \quad (9.1)$$

that minimizes a local fitting criterion within a neighborhood of x_0 .

9.1 Local Polynomial Model

Consider a one-dimensional signal model:

$$y(t) = f(t) + \eta(t) \quad (9.2)$$

where $f(t)$ represents the underlying smooth signal and $\eta(t)$ denotes additive white Gaussian noise with variance σ^2 .

For discrete processing, we work with sampled versions. Let $\mathbf{y} = [y_1, y_2, \dots, y_M]^T$ represent an M -dimensional signal vector extracted from a local neighborhood around a central pixel.

Within a local neighborhood, we assume the signal can be well approximated by a polynomial of degree L :

$$f(t_i) \approx \sum_{j=0}^L \beta_j t_i^j, \quad i = 1, 2, \dots, M \quad (9.3)$$

where $\{\beta_j\}_{j=0}^L$ are the polynomial coefficients to be estimated, and $\{t_i\}_{i=1}^M$ are the spatial locations within the neighborhood.

▲ Degree Constraint**attention 9.1.1**

The polynomial degree L must satisfy $L + 1 \leq M$ to ensure an over-determined system. This constraint prevents overfitting and ensures stable coefficient estimation.

9.1.1 Matrix Formulation

The polynomial approximation can be expressed in matrix form. Define the design matrix $\mathbf{T} \in \mathbb{R}^{M \times (L+1)}$:

$$\mathbf{T} = \begin{pmatrix} 1 & t_1 & t_1^2 & \cdots & t_1^L \\ 1 & t_2 & t_2^2 & \cdots & t_2^L \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & t_M & t_M^2 & \cdots & t_M^L \end{pmatrix} \quad (9.4)$$

The polynomial coefficient vector is $\boldsymbol{\beta} = [\beta_0, \beta_1, \dots, \beta_L]^T \in \mathbb{R}^{L+1}$.

The polynomial approximation becomes:

$$\mathbf{f} \approx \mathbf{T}\boldsymbol{\beta} \quad (9.5)$$

9.1.2 Least Squares Formulation

The optimal polynomial coefficients are obtained by minimizing the squared approximation error:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{L+1}} \|\mathbf{y} - \mathbf{T}\boldsymbol{\beta}\|^2 \quad (9.6)$$

Theorem 9.1.2.1 (Unweighted LPA Solution): The least squares solution to the local polynomial approximation problem is:

$$\hat{\boldsymbol{\beta}} = (\mathbf{T}^T \mathbf{T})^{-1} \mathbf{T}^T \mathbf{y} \quad (9.7)$$

provided that $\mathbf{T}^T \mathbf{T}$ is invertible.

Proof: Taking the derivative of the objective function Equation (9.6) with respect to $\boldsymbol{\beta}$ and setting it to zero:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\beta}} \|\mathbf{y} - \mathbf{T}\boldsymbol{\beta}\|^2 &= \frac{\partial}{\partial \boldsymbol{\beta}} (\mathbf{y} - \mathbf{T}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{T}\boldsymbol{\beta}) \\ &= -2\mathbf{T}^T (\mathbf{y} - \mathbf{T}\boldsymbol{\beta}) = 0 \end{aligned} \quad (9.8)$$

Solving for $\boldsymbol{\beta}$ yields the normal equations:

$$\mathbf{T}^T \mathbf{T} \boldsymbol{\beta} = \mathbf{T}^T \mathbf{y} \quad (9.9)$$

The solution follows directly when $\mathbf{T}^T \mathbf{T}$ is invertible.

9.2 Weighted Local Polynomial Approximation

In many practical scenarios, not all samples within a neighborhood should contribute equally to the polynomial fit. Samples closer to the center of the neighborhood are typically more relevant for estimating the signal at that location. Weighted LPA addresses this by introducing spatially varying weights.

The weighted LPA problem incorporates a diagonal weight matrix $\mathbf{W} \in \mathbb{R}^{M \times M}$:

$$\hat{\beta}_w = \arg \min_{\beta \in \mathbb{R}^{L+1}} \|\mathbf{W}(\mathbf{y} - \mathbf{T}\beta)\|^2 \quad (9.10)$$

where $\mathbf{W} = \text{diag}(w_1, w_2, \dots, w_M)$ with $w_i \geq 0$ for all i .

Theorem 9.2.1 (Weighted LPA Solution): The weighted least squares solution is:

$$\hat{\beta}_w = (\mathbf{T}^T \mathbf{W}^2 \mathbf{T})^{-1} \mathbf{T}^T \mathbf{W}^2 \mathbf{y} \quad (9.11)$$

Proof: The weighted objective function can be written as:

$$\|\mathbf{W}(\mathbf{y} - \mathbf{T}\beta)\|^2 = \|\mathbf{W}\mathbf{y} - \mathbf{W}\mathbf{T}\beta\|^2 \quad (9.12)$$

This is equivalent to solving the unweighted problem:

$$\arg \min_{\beta} \|\tilde{\mathbf{y}} - \tilde{\mathbf{T}}\beta\|^2 \quad (9.13)$$

where $\tilde{\mathbf{y}} = \mathbf{W}\mathbf{y}$ and $\tilde{\mathbf{T}} = \mathbf{W}\mathbf{T}$.

Applying [Theorem 9.1.2.1](#):

$$\begin{aligned} \hat{\beta}_w &= (\tilde{\mathbf{T}}^T \tilde{\mathbf{T}})^{-1} \tilde{\mathbf{T}}^T \tilde{\mathbf{y}} \\ &= ((\mathbf{W}\mathbf{T})^T (\mathbf{W}\mathbf{T}))^{-1} (\mathbf{W}\mathbf{T})^T \mathbf{W}\mathbf{y} \\ &= (\mathbf{T}^T \mathbf{W}^2 \mathbf{T})^{-1} \mathbf{T}^T \mathbf{W}^2 \mathbf{y} \end{aligned} \quad (9.14)$$

Common weight functions include:

1. **Uniform weights:** $w_i = 1$ for all i (reduces to unweighted case)
2. **Gaussian weights:** $w_i = \exp(-(t_i - t_0)^2 / (2\sigma_w^2))$
3. **Binary weights:** $w_i \in \{0, 1\}$ for adaptive support selection

Example 9.2.1 (Binary Weight Example): Consider a signal with a discontinuity. Using binary weights allows selective processing:

- Left-side filter: $\mathbf{w} = [1, 1, 1, 0, 0]^T$
- Right-side filter: $\mathbf{w} = [0, 0, 1, 1, 1]^T$

This prevents blurring across discontinuities while maintaining smoothing within homogeneous regions.

9.3 QR Decomposition Approach

Direct computation of $(\mathbf{T}^T \mathbf{T})^{-1}$ can be numerically unstable, especially when \mathbf{T} is ill-conditioned. The QR decomposition provides a numerically stable alternative while revealing the underlying geometric structure of the problem.

Theorem 9.3.1 (QR Decomposition): Any matrix $\mathbf{T} \in \mathbb{R}^{M \times (L+1)}$ with $M \geq L+1$ and full column rank can be decomposed as:

$$\mathbf{T} = \mathbf{Q}\mathbf{R} \quad (9.15)$$

where $\mathbf{Q} \in \mathbb{R}^{M \times (L+1)}$ has orthonormal columns ($\mathbf{Q}^T \mathbf{Q} = \mathbf{I}_{L+1}$) and $\mathbf{R} \in \mathbb{R}^{(L+1) \times (L+1)}$ is upper triangular.

