

Dartmouth College

ENGS 109

High Dimensional Sensing and Learning

Comprehensive Course Notes

Final Exam Reference Edition

Scope. Compressive sensing, sparse recovery, greedy and convex algorithms, random matrices, matrix completion, robust PCA, dictionary learning, nonnegative matrix factorization, and deep sparse coding.

Exam reference. The opening quick-reference section contains the Lecture 25X final-review guidance and the principal proof templates.

Prepared by: Farhan Sadeek

Term: Spring 2026

Source text: Foucart & Rauhut, *A Mathematical Introduction to Compressive Sensing*

Contents

Preface	1
Exam Quick Reference	2
I Foundations: Sparsity and the Inverse Problem	10
1 An Invitation to Compressive Sensing	10
1.1 The central question	10
1.2 Sparsity and compressibility	10
1.3 Recovery formulations	12
1.4 Linear-program reformulation of P_1 (real case)	12
1.5 Motivating applications	12
1.6 Founding papers	13
2 Sparse Solutions of Underdetermined Systems	14
2.1 When does a sparse solution exist—and is it unique?	14
2.2 Spark	15
2.3 Prony's method: practical recovery from $m = 2s$ Fourier samples	16
2.4 Computational complexity	16
II Algorithms for Sparse Recovery	18
3 Sampling, Least Squares, and Greedy Pursuit	18
3.1 Classical sampling	18
3.2 Least squares geometry	19
3.3 Greedy pursuit algorithms	19
3.4 Thresholding-based algorithms	20
3.5 Phase transitions	21
4 ℓ_0 Minimization and Greedy Performance	21
4.1 Why ℓ_0 is hard	21
4.2 Complexity comparison	21
4.3 Performance theorems	22

4.4 Stopping rules 23

III Theoretical Recovery Conditions 24

5 Spark, Null Space Property, and Basis Pursuit 24

5.1 The Null Space Property 24

5.2 Stable NSP and approximately sparse signals 25

5.3 Robust NSP and noisy measurements 26

5.4 Recovery of individual vectors via dual certificates 27

5.5 Tangent cone characterization 28

5.6 Low-rank matrix recovery 29

6 Coherence 30

6.1 Mutual coherence 30

6.2 The Welch bound 30

6.3 Coherence-based recovery guarantees 31

6.4 The famous coherence bound 33

7 Restricted Isometry Property 33

7.1 The RIP 33

7.2 $RIP \Rightarrow NSP \Rightarrow BP$ recovery 34

7.3 RIP-based guarantees for thresholding and greedy algorithms 35

7.4 Information-theoretic lower bound 37

IV Random Matrices, Sensing, and Algorithms 39

8 Random Matrices and Sensing-Matrix Design 39

8.1 Subgaussian random matrices 39

8.2 Gaussian matrices: explicit constants 40

8.3 Singular-value concentration of Gaussian matrices 41

8.4 Johnson–Lindenstrauss embeddings and the $JL \Leftrightarrow RIP$ duality 42

8.5 Fast and structured constructions 43

8.6 Non-uniform recovery 44

V	Sparse Representation Classification	46
9	Sparse Representation-based Classification	46
9.1	Setup	46
9.2	The SRC algorithm (Wright et al. 2009)	46
9.3	Robust SRC under occlusion	46
9.4	Sparsity Concentration Index (SCI)	47
9.5	Theoretical correctness	47
9.6	Comparison with classical classifiers	48
9.7	Applications	48
VI	ℓ_1-Minimization Algorithms	49
10	Convex Algorithms for ℓ_1 Minimization	49
10.1	The toolkit	49
10.2	Geometry of why ℓ_1 promotes sparsity	49
10.3	Soft thresholding and proximal operators	49
10.4	ISTA and FISTA	50
10.5	ADMM for BP	51
10.6	Primal–dual splitting (Chambolle–Pock 2011)	51
10.7	Iteratively Reweighted Least Squares (IRLS)	52
10.8	LARS / Homotopy method	54
10.9	Iteratively reweighted ℓ_1	55
10.10	Robust PCA and outliers	55
10.11	Software	55
	Appendices	56
A	Probability Tools	56
A.1	Markov, Chebyshev, Chernoff	56
A.2	Cramér and Hoeffding	56
A.3	Subgaussian random variables	57
A.4	Bernstein inequalities	58
A.5	Concentration of measure	60

A.6	Khintchine, Slepian, Gordon, Dudley	60
B	Convex Analysis	63
B.1	Convex sets and functions	63
B.2	Convex conjugate	64
B.3	Subdifferential	64
B.4	Lagrangian duality and the BP dual certificate	65
C	Matrix Analysis	66
C.1	Norms and the SVD	66
C.2	Pseudoinverse and least squares	67
C.3	Vandermonde matrices and the determinant	68
C.4	Norm equivalences	69
C.5	Schatten norms and matrix functions	70
D	Multiple Measurement Vectors and Joint Sparsity	71
D.1	The MMV model and row sparsity	71
D.2	Uniqueness for the MMV problem	71
D.3	Row-NSP: a structural characterization	72
D.4	$\ell_{1,2}$ -recovery via coherence	73
D.5	Simultaneous OMP (SOMP)	74
D.6	$\ell_{1,2}$ -minimization via ALM	75
D.7	Group sparse and hierarchical models	75
D.8	Heterogeneous dictionaries and gross noise	75
D.9	Applications	76
E	Matrix Completion and Robust PCA	76
E.1	Matrix Completion: setup	76
E.2	The rank null-space property	77
E.3	Conditions for exact recovery	78
E.4	Robust PCA	78
E.5	Algorithms: SVT and ALM	80
E.6	Applications	81
F	Dictionary Learning	81
F.1	The dictionary learning problem	81

F.2	K-means (the simplest dictionary)	82
F.3	Method of Optimal Directions (MOD)	82
F.4	K-SVD: rank-1 atom updates	83
F.5	Online dictionary learning	84
F.6	Locality-constrained linear coding (LLC)	85
F.7	Sparse Subspace Clustering (SSC)	85
G	Nonnegative Matrix Factorization	86
G.1	Why parts-based?	86
G.2	Multiplicative update rules (Lee–Seung)	87
G.3	Variants	87
G.4	Pros and cons	88
G.5	Separable NMF	88
G.6	Sparsity from entropy	89
G.7	Applications	90
H	Deep Sparse Coding Networks	90
H.1	Motivation	90
H.2	Composite sparse coding module	90
H.3	Inference and training	91
H.4	Performance	92
H.5	Why SCN works: feature-map clustering	92
H.6	Pros and cons	93
	Course Synthesis	94

Preface

These notes cover ENGS 109 *High Dimensional Sensing and Learning*, a graduate course on **compressive sensing** (CS), sparse representation theory, and applications in machine learning and signal processing. They are written to be self-contained and mathematically rigorous, with every major theorem followed by a proof or proof sketch. The reference text is Foucart & Rauhut, *A Mathematical Introduction to Compressive Sensing* (Birkhäuser/Springer, 2013), abbreviated [FR] below; chapter and theorem numbers refer to that monograph.

The course studies the underdetermined inverse problem

$$\mathbf{y} = A\mathbf{x}, \quad A \in \mathbb{C}^{m \times N}, \quad m < N,$$

and identifies the conditions on A under which a sparse \mathbf{x} can be recovered exactly or stably from \mathbf{y} . The theoretical core is the implication chain

$$\boxed{\text{random matrix}} \implies \boxed{\text{RIP}} \implies \boxed{\text{NSP}} \implies \boxed{\ell_1 \text{ recovers } \ell_0},$$

augmented by the coherence-based theory and the specialized analyses of greedy and thresholding algorithms.

Notation.

Vectors are bold lowercase ($\mathbf{x}, \mathbf{y}, \mathbf{v}, \mathbf{a}_j$); matrices are uppercase (A, V, U, Σ). \mathbf{a}_j denotes the j -th column of A , and A_S the submatrix of columns indexed by $S \subseteq [N] := \{1, \dots, N\}$. The ℓ_p -(quasi)norm is $\|\mathbf{x}\|_p = (\sum_j |x_j|^p)^{1/p}$ for $0 < p < \infty$, $\|\mathbf{x}\|_\infty = \max_j |x_j|$, and $\|\mathbf{x}\|_0 := |\text{supp}(\mathbf{x})|$. The set of s -sparse vectors is

$$\Sigma_s := \{\mathbf{x} \in \mathbb{C}^N : \|\mathbf{x}\|_0 \leq s\}.$$

Best s -term approximation error in the ℓ_p -norm is $\sigma_s(\mathbf{x})_p := \inf_{\mathbf{z} \in \Sigma_s} \|\mathbf{x} - \mathbf{z}\|_p$. The complement of S in $[N]$ is denoted \bar{S} or S^c . We write $\mathbf{v}_S \in \mathbb{C}^N$ for the vector with $v_S(j) = v_j$ if $j \in S$ and 0 otherwise. The Hermitian transpose of A is A^* (or A^\top in the real case); $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{v}^* \mathbf{u}$ is the Euclidean inner product.

Exam Quick Reference

This section is a compact exam reference for the course. It is intended for rapid lookup: identify the model, choose the structural prior, then quote the right condition or algorithm.

Lecture 25X final-exam guidance

Final-review logistics and exam profile

- **Exam time.** Sunday, June 7, 2026, 8:00 am, regular classroom.
- **Allowed materials.** Unlimited *paper* notes and a calculator. The Lecture 25X review explicitly says no phone, laptop, tablet, or smartwatch. Use this as the controlling logistics note for permitted materials.
- **Difficulty.** Comparable to homework proofs. Minh emphasized proof-and-construction problems rather than long computations.
- **Likely format.** Expect one or two questions of the form: “prove statement (A) is equivalent to statement (B),” especially in the null-space, spark, and uniqueness family.
- **Problem-solving approach.** The problem is usually not pure recall: it needs one intermediate construction, such as splitting a support into two pieces or forming the difference of two feasible solutions.

Likely exam archetypes

1. **Equivalence proof.** Show uniqueness of sparse recovery iff a kernel/spark/NSP condition holds. One direction is usually by contradiction; the other is usually by direct construction.
2. **Compute a certificate.** For a small explicit matrix, compute or bound $\text{spark}(A)$, $\mu(A)$, or an RIP-related quantity, then state the largest guaranteed sparsity.
3. **Trace a greedy algorithm.** Run MP/OMP by hand: compute correlations $|A^*r|$, select the largest, solve least squares on the selected support, and update the residual.
4. **Derive an update rule.** Write the loss, take a gradient, then either set it to zero (ALS/MOD) or use a chosen step size (NMF multiplicative updates).
5. **Match structure to method.** Sparse vector $\rightarrow \ell_1$ /OMP; low rank \rightarrow nuclear norm/ALS; low rank plus sparse \rightarrow RPCA; nonnegative parts \rightarrow NMF; learned sparse dictionary \rightarrow K-SVD/MOD.

6. **Short conceptual.** Be ready to explain why ℓ_1 promotes sparsity, why $\|\cdot\|_0$ is not a norm, and why NMF topics are more interpretable than SVD components.

Spark equivalence proof template

This is the worked proof Minh highlighted as exam-level. Learn the construction, not only the theorem statement.

Uniform P_0 uniqueness iff no short kernel vector

For $A \in \mathbb{R}^{m \times N}$ and integer $s \geq 1$, the following are equivalent:

(A) every s -sparse \mathbf{x} is the unique minimizer of $\min_z \|\mathbf{z}\|_0$ s.t. $A\mathbf{z} = A\mathbf{x}$,

and

(B) every nonzero $\mathbf{v} \in \ker A$ has $\|\mathbf{v}\|_0 > 2s$, equivalently $\text{spark}(A) > 2s$.

(A) \Rightarrow (B), contradiction. If (B) fails, choose $\mathbf{0} \neq \mathbf{v} \in \ker A$ with $\|\mathbf{v}\|_0 \leq 2s$. Split $S = \text{supp}(\mathbf{v})$ into disjoint S_1, S_2 with $|S_1|, |S_2| \leq s$. Define

$$\mathbf{x} = \mathbf{v}_{S_1}, \quad \mathbf{z} = -\mathbf{v}_{S_2}.$$

Then \mathbf{x} and \mathbf{z} are distinct s -sparse vectors and $\mathbf{x} - \mathbf{z} = \mathbf{v} \in \ker A$, so $A\mathbf{x} = A\mathbf{z}$. This gives two distinct feasible s -sparse explanations for the same measurement, contradicting uniqueness.

(B) \Rightarrow (A), construction. Let \mathbf{x} be any s -sparse vector and let \mathbf{z} be feasible with $A\mathbf{z} = A\mathbf{x}$ and $\|\mathbf{z}\|_0 \leq \|\mathbf{x}\|_0 \leq s$. Then

$$\mathbf{v} = \mathbf{x} - \mathbf{z} \in \ker A, \quad \|\mathbf{v}\|_0 \leq \|\mathbf{x}\|_0 + \|\mathbf{z}\|_0 \leq 2s.$$

By (B), \mathbf{v} must be zero, hence $\mathbf{z} = \mathbf{x}$. Therefore \mathbf{x} is the unique sparsest feasible vector.

