

# Image Processing

**Rao**

Politecnico di Milano

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

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

# Contents

1	Sparsity and Parsimony .....	3
1.1	Sparsity in Statistics .....	3
1.2	Sparsity in Signal Processing .....	3
2	Signal Processing .....	4
2.1	Discrete Cosine Transform (DCT) .....	4
3	Image Denoising .....	6
3.1	Local Constancy Prior .....	7
3.2	Sparsity-Based Image Prior .....	7
3.3	Noise Standard Deviation Estimation .....	8
3.4	Sliding DCT Algorithm .....	8
3.5	Wiener Filter .....	10
4	Limitations of Sparsity-Based Denoising .....	11
4.1	The Sparsity Problem .....	11
4.2	Experimental Demonstration .....	11
4.3	Overcomplete Dictionaries: The Solution .....	11
4.4	Regularization and Sparse Recovery .....	13
4.5	Towards Sparsity: $\ell_0$ and $\ell_1$ Regularization .....	14
5	Appendix: Fundamental Concepts in Linear Algebra .....	16
5.1	Vector Spaces and Linear Combinations .....	16
5.2	Linear Independence and Basis .....	16
5.3	Orthogonal Vectors .....	17

# 1 Sparsity and Parsimony

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.

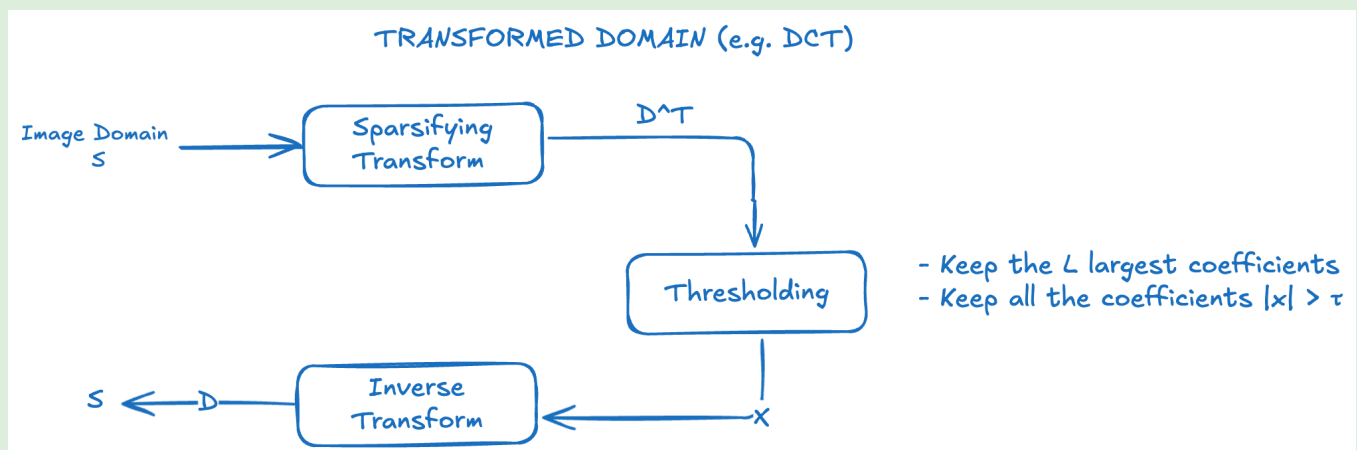


Figure 1.1: Enforce sparsity in signal processing

## 2 Signal Processing

### 2.1 Discrete Cosine Transform (DCT)

#### 2.1.1 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:

$$\text{DCT}_{k(n)} = c_k \cos\left(k\pi \frac{2n+1}{2M}\right) \quad n, k = 0, 1, \dots, M-1 \quad (2.1)$$

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

$$\text{DCT}_{k(n)} = \cos\left(k\pi \frac{2n+1}{2M}\right) \quad (2.2)$$

at  $n = 0, \dots, M-1$ , obtain a vector. Ignore the normalization coefficient. Divide each vector by its  $\ell_2$  norm.

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

$$x = \text{dct2}(s) = D^T s \quad (2.3)$$

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  $x$ :

$$s = \text{idct2}(x) = Dx \quad (2.4)$$

#### 2.1.2 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:  $X^T = D^T Z^T \rightarrow X = D^T S D$

**e.g. JPEG Compression****example 2.1.1**

The image is divided into non-overlapping  $8 \times 8$  blocks. Each block is treated separately during the compression process.