Using the QR decomposition, the LPA solution becomes:

$$\begin{aligned} \hat{\beta} &= (\mathbf{T}^T \mathbf{T})^{-1} \mathbf{T}^T \mathbf{y} \\ &= ((\mathbf{Q}\mathbf{R})^T (\mathbf{Q}\mathbf{R}))^{-1} (\mathbf{Q}\mathbf{R})^T \mathbf{y} \\ &= (\mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R})^{-1} \mathbf{R}^T \mathbf{Q}^T \mathbf{y} \\ &= (\mathbf{R}^T \mathbf{R})^{-1} \mathbf{R}^T \mathbf{Q}^T \mathbf{y} \\ &= \mathbf{R}^{-1} \mathbf{Q}^T \mathbf{y} \end{aligned} \quad (9.16)$$

The signal estimate is:

$$\hat{f} = \mathbf{T} \hat{\beta} = \mathbf{Q} \mathbf{R} \mathbf{R}^{-1} \mathbf{Q}^T \mathbf{y} = \mathbf{Q} \mathbf{Q}^T \mathbf{y} \quad (9.17)$$

Proposition 9.3.1 (Projection Interpretation): The matrix $\mathbf{P} = \mathbf{Q} \mathbf{Q}^T$ represents the orthogonal projection onto the column space of \mathbf{T} . The LPA estimate is the orthogonal projection of the noisy signal onto the space of polynomials of degree L .

For weighted LPA, we apply QR decomposition to the weighted design matrix $\mathbf{W}\mathbf{T}$:

$$\mathbf{W}\mathbf{T} = \tilde{\mathbf{Q}}\tilde{\mathbf{R}} \quad (9.18)$$

The weighted solution becomes:

$$\begin{aligned} \hat{\beta}_w &= \tilde{\mathbf{R}}^{-1} \tilde{\mathbf{Q}}^T \mathbf{W}\mathbf{y} \\ \hat{f}_w &= \tilde{\mathbf{Q}} \tilde{\mathbf{Q}}^T \mathbf{W}\mathbf{y} \end{aligned} \quad (9.19)$$

9.4 Convolution Implementation

A key insight in LPA is that the estimation can be implemented as a convolution operation, enabling efficient computation across entire signals or images.

Consider the estimation of the signal value at the center of the neighborhood. Let i_c denote the central index. The estimate at this location is:

$$\hat{f}(t_{i_c}) = \mathbf{e}_{i_c}^T \mathbf{Q} \mathbf{Q}^T \mathbf{y} \quad (9.20)$$

where \mathbf{e}_{i_c} is the unit vector with 1 in the i_c -th position.

Theorem 9.4.1 (Convolution Kernel Formula): The LPA estimation at the central location can be expressed as:

$$\hat{f}(t_{i_c}) = \sum_{i=1}^M h_i y_i = \mathbf{h}^T \mathbf{y} \quad (9.21)$$

where the convolution kernel is:

$$\mathbf{h} = \mathbf{Q} \mathbf{Q}^T \mathbf{e}_{i_c} \quad (9.22)$$

The convolution kernel can be computed explicitly as:

$$\mathbf{h} = \sum_{j=0}^L \beta_j \mathbf{q}_j \quad (9.23)$$

where \mathbf{q}_j are the columns of \mathbf{Q} and:

$$\beta_j = \mathbf{q}_j^T \mathbf{e}_{i_c} = q_{j,i_c} \quad (9.24)$$

9.5 Special Cases

Example 9.5.1 (Zero-Order Polynomial (Moving Average)): For $L = 0$, the design matrix is $\mathbf{T} = \mathbf{1} = [1, 1, \dots, 1]^T$.

The QR decomposition gives:

$$\begin{aligned} \mathbf{Q} &= \frac{1}{\sqrt{M}} \mathbf{1} \\ \mathbf{R} &= \sqrt{M} \end{aligned} \quad (9.25)$$

The convolution kernel becomes:

$$\mathbf{h} = \frac{1}{M} \mathbf{1} \quad (9.26)$$

This is the standard moving average filter.

Example 9.5.2 (Weighted Zero-Order Polynomial): For weighted zero-order polynomial with normalized weights $\sum_{i=1}^M w_i^2 = 1$:

$$\begin{aligned} \tilde{\mathbf{Q}} &= \mathbf{w} \\ \tilde{\mathbf{R}} &= 1 \end{aligned} \quad (9.27)$$

The convolution kernel is:

$$\mathbf{h} = \mathbf{w} \quad (9.28)$$

This shows that Gaussian smoothing corresponds to weighted zero-order polynomial fitting.

9.6 Statistical Properties and Performance Analysis

9.6.1 Bias-Variance Decomposition

The mean squared error (MSE) of the LPA estimator can be decomposed into bias and variance components:

$$\text{MSE}(\hat{f}(t_0)) = \text{Bias}^2(\hat{f}(t_0)) + \text{Var}(\hat{f}(t_0)) \quad (9.29)$$

9.6.2 Bias Analysis

Theorem 9.6.2.1 (Bias of LPA Estimator): For a signal $f(t)$ that is $(L + 1)$ -times differentiable, the bias of the LPA estimator is:

$$\text{Bias}(\hat{f}(t_0)) = \frac{f^{(L+1)}(t_0)}{(L+1)!} \sum_{i=1}^M h_i (t_i - t_0)^{L+1} + O(h^{L+2}) \quad (9.30)$$

where h is the neighborhood size.

Theorem 9.6.2.2 (Bias for Polynomial Signals): If the true signal is a polynomial of degree L or less, the LPA estimator is unbiased.

9.6.3 Variance Analysis

Theorem 9.6.3.1 (Variance of LPA Estimator): For additive white Gaussian noise with variance σ^2 , the variance of the LPA estimator is:

$$\text{Var}(\hat{f}(t_0)) = \sigma^2 \|\mathbf{h}\|^2 = \sigma^2 \mathbf{e}_{i_c}^T \mathbf{Q} \mathbf{Q}^T \mathbf{e}_{i_c} \quad (9.31)$$

The variance decreases as the effective number of samples increases, but the bias may increase due to the larger neighborhood size. This creates a fundamental bias-variance tradeoff.

9.6.4 Optimal Neighborhood Selection

The optimal neighborhood size balances bias and variance:

$$h_{\text{opt}} = \arg \min_h [\text{Bias}^2(h) + \text{Var}(h)] \quad (9.32)$$

In practice, this leads to adaptive algorithms that select different neighborhood sizes and shapes based on local signal characteristics.

9.7 Adaptive Neighborhood Selection

Fixed neighborhood sizes and shapes are suboptimal for signals with varying local characteristics. Adaptive methods adjust the approximation parameters based on local signal properties.

9.7.1 Directional Filtering

Binary weights enable directional filtering, which is particularly useful near discontinuities:

Definition 9.7.1.1 (Directional Kernels): A set of directional kernels $\{\mathbf{h}_d\}_{d=1}^D$ provides estimates along different directions or orientations. Each kernel uses binary weights to select samples from a specific spatial direction.

