

# Sparsity-Promoting Algorithms and $\ell_0$ Optimization: A Comprehensive Treatment of Sparse Coding Theory

Lecture Notes on Signal Processing and Sparse Representation

July 11, 2025

## Contents

<b>1</b>	<b>Introduction to Sparsity and Redundant Representations</b>	<b>2</b>
1.1	Motivation: From Orthonormal Bases to Redundant Dictionaries . . . . .	2
1.2	The $\ell_0$ Norm: Definition and Properties . . . . .	3
1.3	Geometric Interpretation of Sparsity Constraints . . . . .	3
<b>2</b>	<b>The Sparse Coding Problem: Formulation and Complexity</b>	<b>4</b>
2.1	Problem Formulation . . . . .	4
2.2	Geometric Interpretation: Union of Subspaces . . . . .	4
2.3	Computational Complexity Analysis . . . . .	4
<b>3</b>	<b>Greedy Algorithms: The Matching Pursuit Framework</b>	<b>6</b>
3.1	The Greedy Paradigm . . . . .	6
3.2	Matching Pursuit Algorithm . . . . .	6
3.3	Mathematical Derivation of Key Formulas . . . . .	8
3.3.1	Projection Formula Derivation . . . . .	8
3.3.2	Error Formula Derivation . . . . .	8
3.4	Algorithm Properties and Limitations . . . . .	9
<b>4</b>	<b>Theoretical Analysis and Performance Guarantees</b>	<b>10</b>
4.1	Convergence Analysis . . . . .	10
4.2	Approximation Quality . . . . .	10
4.3	Computational Complexity . . . . .	10
4.4	Comparison with Optimal Solutions . . . . .	11
<b>5</b>	<b>Examples, Applications, and Extensions</b>	<b>12</b>
5.1	Detailed Algorithmic Example . . . . .	12
5.2	Practical Applications . . . . .	13
5.2.1	Image Processing . . . . .	13
5.2.2	Signal Processing Applications . . . . .	13
5.3	Dictionary Learning Perspective . . . . .	13
5.4	Extensions and Variants . . . . .	13

5.4.1	Orthogonal Matching Pursuit (OMP)	13
5.4.2	Regularized Variants	14
<b>6</b>	<b>Summary and Future Directions</b>	<b>15</b>
6.1	Key Takeaways	15
6.2	Theoretical Foundations Established	15
6.3	Algorithmic Contributions	15
6.4	Limitations and Open Questions	16
6.5	Preview of Advanced Topics	16
<b>A</b>	<b>Mathematical Appendix</b>	<b>17</b>
A.1	Detailed Proof of Residual Update Formula	17
A.2	Alternative Derivation Using Lagrangian Methods	17
A.3	Coherence-Based Performance Analysis	17
<b>B</b>	<b>Implementation Notes</b>	<b>18</b>
B.1	Numerical Considerations	18
B.2	Computational Optimizations	18
B.3	Parameter Selection Guidelines	18
<b>C</b>	<b>Notation and Terminology</b>	<b>19</b>
C.1	Mathematical Notation	19
C.2	Algorithmic Terminology	19
C.3	Problem Classifications	19
<b>D</b>	<b>Conclusion</b>	<b>20</b>

# 1 Introduction to Sparsity and Redundant Representations

In the realm of signal processing and machine learning, the concept of *sparsity* has emerged as a fundamental principle for achieving efficient and meaningful representations of data. This lecture explores the mathematical foundations of sparsity-promoting algorithms, with particular emphasis on the  $\ell_0$  norm optimization problem and its computational challenges.

## 1.1 Motivation: From Orthonormal Bases to Redundant Dictionaries

Traditional signal processing relies heavily on orthonormal bases for signal representation. However, as demonstrated in previous lectures, natural signals often exhibit sparse structure with respect to overcomplete dictionaries rather than standard orthonormal bases.

**Definition 1.1** (Redundant Dictionary). A **redundant dictionary**  $\mathbf{D} \in \mathbb{R}^{m \times n}$  with  $n > m$  is a matrix whose columns  $\{\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_n\}$  span the ambient space  $\mathbb{R}^m$  but are linearly dependent due to the overcompleteness constraint  $n > m$ .