Universal hinge. If two candidates satisfy $A\mathbf{x} = A\mathbf{z}$, then $\mathbf{x} - \mathbf{z} \in \ker A$. If both candidates are s -sparse, their difference is at most $2s$ -sparse. Most uniqueness proofs in this course are this argument with a stronger condition replacing $\text{spark}(A) > 2s$.

NMF / PS9 essentials from Lecture 25X

For a nonnegative matrix $V \in \mathbb{R}_{\geq 0}^{m \times n}$, NMF solves

$$\min_{W, H \geq 0} \frac{1}{2} \|V - WH\|_F^2, \quad W \in \mathbb{R}_{\geq 0}^{m \times r}, \quad H \in \mathbb{R}_{\geq 0}^{r \times n}.$$

The problem is nonconvex jointly in (W, H) , but convex in one factor with the other fixed.

Item	Exam-relevant fact
Lee–Seung updates	$H \leftarrow H \circ \frac{W^T V}{W^T W H + \epsilon}, W \leftarrow W \circ \frac{V H^T}{W H H^T + \epsilon}$. Products are matrix multiplications; \circ and division are entrywise.
Shape check	$W^T V$ and $W^T W H$ are $r \times n$, matching H . $V H^T$ and $W H H^T$ are $m \times r$, matching W .
Update order	Use H first, then update W with the newly updated H (Gauss–Seidel style). This is the standard notebook order.
Why nonnegative	Start with $W, H \geq 0$. Each update multiplies a nonnegative entry by a nonnegative ratio, so no projection step is needed.
Derivation target	$\nabla_H \frac{1}{2} \ V - WH\ _F^2 = W^T W H - W^T V$. Choosing $\eta_{ij} = H_{ij} / (W^T W H)_{ij}$ in a gradient step gives $H_{ij} \leftarrow H_{ij} (W^T V)_{ij} / (W^T W H)_{ij}$.
Common implementation errors	Using matrix multiply instead of entrywise multiply/divide in the final line; forgetting ϵ ; getting an increasing reconstruction error curve; transposing the wrong factor.
Topic modeling	<code>TfidfVectorizer</code> returns documents \times terms. Transpose to terms \times documents if columns of W are supposed to be topics over words. TF-IDF is nonnegative, so it is valid NMF input.
NMF vs SVD	SVD components have mixed signs and can use cancellation; NMF topics are additive nonnegative word lists, so they are usually more interpretable.

Selected worked results

- $\|\cdot\|_0$ **is not a norm**. It fails homogeneity: for nonzero \mathbf{x} , $\|2\mathbf{x}\|_0 = \|\mathbf{x}\|_0 \neq 2\|\mathbf{x}\|_0$. It does satisfy positive definiteness and $\|\mathbf{x} + \mathbf{y}\|_0 \leq \|\mathbf{x}\|_0 + \|\mathbf{y}\|_0$ by the support union bound.
- ℓ_1 vs ℓ_2 **two-variable example**. For $A = [1 \ 2]$ and $y = 2$, feasible vectors satisfy $x_1 + 2x_2 = 2$. The ℓ_1 minimizer is $(0, 1)^\top$; the minimum- ℓ_2 solution is $A^\top (A A^\top)^{-1} y = (2/5, 4/5)^\top$.

- **OMP hand-trace result.** For columns $\mathbf{a}_1 = (1, 0, 0)^\top$, $\mathbf{a}_2 = (0, 1, 0)^\top$, $\mathbf{a}_3 = (1/\sqrt{2}, 1/\sqrt{2}, 0)^\top$, $\mathbf{a}_4 = (0, 0, 1)^\top$ and $\mathbf{y} = (2, 0, 3)^\top$, two-step OMP picks support $\{4, 1\}$ and returns $\mathbf{x} = (2, 0, 0, 3)^\top$.

- **Small spark/coherence example.** For $A = \begin{bmatrix} 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 \end{bmatrix}$, $\mu(A) = 1/\sqrt{2}$, $\text{spark}(A) = 3$, and the guaranteed uniform sparsity is only $s = 1$.

- **Matrix completion convexity.** The feasible set $\{X : X_{ij} = M_{ij} \text{ for } (i, j) \in \Omega\}$ is affine, hence convex; $\|X\|_*$ is a norm, hence convex. Therefore nuclear-norm matrix completion is convex.

Problem dictionary

Problem	Model	Hard prior	Convex / practical tool
Sparse recovery	$\mathbf{y} = \mathbf{A}\mathbf{x}, m \ll N$	$\ \mathbf{x}\ _0 \leq s$	BP/BPDN/LASSO: ℓ_1 minimization
Noisy / compressible recovery	$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}, \ \mathbf{e}\ _2 \leq \eta$	\mathbf{x} close to Σ_s	$\ \mathbf{x} - \mathbf{x}^\#\ _2 \lesssim \sigma_s(\mathbf{x})_1/\sqrt{s} + \eta$
MMV / joint sparsity	$\mathbf{Y} = \mathbf{A}\mathbf{X}$, same row support across columns	\mathbf{X} has at most s nonzero rows	$\ell_{1,2}$ minimization, SOMP, SSP/CoSaMP lifts
Matrix completion	$\mathbf{Y} = P_\Omega(\mathbf{X})$, observe entries only on Ω	$\text{rank}(\mathbf{X}) \leq r$	Nuclear norm $\ \mathbf{X}\ _*$, SVT/ALM/APG
Robust PCA	$\mathbf{Y} = \mathbf{L} + \mathbf{S}$	\mathbf{L} low-rank, \mathbf{S} sparse	$\min \ \mathbf{L}\ _* + \lambda \ \mathbf{S}\ _1$
SRC classification	$\mathbf{y} \approx [\mathbf{A}_1 \dots \mathbf{A}_C]\mathbf{x}$	\mathbf{x} concentrated on one class block	BPDN, residuals r_i , SCI rejection
Dictionary learning	$\mathbf{Y} \approx \mathbf{D}\mathbf{X}$ with \mathbf{D} unknown	columns of \mathbf{X} are sparse	K-means, MOD, K-SVD, online dictionary learning
NMF	$\mathbf{Y} \approx \mathbf{D}\mathbf{X}, \mathbf{D}, \mathbf{X} \geq 0$	additive parts / cluster memberships	Lee–Seung multiplicative updates, NNLS
Deep sparse coding	stacked sparse-coding modules	sparse, nonnegative elastic-net features	FISTA inference + supervised back-propagation

Theorem lookup

- **Exact sparse uniqueness.** Every s -sparse vector is the unique s -sparse solution iff $\ker(\mathbf{A}) \cap \Sigma_{2s} = \{\mathbf{0}\}$ iff $\text{spark}(\mathbf{A}) > 2s$ (Theorem 2.1). This also implies the information lower bound $m \geq 2s$ for uniform P_0 recovery.
- **Basis pursuit exact recovery.** BP recovers every s -sparse vector iff \mathbf{A} satisfies the NSP of order s (Theorem 5.2). In proofs, take any competing feasible \mathbf{z} , set $\mathbf{v} = \mathbf{x} - \mathbf{z} \in \ker \mathbf{A}$, and compare $\|\mathbf{x} + \mathbf{v}\|_1$ to $\|\mathbf{x}\|_1$ on $S = \text{supp}(\mathbf{x})$.
- **Approximate and noisy recovery.** Stable NSP gives $\|\mathbf{x} - \mathbf{x}^\#\|_1 \lesssim \sigma_s(\mathbf{x})_1$; robust NSP gives BPDN stability under noise (Theorems 5.5 and 5.7).

- **RIP route.** RIP of order $2s$ implies robust NSP and BPDN recovery. Use $\delta_{2s} < 4/\sqrt{41}$ as the sharp listed condition, or $\delta_{2s} < \sqrt{2} - 1$ as the simpler textbook sufficient condition (Theorem 7.3).
- **Coherence route.** For normalized columns, $\mu(A) = \max_{i \neq j} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|$ and $\mu(A) \geq \sqrt{(N-m)/(m(N-1))}$ (Welch, Theorem 6.3). The simple sufficient condition

$$\mu(A) < \frac{1}{2s-1} \iff s < \frac{1}{2}(1 + 1/\mu(A))$$

gives OMP and BP recovery through the cumulative coherence theorem (Theorem 6.4), but it usually costs $m \gtrsim s^2$.

- **Random matrices.** Subgaussian matrices satisfy RIP with high probability once

$$m \gtrsim \delta^{-2} s \log(eN/s)$$

(Theorem 8.2). This is the optimal compressive-sensing scaling up to constants.

- **Matrix analogues.** Replace support size by rank and ℓ_1 by nuclear norm. Rank NSP characterizes nuclear-norm recovery; Candès–Recht gives exact matrix completion under incoherence and random sampling; Candès–Li–Ma–Wright gives exact RPCA when rank and sparse corruption are both small.

Algorithm selection

Use case	Algorithm	Exam facts to remember
Small exact sparse recovery	Brute-force P_0 over supports	Try supports in increasing size; for each support solve $A_s^\dagger y$ and check residual. Cost $\binom{N}{s}$, so this is only practical for small N, s .
Fast sparse recovery with known s	OMP / MP / IHT / HTP / CoSaMP	OMP adds one atom and refits by LS; IHT thresholds a gradient step; HTP adds an LS refit; CoSaMP identifies $2s$, merges, refits, then prunes.
Best convex sparse recovery certificate	BP/BPDN/LASSO	Use BP for noiseless equality, BPDN for bounded noise, LASSO for penalized least squares. Soft thresholding is the prox of $\lambda \ \cdot \ _1$.
Large ℓ_1 problems	ISTA/FISTA, ADMM, IRLS, homotopy	FISTA has $O(1/k^2)$ objective convergence; ADMM separates the constraint; IRLS solves weighted least squares and updates weights.
Shared support across tasks	SOMP or $\ell_{1,2}$ minimization	MMV uniqueness improves to $s < (\text{spark}(A) - 1 + \text{rank}(Y))/2$ when $\text{rank}(Y) > 1$.
Low-rank missing data	SVT / nuclear-norm minimization	The prox of $\tau \ \cdot \ _*$ is singular-value shrinkage: $U \text{diag}((\sigma_i - \tau)_+) V^*$.
Low-rank plus sparse outliers	ALM / ADMM for RPCA	Alternate singular-value shrinkage for L and entrywise soft-thresholding for S .
Learned sparse representation	K-SVD / MOD	MOD updates $D = YX^\dagger$; K-SVD updates one atom by the top SVD of the restricted residual, preserving the sparsity pattern.
Parts-based nonnegative features	NMF	Multiplicative updates preserve nonnegativity and monotonically decrease $\ Y - DX\ _F^2$, but the problem is nonconvex and nonunique.

Common exam moves

1. **Difference of two feasible solutions.** If $Ax = Az$, set $\mathbf{v} = \mathbf{x} - \mathbf{z} \in \ker A$. If both are s -sparse, then $\mathbf{v} \in \Sigma_{2s}$.
2. **Top- s support split.** For stable/robust proofs, let S be the indices of the largest s entries, then split S^c into blocks S_1, S_2, \dots of size s . The recurring bound is $\sum_{j \geq 2} \|\mathbf{v}_{S_j}\|_2 \leq \|\mathbf{v}_{S^c}\|_1 / \sqrt{s}$.
3. **Coherence proof pattern.** Bound on-support correlations from below by $1 - \mu_1(s - 1)$ and off-support correlations from above by $\mu_1(s)$. Correct selection follows when $\mu_1(s) + \mu_1(s - 1) < 1$.
4. **RIP proof pattern.** Use RIP on a union support and restricted orthogonality for disjoint supports. The union size determines the RIP order: $2s$ for BP, $3s$ for IHT/HTP, $4s$ for CoSaMP.
5. **Matrix proof pattern.** Replace entry magnitudes by singular values. Rank is $\|\sigma(X)\|_0$, nuclear norm is $\|\sigma(X)\|_1$, and Frobenius norm is $\|\sigma(X)\|_2$.
6. **Implementation checks.** Normalize dictionary columns before coherence, OMP, SRC, and dictionary learning. In SRC, split train/test per class before stacking; in matrix completion, missing entire rows or columns makes recovery impossible; in K-SVD, restrict the SVD update to examples that actually use the atom.

Part I

Foundations: Sparsity and the Inverse Problem

1 An Invitation to Compressive Sensing

1.1 The central question

In many problems of science and technology—especially in signal and image processing—one observes data $\mathbf{y} \in \mathbb{C}^m$ produced by a linear measurement process,

$$\mathbf{A}\mathbf{x} = \mathbf{y}, \quad \mathbf{A} \in \mathbb{C}^{m \times N}, \quad \mathbf{x} \in \mathbb{C}^N. \quad (1)$$

Classical wisdom (the Shannon–Nyquist sampling theorem) tells us that the number of measurements m must equal or exceed the signal length N . In compressive sensing, the assumption $m < N$ makes (1) *underdetermined*: it has either no solution or infinitely many. Yet under the additional prior that \mathbf{x} is *sparse*, exact recovery is possible from $m \sim s \log(N/s)$ measurements, provided \mathbf{A} is suitably designed.

The compressive sensing programme thus identifies two coupled problems:

- (P1) *Matrix design*. Which matrices $\mathbf{A} \in \mathbb{C}^{m \times N}$ enable robust sparse recovery, and at what minimal m ?
- (P2) *Recovery algorithm*. Which algorithm reconstructs \mathbf{x} from $\mathbf{y} = \mathbf{A}\mathbf{x}$ efficiently?