Example 9.7.1.1 (One-Dimensional Directional Kernels): For a 1D signal with neighborhood size $M = 5$:

- \mathbf{h}_{left} : weights $[1, 1, 1, 0, 0]$
- $\mathbf{h}_{\text{right}}$: weights $[0, 0, 1, 1, 1]$
- $\mathbf{h}_{\text{center}}$: weights $[0, 1, 1, 1, 0]$

9.7.2 Intersection of Confidence Intervals

The Intersection of Confidence Intervals (ICI) rule provides a principled approach for adaptive neighborhood selection:

1. Compute estimates $\{\hat{f}_d\}$ and confidence intervals $\{\text{CI}_d\}$ for each directional kernel
2. Find the intersection of all confidence intervals: $\text{CI}_{\text{intersect}} = \bigcap_{d=1}^D \text{CI}_d$
3. Select the estimator with the largest neighborhood whose confidence interval contains $\text{CI}_{\text{intersect}}$

9.8 The Intersection of Confidence Intervals (ICI) Rule

The ICI rule provides a data-driven approach to support selection based on statistical confidence intervals:

Definition 9.8.1 (Confidence Interval for LPA): For a given support size h and confidence parameter $\gamma > 0$, the confidence interval for $\hat{f}_h(x_0)$ is:

$$\text{CI}_h(x_0) = \left[\hat{f}_h(x_0) - \gamma \sqrt{V(h, x_0)}, \hat{f}_h(x_0) + \gamma \sqrt{V(h, x_0)} \right] \quad (9.33)$$

For Gaussian noise with known variance σ^2 , choosing $\gamma = z_{\alpha/2}$ (the $\alpha/2$ quantile of the standard normal distribution) yields a $(1 - \alpha)$ confidence interval.

Theorem 9.8.1 (ICI Principle): Consider a sequence of support sizes $h_1 < h_2 < \dots < h_K$. Define:

$$\mathcal{J}_k(x_0) = \bigcap_{j=1}^k \text{CI}_{h_j}(x_0) \quad (9.34)$$

The optimal scale $k^*(x_0)$ is chosen as:

$$k^*(x_0) = \max\{k : \mathcal{J}_k(x_0) \neq \emptyset\} \quad (9.35)$$

Algorithm 7: Intersection of Confidence Intervals (ICI)

```

1: procedure ICI(signal  $\mathbf{y} \in \mathbb{R}^n$ , polynomial degree  $p$ , scales  $\{h_k\}_{k=1}^K$ , parameter  $\gamma$ )
2:
3:    $\triangleright$  Input: Signal  $\mathbf{y} \in \mathbb{R}^n$ , polynomial degree  $p$ , scales  $\{h_k\}_{k=1}^K$ , parameter  $\gamma$ 
4:
5:    $\triangleright$  Output: Optimal scales  $\{k^*(x_i)\}_{i=1}^n$  and estimates  $\{\hat{f}(x_i)\}_{i=1}^n$ 
6:
7:    $\triangleright$  Initialization:
8:   for  $i = 1, \dots, n$  do
9:      $L_0(x_i) \leftarrow -\infty$ 
10:     $U_0(x_i) \leftarrow +\infty$ 
11:     $k^*(x_i) \leftarrow 0$ 
12:   end
13:
14:    $\triangleright$  For each scale  $k = 1, 2, \dots, K$ :
15:   for  $i = 1, \dots, n$  do
16:      $\triangleright$  (a) Compute LPA estimates  $\hat{f}_{h_k}(x_i)$  for all  $i$ 
17:     for  $i = 1, \dots, n$  do
18:        $\hat{f}_{h_k}(x_i) \leftarrow$  LPA estimate with support  $h_k$ 
19:     end
20:
21:      $\triangleright$  (b) Compute variances  $V(h_k, x_i) = \sigma^2 \sum_j h_{j,k}^2(x_i)$ 
22:     for  $i = 1, \dots, n$  do
23:        $V(h_k, x_i) \leftarrow \sigma^2 \sum_j h_{j,k}^2(x_i)$ 
24:     end
25:
26:      $\triangleright$  (c) Calculate confidence bounds:
27:     for  $i = 1, \dots, n$  do
28:        $l_k(x_i) \leftarrow \hat{f}_{h_k}(x_i) - \gamma \sqrt{V(h_k, x_i)}$ 
29:        $u_k(x_i) \leftarrow \hat{f}_{h_k}(x_i) + \gamma \sqrt{V(h_k, x_i)}$ 
30:     end
31:
32:      $\triangleright$  (d) Update intersection bounds:
33:     for  $i = 1, \dots, n$  do
34:        $L_k(x_i) \leftarrow \max\{L_{k-1}(x_i), l_k(x_i)\}$ 
35:        $U_k(x_i) \leftarrow \min\{U_{k-1}(x_i), u_k(x_i)\}$ 
36:     end
37:
38:      $\triangleright$  (e) Check intersection validity:
39:     for  $i = 1, \dots, n$  do
40:       if  $U_k(x_i) < L_k(x_i) \wedge k^*(x_i) = 0$  then
41:          $k^*(x_i) \leftarrow k - 1$ 
42:       end
43:     end
44:   end
45:
46:    $\triangleright$  Final estimates:
47:   for  $i = 1, \dots, n$  do
48:      $\hat{f}(x_i) \leftarrow \hat{f}_{h_{k^*(x_i)}}(x_i)$ 
49:   end
50:   return  $\{k^*(x_i)\}_{i=1}^n, \{\hat{f}(x_i)\}_{i=1}^n$ 
51: end

```

10 Robust Fitting

The problem of *robust model fitting* arises ubiquitously in computer vision applications where we must estimate analytical expressions derived from geometric constraints in the presence of outliers. Unlike standard noise, outliers represent data points that significantly deviate from the expected model with statistical characteristics that are fundamentally different from the inlier distribution.

Definition 10.1 (Robust Fitting Problem): Given a set of observations $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ where $x_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$, the robust fitting problem seeks to estimate parameters $\theta \in \mathbb{R}^p$ of a model $f(x; \theta)$ such that:

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^n \rho(r_i(\theta)) \quad (10.1)$$

where $r_i(\theta) = y_i - f(x_i; \theta)$ is the residual and $\rho : \mathbb{R} \rightarrow \mathbb{R}^+$ is a robust loss function.

The choice of ρ fundamentally determines the robustness properties of the estimator. Classical least squares uses $\rho(r) = r^2$, which, as we shall demonstrate, lacks robustness to outliers.

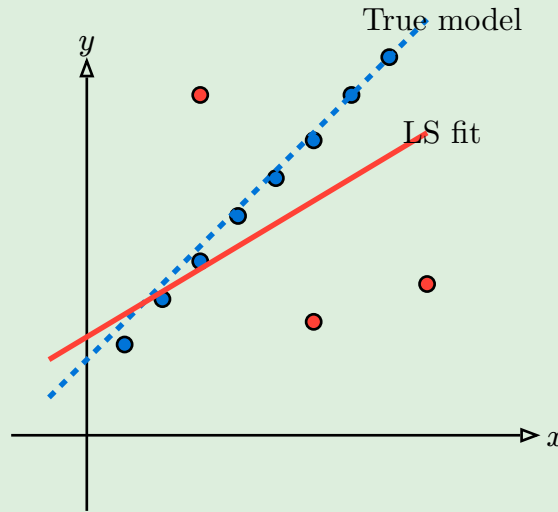


Figure 10.1: Illustration of outlier influence on least squares fitting

10.1 Limitations of Vertical Distance Minimization

We begin with the fundamental problem of fitting a straight line to a set of 2D points. The parametric form of a line is:

$$y = mx + b \quad (10.2)$$

where m is the slope and b is the y-intercept. Given observations $\{(x_i, y_i)\}_{i=1}^n$, ordinary least squares (OLS) seeks:

$$(m^*, b^*) = \arg \min_{m, b} \sum_{i=1}^n (y_i - mx_i - b)^2 \quad (10.3)$$

The ordinary least squares approach measures error along the vertical axis, which introduces several fundamental limitations:

1. **Inability to fit vertical lines:** When the line approaches vertical ($m \rightarrow \infty$), the formulation breaks down as we cannot express x as a function of y .