The fundamental representation problem can be formulated as:

$$\mathbf{y} = \mathbf{D}\mathbf{x} \tag{1}$$

where:

- $\mathbf{y} \in \mathbb{R}^m$  is the observed signal in the ambient space
- $\mathbf{D} \in \mathbb{R}^{m \times n}$  is the redundant dictionary with  $n > m$
- $\mathbf{x} \in \mathbb{R}^n$  is the sparse coefficient vector we seek to determine

**Important:** The key insight is that while traditional orthonormal representations provide unique decompositions, redundant dictionaries offer flexibility at the cost of uniqueness, necessitating additional constraints (such as sparsity) to select meaningful solutions.

## 1.2 The $\ell_0$ Norm: Definition and Properties

**Definition 1.2** ( $\ell_0$  Norm). *The  $\ell_0$  norm (more precisely,  $\ell_0$  pseudo-norm) of a vector  $\mathbf{x} \in \mathbb{R}^n$  is defined as:*

$$\|\mathbf{x}\|_0 := \text{card}\{i : x_i \neq 0\} = \sum_{i=1}^n \mathbf{1}_{x_i \neq 0} \quad (2)$$

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

The  $\ell_0$  norm can be understood as the limit of  $\ell_p$  norms as  $p \rightarrow 0^+$ :

$$\|\mathbf{x}\|_0 = \lim_{p \rightarrow 0^+} \|\mathbf{x}\|_p^p = \lim_{p \rightarrow 0^+} \left( \sum_{i=1}^n |x_i|^p \right) \quad (3)$$

**Proposition 1.3** (Properties of the  $\ell_0$  Norm). *The  $\ell_0$  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  $\lambda \neq 0, \pm 1$

*Proof Sketch.* Properties 1-3 follow directly from the definition as a cardinality function. Property 4 demonstrates why  $\|\cdot\|_0$  is not a true norm: for any  $\mathbf{x} \neq \mathbf{0}$  and  $\lambda \neq 0$ , we have  $\|\lambda \mathbf{x}\|_0 = \|\mathbf{x}\|_0$  regardless of  $|\lambda|$ , violating the homogeneity requirement of a norm.  $\square$

## 1.3 Geometric Interpretation of Sparsity Constraints

The geometric structure imposed by  $\ell_0$  constraints fundamentally differs from traditional linear subspace projections.

**Example 1.4** (Sparsity Sets in Low Dimensions). *In  $\mathbb{R}^3$ , the sets of vectors with different sparsity levels form:*

$$S_0 = \{\mathbf{x} : \|\mathbf{x}\|_0 = 0\} = \{\mathbf{0}\} \quad (\text{origin only}) \quad (4)$$

$$S_1 = \{\mathbf{x} : \|\mathbf{x}\|_0 = 1\} = \text{coordinate axes} \quad (5)$$

$$S_2 = \{\mathbf{x} : \|\mathbf{x}\|_0 = 2\} = \text{coordinate planes} \quad (6)$$

$$S_3 = \{\mathbf{x} : \|\mathbf{x}\|_0 = 3\} = \mathbb{R}^3 \setminus (S_0 \cup S_1 \cup S_2) \quad (7)$$

This geometric perspective reveals that sparsity constraints define a *union of linear subspaces* rather than a single linear subspace, fundamentally altering the optimization landscape.

## 2 The Sparse Coding Problem: Formulation and Complexity

### 2.1 Problem Formulation

The central problem in sparse coding can be formally stated as:

$$\boxed{\begin{array}{ll} \mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}} & \|\mathbf{x}\|_0 \\ \text{subject to} & \mathbf{D}\mathbf{x} = \mathbf{y} \end{array}} \quad (8)$$

This is known as the  $P_0$  **problem** in the sparse coding literature.

**Remark 2.1** (Relationship to Previous Methods). *This formulation extends our previous approach of selecting the largest coefficients in orthonormal decompositions to the overcomplete dictionary setting, where the representation itself must be determined simultaneously with the sparsity constraint.*

### 2.2 Geometric Interpretation: Union of Subspaces

The sparsity constraint fundamentally changes the geometric structure of admissible solutions. Instead of seeking projections onto a single linear subspace (as in Principal Component Analysis), we consider projections onto a *union of low-dimensional subspaces*.