1.2 Sparsity and compressibility

Definition 1.1 (Sparsity, support)

Let $\mathbf{x} \in \mathbb{C}^N$. Its *support* is $\text{supp}(\mathbf{x}) := \{j \in [N] : x_j \neq 0\}$, and

$$\|\mathbf{x}\|_0 := |\text{supp}(\mathbf{x})|$$

is the *sparsity*. Note that $\|\cdot\|_0$ is not a norm (no homogeneity). We call \mathbf{x} *s-sparse* if $\|\mathbf{x}\|_0 \leq s$, and write $\Sigma_s = \{\mathbf{x} \in \mathbb{C}^N : \|\mathbf{x}\|_0 \leq s\}$.

Definition 1.2 (Best s -term approximation)

For $0 < p \leq \infty$, the ℓ_p -error of best s -term approximation is

$$\sigma_s(\mathbf{x})_p := \inf_{\mathbf{z} \in \Sigma_s} \|\mathbf{x} - \mathbf{z}\|_p.$$

The infimum is achieved by the vector \mathbf{x}_s obtained from \mathbf{x} by retaining only the s entries of largest absolute value (and is the same for all p).

Definition 1.3 (Compressibility and weak ℓ_p space)

A vector \mathbf{x} is *compressible* if $\sigma_s(\mathbf{x})_p$ decays quickly in s . The weak ℓ_p -quasinorm is

$$\|\mathbf{x}\|_{p,\infty} := \max_{k \in [N]} k^{1/p} x_k^*$$

where $\mathbf{x}^* \in \mathbb{R}_+^N$ is the non-increasing rearrangement of $|\mathbf{x}|$. The set $\{\mathbf{x} : \|\mathbf{x}\|_p \leq 1\}$ is a non-convex ball for $p < 1$ and serves as a model for compressible signals.

Theorem 1.4 (Stechkin's inequality)

For any $0 < p < q \leq \infty$ and any $\mathbf{x} \in \mathbb{C}^N$,

$$\sigma_s(\mathbf{x})_q \leq \frac{c_{p,q}}{s^{1/p-1/q}} \|\mathbf{x}\|_p, \quad c_{p,q} := \left[\left(\frac{p}{q} \right)^{p/q} \left(1 - \frac{p}{q} \right)^{1-p/q} \right]^{1/p} \leq 1.$$

The choice $p = 1, q = 2$ yields the most-used special case

$$\sigma_s(\mathbf{x})_2 \leq \frac{1}{2\sqrt{s}} \|\mathbf{x}\|_1.$$

Proof. Let \mathbf{x}^* be the non-increasing rearrangement of $|\mathbf{x}|$. The error of best s -term approximation is $\sigma_s(\mathbf{x})_q = (\sum_{j=s+1}^N (x_j^*)^q)^{1/q}$. We wish to maximize this quantity subject to the constraint $\|\mathbf{x}\|_p^p = \sum_{j=1}^N (x_j^*)^p = L$. Let $\alpha_j = (x_j^*)^p$. The problem is to maximize $f(\boldsymbol{\alpha}) = \sum_{j=s+1}^N \alpha_j^{q/p}$ subject to $\sum_{j=1}^N \alpha_j = L$ and $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_N \geq 0$.

Since $q/p > 1$, the function f is convex. The maximum of a convex function over a compact convex set is attained at an extreme point. The extreme points of the constraint set are of the form $\boldsymbol{\alpha}^{(k)} = (L/k, \dots, L/k, 0, \dots, 0)$ for $k \in \{1, \dots, N\}$.

- If $k \leq s$, then $\alpha_j = 0$ for all $j > s$, so $f(\boldsymbol{\alpha}^{(k)}) = 0$.
- If $k > s$, then $f(\boldsymbol{\alpha}^{(k)}) = \sum_{j=s+1}^k (L/k)^{q/p} = (k-s)(L/k)^{q/p} = L^{q/p} \frac{k-s}{k^{q/p}}$.

To find the maximum, we consider the function $g(k) = (k-s)k^{-q/p}$ for $k \in [s, N]$. Setting $g'(k) = k^{-q/p} - \frac{q}{p}k^{-q/p-1}(k-s) = 0$ yields $1 - \frac{q}{p}(1-s/k) = 0$, which implies $k = s \frac{q}{q-p}$. Plugging this value back into $g(k)$ gives the constant $c_{p,q}^q$. Specifically, for $p = 1, q = 2$, the maximum is attained near $k = 2s$, yielding the constant $1/(2\sqrt{s})$.

Note: 2.5 for reference

□

1.3 Recovery formulations

The naive sparse-recovery problem is the ℓ_0 -minimization

$$\min_{z \in \mathbb{C}^N} \|z\|_0 \quad \text{s.t.} \quad Az = \mathbf{y}. \quad (\text{P}_0)$$

With noise it becomes

$$\min_z \|z\|_0 \quad \text{s.t.} \quad \|Az - \mathbf{y}\|_2 \leq \eta. \quad (\text{P}_{0,\eta})$$

The convex relaxation, the *basis pursuit*, replaces $\|\cdot\|_0$ by $\|\cdot\|_1$:

$$\min_z \|z\|_1 \quad \text{s.t.} \quad Az = \mathbf{y}. \quad (\text{P}_1)$$

Variants include

$$\text{BPDN} \quad \min \|z\|_1 \quad \text{s.t.} \quad \|Az - \mathbf{y}\|_2 \leq \eta, \quad (2)$$

$$\text{LASSO} \quad \min \frac{1}{2} \|Az - \mathbf{y}\|_2^2 + \lambda \|z\|_1, \quad (3)$$

$$\text{Dantzig selector} \quad \min \|z\|_1 \quad \text{s.t.} \quad \|A^*(Az - \mathbf{y})\|_\infty \leq \tau. \quad (4)$$

P_0 is NP-hard (Theorem 2.7); P_1 is a convex linear program.

1.4 Linear-program reformulation of P_1 (real case)

Introduce $\mathbf{u} \in \mathbb{R}^N$ with $-u_j \leq x_j \leq u_j$. Then

$$\min_{\mathbf{x}, \mathbf{u}} \mathbf{1}^\top \mathbf{u} \quad \text{s.t.} \quad A\mathbf{x} = \mathbf{y}, \quad -\mathbf{u} \leq \mathbf{x} \leq \mathbf{u}$$

is a linear program equivalent to P_1 . Equivalently, splitting $\mathbf{x} = \mathbf{z}^+ - \mathbf{z}^-$ with $\mathbf{z}^\pm \geq 0$:

$$\min_{\mathbf{z}^\pm \geq 0} \mathbf{1}^\top (\mathbf{z}^+ + \mathbf{z}^-) \quad \text{s.t.} \quad [A \mid -A] \begin{bmatrix} \mathbf{z}^+ \\ \mathbf{z}^- \end{bmatrix} = \mathbf{y}.$$

1.5 Motivating applications

Single-pixel camera (Rice).

A digital micromirror device steers light from random pixel patterns onto a single photodiode. Each measurement $y_\ell = \langle \mathbf{a}_\ell, \mathbf{z} \rangle$ uses a random ± 1 pattern. If $\mathbf{z} = W\mathbf{x}$ with W a wavelet basis and \mathbf{x} sparse, the system becomes $\mathbf{y} = AW\mathbf{x} = A'\mathbf{x}$, a standard CS problem.

Magnetic resonance imaging (MRI).

The continuous measurement is $f(t) = \int_{\mathbb{R}^3} |X(\mathbf{z})| e^{-2\pi i \mathbf{k}(t) \cdot \mathbf{z}} d\mathbf{z}$, the Fourier transform of the magnetization sampled along trajectory $\mathbf{k}(t)$. Discretizing yields $\mathbf{y} = R_K \mathcal{F} \mathbf{x}$, a partial-Fourier system. Sparsity in wavelets or total variation enables sub-Nyquist scans.

Radar.

The reflected signal from s scatterers takes the form $\mathbf{y} = \sum_{(k,\ell)} x_{k\ell} T_k M_\ell \mathbf{g}$ with T_k time-shifts and M_ℓ frequency modulations of a known pulse \mathbf{g} . Few targets \Rightarrow sparse \mathbf{x} .

Robust PCA.

A data matrix M that is the sum of low-rank L and sparse S admits a convex decomposition $\min \|L\|_* + \lambda \|S\|_1$ s.t. $M = L + S$ (Candès–Li–Ma–Wright 2011). Removes “snow on a campus photo,” separates background from foreground.

Sparse face recognition.

Stack training faces of class i as columns of A_i and concatenate $A = [A_1 | \dots | A_C]$. A novel face \mathbf{y} has a near-block-sparse representation $\mathbf{y} \approx A\mathbf{x}$ that simultaneously identifies the class and the contributing training samples (Lecture 8).

1.6 Founding papers

The CS programme was launched by:

- E. Candès, J. Romberg, T. Tao (2006) – introduced the term *compressed sensing*, proved the foundational ℓ_1 recovery results, and identified the RIP.
- D. Donoho (2006) – coined *compressive sensing* and developed the Gelfand-width and combinatorial-geometry perspectives.

Antecedents include Prony (1795, sparse Fourier), Tibshirani (1996, LASSO), and Mallat–Zhang (1993, matching pursuit).

2 Sparse Solutions of Underdetermined Systems

2.1 When does a sparse solution exist—and is it unique?

We examine *uniform* recovery: which matrices A allow recovery of every s -sparse vector? The answer is a beautiful characterization in linear-algebra language.

Theorem 2.1 (Equivalent uniqueness conditions)

For $A \in \mathbb{C}^{m \times N}$, the following are equivalent:

1. Every s -sparse \mathbf{x} is the unique s -sparse solution of $Az = A\mathbf{x}$.
2. $\ker(A) \cap \Sigma_{2s} = \{\mathbf{0}\}$.
3. For every $S \subseteq [N]$ with $|S| \leq 2s$, the submatrix A_S is injective.
4. Every set of $2s$ columns of A is linearly independent.

Proof. (2) \Leftrightarrow (3) \Leftrightarrow (4) are standard results in linear algebra: A_S is injective iff its columns are linearly independent, and $\ker(A) \cap \Sigma_{2s} = \{\mathbf{0}\}$ is simply the statement that no non-trivial linear combination of $2s$ columns can equal zero.

(1) \implies (2): We prove the contrapositive. Suppose there exists $\mathbf{v} \in \ker(A) \cap \Sigma_{2s}$ such that $\mathbf{v} \neq \mathbf{0}$. We can decompose \mathbf{v} as $\mathbf{v} = \mathbf{x} - \mathbf{x}'$, where \mathbf{x} and \mathbf{x}' are s -sparse vectors with disjoint supports (e.g., let \mathbf{x} take the first s non-zero entries of \mathbf{v} and \mathbf{x}' take the negative of the remaining entries). Since $A\mathbf{v} = \mathbf{0}$, we have $A\mathbf{x} = A\mathbf{x}'$. Thus \mathbf{x} is an s -sparse solution to $Az = A\mathbf{x}$, but it is not unique because \mathbf{x}' is another s -sparse solution. This contradicts (1).

(2) \implies (1): Suppose $\ker(A) \cap \Sigma_{2s} = \{\mathbf{0}\}$. Let $\mathbf{x} \in \Sigma_s$ and suppose there exists $\mathbf{x}' \in \Sigma_s$ such that $A\mathbf{x} = A\mathbf{x}'$. Let $\mathbf{v} = \mathbf{x} - \mathbf{x}'$. Then $A\mathbf{v} = \mathbf{0}$. Furthermore, $\|\mathbf{v}\|_0 = \|\mathbf{x} - \mathbf{x}'\|_0 \leq \|\mathbf{x}\|_0 + \|\mathbf{x}'\|_0 \leq 2s$. Thus $\mathbf{v} \in \ker(A) \cap \Sigma_{2s}$. By (2), $\mathbf{v} = \mathbf{0}$, which implies $\mathbf{x} = \mathbf{x}'$. Hence \mathbf{x} is the unique s -sparse solution.

Note: 2.13 for reference

□

Corollary 2.2

Uniform recovery of all s -sparse signals requires $m \geq 2s$.

This bound is achievable in principle:

Theorem 2.3 (Vandermonde construction)

For any $N \geq 2s$, the Vandermonde matrix

$$A = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ t_1 & t_2 & \cdots & t_N \\ \vdots & \vdots & \ddots & \vdots \\ t_1^{2s-1} & t_2^{2s-1} & \cdots & t_N^{2s-1} \end{pmatrix} \in \mathbb{C}^{2s \times N}$$

with distinct nodes t_1, \dots, t_N allows recovery of every s -sparse signal via P_0 . Its determinant on any $2s$ -subset of columns equals $\prod_{k < \ell} (t_{j_\ell} - t_{j_k}) \neq 0$, so condition (3) of Theorem 2.1 holds.

Note: 2.14 for reference

Theorem 2.4 (Non-uniform recovery requires only $m = s + 1$)

For any $N \geq s + 1$ and any fixed s -sparse \mathbf{x} , almost every $A \in \mathbb{C}^{(s+1) \times N}$ recovers \mathbf{x} from $\mathbf{y} = A\mathbf{x}$ via P_0 .