2. **Asymmetric treatment of variables:** The choice of dependent vs. independent variable artificially privileges one coordinate axis.
3. **Scale dependency:** The error metric depends on the coordinate system orientation.

10.2 The Direct Linear Transformation (DLT) Algorithm

To overcome the limitations of parametric representations, we adopt the implicit form:

$$ax + by + c = 0 \quad (10.4)$$

where $(a, b, c) \in \mathbb{R}^3$ are the line parameters. This representation elegantly handles all line orientations, including vertical lines (where $b = 0$).

▲ Scale Ambiguity

attention 10.2.1

The implicit representation introduces a scale ambiguity: the lines defined by (a, b, c) and $(\lambda a, \lambda b, \lambda c)$ for any $\lambda \neq 0$ are identical. This necessitates a normalization constraint.

10.2.1 Algebraic Error Minimization

Rather than minimizing geometric distance (which leads to nonlinear optimization), we minimize the *algebraic error*:

Definition 10.2.1.1 (Algebraic Error): For a point (x_i, y_i) and line parameters $\theta = [a, b, c]^T$, the algebraic error is:

$$r_i^{\text{alg}(\theta)} = ax_i + by_i + c \quad (10.5)$$

Definition 10.2.1.2 (Geometric Error): The geometric error is the perpendicular distance from the point to the line:

$$r_i^{\text{geom}(\theta)} = \frac{|ax_i + by_i + c|}{\sqrt{a^2 + b^2}} \quad (10.6)$$

The relationship between algebraic and geometric errors is:

$$r_i^{\text{geom}(\theta)} = \frac{|r_i^{\text{alg}(\theta)}|}{\sqrt{a^2 + b^2}} \quad (10.7)$$

10.2.2 Constrained Least Squares Formulation

The DLT problem is formulated as:

$$\theta^* = \arg \min_{\theta} \|A\theta\|_2^2 \quad \text{subject to} \quad \|\theta\|_2 = 1 \quad (10.8)$$

where the design matrix is:

$$A = \begin{pmatrix} x_1 & y_1 & 1 \\ x_2 & y_2 & 1 \\ \vdots & \vdots & \vdots \\ x_n & y_n & 1 \end{pmatrix} \in \mathbb{R}^{n \times 3} \quad (10.9)$$

Theorem 10.2.2.1 (DLT Solution via SVD): Let $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ be the singular value decomposition of \mathbf{A} , where:

- $\mathbf{U} \in \mathbb{R}^{n \times n}$ and $\mathbf{V} \in \mathbb{R}^{3 \times 3}$ are orthogonal matrices
- $\mathbf{\Sigma} = \text{diag}(\sigma_1, \sigma_2, \sigma_3)$ with $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq 0$

Then the optimal solution is $\boldsymbol{\theta}^* = \mathbf{v}_3$, the last column of \mathbf{V} .

Proof: Using the orthogonality of \mathbf{U} and the substitution $\mathbf{w} = \mathbf{V}^T \boldsymbol{\theta}$:

$$\begin{aligned} \|\mathbf{A}\boldsymbol{\theta}\|_2^2 &= \|\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \boldsymbol{\theta}\|_2^2 \\ &= \|\mathbf{\Sigma}\mathbf{V}^T \boldsymbol{\theta}\|_2^2 \quad (\text{orthogonal invariance}) \\ &= \|\mathbf{\Sigma}\mathbf{w}\|_2^2 = \sum_{i=1}^3 \sigma_i^2 w_i^2 \end{aligned} \tag{10.10}$$

Since $\|\boldsymbol{\theta}\|_2 = \|\mathbf{w}\|_2 = 1$, minimizing the objective requires placing all weight on the smallest singular value:

$$\mathbf{w}^* = [0, 0, 1]^T \Rightarrow \boldsymbol{\theta}^* = \mathbf{V}\mathbf{w}^* = \mathbf{v}_3 \tag{10.11}$$

Numerical conditioning is crucial for stable DLT computation, especially when dealing with higher-order geometric entities.

Algorithm 8: Direct Linear Transformation (DLT)

```

1: procedure DLT(data points  $\{(x_i, y_i)\}_{i=1}^n$ )
2:    $\triangleright$  Input: Data points  $\{(x_i, y_i)\}_{i=1}^n$ 
3:    $\triangleright$  Output: Line parameters  $\boldsymbol{\theta}^* = [a, b, c]^T$ 
4:
5:    $\triangleright$  Construct design matrix:
6:    $\mathbf{A} \leftarrow \begin{pmatrix} x_1 & y_1 & 1 \\ x_2 & y_2 & 1 \\ \vdots & \vdots & \vdots \\ x_n & y_n & 1 \end{pmatrix} \in \mathbb{R}^{n \times 3}$ 
7:
8:    $\triangleright$  Compute SVD of design matrix:
9:    $\mathbf{U}, \mathbf{\Sigma}, \mathbf{V} \leftarrow \text{SVD}(\mathbf{A})$ 
10:   $\triangleright$  where  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ 
11:
12:   $\triangleright$  Extract solution from smallest singular vector:
13:   $\boldsymbol{\theta}^* \leftarrow \mathbf{v}_3$ 
14:   $\triangleright$  last column of  $\mathbf{V}$  corresponding to smallest singular value
15:
16:   $\triangleright$  Optional: Normalize solution
17:  if normalization desired then
18:     $\boldsymbol{\theta}^* \leftarrow \frac{\boldsymbol{\theta}^*}{\|\boldsymbol{\theta}^*\|_2}$ 
19:  end
20:
21:  return  $\boldsymbol{\theta}^*$ 
22: end
```

10.2.3 Preconditioning for Numerical Stability

Numerical conditioning is crucial for stable DLT computation, especially when dealing with higher-order geometric entities.

Algorithm 9: Normalized DLT

```

1: procedure NORMALIZED-DLT(data points  $\{(x_i, y_i)\}_{i=1}^n$ )
2:   ▷ Input: Data points  $\{(x_i, y_i)\}_{i=1}^n$ 
3:
4:   ▷ Output: Line parameters  $\theta^* = [a, b, c]^T$ 
5:
6:   ▷ Compute data centroid and scale:
7:    $| (x) \leftarrow \frac{1}{n} \sum_{i=1}^n x_i$ 
8:
9:    $| (y) \leftarrow \frac{1}{n} \sum_{i=1}^n y_i$ 
10:
11:  ▷ Apply normalization transformation:
12:
13:   $T \leftarrow \begin{pmatrix} s_x & 0 & -s_x | (x) \\ 0 & s_y & -s_y | (y) \\ 0 & 0 & 1 \end{pmatrix}$ 
14:  ▷ where  $s_x$  and  $s_y$  are chosen such that the normalized data has unit variance
15:
16:  ▷ Solve DLT for normalized points:
17:  for  $i = 1, \dots, n$  do
18:     $\tilde{p}_i \leftarrow T \begin{pmatrix} x_i \\ y_i \\ 1 \end{pmatrix}$ 
19:
20:  end
21:   $\tilde{\theta}^* \leftarrow \text{DLT}(\{\tilde{p}_i\})$ 
22:
23:  ▷ Denormalize the solution:
24:   $\theta^* \leftarrow T^T \tilde{\theta}^*$ 
25:
26:  return  $\theta^*$ 
27: end

```

10.3 Robust Estimation Methods

Definition 10.3.1 (Breakdown Point): The breakdown point ε^* of an estimator is the largest fraction of arbitrarily corrupted observations that the estimator can handle while still producing a bounded estimate close to the true value.