**Definition 2.2** (Sparsity-Constrained Subspaces). *For a given sparsity level  $k$ , the set of  $k$ -sparse vectors representable by dictionary  $\mathbf{D}$  is:*

$$\mathcal{S}_k(\mathbf{D}) = \{\mathbf{D}\mathbf{x} : \mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|_0 \leq k\} \quad (9)$$

*This set forms a union of  $\binom{n}{k}$  distinct  $k$ -dimensional subspaces.*

**Example 2.3** (Two-Dimensional Illustration). *Consider  $\mathbf{D} \in \mathbb{R}^{2 \times 4}$  with columns  $\{\mathbf{d}_1, \mathbf{d}_2, \mathbf{d}_3, \mathbf{d}_4\}$ . The 1-sparse representable vectors form the union:*

$$\mathcal{S}_1(\mathbf{D}) = \operatorname{span}\{\mathbf{d}_1\} \cup \operatorname{span}\{\mathbf{d}_2\} \cup \operatorname{span}\{\mathbf{d}_3\} \cup \operatorname{span}\{\mathbf{d}_4\} \quad (10)$$

*This creates four distinct lines through the origin, and any signal  $\mathbf{y}$  would be projected onto the closest of these four lines.*

### 2.3 Computational Complexity Analysis

**Theorem 2.4** (NP-Hardness of  $\ell_0$  Minimization). *The sparse coding problem (8) is NP-hard in general.*

*Proof Outline.* The proof follows by reduction from the NP-complete subset selection problem. For any desired sparsity level  $k$ , one must potentially examine all  $\binom{n}{k}$  possible support sets, leading to combinatorial explosion.  $\square$

**Brute Force Complexity Analysis** A brute force approach would involve:

1. Testing all possible sparsity levels  $s = 1, 2, \dots, k$
2. For each sparsity level  $s$ , examining all  $\binom{n}{s}$  possible support sets
3. Solving a least-squares problem for each support set

The total computational complexity becomes:

$$\mathcal{O} \left( \sum_{s=1}^k \binom{n}{s} \cdot s^3 \right) \approx \mathcal{O} \left( \binom{n}{k} \cdot k^3 \right) \quad (11)$$

**Example 2.5** (Computational Intractability). *Consider a realistic scenario with:*

- *Signal dimension:  $m = 500$  (e.g., a  $22 \times 22$  image patch)*
- *Dictionary size:  $n = 1000$  (4E overcomplete)*
- *Target sparsity:  $k = 20$  (4% sparse)*
- *Linear system solution time:  $10^{-9}$  seconds per system*

*The number of combinations to test is:*

$$\binom{1000}{20} \approx 10^{51} \quad (12)$$

*Even with supercomputer capabilities, this would require approximately  $10^{31}$  years far exceeding the age of the universe.*

### 3 Greedy Algorithms: The Matching Pursuit Framework

Given the computational intractability of exact  $\ell_0$  minimization, we turn to greedy approximation algorithms that provide computationally feasible solutions.

#### 3.1 The Greedy Paradigm

**Definition 3.1** (Greedy Algorithm Principle). A **greedy algorithm** for sparse coding makes locally optimal choices at each iteration without reconsidering previous decisions, building up the solution incrementally by adding one dictionary atom at a time.

**Example 3.2** (Coin Change Analogy). The greedy approach mirrors the coin change problem:

- **Goal:** Minimize the number of coins to make change
- **Greedy strategy:** Always use the largest denomination possible
- **Limitation:** Optimal only for specially designed coin systems

For standard currency systems (e.g.,  $\{1, 2, 5, 10, 20, 50\}$ ), greedy gives optimal solutions. However, for pathological systems (e.g.,  $\{1, 3, 4\}$ ), greedy fails: making change for 6 units gives greedy solution  $4+1+1$  (3 coins) vs. optimal  $3+3$  (2 coins).