Proof. Fix an s -sparse vector \mathbf{x} with support S . Recovery fails via P_0 if there exists $\mathbf{z} \in \Sigma_s$ such that $A\mathbf{z} = A\mathbf{x}$ and $\mathbf{z} \neq \mathbf{x}$. This occurs if $A(\mathbf{x} - \mathbf{z}) = \mathbf{0}$ for some $\mathbf{z} \in \Sigma_s \setminus \{\mathbf{x}\}$. Let $\mathbf{v} = \mathbf{x} - \mathbf{z}$. Then \mathbf{v} is supported on $S \cup T$ for some T with $|T| \leq s$. The condition $A\mathbf{x} = A\mathbf{z}$ is equivalent to $A\mathbf{x} \in \text{Range}(A_T)$.

For a fixed T with $|T| \leq s$ and $T \neq S$, the set of matrices A such that $A\mathbf{x} \in \text{Range}(A_T)$ is defined by the vanishing of the determinant of the matrix $[A_T | A\mathbf{x}]$. Since $A\mathbf{x} = \sum_{j \in S} x_j \mathbf{a}_j$, the columns of $[A_T | A\mathbf{x}]$ are linear combinations of the columns of A . The determinant is a non-trivial polynomial in the entries of A (it is not identically zero because one can construct a matrix where $A\mathbf{x}$ is linearly independent of A_T columns, given $m = s + 1$).

The zero set of a non-trivial polynomial in $\mathbb{C}^{(s+1) \times N}$ has Lebesgue measure zero. Since there are only finitely many subsets $T \subseteq [N]$ with $|T| \leq s$, the exceptional set of matrices is a finite union of measure-zero sets, and thus itself has measure zero.

Note: 2.16 for reference

□

2.2 Spark

The number “ $2s$ ” appearing in Theorem 2.1 motivates:

Definition 2.5 (Spark, Donoho–Elad 2003)

$$\text{spark}(A) := \min\{\|\mathbf{z}\|_0 : \mathbf{z} \in \ker A \setminus \{\mathbf{0}\}\}.$$

Theorem 2.1 can be restated: P_0 has a unique solution for every s -sparse \mathbf{x} iff $\text{spark}(A) > 2s$. Spark

ranges in $\{2, \dots, m+1\}$ for full-row-rank A , and equals $m+1$ a.s. for matrices with i.i.d. continuous entries.

2.3 Prony's method: practical recovery from $m = 2s$ Fourier samples

Theorem 2.6 (Prony's method)

For $N \geq 2s$, every s -sparse $\mathbf{x} \in \mathbb{C}^N$ can be recovered explicitly from its first $2s$ DFT coefficients $\hat{x}(0), \dots, \hat{x}(2s-1)$.

Proof. Let $S = \{j_1, \dots, j_s\} \subseteq \{0, \dots, N-1\}$ be the support of \mathbf{x} . Define the *annihilating polynomial* $p(z) = \prod_{k=1}^s (z - e^{2\pi i j_k / N})$. The coefficients h_k of $p(z) = \sum_{k=0}^s h_k z^k$ (with $h_s = 1$) satisfy:

$$\sum_{k=0}^s h_k e^{2\pi i j_\ell k / N} = 0 \quad \text{for each } \ell \in \{1, \dots, s\}.$$

The DFT coefficients are $\hat{x}(n) = \sum_{\ell=1}^s x_{j_\ell} e^{-2\pi i j_\ell n / N}$. Consider the convolution of the sequence h with \hat{x} :

$$\sum_{k=0}^s h_k \hat{x}(n-k) = \sum_{k=0}^s h_k \sum_{\ell=1}^s x_{j_\ell} e^{-2\pi i j_\ell (n-k) / N} = \sum_{\ell=1}^s x_{j_\ell} e^{-2\pi i j_\ell n / N} \underbrace{\sum_{k=0}^s h_k e^{2\pi i j_\ell k / N}}_{=0} = 0.$$

This gives a system of s linear equations for the s unknown coefficients h_0, \dots, h_{s-1} :

$$\begin{pmatrix} \hat{x}(s-1) & \hat{x}(s-2) & \cdots & \hat{x}(0) \\ \hat{x}(s) & \hat{x}(s-1) & \cdots & \hat{x}(1) \\ \vdots & \vdots & \ddots & \vdots \\ \hat{x}(2s-2) & \hat{x}(2s-3) & \cdots & \hat{x}(s-1) \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \\ \vdots \\ h_s \end{pmatrix} = - \begin{pmatrix} \hat{x}(s) \\ \hat{x}(s+1) \\ \vdots \\ \hat{x}(2s-1) \end{pmatrix}.$$

Once h is found, the roots of $p(z)$ reveal the support S . Finally, \mathbf{x}_S is recovered by solving the $2s \times s$ Vandermonde system $\mathbf{V}_S \mathbf{x}_S = \mathbf{y}$.

Note: 2.15 for reference

□

2.4 Computational complexity

Theorem 2.7 (NP-hardness of P_0 , Natarajan 1995)

For any $\eta \geq 0$, the problem $(P_{0,\eta})$ is NP-hard.

Reduction from Exact Cover by 3-Sets. Given a collection $\{C_i\}_{i=1}^N$ of 3-element subsets of $[m]$, set $(\mathbf{a}_i)_j = 1$ if $j \in C_i$, 0 otherwise; let $A = [\mathbf{a}_1 \mid \dots \mid \mathbf{a}_N]$ and $\mathbf{y} = \mathbf{1} \in \mathbb{R}^m$. Any \mathbf{z} feasible for $\|A\mathbf{z} - \mathbf{y}\|_2 \leq \eta < 1$ has all entries of $A\mathbf{z}$ nonzero; since each column of A has exactly 3 nonzeros, $\|\mathbf{z}\|_0 \geq m/3$. The output \mathbf{x}^*

of $P_{0,\eta}$ has $\|\mathbf{x}^*\|_0 = m/3$ iff $\{C_j\}_{j \in \text{supp}(\mathbf{x}^*)}$ is an exact 3-cover. Thus P_0 would solve EXACT-COVER-3, which is NP-complete.

Note: 2.17 for reference

□

Remark 2.8

NP-hardness pertains to general A and \mathbf{y} ; for the structured A used in practice (Gaussian, partial Fourier, ETFs), tractable algorithms (greedy, ℓ_1 , thresholding) succeed. The trade-off: stronger structural assumptions on A buy polynomial-time recovery.

Part II

Algorithms for Sparse Recovery

3 Sampling, Least Squares, and Greedy Pursuit

3.1 Classical sampling

Theorem 3.1 (Shannon–Nyquist)

A function $f \in L^2(\mathbb{R})$ bandlimited to $[-B, B]$ is exactly reconstructed from samples $\{f(n/(2B))\}_{n \in \mathbb{Z}}$ via

$$f(t) = \sum_{n \in \mathbb{Z}} f(n/(2B)) \operatorname{sinc}(2Bt - n).$$

Sampling below the Nyquist rate $f_s = 2B$ produces aliasing (overlapping spectral copies in the periodic Fourier representation).

Proof. Let $f \in L^2(\mathbb{R})$ have Fourier transform $\hat{f}(\xi) = \int_{\mathbb{R}} f(t)e^{-2\pi i \xi t} dt$, supported on $[-B, B]$. We can represent \hat{f} by its Fourier series on $[-B, B]$. The periodized version is $F(\xi) = \sum_{k \in \mathbb{Z}} \hat{f}(\xi - 2Bk)$, which is $2B$ -periodic. Its Fourier coefficients c_n are:

$$c_n = \frac{1}{2B} \int_{-B}^B F(\xi) e^{-2\pi i n \xi / (2B)} d\xi = \frac{1}{2B} \int_{-B}^B \hat{f}(\xi) e^{-2\pi i n \xi / (2B)} d\xi.$$

By the inverse Fourier transform formula, $f(t) = \int_{-B}^B \hat{f}(\xi) e^{2\pi i \xi t} d\xi$. Comparing this with the expression for c_n , we see that $c_n = \frac{1}{2B} f(-n/(2B))$. Thus $F(\xi) = \sum_{n \in \mathbb{Z}} \frac{1}{2B} f(n/(2B)) e^{2\pi i n \xi / (2B)}$ on $[-B, B]$. Since $\hat{f}(\xi) = F(\xi) \cdot \operatorname{rect}(\xi/(2B))$, taking the inverse Fourier transform:

$$\begin{aligned} f(t) &= \int_{-B}^B \left(\sum_{n \in \mathbb{Z}} \frac{1}{2B} f(n/(2B)) e^{2\pi i n \xi / (2B)} \right) e^{2\pi i \xi t} d\xi \\ &= \sum_{n \in \mathbb{Z}} f(n/(2B)) \int_{-B}^B \frac{1}{2B} e^{2\pi i \xi (t+n/(2B))} d\xi \\ &= \sum_{n \in \mathbb{Z}} f(n/(2B)) \operatorname{sinc}(2Bt + n) = \sum_{n \in \mathbb{Z}} f(n/(2B)) \operatorname{sinc}(2Bt - n). \quad \square \end{aligned}$$

CS replaces the bandlimit assumption by sparsity, enabling sub-Nyquist sampling when the signal is sparse in some basis.

3.2 Least squares geometry

For $A \in \mathbb{C}^{m \times N}$ of full column rank, the unique least-squares solution to $Ax = y$ is

$$x_{ls} = \arg \min_x \|Ax - y\|_2^2 = (A^*A)^{-1}A^*y = A^\dagger y,$$

characterized by the normal equations $A^*Ax_{ls} = A^*y$ and the orthogonality $A^*(y - Ax_{ls}) = 0$.

For underdetermined A ($m < N$, full row rank) the least-norm solution is $x_2 = A^\dagger y = A^*(AA^*)^{-1}y$. From the SVD $A = U\Sigma V^*$, $A^\dagger = V\Sigma^\dagger U^*$ where Σ^\dagger inverts the non-zero singular values (Appendix A).

3.3 Greedy pursuit algorithms

The combinatorial P_0 is approximated greedily: select one column (*atom*) per iteration that best correlates with the residual.

Algorithm 1 Matching Pursuit (MP)

Require: $A \in \mathbb{C}^{m \times N}$ with ℓ_2 -normalized columns, $y \in \mathbb{C}^m$, sparsity s

- 1: $r_0 \leftarrow y, x \leftarrow \mathbf{0}$
 - 2: **for** $n = 0, \dots, s - 1$ **do**
 - 3: $j_{n+1} \leftarrow \arg \max_j |\langle a_j, r_n \rangle|$
 - 4: $c \leftarrow \langle a_{j_{n+1}}, r_n \rangle$
 - 5: $x_{j_{n+1}} \leftarrow x_{j_{n+1}} + c; \quad r_{n+1} \leftarrow r_n - c a_{j_{n+1}}$
 - 6: **end for**
 - 7: **return** x
-

Algorithm 2 Orthogonal Matching Pursuit (OMP)

Require: $A \in \mathbb{C}^{m \times N}$, $y \in \mathbb{C}^m$, sparsity s

- 1: $S_0 \leftarrow \emptyset, x^0 \leftarrow \mathbf{0}$
 - 2: **for** $n = 0, \dots, s - 1$ **do**
 - 3: $j_{n+1} \leftarrow \arg \max_{j \notin S_n} |\langle a_j, y - Ax^n \rangle|$
 - 4: $S_{n+1} \leftarrow S_n \cup \{j_{n+1}\}$
 - 5: $x^{n+1} \leftarrow \arg \min_{\text{supp}(z) \subseteq S_{n+1}} \|y - Az\|_2$ ▷ LS re-fit
 - 6: **end for**
 - 7: **return** x^s
-

Algorithm 3 Weak Matching Pursuit (WMP, parameter $\rho \in (0, 1]$)

1: Replace OMP's selection by any j_{n+1} with $|\langle \mathbf{a}_{j_{n+1}}, \mathbf{r}_n \rangle| \geq \rho \max_j |\langle \mathbf{a}_j, \mathbf{r}_n \rangle|$

The crucial difference between MP and OMP: OMP re-projects orthogonally after every selection, so $A_{S_n}^* \mathbf{r}_n = 0$ and no atom is ever selected twice.

3.4 Thresholding-based algorithms

Definition 3.2 (Hard thresholding operator)

For $\mathbf{z} \in \mathbb{C}^N$, $H_s(\mathbf{z})$ is the vector that retains the s largest-magnitude entries of \mathbf{z} and zeros the rest. The index set of these s largest entries is denoted $L_s(\mathbf{z})$.