For ordinary least squares, the breakdown point is $\varepsilon^* = 0$, meaning a single outlier with sufficient leverage can arbitrarily corrupt the estimate. This motivates the development of robust alternatives.

M-estimators generalize maximum likelihood estimation by replacing the squared loss with robust alternatives:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n \rho(r_i(\boldsymbol{\theta})) \quad (10.12)$$

Common choices for ρ include:

Name	$\rho(r)$	Influence Function $\psi(r) = \rho'(r)$
L2 (Least Squares)	r^2	$2r$
L1 (Least Absolute)	$ r $	$\text{sign}(r)$
Huber	$\begin{cases} r^2/2 & r \leq k \\ k r - k^2/2 & r > k \end{cases}$	$\begin{cases} r & r \leq k \\ k \text{sign}(r) & r > k \end{cases}$
Tukey's Bisquare	$\begin{cases} k^2/6 [1 - (1 - r^2/k^2)^3] & r \leq k \\ k^2/6 & r > k \end{cases}$	$\begin{cases} r(1 - r^2/k^2)^2 & r \leq k \\ 0 & r > k \end{cases}$

Table 2: Common M-estimator loss functions and their derivatives

10.4 RANSAC: Random Sample Consensus

RANSAC takes a fundamentally different approach by explicitly modeling the presence of outliers through a consensus framework.

10.4.1 The Consensus Set

Definition 10.4.1.1 (Consensus Set): Given parameters θ and threshold ε , the consensus set is:

$$\mathcal{C}(\theta) = \{i : |r_i(\theta)| \leq \varepsilon\} \quad (10.13)$$

RANSAC maximizes the cardinality of the consensus set:

$$\theta^* = \arg \max_{\theta} |\mathcal{C}(\theta)| \quad (10.14)$$

10.4.2 The RANSAC Algorithm

Algorithm 10: RANSAC

```

1: procedure RANSAC(Data points  $\mathcal{D}$ , inlier threshold  $\varepsilon$ , confidence  $p$ )
2:   ▷ Input: Data points  $\mathcal{D}$ , inlier threshold  $\varepsilon$ , confidence  $p$ 
3:   ▷ Output: Model parameters  $\theta^*$ 
4:
5:   ▷ Initialize:
6:    $\mathcal{C}^* \leftarrow \emptyset$ 
7:    $k \leftarrow 0$ 
8:
9:   while  $k < N$  do
10:    ▷ (a) Randomly sample minimal set  $\mathcal{S} \subset \mathcal{D}$  with  $|\mathcal{S}| = m$ 
11:     $\mathcal{S} \leftarrow$  randomly sample  $m$  points from  $\mathcal{D}$ 
12:
13:    ▷ (b) Fit model:
14:     $\theta_k \leftarrow \text{FitModel}(\mathcal{S})$ 
15:
16:    ▷ (c) Evaluate consensus:
17:     $\mathcal{C}_k \leftarrow \{i : |r_i(\theta_k)| \leq \varepsilon\}$ 
18:
19:    ▷ (d) if  $|\mathcal{C}_k| > |\mathcal{C}^*|$  then update best model
20:    if  $|\mathcal{C}_k| > |\mathcal{C}^*|$  then
21:       $\mathcal{C}^* \leftarrow \mathcal{C}_k$ 
22:       $\theta^* \leftarrow \theta_k$ 
23:    end
24:
25:     $k \leftarrow k + 1$ 
26:  end
27:
28:  ▷ Refine:
29:   $\theta^* \leftarrow \text{LeastSquares}(\mathcal{C}^*)$ 
30:
31:  return  $\theta^*$ 
32: end

```

10.4.3 Determining the Number of Iterations

The number of iterations N is determined probabilistically to ensure finding at least one outlier-free sample with confidence p :

Theorem 10.4.3.1 (RANSAC Iteration Count): Given inlier ratio w and desired confidence p , the required number of iterations is:

$$N = \frac{\log(1-p)}{\log(1-w^m)} \quad (10.15)$$

where m is the minimal sample size.

Proof: The probability of selecting an all-inlier minimal sample is w^m . The probability of failing to select such a sample in N attempts is $(1-w^m)^N$. Setting this equal to $1-p$ and solving for N yields the result.

Example 10.4.3.1 (RANSAC Iteration Example): For line fitting ($m = 2$) with 50% inliers ($w = 0.5$) and 99% confidence ($p = 0.99$):

$$N = \frac{\log(0.01)}{\log(1-0.5^2)} = \frac{\log(0.01)}{\log(0.75)} \approx 17 \quad (10.16)$$

10.5 MSAC and MLESAC Variants

10.5.1 MSAC: M-estimator Sample Consensus

MSAC refines RANSAC by considering the magnitude of residuals within the inlier band:

$$L_{\text{MSAC}}(\theta) = \sum_{i=1}^n \min(|r_i(\theta)|, \varepsilon) \quad (10.17)$$

This formulation corresponds to a truncated L1 loss, providing better discrimination between competing models with similar consensus set sizes.

10.5.2 MLESAC: Maximum Likelihood Sample Consensus

MLESAC models the error distribution explicitly as a mixture:

$$p(r_i) = \gamma \mathcal{N}(0, \sigma^2) + (1-\gamma) \mathcal{U}(-v, v) \quad (10.18)$$

where γ is the inlier ratio, and maximizes the likelihood of the observed residuals.

11 Multiple Model Fitting

The extension from single-model to multi-model fitting introduces fundamental challenges that require sophisticated approaches. The core difficulty lies in the chicken-and-egg problem: clustering points requires knowledge of models, while fitting models requires knowledge of point clusters.

Definition 11.1 (Multi-Model Fitting Problem): Given a dataset $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$ and a model family \mathcal{M} , find:

$$\Theta = \{\theta_1, \theta_2, \dots, \theta_K\} \quad (11.1)$$

$$\mathcal{P} = \{\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_K\} \quad (11.2)$$

where Θ represents the set of model parameters and \mathcal{P} represents a partition of the data such that:

$$\bigcup_{k=1}^K \mathcal{P}_k = \mathcal{D} \setminus \mathcal{O} \quad (11.3)$$

$$\mathcal{P}_i \cap \mathcal{P}_j = \emptyset \quad \forall i \neq j \quad (11.4)$$

and each θ_k optimally fits the points in \mathcal{P}_k .

The fundamental challenge in multi-model fitting can be formalized as follows:

1. **Clustering Requirement:** To fit model θ_k , we need to identify the subset $\mathcal{P}_k \subset \mathcal{D}$ of points that belong to model k
2. **Model Requirement:** To cluster point \mathbf{x}_i into partition \mathcal{P}_k , we need to know model θ_k to compute the distance $d(\mathbf{x}_i, \theta_k)$

This circular dependency necessitates iterative or simultaneous approaches that can break the cycle through:

- *Sequential strategies:* Fit models one at a time, removing inliers after each fit
- *Simultaneous strategies:* Jointly optimize over all models and partitions
- *Voting strategies:* Use parameter space voting to identify multiple models

11.1 Outlier Complexity in Multi-Model Settings

In multi-model fitting, the effective outlier ratio increases significantly. Consider fitting model θ_k in the presence of K models:

$$\varepsilon_{\text{effective}} = \varepsilon_{\text{true}} + \frac{K-1}{K}(1 - \varepsilon_{\text{true}}) \quad (11.5)$$

where $\varepsilon_{\text{true}}$ is the true outlier ratio and $\varepsilon_{\text{effective}}$ is the effective outlier ratio seen by each individual model.