#### 3.2 Matching Pursuit Algorithm

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

---

**Input:** Signal  $\mathbf{y}$ , dictionary  $\mathbf{D}$ , stopping criterion

**Output:** Sparse representation  $\mathbf{x}$

1. **Initialize:**

$$\mathbf{x}^{(0)} = \mathbf{0} \quad (\text{coefficient vector}) \quad (13)$$

$$\mathbf{r}^{(0)} = \mathbf{y} \quad (\text{residual}) \quad (14)$$

$$\Omega^{(0)} = \emptyset \quad (\text{active set}) \quad (15)$$

$$k = 0 \quad (\text{iteration counter}) \quad (16)$$

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

$$E_j^{(k)} = \left\| \mathbf{r}^{(k)} - \frac{(\mathbf{r}^{(k)})^T \mathbf{d}_j}{\|\mathbf{d}_j\|_2^2} \mathbf{d}_j \right\|_2^2 \quad (17)$$

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

$$j^* = \operatorname{argmin}_{j=1, \dots, n} E_j^{(k)} \quad (18)$$

Equivalently (by maximizing correlation):

$$j^* = \operatorname{argmax}_{j=1, \dots, n} \frac{|(\mathbf{r}^{(k)})^T \mathbf{d}_j|^2}{\|\mathbf{d}_j\|_2^2} \quad (19)$$

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

$$z_{j^*}^{(k)} = \frac{(\mathbf{r}^{(k)})^T \mathbf{d}_{j^*}}{\|\mathbf{d}_{j^*}\|_2^2} \quad (20)$$

5. **Solution Update:**

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + z_{j^*}^{(k)} \mathbf{e}_{j^*} \quad (21)$$

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

6. **Residual Update:**

$$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - z_{j^*}^{(k)} \mathbf{d}_{j^*} \quad (22)$$

7. **Active Set Update:**

$$\Omega^{(k+1)} = \Omega^{(k)} \cup \{j^*\} \quad (23)$$

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

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

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



### 3.3 Mathematical Derivation of Key Formulas

#### 3.3.1 Projection Formula Derivation

The projection of residual  $\mathbf{r}^{(k)}$  onto atom  $\mathbf{d}_j$  is obtained by solving:

$$\min_{z_j} \|\mathbf{r}^{(k)} - z_j \mathbf{d}_j\|_2^2 \quad (24)$$

$$\frac{d}{dz_j} \|\mathbf{r}^{(k)} - z_j \mathbf{d}_j\|_2^2 = \frac{d}{dz_j} [(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)} - 2z_j (\mathbf{r}^{(k)})^T \mathbf{d}_j + z_j^2 \mathbf{d}_j^T \mathbf{d}_j] \quad (25)$$

$$= -2(\mathbf{r}^{(k)})^T \mathbf{d}_j + 2z_j \|\mathbf{d}_j\|_2^2 \quad (26)$$

Setting the derivative to zero:

$$z_j^* = \frac{(\mathbf{r}^{(k)})^T \mathbf{d}_j}{\|\mathbf{d}_j\|_2^2} \quad (27)$$

This confirms equation (20).

#### 3.3.2 Error Formula Derivation

The approximation error becomes:

$$E_j^{(k)} = \|\mathbf{r}^{(k)} - z_j^* \mathbf{d}_j\|_2^2 \quad (28)$$

$$= (\mathbf{r}^{(k)})^T \mathbf{r}^{(k)} - 2z_j^* (\mathbf{r}^{(k)})^T \mathbf{d}_j + (z_j^*)^2 \|\mathbf{d}_j\|_2^2 \quad (29)$$

Substituting the optimal  $z_j^*$ :

$$E_j^{(k)} = \|\mathbf{r}^{(k)}\|_2^2 - 2 \frac{((\mathbf{r}^{(k)})^T \mathbf{d}_j)^2}{\|\mathbf{d}_j\|_2^2} + \frac{((\mathbf{r}^{(k)})^T \mathbf{d}_j)^2}{\|\mathbf{d}_j\|_2^2} \quad (30)$$

$$= \|\mathbf{r}^{(k)}\|_2^2 - \frac{((\mathbf{r}^{(k)})^T \mathbf{d}_j)^2}{\|\mathbf{d}_j\|_2^2} \quad (31)$$