Algorithm 4 Basic Thresholding (BT)

1: $S^* \leftarrow L_s(A^* \mathbf{y})$
 2: $\mathbf{x} \leftarrow \arg \min_{\text{supp}(\mathbf{z}) \subseteq S^*} \|\mathbf{y} - A\mathbf{z}\|_2$

Algorithm 5 Iterative Hard Thresholding (IHT)

1: $\mathbf{x}^0 \leftarrow \mathbf{0}$
 2: **for** $n = 0, 1, \dots$ **do**
 3: $\mathbf{x}^{n+1} \leftarrow H_s(\mathbf{x}^n + A^*(\mathbf{y} - A\mathbf{x}^n))$
 4: **end for**

Algorithm 6 Hard Thresholding Pursuit (HTP)

1: $\mathbf{x}^0 \leftarrow \mathbf{0}$
 2: **for** $n = 0, 1, \dots$ **do**
 3: $S^{n+1} \leftarrow L_s(\mathbf{x}^n + A^*(\mathbf{y} - A\mathbf{x}^n))$
 4: $\mathbf{x}^{n+1} \leftarrow \arg \min_{\text{supp}(\mathbf{z}) \subseteq S^{n+1}} \|\mathbf{y} - A\mathbf{z}\|_2$
 5: **end for**

Algorithm 7 CoSaMP (Compressive Sampling Matching Pursuit)

```

1:  $\mathbf{x}^0 \leftarrow \mathbf{0}$ 
2: for  $n = 0, 1, \dots$  do
3:    $U^{n+1} \leftarrow \text{supp}(\mathbf{x}^n) \cup L_{2s}(A^*(\mathbf{y} - A\mathbf{x}^n))$ 
4:    $\mathbf{u}^{n+1} \leftarrow \arg \min_{\text{supp}(z) \subseteq U^{n+1}} \|\mathbf{y} - Az\|_2$ 
5:    $\mathbf{x}^{n+1} \leftarrow H_s(\mathbf{u}^{n+1})$ 
6: end for

```

Algorithm 8 Subspace Pursuit (SP)

```

1:  $S_0 \leftarrow \emptyset, \mathbf{x}^0 \leftarrow \mathbf{0}$ 
2: for  $n = 0, 1, \dots$  do
3:    $U^{n+1} \leftarrow S_n \cup L_s(A^*(\mathbf{y} - A\mathbf{x}^n))$ 
4:    $\mathbf{u}^{n+1} \leftarrow \arg \min_{\text{supp}(z) \subseteq U^{n+1}} \|\mathbf{y} - Az\|_2$ 
5:    $S_{n+1} \leftarrow L_s(\mathbf{u}^{n+1}); \mathbf{x}^{n+1} \leftarrow \arg \min_{\text{supp}(z) \subseteq S_{n+1}} \|\mathbf{y} - Az\|_2$ 
6: end for

```

3.5 Phase transitions

Empirical performance of these algorithms is summarized by phase-transition curves in the (ρ, δ) -plane, $\rho = s/m$, $\delta = m/N$. Below each algorithm's curve, recovery succeeds with high probability; above, it fails. The Donoho–Tanner curve gives the threshold for ℓ_1 -minimization in the asymptotic Gaussian setting, and is the gold standard.

4 ℓ_0 Minimization and Greedy Performance

4.1 Why ℓ_0 is hard

By Theorem 2.7, P_0 is NP-hard. Exhaustive search over supports of size s costs $\binom{N}{s}$, prohibitive even for moderate N, s . Greedy and convex-relaxation methods bypass this combinatorial barrier under structural assumptions on A .

4.2 Complexity comparison

Per iteration:

- MP: $O(mN)$ correlation, $O(m)$ residual update; total $O(smN)$.

- OMP: $O(mN)$ correlation plus an LS re-fit. Maintaining a Cholesky factor of $A_{S_n}^* A_{S_n}$ allows incremental update in $O(n^2)$, giving total $O(smN + s^3)$.
- IHT: $O(mN)$ matrix-vector product plus $O(N)$ thresholding per iteration; typically $O(\log(1/\epsilon))$ iterations.
- CoSaMP, SP: $O(mN + s^3)$ per iteration with linear convergence.

4.3 Performance theorems

Theorem 4.1 (OMP under coherence, see Lecture 5 / Section 6.3)

If the coherence $\mu(A) < \frac{1}{2s-1}$ and \mathbf{x} is s -sparse, OMP recovers \mathbf{x} exactly in s iterations.

Theorem 4.2 (Geometric residual decay)

Under the same hypothesis, the OMP residuals satisfy $\|\mathbf{r}_{n+1}\|_2^2 \leq (1 - c)\|\mathbf{r}_n\|_2^2$ for a constant $c = c(\mu, s) > 0$.

Proof. Let S be the true support and S_n be the support after n iterations. The residual \mathbf{r}_n is orthogonal to A_{S_n} and lies in the span of A_S . Specifically, $\mathbf{r}_n = A_S(\mathbf{x}_S - \mathbf{x}_S^n) = \sum_{j \in S \setminus S_n} \tilde{x}_j \mathbf{a}_j$. The norm of the residual satisfies:

$$\|\mathbf{r}_n\|_2^2 = \langle \mathbf{r}_n, \sum_{j \in S \setminus S_n} \tilde{x}_j \mathbf{a}_j \rangle = \sum_{j \in S \setminus S_n} \tilde{x}_j \langle \mathbf{a}_j, \mathbf{r}_n \rangle \leq \|\tilde{\mathbf{x}}\|_1 \max_{j \in S \setminus S_n} |\langle \mathbf{a}_j, \mathbf{r}_n \rangle|.$$

Using the bound $\|\tilde{\mathbf{x}}\|_1 \leq \sqrt{s - |S_n|} \|\tilde{\mathbf{x}}\|_2$ and the fact that $\|\mathbf{r}_n\|_2^2 \geq (1 - \mu_1(s - 1)) \|\tilde{\mathbf{x}}\|_2^2$ (from Gerschgorin), we have $\|\tilde{\mathbf{x}}\|_2 \leq \|\mathbf{r}_n\|_2 / \sqrt{1 - \mu_1(s - 1)}$. Substituting:

$$\|\mathbf{r}_n\|_2^2 \leq \frac{\sqrt{s - |S_n|}}{\sqrt{1 - \mu_1(s - 1)}} \|\mathbf{r}_n\|_2 \max_{j \in S} |\langle \mathbf{a}_j, \mathbf{r}_n \rangle|.$$

This gives $\max_{j \in S} |\langle \mathbf{a}_j, \mathbf{r}_n \rangle| \geq \frac{\sqrt{1 - \mu_1(s - 1)}}{\sqrt{s - |S_n|}} \|\mathbf{r}_n\|_2$. Since OMP picks $j_{n+1} = \arg \max_j |\langle \mathbf{a}_j, \mathbf{r}_n \rangle|$, and $\|\mathbf{r}_{n+1}\|_2^2 \leq \|\mathbf{r}_n - \langle \mathbf{a}_{j_{n+1}}, \mathbf{r}_n \rangle \mathbf{a}_{j_{n+1}}\|_2^2 = \|\mathbf{r}_n\|_2^2 - |\langle \mathbf{a}_{j_{n+1}}, \mathbf{r}_n \rangle|^2$:

$$\|\mathbf{r}_{n+1}\|_2^2 \leq \|\mathbf{r}_n\|_2^2 \left(1 - \frac{1 - \mu_1(s - 1)}{s - |S_n|} \right).$$

For μ small enough, this yields the geometric decay with $c = (1 - \mu_1(s - 1))/s$. □

Theorem 4.3 (CoSaMP near-optimal recovery, Needell–Tropp 2009)

If $\delta_{4s}(A) \leq 0.1$, then CoSaMP applied to $\mathbf{y} = A\mathbf{x} + \mathbf{e}$ satisfies, for every $\mathbf{x} \in \mathbb{C}^N$ and noise $\|\mathbf{e}\|_2 \leq \eta$,

$$\|\mathbf{x} - \mathbf{x}^n\|_2 \leq 2^{-n} \|\mathbf{x}\|_2 + C \left(\frac{\sigma_s(\mathbf{x})_1}{\sqrt{s}} + \eta \right).$$

Proof. Let S be the support of the best s -term approximation of \mathbf{x} . CoSaMP iteratively updates the support Ω^{n+1} by combining the previous support with the $2s$ largest components of the proxy signal $A^*(\mathbf{y} - A\mathbf{x}^n)$. The union support $T = S \cup \Omega^{n+1} \cup \text{supp}(\mathbf{x}^n)$ has size at most $4s$.

Restrict analysis to T . Let $\mathbf{v} = \mathbf{x} - \mathbf{x}^n$. The proxy step identifies a set J of $2s$ atoms. Using the RIP of order $4s$, one can show that the energy of \mathbf{v} on $S \setminus \Omega^{n+1}$ is small:

$$\|\mathbf{x} - \mathbf{u}^{n+1}\|_2 \leq \text{coeff} \cdot \delta_{4s} \|\mathbf{x} - \mathbf{x}^n\|_2 + \text{coeff} \cdot \|\mathbf{e}\|_2.$$

The least-squares step on the enlarged support Ω^{n+1} (size $3s$) further reduces the error. Specifically, the solution \mathbf{u}^{n+1} satisfies $\|\mathbf{u}^{n+1} - \mathbf{x}\|_2 \leq \frac{1+\delta_{4s}}{1-\delta_{4s}} \|\mathbf{x}_{T \setminus \Omega^{n+1}}\|_2$. Finally, the pruning step $\mathbf{x}^{n+1} = H_s(\mathbf{u}^{n+1})$ introduces at most a factor of 2 in the error relative to the best s -term approximation of \mathbf{u}^{n+1} . Combining these estimates, and using $\delta_{4s} \leq 0.1$, the error satisfies a contraction $\epsilon_{n+1} \leq \frac{1}{2}\epsilon_n + C\eta$. Summing the geometric series yields the result. \square

4.4 Stopping rules

For known s , run s iterations. Otherwise stop when $\|\mathbf{r}_n\|_2 < \eta$ (matches noise level) or $|\|\mathbf{r}_n\|_2 - \|\mathbf{r}_{n-1}\|_2|$ falls below tolerance.

Part III

Theoretical Recovery Conditions

5 Spark, Null Space Property, and Basis Pursuit

5.1 The Null Space Property

Definition 5.1 (NSP)

A matrix $A \in \mathbb{C}^{m \times N}$ satisfies the NSP relative to $S \subseteq [N]$ if

$$\|v_S\|_1 < \|v_{\bar{S}}\|_1 \quad \forall v \in \ker A \setminus \{0\}.$$

A satisfies the NSP of order s if it satisfies the NSP relative to every $|S| \leq s$.

Note: 4.1 for reference

Equivalent formulations:

$$(i) \quad 2\|v_S\|_1 < \|v\|_1 \quad \forall v \in \ker A \setminus \{0\}, \tag{5}$$

$$(ii) \quad \|v\|_1 < 2\sigma_s(v)_1 \quad \forall v \in \ker A \setminus \{0\}. \tag{6}$$

Theorem 5.2 (NSP characterizes BP exact recovery)

A satisfies the NSP of order s if and only if every s -sparse $x \in \mathbb{C}^N$ is the unique minimizer of P_1 with $y = Ax$.

Proof. (\Leftarrow) Suppose every s -sparse vector is the unique ℓ_1 -minimizer. Take $v \in \ker A \setminus \{0\}$ and any $|S| \leq s$. The vector $x = v_S$ is s -sparse and $Ax = Ax - Av = A(-v_{\bar{S}})$, so $-v_{\bar{S}}$ is feasible. By unique optimality of $x = v_S$, $\|v_S\|_1 < \|v_{\bar{S}}\|_1$.

(\Rightarrow) Suppose NSP holds. Let x be s -sparse with $S = \text{supp}(x)$, and let $z \neq x$ satisfy $Az = Ax$. Then $v = x - z \in \ker A \setminus \{0\}$, so by NSP $\|v_S\|_1 < \|v_{\bar{S}}\|_1$. Hence

$$\begin{aligned} \|x\|_1 &= \|x_S\|_1 = \|v_S + z_S\|_1 \leq \|v_S\|_1 + \|z_S\|_1 \\ &< \|v_{\bar{S}}\|_1 + \|z_S\|_1 = \|z_{\bar{S}}\|_1 + \|z_S\|_1 = \|z\|_1, \end{aligned}$$

since $x_{\bar{S}} = 0$ implies $v_{\bar{S}} = -z_{\bar{S}}$. Thus x is the strict ℓ_1 -minimizer. □

Note: 4.4 for reference

Theorem 5.3 (Real = complex NSP)

For $A \in \mathbb{R}^{m \times N}$, the real and complex versions of the NSP are equivalent.

Proof. (Complex NSP \Rightarrow Real NSP): A real $\mathbf{v} \in \ker A$ is also a complex element of $\ker A$, so the complex NSP applied to it gives the real NSP inequality immediately.

(Real NSP \Rightarrow Complex NSP): Let $\mathbf{v} \in \ker A \setminus \{0\}$ be complex. Write $v_j = r_j e^{i\theta_j}$ with $r_j = |v_j| \geq 0$. For any $\theta \in \mathbb{R}$ the rotated real vector $\mathbf{w}(\theta) = \Re(e^{-i\theta} \mathbf{v})$ has entries $w_j(\theta) = r_j \cos(\theta_j - \theta)$ and satisfies $A\mathbf{w}(\theta) = \mathbf{0}$ (since A is real, the real and imaginary parts of \mathbf{v} separately lie in $\ker A$, so any real linear combination is too). The real NSP applied to $\mathbf{w}(\theta)$ on S gives $\sum_{j \in S} r_j |\cos(\theta_j - \theta)| < \sum_{j \in \bar{S}} r_j |\cos(\theta_j - \theta)|$ whenever $\mathbf{w}(\theta) \neq \mathbf{0}$. Integrate both sides over $\theta \in [0, 2\pi)$: by translation invariance,

$$\int_0^{2\pi} |\cos(\theta_j - \theta)| d\theta = 4 \quad (\text{independent of } \theta_j).$$

So $4 \sum_{j \in S} r_j \leq 4 \sum_{j \in \bar{S}} r_j$, i.e. $\|\mathbf{v}_S\|_1 \leq \|\mathbf{v}_{\bar{S}}\|_1$. Strict inequality follows because the integrand inequality is strict on a set of positive measure (whenever $\mathbf{w}(\theta)$ is non-zero, which holds for almost every θ since $\mathbf{v} \neq \mathbf{0}$).