Example: With 10 equally represented models and 10% true outliers:

$$\varepsilon_{\text{effective}} = 0.1 + \frac{9}{10}(1 - 0.1) = 0.1 + 0.9 \cdot 0.9 = 0.91 \quad (11.6)$$

This dramatic increase in effective outlier ratio severely impacts the performance of standard RANSAC, motivating specialized multi-model approaches.

11.2 Sequential Multi-Model Fitting

Sequential approaches address multi-model fitting by iteratively applying single-model fitting techniques, removing inliers after each successful fit. While conceptually simple, these methods can suffer from order dependency and accumulation of errors.

Algorithm 11: Sequential RANSAC

```

1: procedure SEQUENTIAL-RANSAC(dataset  $\mathcal{D}$ , minimum consensus  $\tau_{\min}$ , maximum models  $K_{\max}$ )
2:    $\triangleright$  Input: Dataset  $\mathcal{D}$ , minimum consensus  $\tau_{\min}$ , maximum models  $K_{\max}$ 
3:    $\triangleright$  Output: Set of models  $\Theta$ 
4:
5:    $\triangleright$  Initialize:
6:    $\mathcal{D}_{\text{current}} \leftarrow \mathcal{D}$ 
7:    $\Theta \leftarrow \emptyset$ 
8:    $k \leftarrow 0$ 
9:
10:   $\triangleright$  While  $k < K_{\max}$  and  $|\mathcal{D}_{\text{current}}| \geq \tau_{\min}$ :
11:  while  $k < K_{\max} \wedge |\mathcal{D}_{\text{current}}| \geq \tau_{\min}$  do
12:     $\triangleright$  (a) Apply RANSAC to  $\mathcal{D}_{\text{current}}$  to find model  $\theta_k$ 
13:     $\theta_k \leftarrow \text{RANSAC}(\mathcal{D}_{\text{current}})$ 
14:
15:     $\triangleright$  (b) Compute consensus set
16:     $\mathcal{C}_k \leftarrow \{x_i \in \mathcal{D}_{\text{current}} : d(x_i, \theta_k) \leq \tau\}$ 
17:
18:     $\triangleright$  (c) If  $|\mathcal{C}_k| < \tau_{\min}$ , terminate
19:    if  $|\mathcal{C}_k| < \tau_{\min}$  then
20:       $\triangleright$  Terminate - insufficient consensus
21:      BREAK
22:    end
23:
24:     $\triangleright$  (d) Refine  $\theta_k$  using least squares on  $\mathcal{C}_k$ 
25:     $\theta_k \leftarrow \text{LeastSquares}(\mathcal{C}_k)$ 
26:
27:     $\triangleright$  (e) Add  $\theta_k$  to  $\Theta$ 
28:     $\Theta \leftarrow \Theta \cup \{\theta_k\}$ 
29:
30:     $\triangleright$  (f) Update: remove inliers from current dataset
31:     $\mathcal{D}_{\text{current}} \leftarrow \mathcal{D}_{\text{current}} \setminus \mathcal{C}_k$ 
32:
33:     $\triangleright$  (g) Increment model counter
34:     $k \leftarrow k + 1$ 
35:  end
36:
37:  return  $\Theta$ 
38: end

```

Theorem 11.2.1 (Sequential RANSAC Convergence): Under the assumption that models are well-separated and the minimum consensus threshold is appropriately chosen, Sequential RANSAC will find all models with probability:

$$P_{\text{success}} = \prod_{k=1}^K \left(1 - (1 - (1 - \varepsilon_k)^m)^{T_k} \right) \quad (11.7)$$

where ε_k is the outlier ratio when fitting model k and T_k is the number of RANSAC iterations for model k .

Limitations of Sequential RANSAC:

1. *Order Dependency*: The sequence in which models are found depends on their relative support and may not reflect the true underlying structure
2. *Error Accumulation*: Misclassified points in early iterations affect subsequent model fitting
3. *Threshold Sensitivity*: Performance heavily depends on the choice of inlier threshold and minimum consensus
4. *Suboptimal Solutions*: Greedy selection may lead to locally optimal but globally suboptimal solutions

Sequential RANSAC can be understood through the lens of preference matrices, which provide insight into simultaneous multi-model approaches.

Definition 11.2.1 (Preference Matrix): The preference matrix $\mathbf{P} \in \mathbb{R}^{N \times M}$ is defined as:

$$P_{ij} = \begin{cases} 1 & \text{if } d(\mathbf{x}_i, \theta_j) \leq \tau \\ 0 & \text{otherwise} \end{cases} \quad (11.8)$$

where N is the number of data points and M is the number of model hypotheses.

Alternatively, the preference matrix can store residuals:

$$P_{ij} = d(\mathbf{x}_i, \theta_j) \quad (11.9)$$

Conceptual Insights:

- Each row represents a data point's affinity to all model hypotheses
- Each column represents a model's support across all data points
- Sequential RANSAC selects the column with maximum support, then removes corresponding rows
- Simultaneous approaches can exploit the full matrix structure

11.3 Simultaneous Multi-Model Fitting

Simultaneous approaches to multi-model fitting attempt to identify all models jointly, avoiding the limitations of sequential methods. The Hough transform represents a classical voting-based approach that has been widely successful in computer vision applications.

Instead of sequentially selecting columns from the preference matrix, simultaneous approaches seek to:

1. *Cluster rows*: Group data points with similar preference patterns
2. *Regularize solutions*: Incorporate sparsity constraints to prefer simpler models
3. *Optimize jointly*: Minimize a global objective function over all models and assignments

Definition 11.3.1 (Joint Optimization Problem): The simultaneous multi-model fitting problem can be formulated as:

$$\min_{\Theta, \mathcal{P}} \sum_{k=1}^K \sum_{\mathbf{x}_i \in \mathcal{P}_k} \rho(d(\mathbf{x}_i, \theta_k)) + \lambda \|\Theta\|_0 \quad (11.10)$$

$$\text{subject to } \mathcal{P}_i \cap \mathcal{P}_j = \emptyset \quad \forall i \neq j \quad (11.11)$$

$$\bigcup_{k=1}^K \mathcal{P}_k \subset \mathcal{D} \quad (11.12)$$

where $\rho(\cdot)$ is a robust loss function and λ controls model complexity.

11.3.1 The Hough Transform

The Hough transform provides an elegant solution to multi-model fitting by transforming the point clustering problem into a peak detection problem in parameter space.

The key insight of the Hough transform is the duality between point space and parameter space:

- Each point in data space corresponds to a curve in parameter space
- Each model in parameter space corresponds to a point in data space
- Points lying on the same model generate curves that intersect at the model's parameters

11.3.2 Line Detection via Hough Transform

For line detection, we parameterize lines using the normal form:

$$\rho = x \cos \theta + y \sin \theta \quad (11.13)$$

where ρ is the perpendicular distance from the origin to the line, and θ is the angle of the normal vector.