**Proposition 3.3** (Error Positivity). *The approximation error  $E_j^{(k)}$  is always non-negative, with equality if and only if  $\mathbf{r}^{(k)}$  is parallel to  $\mathbf{d}_j$ .*

*Proof.* By the Cauchy-Schwarz inequality:

$$((\mathbf{r}^{(k)})^T \mathbf{d}_j)^2 \leq \|\mathbf{r}^{(k)}\|_2^2 \|\mathbf{d}_j\|_2^2 \quad (32)$$

Therefore:

$$\frac{((\mathbf{r}^{(k)})^T \mathbf{d}_j)^2}{\|\mathbf{d}_j\|_2^2} \leq \|\mathbf{r}^{(k)}\|_2^2 \quad (33)$$

Equality holds if and only if  $\mathbf{r}^{(k)}$  and  $\mathbf{d}_j$  are linearly dependent.  $\square$

### 3.4 Algorithm Properties and Limitations

**Proposition 3.4** (Residual Monotonicity). *The Matching Pursuit algorithm produces a monotonically decreasing sequence of residual norms:*

$$\|\mathbf{r}^{(k+1)}\|_2 \leq \|\mathbf{r}^{(k)}\|_2 \quad (34)$$

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

**Remark 3.5** (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*

## 4 Theoretical Analysis and Performance Guarantees

### 4.1 Convergence Analysis

**Theorem 4.1** (Convergence of Matching Pursuit). *For any finite dictionary  $\mathbf{D}$  and signal  $\mathbf{y}$ , the Matching Pursuit algorithm converges in the sense that:*

$$\lim_{k \rightarrow \infty} \|\mathbf{r}^{(k)}\|_2 = \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2 \quad (35)$$

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

*Proof Sketch.* The proof relies on the fact that the residual energy decreases monotonically and is bounded below by zero. The key insight is that if the algorithm fails to converge to the optimal approximation, there must exist some atom that maintains significant correlation with the residual, contradicting the optimality of atom selection.  $\square$

### 4.2 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.

**Definition 4.2** (Dictionary Coherence). *The **coherence** of a dictionary  $\mathbf{D}$  with normalized columns is:*

$$\mu(\mathbf{D}) = \max_{i \neq j} |\mathbf{d}_i^T \mathbf{d}_j| \quad (36)$$

**Theorem 4.3** (Approximation Bound). *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.*

### 4.3 Computational Complexity

**Proposition 4.4** (MP 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.

## 4.4 Comparison with Optimal Solutions

**Example 4.5** (Suboptimality Illustration). *Consider a simple 2D case where the true 1-sparse representation uses atom  $\mathbf{d}_3$ , but the signal has small noise components that align better with  $\mathbf{d}_1$ . Matching Pursuit might select  $\mathbf{d}_1$  first, then require additional atoms to approximate the remaining signal components, resulting in a less sparse solution than optimal.*

This limitation motivates the development of orthogonal variants and more sophisticated algorithms that will be covered in subsequent lectures.

## 5 Examples, Applications, and Extensions

### 5.1 Detailed Algorithmic Example

**Example 5.1** (Complete MP Execution). *Consider a 2D signal  $\mathbf{y} = [3, 1]^T$  and dictionary:*

$$\mathbf{D} = \begin{bmatrix} 1 & 0 & 0.6 & 0.8 \\ 0 & 1 & 0.8 & 0.6 \end{bmatrix} \quad (37)$$

**Iteration 1:**

- *Initial residual:*  $\mathbf{r}^{(0)} = [3, 1]^T$
- *Correlations:*

$$(\mathbf{r}^{(0)})^T \mathbf{d}_1 = 3 \quad (38)$$

$$(\mathbf{r}^{(0)})^T \mathbf{d}_2 = 1 \quad (39)$$

$$(\mathbf{r}^{(0)})^T \mathbf{d}_3 = 3(0.6) + 1(0.8) = 2.6 \quad (40)$$

$$(\mathbf{r}^{(0)})^T \mathbf{d}_4 = 3(0.8) + 1(0.6) = 3.0 \quad (41)$$

- *Normalized correlations (all atoms have unit norm):*

$$\frac{|(\mathbf{r}^{(0)})^T \mathbf{d}_1|^2}{\|\mathbf{d}_1\|_2^2} = 9 \quad (42)$$