Note: 4.7 for reference

□

5.2 Stable NSP and approximately sparse signals

Definition 5.4 (Stable NSP)

A satisfies the *stable NSP* of order s with constant $0 < \rho < 1$ if

$$\|\mathbf{v}_S\|_1 \leq \rho \|\mathbf{v}_{\bar{S}}\|_1 \quad \forall \mathbf{v} \in \ker A, |S| \leq s.$$

Note: 4.10 for reference

Theorem 5.5 (Stability of BP)

If A satisfies the stable NSP of order s with constant $\rho < 1$, then for every $\mathbf{x} \in \mathbb{C}^N$ and the BP minimizer $\mathbf{x}^\#$,

$$\|\mathbf{x} - \mathbf{x}^\#\|_1 \leq \frac{2(1 + \rho)}{1 - \rho} \sigma_s(\mathbf{x})_1.$$

Proof. Let S index the s largest entries of $|\mathbf{x}|$, so $\|\mathbf{x}_{\bar{S}}\|_1 = \sigma_s(\mathbf{x})_1$. Set $\mathbf{v} = \mathbf{x}^\# - \mathbf{x} \in \ker A$.

Optimality of $\mathbf{x}^\#$ gives $\|\mathbf{x}^\#\|_1 \leq \|\mathbf{x}\|_1$. Splitting both sides by support S and the triangle inequality,

$$\begin{aligned} \|\mathbf{x}^\#\|_1 &= \|\mathbf{x}_S^\#\|_1 + \|\mathbf{x}_{\bar{S}}^\#\|_1 \geq \|\mathbf{x}_S\|_1 - \|\mathbf{v}_S\|_1 + \|\mathbf{v}_{\bar{S}}\|_1 - \|\mathbf{x}_{\bar{S}}\|_1 \\ &= \|\mathbf{x}\|_1 - 2\|\mathbf{x}_{\bar{S}}\|_1 - \|\mathbf{v}_S\|_1 + \|\mathbf{v}_{\bar{S}}\|_1. \end{aligned}$$

Combining with $\|\mathbf{x}^\#\|_1 \leq \|\mathbf{x}\|_1$ and rearranging,

$$\|\mathbf{v}_{\bar{S}}\|_1 \leq \|\mathbf{v}_S\|_1 + 2\sigma_s(\mathbf{x})_1.$$

Now apply the stable NSP to $\mathbf{v} \in \ker A$ on the set S : $\|\mathbf{v}_S\|_1 \leq \rho\|\mathbf{v}_{\bar{S}}\|_1$. Substituting, $\|\mathbf{v}_{\bar{S}}\|_1 \leq \rho\|\mathbf{v}_{\bar{S}}\|_1 + 2\sigma_s(\mathbf{x})_1$, so $\|\mathbf{v}_{\bar{S}}\|_1 \leq \frac{2}{1-\rho}\sigma_s(\mathbf{x})_1$. Then $\|\mathbf{v}\|_1 = \|\mathbf{v}_S\|_1 + \|\mathbf{v}_{\bar{S}}\|_1 \leq (1+\rho)\|\mathbf{v}_{\bar{S}}\|_1 \leq \frac{2(1+\rho)}{1-\rho}\sigma_s(\mathbf{x})_1$.

Note: 4.11 for reference

□

5.3 Robust NSP and noisy measurements

Definition 5.6 (Robust NSP)

A satisfies the ℓ_q -robust NSP of order s with constants $0 < \rho < 1$, $\tau > 0$ if for every $\mathbf{v} \in \mathbb{C}^N$ and every $|S| \leq s$,

$$\|\mathbf{v}_S\|_q \leq \frac{\rho}{s^{1-1/q}} \|\mathbf{v}_{\bar{S}}\|_1 + \tau \|\mathbf{A}\mathbf{v}\|_2.$$

Note: 4.16/4.20 for reference

Theorem 5.7 (BPDN under robust NSP)

Suppose A satisfies the ℓ_2 -robust NSP of order s with constants $\rho < 1$, $\tau > 0$. For every $\mathbf{x} \in \mathbb{C}^N$, \mathbf{e} with $\|\mathbf{e}\|_2 \leq \eta$, the BPDN minimizer $\mathbf{x}^\#$ of $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}$ satisfies, for $1 \leq p \leq 2$,

$$\|\mathbf{x} - \mathbf{x}^\#\|_p \leq \frac{C}{s^{1-1/p}} \sigma_s(\mathbf{x})_1 + D s^{1/p-1/2} \eta,$$

where C, D depend only on ρ, τ .

Proof. Let $\mathbf{v} = \mathbf{x}^\# - \mathbf{x}$. Since $\mathbf{x}^\#$ is the BPDN minimizer and \mathbf{x} is feasible, $\|\mathbf{x}^\#\|_1 \leq \|\mathbf{x}\|_1$. Let S be the support of the s largest entries of \mathbf{x} . Then $\|\mathbf{x}_S + \mathbf{v}_S\|_1 + \|\mathbf{x}_{\bar{S}} + \mathbf{v}_{\bar{S}}\|_1 \leq \|\mathbf{x}_S\|_1 + \|\mathbf{x}_{\bar{S}}\|_1$. By the triangle inequality, $\|\mathbf{x}_S\|_1 - \|\mathbf{v}_S\|_1 + \|\mathbf{v}_{\bar{S}}\|_1 - \|\mathbf{x}_{\bar{S}}\|_1 \leq \|\mathbf{x}_S\|_1 + \|\mathbf{x}_{\bar{S}}\|_1$, which rearranges to $\|\mathbf{v}_{\bar{S}}\|_1 \leq \|\mathbf{v}_S\|_1 + 2\sigma_s(\mathbf{x})_1$.

Now apply the ℓ_2 -robust NSP to \mathbf{v} :

$$\|\mathbf{v}_S\|_2 \leq \frac{\rho}{\sqrt{s}} \|\mathbf{v}_{\bar{S}}\|_1 + \tau \|\mathbf{A}\mathbf{v}\|_2 \leq \frac{\rho}{\sqrt{s}} (\|\mathbf{v}_S\|_1 + 2\sigma_s(\mathbf{x})_1) + \tau(2\eta),$$

where $\|A\mathbf{v}\|_2 \leq \|A\mathbf{x}^\# - \mathbf{y}\|_2 + \|\mathbf{y} - A\mathbf{x}\|_2 \leq 2\eta$. Using $\|\mathbf{v}_S\|_1 \leq \sqrt{S}\|\mathbf{v}_S\|_2$, we obtain

$$(1 - \rho)\|\mathbf{v}_S\|_2 \leq \frac{2\rho}{\sqrt{S}}\sigma_s(\mathbf{x})_1 + 2\tau\eta.$$

This bounds $\|\mathbf{v}_S\|_2$. For $\rho = 1$, $\|\mathbf{v}\|_1 = \|\mathbf{v}_S\|_1 + \|\mathbf{v}_{\bar{S}}\|_1 \leq 2\|\mathbf{v}_S\|_1 + 2\sigma_s(\mathbf{x})_1 \leq 2\sqrt{S}\|\mathbf{v}_S\|_2 + 2\sigma_s(\mathbf{x})_1$. For $\rho = 2$, $\|\mathbf{v}_{\bar{S}}\|_2 \leq \|\mathbf{v}_{\bar{S}}\|_1/\sqrt{S}$ by the sorting property of approximately sparse vectors. Combining these via Hölder's inequality $\|\mathbf{v}\|_p \leq \|\mathbf{v}\|_1^{2/p-1}\|\mathbf{v}\|_2^{2-2/p}$ gives the result.

Note: 4.22 for reference

□

5.4 Recovery of individual vectors via dual certificates

Theorem 5.8 (Dual certificate for BP)

$\mathbf{x} \in \mathbb{C}^N$ with support S is the unique BP minimizer iff

- A_S is injective, and
- there exists $\mathbf{h} \in \mathbb{C}^m$ such that $(A^*\mathbf{h})_j = \text{sgn}(x_j)$ for $j \in S$ and $|(A^*\mathbf{h})_\ell| < 1$ for $\ell \notin S$.

The vector \mathbf{h} is called a *dual certificate*.

Proof. The optimization problem $\min \|\mathbf{z}\|_1$ s.t. $A\mathbf{z} = A\mathbf{x}$ is convex. A feasible \mathbf{x} is a minimizer if and only if there exists a dual vector $\mathbf{h} \in \mathbb{C}^m$ such that $A^*\mathbf{h}$ lies in the subdifferential $\partial\|\mathbf{x}\|_1$. The subdifferential is the set of all vectors \mathbf{u} such that $u_j = \text{sgn}(x_j)$ for $j \in S$ and $|u_j| \leq 1$ for $j \notin S$.

(\Leftarrow) Suppose such an \mathbf{h} exists with $|(A^*\mathbf{h})_j| < 1$ for $j \notin S$, and A_S is injective. Let \mathbf{z} be another feasible point, $A\mathbf{z} = A\mathbf{x}$, and let $\mathbf{v} = \mathbf{z} - \mathbf{x}$. Then $A\mathbf{v} = \mathbf{0}$ and $\mathbf{v} \neq \mathbf{0}$. Since A_S is injective, \mathbf{v} cannot be supported on S alone, so $\mathbf{v}_{\bar{S}} \neq \mathbf{0}$. Then

$$\begin{aligned} \|\mathbf{z}\|_1 &= \|\mathbf{x} + \mathbf{v}\|_1 \geq \text{Re}\langle A^*\mathbf{h}, \mathbf{x} + \mathbf{v} \rangle + \sum_{j \notin S} (1 - |(A^*\mathbf{h})_j|)|v_j| \\ &= \text{Re}\langle \mathbf{h}, A\mathbf{x} \rangle + \underbrace{\text{Re}\langle \mathbf{h}, A\mathbf{v} \rangle}_{=0} + \sum_{j \notin S} (1 - |(A^*\mathbf{h})_j|)|v_j| \\ &= \|\mathbf{x}\|_1 + \sum_{j \notin S} (1 - |(A^*\mathbf{h})_j|)|v_j|. \end{aligned}$$

Since $|(A^*\mathbf{h})_j| < 1$ and $\mathbf{v}_{\bar{S}} \neq \mathbf{0}$, the last term is strictly positive, so $\|\mathbf{z}\|_1 > \|\mathbf{x}\|_1$. Thus \mathbf{x} is the unique minimizer.

Note: 4.25 for reference

□

Theorem 5.9 (Inexact dual certificate)

Let \mathbf{x} be supported on S , $\|(A_S^* A_S)^{-1}\|_{2 \rightarrow 2} \leq \alpha$ and $\max_{\ell \notin S} \|A_S^* \mathbf{a}_\ell\|_2 \leq \beta$. If there exists $\mathbf{u} = A^* \mathbf{h}$ with $\|\mathbf{u}_S - \text{sgn}(\mathbf{x}_S)\|_2 \leq \gamma$ and $\|\mathbf{u}_{\bar{S}}\|_\infty \leq \theta$, and if $\theta + \alpha\beta\gamma < 1$, then \mathbf{x} is the unique BP minimizer.

Proof. Let $\mathbf{v} \in \ker A \setminus \{0\}$. We show $\|\mathbf{v}_S\|_1 < \|\mathbf{v}_{\bar{S}}\|_1$ by first bounding $\|\mathbf{v}_S\|_2$. From $A\mathbf{v} = \mathbf{0}$, we have $A_S \mathbf{v}_S = -A_{\bar{S}} \mathbf{v}_{\bar{S}}$. Multiplying by A_S^* :

$$A_S^* A_S \mathbf{v}_S = -A_S^* A_{\bar{S}} \mathbf{v}_{\bar{S}} = -\sum_{\ell \notin S} v_\ell A_S^* \mathbf{a}_\ell.$$

Thus $\mathbf{v}_S = -(A_S^* A_S)^{-1} \sum_{\ell \notin S} v_\ell A_S^* \mathbf{a}_\ell$. Taking the ℓ_2 norm:

$$\|\mathbf{v}_S\|_2 \leq \|(A_S^* A_S)^{-1}\|_{2 \rightarrow 2} \sum_{\ell \notin S} |v_\ell| \|A_S^* \mathbf{a}_\ell\|_2 \leq \alpha\beta \|\mathbf{v}_{\bar{S}}\|_1.$$

Now using the dual vector $\mathbf{u} = A^* \mathbf{h}$, we have $\langle \mathbf{u}, \mathbf{v} \rangle = \langle \mathbf{h}, A\mathbf{v} \rangle = 0$. So $\langle \mathbf{u}_S, \mathbf{v}_S \rangle = -\langle \mathbf{u}_{\bar{S}}, \mathbf{v}_{\bar{S}} \rangle$. Then:

$$\begin{aligned} |\langle \text{sgn}(\mathbf{x}_S), \mathbf{v}_S \rangle| &= |\langle \text{sgn}(\mathbf{x}_S) - \mathbf{u}_S, \mathbf{v}_S \rangle + \langle \mathbf{u}_S, \mathbf{v}_S \rangle| \\ &= |\langle \text{sgn}(\mathbf{x}_S) - \mathbf{u}_S, \mathbf{v}_S \rangle - \langle \mathbf{u}_{\bar{S}}, \mathbf{v}_{\bar{S}} \rangle| \\ &\leq \|\text{sgn}(\mathbf{x}_S) - \mathbf{u}_S\|_2 \|\mathbf{v}_S\|_2 + \|\mathbf{u}_{\bar{S}}\|_\infty \|\mathbf{v}_{\bar{S}}\|_1 \\ &\leq \gamma(\alpha\beta \|\mathbf{v}_{\bar{S}}\|_1) + \theta \|\mathbf{v}_{\bar{S}}\|_1 = (\theta + \alpha\beta\gamma) \|\mathbf{v}_{\bar{S}}\|_1. \end{aligned}$$

Since $\theta + \alpha\beta\gamma < 1$, we have $|\langle \text{sgn}(\mathbf{x}_S), \mathbf{v}_S \rangle| < \|\mathbf{v}_{\bar{S}}\|_1$. Since $|\langle \text{sgn}(\mathbf{x}_S), \mathbf{v}_S \rangle| \leq \|\mathbf{v}_S\|_1$, and A_S is injective, this implies the unique recovery condition.