Algorithm 12: Hough Transform for Line Detection

```

1: procedure HOUGH-TRANSFORM(edge points  $\{(x_i, y_i)\}_{i=1}^N$ )
2:    $\triangleright$  Input: Edge points  $\{(x_i, y_i)\}_{i=1}^N$ 
3:    $\triangleright$  Output: Line parameters  $\{(\rho_k, \theta_k)\}$ 
4:
5:    $\triangleright$  Initialize accumulator array  $A[\rho, \theta]$  with appropriate discretization
6:    $A[\rho, \theta] \leftarrow 0$  for all  $\rho, \theta$ 
7:
8:    $\triangleright$  For each edge point  $(x_i, y_i)$ :
9:   for  $i = 1, \dots, N$  do
10:     $\triangleright$  (a) For  $\theta = -\frac{\pi}{2}$  to  $\frac{\pi}{2}$  with step  $\Delta\theta$ :
11:    for  $\theta = -\frac{\pi}{2}, \dots, \frac{\pi}{2}$ , Step:  $\Delta\theta$  do
12:       $\triangleright$  (i) Compute  $\rho = x_i \cos \theta + y_i \sin \theta$ 
13:       $\rho \leftarrow x_i \cos \theta + y_i \sin \theta$ 
14:
15:       $\triangleright$  (ii) Increment  $A[\rho, \theta]$ 
16:       $A[\rho, \theta] \leftarrow A[\rho, \theta] + 1$ 
17:    end
18:  end
19:
20:   $\triangleright$  Find local maxima in  $A[\rho, \theta]$  above threshold
21:  Peaks  $\leftarrow$  FindLocalMaxima( $A$ , threshold)
22:
23:   $\triangleright$  Each maximum corresponds to a line with parameters  $(\rho, \theta)$ 
24:  for each peak  $(\rho_k, \theta_k)$  in Peaks do
25:     $\triangleright$  Line detected with parameters  $(\rho_k, \theta_k)$ 
26:  end
27:
28:  return "Peaks"
29: end

```

Advantages of Normal Parameterization:

1. *Bounded Parameter Space:* $\rho \in [0, \rho_{\max}]$ and $\theta \in [-\frac{\pi}{2}, \frac{\pi}{2}]$
2. *No Singularities:* Unlike slope-intercept form, vertical lines are handled naturally
3. *Uniform Discretization:* Parameter space can be uniformly discretized

11.3.3 Extension to Circle Detection

Circle detection requires a three-dimensional parameter space (x_c, y_c, r) where (x_c, y_c) is the center and r is the radius. The circle equation is:

$$(x - x_c)^2 + (y - y_c)^2 = r^2 \quad (11.14)$$

Algorithm 13: Hough Transform for Circle Detection

```

1: procedure HOUGH-TRANSFORM-CIRCLE(edge points  $\{(x_i, y_i)\}_{i=1}^N$ , known radius  $r$ )
2:   ▷ Input: Edge points  $\{(x_i, y_i)\}_{i=1}^N$ , known radius  $r$ 
3:   ▷ Output: Circle centers  $\{(x_{c,k}, y_{c,k})\}$ 
4:
5:   ▷ Initialize 2D accumulator  $A[x_c, y_c]$ 
6:    $A[x_c, y_c] \leftarrow 0$  for all  $x_c, y_c$ 
7:
8:   ▷ For each edge point  $(x_i, y_i)$ :
9:   for  $i = 1, \dots, N$  do
10:    ▷ (a) For  $\theta = 0$  to  $2\pi$  with step  $\Delta\theta$ :
11:    for  $\theta = 0, \dots, 2\pi$  do
12:      ▷ (i) Compute center coordinates
13:       $x_c \leftarrow x_i + r \cos \theta$ 
14:       $y_c \leftarrow y_i + r \sin \theta$ 
15:
16:      ▷ (ii) Increment  $A[x_c, y_c]$ 
17:       $A[x_c, y_c] \leftarrow A[x_c, y_c] + 1$ 
18:    end
19:  end
20:
21:  ▷ Find peaks in  $A[x_c, y_c]$ 
22:  Peaks  $\leftarrow$  FindLocalMaxima( $A$ , threshold)
23:
24:  return "Peaks"
25: end

```

For unknown radius, a 3D accumulator is required, significantly increasing computational cost.

12 Appendix: Fundamental Concepts in Linear Algebra

12.1 Vector Spaces and Linear Combinations

Definition 12.1.1 (Span of Vectors): Given a set of vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\} \subset \mathbb{R}^m$, the *span* of these vectors is defined as:

$$\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\} = \left\{ \sum_{i=1}^n \lambda_i \mathbf{v}_i : \lambda_i \in \mathbb{R} \right\} \quad (12.1)$$

The span represents the set of all possible linear combinations of the given vectors, forming a vector subspace of \mathbb{R}^m . This concept is fundamental to understanding how different sets of vectors can generate different subspaces.

12.2 Linear Independence and Basis

Definition 12.2.1 (Linear Independence): A set of vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ is *linearly independent* if and only if:

$$\sum_{i=1}^n \lambda_i \mathbf{v}_i = \mathbf{0} \Rightarrow \lambda_i = 0 \text{ for all } i = 1, 2, \dots, n \quad (12.2)$$

This definition captures the fundamental property that no vector in the set can be expressed as a linear combination of the others. The importance of linear independence becomes clear when we consider the uniqueness of representations.

Theorem 12.2.1 (Uniqueness of Representation): Let $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\} \subset \mathbb{R}^m$ be a linearly independent set of vectors. If $\mathbf{s} \in \text{span}\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$, then there exists a unique representation:

$$\mathbf{s} = \sum_{i=1}^n x_i \mathbf{e}_i \quad (12.3)$$

where the coefficients $x_i \in \mathbb{R}$ are uniquely determined.

Proof: Suppose \mathbf{s} admits two different representations:

$$\begin{aligned} \mathbf{s} &= \sum_{i=1}^n x_i \mathbf{e}_i \\ \mathbf{s} &= \sum_{i=1}^n y_i \mathbf{e}_i \end{aligned} \quad (12.4)$$

Subtracting these equations:

$$\mathbf{0} = \sum_{i=1}^n (x_i - y_i) \mathbf{e}_i \quad (12.5)$$

By linear independence, $(x_i - y_i) = 0$ for all i , implying $x_i = y_i$ for all i . Therefore, the representation is unique.

12.3 Orthogonal Vectors

Definition 12.3.1 (Orthonormal Basis): A set of vectors $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\} \subset \mathbb{R}^n$ forms an *orthonormal basis* if:

1. $\langle \mathbf{e}_i, \mathbf{e}_j \rangle = \delta_{ij}$ (orthonormality condition)
2. $\text{span}\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\} = \mathbb{R}^n$ (spanning condition)

where δ_{ij} is the Kronecker delta. Kronecker delta is defined as:

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \quad (12.6)$$

The power of orthonormal bases lies in their computational convenience. For any signal $\mathbf{s} \in \mathbb{R}^n$ and orthonormal basis $\mathbf{D} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n]$, the coefficient computation is straightforward:

$$\mathbf{x} = \mathbf{D}^T \mathbf{s} \quad (12.7)$$

where $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$ and $x_i = \langle \mathbf{e}_i, \mathbf{s} \rangle$.

12.4 The ℓ_0 Norm

The ℓ_0 norm (more precisely, ℓ_0 pseudo-norm) of a vector $\mathbf{x} \in \mathbb{R}^n$ is defined as:

$$\|\mathbf{x}\|_0 := |\{i : x_i \neq 0\}| = \sum_{i=1}^n \mathbf{1}_{x_i \neq 0} \quad (12.8)$$

where $\mathbf{1}_{\{x_i \neq 0\}}$ is the indicator function that equals 1 if $x_i \neq 0$ and 0 otherwise.