$$\frac{|(\mathbf{r}^{(0)})^T \mathbf{d}_2|^2}{\|\mathbf{d}_2\|_2^2} = 1 \quad (43)$$

$$\frac{|(\mathbf{r}^{(0)})^T \mathbf{d}_3|^2}{\|\mathbf{d}_3\|_2^2} = 6.76 \quad (44)$$

$$\frac{|(\mathbf{r}^{(0)})^T \mathbf{d}_4|^2}{\|\mathbf{d}_4\|_2^2} = 9 \quad (45)$$

- *Both  $\mathbf{d}_1$  and  $\mathbf{d}_4$  achieve maximum correlation. Choose  $\mathbf{d}_1$  (by convention).*
- *Update:*  $x_1^{(1)} = 3$ ,  $\mathbf{r}^{(1)} = [0, 1]^T$

**Iteration 2:**

- *Current residual:*  $\mathbf{r}^{(1)} = [0, 1]^T$
- *Select  $\mathbf{d}_2$  (perfect alignment)*
- *Update:*  $x_2^{(1)} = 1$ ,  $\mathbf{r}^{(2)} = [0, 0]^T$

*Final sparse representation:*  $\mathbf{x} = [3, 1, 0, 0]^T$  *with sparsity 2.*

## 5.2 Practical Applications

### 5.2.1 Image Processing

In image processing, sparse coding with overcomplete dictionaries enables:

- **Denoising:** Natural images are often sparse in learned dictionaries, allowing separation of signal from noise
- **Compression:** Sparse representations require storage of only non-zero coefficients and their locations
- **Inpainting:** Missing image regions can be reconstructed using sparse priors from remaining data
- **Super-resolution:** High-resolution details can be recovered using sparse models learned from training data

### 5.2.2 Signal Processing Applications

- **Audio processing:** Speech and music signals exhibit sparsity in time-frequency dictionaries
- **Biomedical signals:** ECG, EEG, and other physiological signals benefit from adaptive sparse representations
- **Radar and communications:** Sparse channel estimation and signal detection in wireless systems

## 5.3 Dictionary Learning Perspective

**Remark 5.2** (Adaptive Dictionaries). *While this lecture focuses on sparse coding with fixed dictionaries, practical applications often involve **dictionary learning** the joint optimization of both the dictionary  $\mathbf{D}$  and sparse codes  $\mathbf{X}$  from training data:*

$$\min_{\mathbf{D}, \mathbf{X}} \|\mathbf{Y} - \mathbf{DX}\|_F^2 + \lambda \|\mathbf{X}\|_0 \quad (46)$$

where  $\mathbf{Y}$  contains training signals as columns and  $\mathbf{X}$  contains corresponding sparse codes.

## 5.4 Extensions and Variants

### 5.4.1 Orthogonal Matching Pursuit (OMP)

A key limitation of standard Matching Pursuit is the potential for atom reselection. **Orthogonal Matching Pursuit** addresses this by:

1. Maintaining orthogonality of residuals to previously selected atoms
2. Solving least-squares problems over the active set at each iteration
3. Guaranteeing that each atom is selected at most once

### 5.4.2 Regularized Variants

Alternative formulations replace the combinatorial  $\ell_0$  constraint with tractable regularizers:

- **$\ell_1$  regularization (LASSO):**  $\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$
- **Elastic net:** Combines  $\ell_1$  and  $\ell_2$  penalties
- **Group sparsity:** Encourages sparsity at the group level

## 6 Summary and Future Directions

### 6.1 Key Takeaways

**Important: Fundamental Insights:**

1. Redundant dictionaries provide representational flexibility at the cost of uniqueness
2. The  $\ell_0$  norm captures sparsity but leads to combinatorially hard optimization problems
3. Greedy algorithms like Matching Pursuit provide computationally tractable approximations
4. The geometric structure of sparsity involves unions of low-dimensional subspaces rather than single linear subspaces

### 6.2 Theoretical Foundations Established

This lecture has established several crucial theoretical foundations:

1. **Problem formulation:** The sparse coding problem as  $\ell_0$ -constrained optimization
2. **Complexity analysis:** NP-hardness of exact solutions and polynomial-time greedy approximations
3. **Algorithmic framework:** The Matching Pursuit paradigm and its mathematical derivation
4. **Geometric insight:** Sparsity as projection onto unions of subspaces

### 6.3 Algorithmic Contributions

The Matching Pursuit algorithm provides:

- **Computational tractability:**  $\mathcal{O}(kmn)$  complexity vs. exponential for exact methods
- **Convergence guarantees:** Monotonic residual decrease and asymptotic optimality
- **Implementation simplicity:** Clear iterative structure suitable for practical deployment
- **Theoretical foundation:** Mathematical framework for understanding greedy sparse approximation



## 6.4 Limitations and Open Questions

Several important limitations motivate further research:

1. **Suboptimality:** Greedy selection may miss globally optimal sparse solutions
2. **Atom reselection:** Standard MP may select atoms multiple times, reducing efficiency
3. **Dictionary dependence:** Performance heavily depends on dictionary structure and coherence
4. **Parameter selection:** Choice of stopping criteria affects solution quality and computational cost

## 6.5 Preview of Advanced Topics

Future lectures will address these limitations through:

- **Orthogonal Matching Pursuit:** Eliminating atom reselection through orthogonal projections
- **Stagewise Orthogonal Matching Pursuit:** Improved atom selection strategies
- **Convex relaxations:**  $\ell_1$  minimization and its theoretical guarantees
- **Dictionary learning:** Joint optimization of dictionaries and sparse codes
- **Compressed sensing:** Theoretical foundations for sparse signal recovery

## A Mathematical Appendix

### A.1 Detailed Proof of Residual Update Formula

**Lemma A.1** (Residual Update Correctness). *The residual update in Matching Pursuit satisfies:*

$$\mathbf{r}^{(k+1)} = \mathbf{y} - \mathbf{D}\mathbf{x}^{(k+1)} \quad (47)$$

where  $\mathbf{x}^{(k+1)}$  is the updated coefficient vector.

*Proof.* Starting from the coefficient update:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + z_{j^*}^{(k)} \mathbf{e}_{j^*} \quad (48)$$

The residual becomes:

$$\mathbf{r}^{(k+1)} = \mathbf{y} - \mathbf{D}\mathbf{x}^{(k+1)} \quad (49)$$

$$= \mathbf{y} - \mathbf{D}(\mathbf{x}^{(k)} + z_{j^*}^{(k)} \mathbf{e}_{j^*}) \quad (50)$$

$$= \mathbf{y} - \mathbf{D}\mathbf{x}^{(k)} - z_{j^*}^{(k)} \mathbf{D}\mathbf{e}_{j^*} \quad (51)$$

$$= \mathbf{r}^{(k)} - z_{j^*}^{(k)} \mathbf{d}_{j^*} \quad (52)$$

This confirms the residual update formula (22).  $\square$

### A.2 Alternative Derivation Using Lagrangian Methods

The sparse coding problem can also be approached using constrained optimization theory. Consider the Lagrangian relaxation:

$$L(\mathbf{x}, \lambda) = \|\mathbf{x}\|_0 + \lambda \|\mathbf{D}\mathbf{x} - \mathbf{y}\|_2^2 \quad (53)$$

While this approach doesn't directly solve the combinatorial problem, it provides insight into the trade-off between sparsity and reconstruction fidelity that underlies practical algorithms.

### A.3 Coherence-Based Performance Analysis

**Theorem A.2** (Exact Recovery Conditions). *Let  $\mathbf{x}^*$  be a  $k$ -sparse vector and  $\mathbf{y} = \mathbf{D}\mathbf{x}^*$ . If the dictionary  $\mathbf{D}$  satisfies:*

$$\mu(\mathbf{D}) < \frac{1}{2k-1} \quad (54)$$

*then Matching Pursuit exactly recovers  $\mathbf{x}^*$  in  $k$  iterations.*

*Proof Outline.* The proof relies on showing that under the coherence condition, the correlations between the residual and atoms in the true support always exceed correlations with atoms outside the support. This ensures correct atom selection at each iteration.  $\square$