Note: 4.31 for reference

□

5.5 Tangent cone characterization

Theorem 5.10 (Tangent cone)

For $\mathbf{x} \in \mathbb{R}^N$ and $A \in \mathbb{R}^{m \times N}$, define $T(\mathbf{x}) = \text{cone}\{\mathbf{z} - \mathbf{x} : \|\mathbf{z}\|_1 \leq \|\mathbf{x}\|_1\}$. Then \mathbf{x} is the unique BP minimizer iff $\ker A \cap T(\mathbf{x}) = \{0\}$.

Proof. Recall $T(\mathbf{x}) = \text{cone}\{\mathbf{z} - \mathbf{x} : \|\mathbf{z}\|_1 \leq \|\mathbf{x}\|_1\}$ consists of all positive scalar multiples of differences $\mathbf{z} - \mathbf{x}$ where \mathbf{z} has ℓ_1 -norm at most that of \mathbf{x} .

(\Rightarrow) Suppose \mathbf{x} is the unique BP minimizer. If $\mathbf{v} \in \ker A \cap T(\mathbf{x})$ is non-zero, then $\mathbf{v} = \lambda(\mathbf{z} - \mathbf{x})$ for some $\lambda > 0$ and some \mathbf{z} with $\|\mathbf{z}\|_1 \leq \|\mathbf{x}\|_1$. Since $\lambda > 0$ and $\mathbf{v} \neq 0$, $\mathbf{z} \neq \mathbf{x}$. Now $\mathbf{x} + \mathbf{v}/\lambda = \mathbf{z}$ has $A\mathbf{z} = A\mathbf{x} + \lambda^{-1}A\mathbf{v} = A\mathbf{x}$, so \mathbf{z} is feasible. Optimality of \mathbf{x} gives $\|\mathbf{z}\|_1 > \|\mathbf{x}\|_1$ (strict, by uniqueness), contradicting the construction. Hence $\ker A \cap T(\mathbf{x}) = \{0\}$.

(\Leftarrow) Suppose $\ker A \cap T(\mathbf{x}) = \{0\}$. Let $\mathbf{z} \neq \mathbf{x}$ be any feasible point, $A\mathbf{z} = A\mathbf{x} = \mathbf{y}$. Then $\mathbf{v} := \mathbf{z} - \mathbf{x} \in$

$\ker A$ and $\mathbf{v} \neq 0$. By hypothesis $\mathbf{v} \notin T(\mathbf{x})$. The defining characterization of the descent cone says $\mathbf{v} \in T(\mathbf{x})$ iff $\|\mathbf{x} + t\mathbf{v}\|_1 \leq \|\mathbf{x}\|_1$ for some $t > 0$ (equivalently, \mathbf{v} is a feasible descent direction for $\|\cdot\|_1$ at \mathbf{x}). Since $\mathbf{v} \notin T(\mathbf{x})$, $\|\mathbf{x} + t\mathbf{v}\|_1 > \|\mathbf{x}\|_1$ for every $t > 0$; in particular at $t = 1$, $\|\mathbf{z}\|_1 = \|\mathbf{x} + \mathbf{v}\|_1 > \|\mathbf{x}\|_1$. Hence \mathbf{x} is the unique BP minimizer.

Note: 4.34 for reference

□

Theorem 5.11 (Tangent-cone stability)

With $\mathbf{y} = A\mathbf{x} + \mathbf{e}$, $\|\mathbf{e}\|_2 \leq \eta$, if $\inf_{\mathbf{z} \in T(\mathbf{x}), \|\mathbf{z}\|_2=1} \|A\mathbf{z}\|_2 \geq \tau > 0$, then the BPDN minimizer satisfies $\|\mathbf{x} - \mathbf{x}^\# \|_2 \leq 2\eta/\tau$.

Proof. Let $\mathbf{x}^\#$ be the BPDN minimizer. Since \mathbf{x} is feasible for the constraint $\|A\mathbf{z} - \mathbf{y}\|_2 \leq \eta$ (assuming the noise level is correctly specified), we have $\|\mathbf{x}^\#\|_1 \leq \|\mathbf{x}\|_1$. By definition of the tangent cone $T(\mathbf{x})$, the difference $\mathbf{v} = \mathbf{x}^\# - \mathbf{x}$ lies in $T(\mathbf{x})$.

The distance between images satisfies:

$$\|A\mathbf{v}\|_2 = \|A\mathbf{x}^\# - A\mathbf{x}\|_2 = \|(A\mathbf{x}^\# - \mathbf{y}) - (A\mathbf{x} - \mathbf{y})\|_2 \leq \|A\mathbf{x}^\# - \mathbf{y}\|_2 + \|\mathbf{y} - A\mathbf{x}\|_2 \leq 2\eta.$$

By the hypothesis $\inf_{\mathbf{z} \in T(\mathbf{x}), \|\mathbf{z}\|_2=1} \|A\mathbf{z}\|_2 \geq \tau$, it follows that for any $\mathbf{v} \in T(\mathbf{x})$, $\|A\mathbf{v}\|_2 \geq \tau\|\mathbf{v}\|_2$. Combining these inequalities:

$$\tau\|\mathbf{v}\|_2 \leq \|A\mathbf{v}\|_2 \leq 2\eta \implies \|\mathbf{v}\|_2 \leq \frac{2\eta}{\tau}.$$

□

Note: 4.36 for reference

5.6 Low-rank matrix recovery

Definition 5.12 (Nuclear norm)

For $X \in \mathbb{C}^{n_1 \times n_2}$, $\|X\|_* = \sum_j \sigma_j(X)$, the ℓ_1 -norm of the singular value vector.

Theorem 5.13 (Rank NSP)

A linear map $\mathcal{A} : \mathbb{C}^{n_1 \times n_2} \rightarrow \mathbb{C}^m$ recovers every rank- r matrix via nuclear-norm minimization $\min \|Z\|_*$ s.t. $\mathcal{A}(Z) = \mathcal{A}(X)$ iff for every nonzero $M \in \ker \mathcal{A}$ with singular values $\sigma_1 \geq \dots \geq \sigma_n$,

$$\sum_{j=1}^r \sigma_j(M) < \sum_{j=r+1}^n \sigma_j(M).$$

Proof. This is the matrix analogue of Theorem 5.2. Let Σ_r be the set of matrices of rank at most r .

(\implies) Suppose every rank- r matrix is the unique nuclear-norm minimizer. Let $M \in \ker \mathcal{A} \setminus \{0\}$ and let

$M = U\Sigma V^*$ be its SVD. Let $M_1 = U\Sigma_r V^*$ where Σ_r contains only the r largest singular values. Then M_1 has rank at most r . Let $M_2 = M - M_1$. Since $\mathcal{A}(M) = \mathcal{A}(M_1 + M_2) = 0$, we have $\mathcal{A}(M_1) = \mathcal{A}(-M_2)$. By the unique recovery property, $\|M_1\|_* < \| -M_2\|_* = \|M_2\|_*$. Since $\|M_1\|_* = \sum_{j=1}^r \sigma_j(M)$ and $\|M_2\|_* = \sum_{j=r+1}^n \sigma_j(M)$, the condition holds.

(\Leftarrow) Suppose the condition holds. Let X be a rank- r matrix and let Z be another matrix such that $\mathcal{A}(Z) = \mathcal{A}(X)$. Let $M = Z - X \in \ker \mathcal{A}$. By the triangle inequality for singular values (Ky Fan norms) or directly from the properties of the nuclear norm, one can show that if the condition holds, $\|X + M\|_* > \|X\|_*$ for any $M \in \ker \mathcal{A} \setminus \{0\}$. This follows from decomposing M into parts aligned and misaligned with the singular spaces of X , analogue to the support decomposition in the vector case.

Note: 4.40 for reference

□

6 Coherence

6.1 Mutual coherence

Definition 6.1 (Coherence)

For $A \in \mathbb{C}^{m \times N}$ with ℓ_2 -normalized columns,

$$\mu = \mu(A) = \max_{1 \leq i \neq j \leq N} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|.$$

Note: 5.1 for reference

Definition 6.2 (ℓ_1 -coherence function)

For $1 \leq s \leq N - 1$,

$$\mu_1(s) := \max_{i \in [N]} \max_{|S|=s, i \notin S} \sum_{j \in S} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|.$$

Note: 5.2 for reference

Basic bounds: $\mu \leq \mu_1(s) \leq s\mu$, and $\mu_1(s) \leq \mu_1(s - 1) + \mu$.

6.2 The Welch bound

Theorem 6.3 (Welch bound)

For $A \in \mathbb{C}^{m \times N}$ with ℓ_2 -normalized columns ($m < N$),

$$\mu(A) \geq \sqrt{\frac{N-m}{m(N-1)}}.$$

Equality holds iff the columns form an Equiangular Tight Frame (ETF).

Proof. Let $G = A^*A \in \mathbb{C}^{N \times N}$. Then $\text{tr}(G) = \sum_i \|\mathbf{a}_i\|_2^2 = N$. By Cauchy–Schwarz on traces,

$$N^2 = (\text{tr } G)^2 \leq m \text{tr}(G^2),$$

since G has rank $\leq m$. Now

$$\text{tr}(G^2) = \sum_{i,j} |G_{ij}|^2 = \sum_i \|\mathbf{a}_i\|_2^4 + \sum_{i \neq j} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|^2 = N + \sum_{i \neq j} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|^2.$$

Since each off-diagonal term $\leq \mu^2$ and there are $N(N-1)$ of them, $\text{tr}(G^2) \leq N + N(N-1)\mu^2$. Combining,

$$N^2 \leq m(N + N(N-1)\mu^2) \implies \mu^2 \geq \frac{N-m}{m(N-1)}. \quad \square$$

For $m \ll N$ the Welch bound gives $\mu \gtrsim 1/\sqrt{m}$, achieved (up to constants) by:

- **Random partial DFT:** $\mu \approx 1/\sqrt{m}$.
- **Alltop sequences (prime $m \geq 5$):** $\mu = 1/\sqrt{m}$ for an $m \times m^2$ matrix.
- **Gabor frames** from time-frequency shifts of a window function.
- **ETFs:** achieve the Welch bound with equality (rare; require divisibility conditions).

6.3 Coherence-based recovery guarantees

Theorem 6.4 (OMP via coherence)

If $\mu_1(s) + \mu_1(s-1) < 1$, then OMP recovers every s -sparse \mathbf{x} from $\mathbf{y} = A\mathbf{x}$ in exactly s iterations. In particular, the simpler condition $\mu(A) < \frac{1}{2s-1}$ (which implies the above) suffices.

Proof. Let $S = \text{supp}(\mathbf{x})$ with $|S| = s$, and let $S_n \subseteq S$ be the indices selected after n iterations of OMP (we prove correctness by induction on n). At iteration n the residual is $\mathbf{r}_n = A\mathbf{x} - A\mathbf{x}^n$ where \mathbf{x}^n is the LS-fit on S_n . Since OMP's LS step makes \mathbf{r}_n orthogonal to $\{\mathbf{a}_j : j \in S_n\}$, we have $A_{S_n}^* \mathbf{r}_n = \mathbf{0}$. Write $\mathbf{r}_n = A(\mathbf{x} - \mathbf{x}^n) = \sum_{j \in S \setminus S_n} \tilde{x}_j \mathbf{a}_j$ for some coefficients \tilde{x}_j (the residual lives in the column span of the un-selected support).

For $j \in S \setminus S_n$,

$$\langle \mathbf{a}_j, \mathbf{r}_n \rangle = \tilde{x}_j + \sum_{k \in S \setminus S_n, k \neq j} \tilde{x}_k \langle \mathbf{a}_j, \mathbf{a}_k \rangle.$$

Pick $j^* = \arg \max_{j \in S \setminus S_n} |\tilde{x}_j|$. Then $|\langle \mathbf{a}_{j^*}, \mathbf{r}_n \rangle| \geq |\tilde{x}_{j^*}|(1 - \mu_1(s - |S_n| - 1)) \geq |\tilde{x}_{j^*}|(1 - \mu_1(s - 1))$ by the definition of μ_1 (with $|S \setminus S_n| - 1$ summed inner products of magnitude at most $\mu_1(s - 1)$).