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

The compression ratio is defined as:

$$\text{Comp Ratio} = 1 - \frac{\# \text{Non-zero coefficients}}{\# \text{Pixels in the image}} \quad (2.5)$$

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

$$\text{PSNR} = 10 \log_{10} \left( \frac{1}{\text{MSE}(Y, \hat{Y})} \right) \quad (2.6)$$

### 3 Image Denoising

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

$$z(x) = y(x) + \eta(x) \quad (3.1)$$

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):

$$\text{MSE}(\hat{y}) = \mathbb{E}[\|\hat{y} - y\|_2^2] \quad (3.2)$$

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

#### Def Noise Prior

#### definition 3.0.1

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  $z(x)$ .

#### Def Image Prior

#### definition 3.0.2

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.

### 3.1 Local Constancy Prior

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

$$\hat{y} = \frac{1}{M} \sum_{i=1}^M z(x_i) \quad (3.3)$$

Properties of the this estimator:

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

**Limitations:** Local averaging introduces *bias* at edges:

$$\text{Bias} = |\mathbb{E}[\hat{y}] - y| \gg 0 \text{ at discontinuities}$$

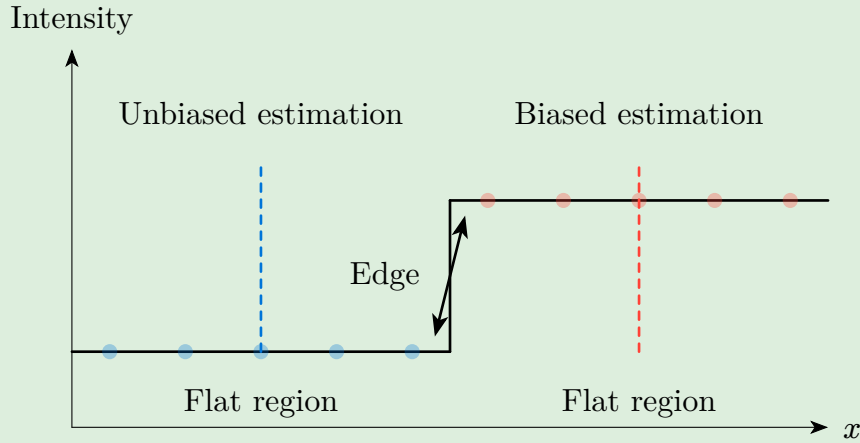


Figure 3.1: Bias-variance tradeoff in local averaging

### 3.2 Sparsity-Based Image Prior

#### 3.2.1 Motivation for Sparsity

Natural images have **sparse representations** in certain transform domains (e.g., DCT), as evidenced by the success of JPEG compression. **Key Insight:** If images can be sparsely represented for compression, this same property can be leveraged for denoising.

#### 3.2.2 DCT-Based Denoising Pipeline

**Step 1 : Analysis**

$$X = D^T S \quad (3.4)$$

where:

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

**Step 2: Enforce Sparsity (Thresholding)**

$$\hat{X}_i = \begin{cases} X_i & \text{if } |X_i| \geq \gamma \\ 0 & \text{if } |X_i| < \gamma \end{cases} \quad (3.5)$$

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

**Step 3: Synthesis**

$$\hat{S} = D\hat{X} \quad (3.6)$$

#### 3.2.3 Threshold Selection

**Thm Universal Thresholding****theorem 3.2.1**

For AWGN with variance  $\sigma^2$ , the optimal threshold for denoising is given by:

$$\gamma = \sigma \sqrt{2 \log(n)} \quad (3.7)$$

where:

- $\sigma$  is the noise standard deviation
- $n$  is the dimension of coefficients vector

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

### 3.3 Noise Standard Deviation Estimation

To use the universal thresholding, we need to estimate  $\sigma$  from the noisy image itself.

**Thm Robust Estimation of Noise Standard Deviation****theorem 3.3.1**

Given a noisy image  $Z$ , the noise standard deviation can be estimated using the following robust method:

$$\hat{\sigma} = \frac{\text{MAD}}{0.6745 \times \sqrt{2}} \quad (3.8)$$

where MAD = Median Absolute Deviation, defined as:

$$\text{MAD}(D) = \text{median}(|D - \text{median}(D)|) \quad (3.9)$$

where  $D$  denotes the horizontal differences of the image.

### 3.4 Sliding DCT Algorithm

#### 3.4.1 Non-Overlapping Tiles (No Aggregation)

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

$$X_{i,j} = \text{DCT}_2(B_{i,j}) \quad (3.10)$$

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

$$\hat{X}_{i,j}(u, v) = \begin{cases} X_{i,j}(u, v) & \text{if } |X_{i,j}(u, v)| \geq \gamma \\ 0 & \text{otherwise} \end{cases} \quad (3.11)$$

Reconstruct each block:

$$\hat{S}_{i,j} = \text{IDCT}_2(\hat{X}_{i,j}) \quad (3.12)$$

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

$$\hat{Y} = \bigcup_{i,j} \hat{B}_{i,j} \quad (3.13)$$

#### Properties:

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

#### 3.4.2 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:

$$\hat{I}(m, n) = \frac{\sum_{i,j \in \Omega_{m,n}} w \cdot \hat{S}_{i,j}(m - i, n - j)}{\sum_{i,j \in \Omega_{m,n}} w + \varepsilon} \quad (3.14)$$

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.