## B Implementation Notes

### B.1 Numerical Considerations

- **Dictionary normalization:** Ensure all atoms have unit norm to avoid bias in correlation computations
- **Numerical precision:** Use appropriate stopping criteria to handle floating-point arithmetic limitations
- **Tie-breaking:** Implement consistent tie-breaking rules for atoms with equal correlations

### B.2 Computational Optimizations

- **Precomputed norms:** Store  $\|\mathbf{d}_j\|_2^2$  values to avoid repeated computation
- **Matrix-vector products:** Utilize optimized BLAS routines for correlation computations
- **Early stopping:** Monitor convergence criteria to avoid unnecessary iterations

### B.3 Parameter Selection Guidelines

- **Maximum sparsity:** Set based on expected signal characteristics and computational constraints
- **Residual threshold:** Choose based on noise level and desired approximation quality
- **Dictionary size:** Balance representational power with computational complexity

## C Notation and Terminology

### C.1 Mathematical Notation

Symbol	Definition
$\mathbf{y} \in \mathbb{R}^m$	Observed signal vector
$\mathbf{D} \in \mathbb{R}^{m \times n}$	Dictionary matrix with $n > m$
$\mathbf{x} \in \mathbb{R}^n$	Sparse coefficient vector
$\mathbf{d}_j$	$j$ -th column (atom) of dictionary $\mathbf{D}$
$\ \mathbf{x}\ _0$	$\ell_0$ pseudo-norm (cardinality)
$\ \mathbf{x}\ _p$	$\ell_p$ norm: $(\sum_i  x_i ^p)^{1/p}$
$\mathbf{r}^{(k)}$	Residual at iteration $k$
$\mathbf{x}^{(k)}$	Coefficient estimate at iteration $k$
$\Omega^{(k)}$	Active set (support) at iteration $k$
$\mu(\mathbf{D})$	Coherence of dictionary $\mathbf{D}$
$\text{supp}(\mathbf{x})$	Support of vector $\mathbf{x}$
$\text{span}(\mathcal{S})$	Linear span of set $\mathcal{S}$
$\text{argmin}_x f(x)$	Argument minimizing function $f$
$\text{card}(\mathcal{S})$	Cardinality of set $\mathcal{S}$

### C.2 Algorithmic Terminology

**Atom** A single column of the dictionary matrix

**Active set** The set of indices corresponding to non-zero coefficients

**Support** The locations of non-zero elements in a sparse vector

**Residual** The approximation error at each iteration

**Sweep stage** The process of testing all atoms for best correlation

**Greedy selection** Locally optimal choice without global consideration

**Sparsity level** The number of non-zero coefficients ( $\ell_0$  norm)

**Coherence** Maximum absolute correlation between distinct dictionary atoms

### C.3 Problem Classifications

**$P_0$  problem** Exact  $\ell_0$  minimization (NP-hard)

**$P_1$  problem**  $\ell_1$  relaxation (convex, tractable)

**Matching Pursuit** Greedy  $\ell_0$  approximation algorithm

**Orthogonal MP** Greedy algorithm with orthogonality constraints

**Basis Pursuit**  $\ell_1$  minimization approach to sparse coding

**LASSO**  $\ell_1$ -regularized least squares

## D Conclusion

This comprehensive treatment of  $\ell_0$  optimization and the Matching Pursuit algorithm establishes the mathematical foundations for understanding sparse coding in overcomplete dictionaries. The key insights that sparsity induces a union of subspaces structure, that exact optimization is computationally intractable, and that greedy algorithms provide practical approximations form the basis for more advanced sparse coding techniques.

The theoretical framework developed here, including the geometric interpretation of sparsity constraints, the complexity analysis of exact methods, and the mathematical derivation of the Matching Pursuit algorithm, provides the necessary background for understanding modern sparse representation theory and its applications across signal processing, machine learning, and computational mathematics.

Future developments in this field continue to build upon these foundational concepts, addressing the limitations of basic greedy methods through orthogonal projections, convex relaxations, and adaptive dictionary learning approaches that will be explored in subsequent lectures.