The ℓ_0 norm can be understood as the limit of ℓ_p norms as $p \rightarrow 0^+$:

$$\|\mathbf{x}\|_0 = \lim_{p \rightarrow 0^+} \|\mathbf{x}\|_p^p = \lim_{p \rightarrow 0^+} \left(\sum_{i=1}^n |x_i|^p \right) \quad (12.9)$$

The ℓ_0 norm satisfies the following properties:

1. **Non-negativity:** $\|\mathbf{x}\|_0 \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$
2. **Zero property:** $\|\mathbf{x}\|_0 = 0$ if and only if $\mathbf{x} = \mathbf{0}$
3. **Triangle inequality:** $\|\mathbf{x} + \mathbf{y}\|_0 \leq \|\mathbf{x}\|_0 + \|\mathbf{y}\|_0$
4. **Failure of homogeneity:** $\|\lambda \mathbf{x}\|_0 \neq |\lambda| \|\mathbf{x}\|_0$ for $\lambda \neq 0, \pm 1$

Let us now interpret ℓ_0 -sparsity geometrically in \mathbb{R}^3 .

- $\|\boldsymbol{\alpha}\|_0 = 0$: Only the origin $(0, 0, 0)$
- $\|\boldsymbol{\alpha}\|_0 = 1$: Points on coordinate axes, e.g., $(7, 0, 0)$, $(0, 3, 0)$
- $\|\boldsymbol{\alpha}\|_0 = 2$: Points lying in coordinate planes, e.g., $(5, 2, 0)$
- $\|\boldsymbol{\alpha}\|_0 = 3$: All other points in \mathbb{R}^3

12.5 Matrix Spark

The concept of matrix spark provides the theoretical foundation for understanding when sparse solutions are unique.

Definition 12.5.1 (Matrix Spark): For a matrix $\mathbf{D} \in \mathbb{R}^{m \times n}$, the spark of \mathbf{D} , denoted $\text{spark}(\mathbf{D})$, is defined as:

$$\text{spark}(\mathbf{D}) = \min\{|\mathcal{S}| : \mathcal{S} \subseteq \{1, 2, \dots, n\}, \mathbf{D}_{\mathcal{S}} \text{ is linearly dependent}\} \quad (12.10)$$

where $|\mathcal{S}|$ denotes the cardinality of the set \mathcal{S} and $\mathbf{D}_{\mathcal{S}}$ represents the submatrix of \mathbf{D} formed by columns indexed by \mathcal{S} .

Equivalently, the spark can be defined in terms of the ℓ_0 norm as:

$$\text{spark}(\mathbf{D}) = \min\{\|\mathbf{x}\|_0 : \mathbf{D}\mathbf{x} = \mathbf{0}, \mathbf{x} \neq \mathbf{0}\} \quad (12.11)$$

The spark and rank of a matrix are related but distinct concepts:

Proposition 12.5.1 (Spark-Rank Relationship): For any matrix $\mathbf{D} \in \mathbb{R}^{m \times n}$ with $\text{rank}(\mathbf{D}) = r$:

$$1 \leq \text{spark}(\mathbf{D}) \leq r + 1 \quad (12.12)$$

Proof: The lower bound follows from the definition. For the upper bound, consider that any $r + 1$ columns must be linearly dependent in an r -dimensional space, hence $\text{spark}(\mathbf{D}) \leq r + 1$.

Lemma 12.5.1 (Spark and Homogeneous Solutions): If $\mathbf{D}\mathbf{x} = \mathbf{0}$ has a solution $\mathbf{x} \neq \mathbf{0}$, then:

$$\text{spark}(\mathbf{D}) \leq \|\mathbf{x}\|_0 \quad (12.13)$$

Proof: If $\mathbf{D}\mathbf{x} = \mathbf{0}$ with $\mathbf{x} \neq \mathbf{0}$, then $\sum_{i \in \text{supp}(\mathbf{x})} x_i \mathbf{d}_i = \mathbf{0}$, showing that the columns $\{\mathbf{d}_i : i \in \text{supp}(\mathbf{x})\}$ are linearly dependent. By definition of spark, $\text{spark}(\mathbf{D}) \leq |\text{supp}(\mathbf{x})| = \|\mathbf{x}\|_0$.

13 Appendix: Optimization Theory for Non-Differentiable Functions

13.1 Lipschitz Continuity

Definition 13.1.1 (Lipschitz Continuous Gradient): A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ has an **L -Lipschitz continuous gradient** if there exists a constant $L > 0$ such that:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2 \quad (13.1)$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. The smallest such constant L is called the Lipschitz constant of the gradient.

Definition 13.1.2 (Lipschitz Continuity of Gradient): A function f has L -Lipschitz continuous gradient if:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n \quad (13.2)$$

For the quadratic function $f(\mathbf{x}) = \frac{1}{2}\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$:

Proposition 13.1.1 (Lipschitz Constant for Quadratic Functions): For the quadratic function $f(\mathbf{x}) = \frac{1}{2}\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$, the Lipschitz constant can be derived as follows:

$$L = \|\mathbf{A}^T \mathbf{A}\|_2 = \lambda_{\max}(\mathbf{A}^T \mathbf{A}) \quad (13.3)$$

where λ_{\max} denotes the largest eigenvalue.

Proof: The gradient is $\nabla f(\mathbf{x}) = \mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{b})$. Thus:

$$\begin{aligned} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 &= \|\mathbf{A}^T \mathbf{A}(\mathbf{x} - \mathbf{y})\|_2 \\ &\leq \|\mathbf{A}^T \mathbf{A}\|_2 \|\mathbf{x} - \mathbf{y}\|_2 \\ &= \lambda_{\max}(\mathbf{A}^T \mathbf{A}) \|\mathbf{x} - \mathbf{y}\|_2 \end{aligned} \quad (13.4)$$

where we used the fact that the spectral norm equals the largest eigenvalue for symmetric positive semidefinite matrices.

13.2 Majorization-Minimization

For smooth convex functions, we can construct quadratic majorizers that facilitate optimization (subgradient methods).

Lemma 13.2.1 (Descent Lemma): Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex, differentiable function with Lipschitz continuous gradient. Then for any $\mathbf{x}_k \in \mathbb{R}^n$, there exists $L > 0$ such that:

$$f(\mathbf{x}) \leq Q_L(\mathbf{x}; \mathbf{x}_k) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \quad (13.5)$$

for all $\mathbf{x} \in \mathbb{R}^n$.

Definition 13.2.1 (Majorization-Minimization Algorithm): Given a convex function f , the majorization-minimization approach generates a sequence $\{\mathbf{x}_k\}$ by:

$$\begin{aligned} \mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} Q_L(\mathbf{x}; \mathbf{x}_k) \\ &= \arg \min_{\mathbf{x}} \left[f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \right] \end{aligned} \quad (13.6)$$

Minimizing the majorizer $Q_L(\mathbf{x}; \mathbf{x}_k)$ with respect to \mathbf{x} :

$$\begin{aligned} \nabla_{\mathbf{x}} Q_L(\mathbf{x}; \mathbf{x}_k) &= \nabla f(\mathbf{x}_k) + L(\mathbf{x} - \mathbf{x}_k) = \mathbf{0} \\ \Rightarrow \mathbf{x}_{k+1} &= \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k) \end{aligned} \quad (13.7)$$

This recovers the standard gradient descent update with step size $\gamma = \frac{1}{L}$.

For a convex, differentiable function f with Lipschitz continuous gradient, the gradient descent algorithm converges to the global minimum.