### 3.4.3 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:

$$w_{i,j} = \text{nnz}(\hat{X}_{i,j}) \quad (3.15)$$

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:

$$\hat{I}(m, n) = \frac{\sum_{i,j \in \Omega_{m,n}} w_{i,j} \cdot \hat{S}_{i,j}(m - i, n - j)}{\sum_{i,j \in \Omega_{m,n}} w_{i,j} + \varepsilon} \quad (3.16)$$

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

### 3.5 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.

#### 3.5.1 Empirical Wiener Filter

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

$$\hat{\mathbf{x}}^{\text{HT}} = D^T \hat{\mathbf{y}}^{\text{HT}} \quad (3.17)$$

The empirical Wiener filter attenuates the DCT coefficients as:

$$\hat{x}_i^{\text{Wie}} = \frac{(\hat{x}_i^{\text{HT}})^2}{(\hat{x}_i^{\text{HT}})^2 + \sigma^2} x_i \quad (3.18)$$

The empirical Wiener estimate is thus:

$$\hat{\mathbf{y}}^{\text{HT}} = D \hat{\mathbf{x}}^{\text{Wie}} \quad (3.19)$$

#### 3.5.2 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  $(i, j)$ :

$$\mathbf{P}_{i,j} = \mathbf{Y} \left[ i - \left\lfloor \frac{p}{2} \right\rfloor : i + \left\lfloor \frac{p}{2} \right\rfloor, j - \left\lfloor \frac{p}{2} \right\rfloor : j + \left\lfloor \frac{p}{2} \right\rfloor \right] \quad (3.20)$$

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  $\hat{\mathbf{p}}_{i,j} = \mathbf{T}^{-1} \hat{\mathbf{p}}_{i,j}$

## 4 Limitations of Sparsity-Based Denoising

### 4.1 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.*

#### e.g. DCT Basis Limitation

#### example 4.1.1

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

$$\mathbf{s}_0 = \mathbf{D}_{\text{DCT}} \mathbf{x}_0 \quad (4.1)$$

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

Now consider the modified signal:

$$\mathbf{s} = \mathbf{s}_0 + \lambda \mathbf{e}_j \quad (4.2)$$

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:

$$\mathbf{x} = \mathbf{D}_{\text{DCT}}^T \mathbf{s} = \mathbf{D}_{\text{DCT}}^T \mathbf{s}_0 + \lambda \mathbf{D}_{\text{DCT}}^T \mathbf{e}_j = \mathbf{x}_0 + \lambda \mathbf{D}_{\text{DCT}}^T \mathbf{e}_j \quad (4.3)$$

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.

### 4.2 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.

### 4.3 Overcomplete Dictionaries: The Solution

#### 4.3.1 Motivation for Redundancy

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  $m > n$ .

#### Def Overcomplete Dictionary

#### definition 4.3.1

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

$$\text{span}\{\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_m\} = \mathbb{R}^n \quad (4.4)$$

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

### 4.3.2 Construction of Overcomplete Dictionaries

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

$$\mathbf{D} = [\mathbf{D}_{\text{DCT}} \ \mathbf{I}] \in \mathbb{R}^{n \times 2n} \quad (4.5)$$

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

#### e.g. Sparse Representation with Overcomplete Dictionary

#### example 4.3.1

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:

$$\mathbf{s} = \mathbf{D} \begin{bmatrix} \mathbf{x}_0 \\ \lambda \mathbf{e}_j \end{bmatrix} \quad (4.6)$$

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

### 4.3.3 Theoretical Properties of Overcomplete Systems

#### Thm Rouché-Capelli Theorem

#### theorem 4.3.1

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:

$$\text{rank}(\mathbf{D}) = \text{rank}([\mathbf{D} \ \mathbf{s}]) \quad (4.7)$$

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

#### Thm Solution Space Dimension

#### theorem 4.3.2

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  $m - n$ .

## 4.4 Regularization and Sparse Recovery

### 4.4.1 The Ill-Posed Nature of Overcomplete Systems

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

#### Def Regularization

#### definition 4.4.1

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

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} J(\mathbf{x}) \quad \text{subject to} \quad D\mathbf{x} = \mathbf{s} \quad (4.8)$$

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

### 4.4.2 $\ell_2$ Regularization: Ridge Regression

The most mathematically tractable regularization is the  $\ell_2$  norm:

$$J(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2 = \frac{1}{2} \sum_{i=1}^m x_i^2 \quad (4.9)$$

This leads to the constrained optimization problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x}\|_2^2 \quad \text{subject to} \quad D\mathbf{x} = \mathbf{s} \quad (4.10)$$

Alternatively, we can formulate the unconstrained version:

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

### 4.4.3 Analytical Solution via Matrix Calculus

The unconstrained  $\ell_2$  regularization problem admits a closed-form solution. To derive this, we use matrix calculus.

**Thm Ridge Regression Solution****theorem 4.4.1**

The solution to the ridge regression problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{D}\mathbf{x} - \mathbf{s}\|_2^2 + \frac{\lambda}{2} \|\mathbf{x}\|_2^2 \quad (4.12)$$

is given by:

$$\hat{\mathbf{x}} = (\mathbf{D}^T \mathbf{D} + \lambda \mathbf{I})^{-1} \mathbf{D}^T \mathbf{s} \quad (4.13)$$

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

**Proof:** Define the objective function:

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{D}\mathbf{x} - \mathbf{s}\|_2^2 + \frac{\lambda}{2} \|\mathbf{x}\|_2^2 \quad (4.14)$$

Expanding the squared norms:

$$\begin{aligned} f(\mathbf{x}) &= \frac{1}{2} (\mathbf{D}\mathbf{x} - \mathbf{s})^T (\mathbf{D}\mathbf{x} - \mathbf{s}) + \frac{\lambda}{2} \mathbf{x}^T \mathbf{x} \\ &= \frac{1}{2} [\mathbf{x}]^T \mathbf{D}^T \mathbf{D} [\mathbf{x}] - [\mathbf{s}]^T \mathbf{D} [\mathbf{x}] + \frac{1}{2} [\mathbf{s}]^T [\mathbf{s}] + \frac{\lambda}{2} [\mathbf{x}]^T [\mathbf{x}] \end{aligned} \quad (4.15)$$

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

$$\nabla f(\mathbf{x}) = \mathbf{D}^T \mathbf{D} \mathbf{x} - \mathbf{D}^T \mathbf{s} + \lambda \mathbf{x} \quad (4.16)$$

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

$$(\mathbf{D}^T \mathbf{D} + \lambda \mathbf{I}) \mathbf{x} = \mathbf{D}^T \mathbf{s} \quad (4.17)$$

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

#### 4.4.4 Limitations of $\ell_2$ Regularization

While  $\ell_2$  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. $\ell_2$ Regularization

#### attention 4.4.1

The  $\ell_2$  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  $\ell_1$  norm or  $\ell_0$  pseudo-norm.

### 4.5 Towards Sparsity: $\ell_0$ and $\ell_1$ Regularization

#### 4.5.1 The $\ell_0$ “Norm” and True Sparsity

The most natural regularization for sparse recovery is the  $\ell_0$  “norm” (technically a pseudo-norm):

$$\|\mathbf{x}\|_0 = |\{i : x_i \neq 0\}| \quad (4.18)$$

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

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_0 \quad \text{subject to} \quad \mathbf{D}\mathbf{x} = \mathbf{s} \quad (4.19)$$

directly seeks the sparsest representation.

#### 4.5.2 Computational Challenges

The  $\ell_0$  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 Appendix: Fundamental Concepts in Linear Algebra

### 5.1 Vector Spaces and Linear Combinations

#### Def Span of Vectors

#### definition 5.1.1

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:

$$\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\} = \left\{ \sum_{i=1}^n \lambda_i \mathbf{v}_i : \lambda_i \in \mathbb{R} \right\} \quad (5.1)$$

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.

### 5.2 Linear Independence and Basis

#### Def Linear Independence

#### definition 5.2.1

A set of vectors  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  is *linearly independent* if and only if:

$$\sum_{i=1}^n \lambda_i \mathbf{v}_i = \mathbf{0} \Rightarrow \lambda_i = 0 \text{ for all } i = 1, 2, \dots, n \quad (5.2)$$

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.

#### Thm Uniqueness of Representation

#### theorem 5.2.1

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:

$$\mathbf{s} = \sum_{i=1}^n x_i \mathbf{e}_i \quad (5.3)$$

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

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

$$\begin{aligned} \mathbf{s} &= \sum_{i=1}^n x_i \mathbf{e}_i \\ \mathbf{s} &= \sum_{i=1}^n y_i \mathbf{e}_i \end{aligned} \quad (5.4)$$

Subtracting these equations:

$$\mathbf{0} = \sum_{i=1}^n (x_i - y_i) \mathbf{e}_i \quad (5.5)$$

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



## 5.3 Orthogonal Vectors

### Def Orthonormal Basis

### definition 5.3.1

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  $\delta_{ij}$  is the Kronecker delta.

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:

$$\mathbf{x} = \mathbf{D}^T \mathbf{s} \tag{5.6}$$

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