For $\ell \notin S$, $|\langle \mathbf{a}_\ell, \mathbf{r}_n \rangle| = |\sum_{k \in S \setminus S_n} \tilde{x}_k \langle \mathbf{a}_\ell, \mathbf{a}_k \rangle| \leq \mu_1(s - |S_n|) \max_k |\tilde{x}_k| \leq \mu_1(s) |\tilde{x}_{j^*}|$.

Hence $\max_{j \in S} |\langle \mathbf{a}_j, \mathbf{r}_n \rangle| \geq (1 - \mu_1(s - 1)) |\tilde{x}_{j^*}| > \mu_1(s) |\tilde{x}_{j^*}| \geq \max_{\ell \notin S} |\langle \mathbf{a}_\ell, \mathbf{r}_n \rangle|$ whenever $\mu_1(s) + \mu_1(s - 1) < 1$. Thus OMP selects from S at iteration $n + 1$. By induction $S_n \subseteq S$ for all n , so after s iterations $S_n = S$ and the LS-fit recovers \mathbf{x} exactly.

Note: 5.14 for reference

□

Theorem 6.5 (BP via coherence)

The same condition $\mu_1(s) + \mu_1(s - 1) < 1$ implies the NSP of order s , hence ℓ_1 -recovery succeeds for every s -sparse signal.

Note: 5.15 for reference

Theorem 6.6 (BT via coherence)

If $\mu_1(s) + \mu_1(s - 1) < \min_{j \in S} |x_j| / \max_{j \in S} |x_j|$, basic thresholding recovers \mathbf{x} exactly.

Proof. Let S be the support of \mathbf{x} . Basic thresholding identifies S correctly if $\min_{j \in S} |(A^* \mathbf{y})_j| > \max_{\ell \notin S} |(A^* \mathbf{y})_\ell|$. Since $\mathbf{y} = \sum_{k \in S} x_k \mathbf{a}_k$:

$$(A^* \mathbf{y})_j = \langle \mathbf{a}_j, \sum_{k \in S} x_k \mathbf{a}_k \rangle = x_j + \sum_{k \in S, k \neq j} x_k \langle \mathbf{a}_j, \mathbf{a}_k \rangle.$$

For $j \in S$:

$$|(A^* \mathbf{y})_j| \geq |x_j| - \sum_{k \in S, k \neq j} |x_k| |\langle \mathbf{a}_j, \mathbf{a}_k \rangle| \geq \min_{k \in S} |x_k| - \mu_1(s - 1) \max_{k \in S} |x_k|.$$

For $\ell \notin S$:

$$|(A^* \mathbf{y})_\ell| = |\sum_{k \in S} x_k \langle \mathbf{a}_\ell, \mathbf{a}_k \rangle| \leq \mu_1(s) \max_{k \in S} |x_k|.$$

Correct identification holds if:

$$\min_{k \in S} |x_k| - \mu_1(s - 1) \max_{k \in S} |x_k| > \mu_1(s) \max_{k \in S} |x_k| \iff \mu_1(s) + \mu_1(s - 1) < \frac{\min_{k \in S} |x_k|}{\max_{k \in S} |x_k|}.$$

Once S is identified, the LS step $\mathbf{x} = A_S^\dagger \mathbf{y}$ recovers the true coefficients exactly.

Note: 5.20 for reference

□

6.4 The famous coherence bound

Combining $\mu_1(s) \leq s\mu$ with the conditions above gives the celebrated

$$s < \frac{1}{2} \left(1 + \frac{1}{\mu(A)} \right) \iff \mu < \frac{1}{2s-1}.$$

Combined with the Welch bound $\mu \gtrsim 1/\sqrt{m}$, this yields the quadratic bottleneck:

$$m \gtrsim s^2.$$

RIP-based theory will improve this to the optimal $m \gtrsim s \log(N/s)$.

7 Restricted Isometry Property

7.1 The RIP

Definition 7.1 (Restricted Isometry Property (RIP))

The s -th restricted isometry constant $\delta_s = \delta_s(A)$ is the smallest $\delta \geq 0$ such that

$$(1 - \delta)\|\mathbf{x}\|_2^2 \leq \|A\mathbf{x}\|_2^2 \leq (1 + \delta)\|\mathbf{x}\|_2^2 \quad \forall \mathbf{x} \in \Sigma_s.$$

Equivalently,

$$\delta_s = \max_{|S| \leq s} \|A_S^* A_S - I\|_{2 \rightarrow 2}.$$

Note: 6.1 for reference

Monotonicity: $\delta_1 \leq \delta_2 \leq \dots \leq \delta_N$. Connection to coherence: $\delta_s \leq (s-1)\mu(A)$.

Proposition 7.2 (Restricted orthogonality)

For disjointly supported $\mathbf{u} \in \Sigma_s, \mathbf{v} \in \Sigma_t$,

$$|\langle A\mathbf{u}, A\mathbf{v} \rangle| \leq \delta_{s+t} \|\mathbf{u}\|_2 \|\mathbf{v}\|_2.$$

Proof. Suppose $\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1$. The vectors $\mathbf{u} \pm \mathbf{v}$ are $(s+t)$ -sparse because their supports are contained in the union of the disjoint supports of \mathbf{u} and \mathbf{v} . Furthermore, $\|\mathbf{u} \pm \mathbf{v}\|_2^2 = \|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2 = 2$. Applying the RIP of order $s+t$:

$$(1 - \delta_{s+t}) \cdot 2 \leq \|A(\mathbf{u} \pm \mathbf{v})\|_2^2 \leq (1 + \delta_{s+t}) \cdot 2.$$

By the parallelogram identity:

$$\begin{aligned} 4\operatorname{Re}\langle A\mathbf{u}, A\mathbf{v} \rangle &= \|A(\mathbf{u} + \mathbf{v})\|_2^2 - \|A(\mathbf{u} - \mathbf{v})\|_2^2 \\ &\leq 2(1 + \delta_{s+t}) - 2(1 - \delta_{s+t}) = 4\delta_{s+t}. \end{aligned}$$

Similarly, $4\operatorname{Re}\langle A\mathbf{u}, A\mathbf{v} \rangle \geq -4\delta_{s+t}$. The same argument applied to $i\mathbf{v}$ bounds the imaginary part. For general norms, rescale by $\|\mathbf{u}\|_2\|\mathbf{v}\|_2$. \square

7.2 RIP \Rightarrow NSP \Rightarrow BP recovery

Theorem 7.3 (BPDN recovery under RIP, Candès 2008 / Foucart 2010)

If $\delta_{2s}(A) < \frac{4}{\sqrt{41}} \approx 0.6246$, then for every $\mathbf{x} \in \mathbb{C}^N$, \mathbf{e} with $\|\mathbf{e}\|_2 \leq \eta$, the BPDN solution satisfies, for $1 \leq \rho \leq 2$,

$$\|\mathbf{x} - \mathbf{x}^\#\|_\rho \leq \frac{C \sigma_s(\mathbf{x})_1}{s^{1-1/\rho}} + D s^{1/\rho-1/2} \eta,$$

with C, D explicit functions of δ_{2s} . The simpler textbook condition $\delta_{2s} < \sqrt{2} - 1 \approx 0.414$ suffices.

Proof. We show that RIP of order $2s$ implies the robust NSP. Let $\mathbf{v} \in \mathbb{C}^N$ and let S be the index set of the s largest entries of \mathbf{v} . Split the remaining indices \bar{S} into blocks S_1, S_2, \dots of size s each, such that S_1 contains the s largest entries of $\mathbf{v}_{\bar{S}}$, S_2 the next s , and so on. By the monotonicity of sorted entries, for $j \geq 1$ and $k \in S_{j+1}$, $|v_k| \leq \|\mathbf{v}_{S_j}\|_1/s$. Summing over $k \in S_{j+1}$ gives $\|\mathbf{v}_{S_{j+1}}\|_2^2 \leq s(\|\mathbf{v}_{S_j}\|_1/s)^2$, so $\|\mathbf{v}_{S_{j+1}}\|_2 \leq \|\mathbf{v}_{S_j}\|_1/\sqrt{s}$. Summing over $j \geq 1$:

$$\sum_{j \geq 2} \|\mathbf{v}_{S_j}\|_2 \leq \frac{1}{\sqrt{s}} \sum_{j \geq 1} \|\mathbf{v}_{S_j}\|_1 = \frac{1}{\sqrt{s}} \|\mathbf{v}_{\bar{S}}\|_1.$$

Let $S_{01} = S \cup S_1$. This set has size $2s$. Applying RIP to $\mathbf{v}_{S_{01}}$:

$$(1 - \delta_{2s})\|\mathbf{v}_{S_{01}}\|_2^2 \leq \|A\mathbf{v}_{S_{01}}\|_2^2 = \langle A\mathbf{v}_{S_{01}}, A\mathbf{v} - \sum_{j \geq 2} A\mathbf{v}_{S_j} \rangle.$$

Using the restricted orthogonality (Proposition 7.1) between S_{01} and each S_j :

$$(1 - \delta_{2s})\|\mathbf{v}_{S_{01}}\|_2^2 \leq \|A\mathbf{v}_{S_{01}}\|_2 \|A\mathbf{v}\|_2 + \sum_{j \geq 2} \delta_{2s} \|\mathbf{v}_{S_{01}}\|_2 \|\mathbf{v}_{S_j}\|_2.$$

Dividing by $\|\mathbf{v}_{S_{01}}\|_2$ and using $\|\mathbf{v}_S\|_2 \leq \|\mathbf{v}_{S_{01}}\|_2$:

$$(1 - \delta_{2s})\|\mathbf{v}_S\|_2 \leq \sqrt{1 + \delta_{2s}} \|A\mathbf{v}\|_2 + \frac{\delta_{2s}}{\sqrt{s}} \|\mathbf{v}_{\bar{S}}\|_1.$$

This is the robust NSP inequality. The condition $\rho = \delta_{2s}/(1 - \delta_{2s}) < 1$ is satisfied if $\delta_{2s} < 1/2$.

Note: 6.9 for reference

\square

7.3 RIP-based guarantees for thresholding and greedy algorithms

Theorem 7.4 (IHT)

If $\delta_{3s}(A) < 1/2$, IHT applied to $\mathbf{y} = A\mathbf{x} + \mathbf{e}$ with sparsity parameter s converges to \mathbf{x} at rate

$$\|\mathbf{x} - \mathbf{x}^n\|_2 \leq \rho^n \|\mathbf{x}\|_2 + \tau \|\mathbf{e}\|_2, \quad \rho = 2\delta_{3s} < 1.$$

Proof. Let $\mathbf{e}^n = \mathbf{x}^n - \mathbf{x}$. The IHT update is $\mathbf{x}^{n+1} = H_s(\mathbf{x}^n + A^*(\mathbf{y} - A\mathbf{x}^n))$. Using $\mathbf{y} = A\mathbf{x} + \mathbf{e}$, we have:

$$\mathbf{x}^{n+1} = H_s(\mathbf{x} + (I - A^*A)\mathbf{e}^n + A^*\mathbf{e}).$$

Let $\mathbf{u} = \mathbf{x} + (I - A^*A)\mathbf{e}^n + A^*\mathbf{e}$. Both \mathbf{x} and \mathbf{x}^{n+1} are s -sparse. Let $T = \text{supp}(\mathbf{x}) \cup \text{supp}(\mathbf{x}^n) \cup \text{supp}(\mathbf{x}^{n+1})$. The size of T is at most $3s$. Since $\mathbf{x}^{n+1} = H_s(\mathbf{u})$, it is the best s -term approximation to \mathbf{u} . Thus $\|\mathbf{u} - \mathbf{x}^{n+1}\|_2 \leq \|\mathbf{u} - \mathbf{x}\|_2$. Using the triangle inequality:

$$\|\mathbf{x}^{n+1} - \mathbf{x}\|_2 \leq \|\mathbf{x}^{n+1} - \mathbf{u}\|_2 + \|\mathbf{u} - \mathbf{x}\|_2 \leq 2\|\mathbf{u} - \mathbf{x}\|_2.$$

Restricting to the support T :

$$\begin{aligned} \|\mathbf{e}^{n+1}\|_2 &\leq 2\|((I - A^*A)\mathbf{e}^n + A^*\mathbf{e})_T\|_2 \\ &\leq 2\|(I_T - A_T^*A_T)\mathbf{e}^n\|_2 + 2\|A_T^*\mathbf{e}\|_2 \\ &\leq 2\delta_{3s}\|\mathbf{e}^n\|_2 + 2\sqrt{1 + \delta_{3s}}\|\mathbf{e}\|_2. \end{aligned}$$

If $\delta_{3s} < 1/2$, then $\rho = 2\delta_{3s} < 1$, and the error converges geometrically to a noise-dependent floor.

Note: 6.18 for reference

□

Theorem 7.5 (HTP)

If $\delta_{3s}(A) < 1/\sqrt{3} \approx 0.577$, HTP recovers every s -sparse signal from noiseless measurements in at most s iterations, and is robust under noise.

Proof. Let S be the true support and S^{n+1} be the support identified at iteration n . The HTP algorithm combines a thresholding step to identify a support and a least squares step to estimate coefficients on that support. In the noiseless case:

$$\mathbf{x}^{n+1} = \operatorname{argmin}_{\text{supp}(\mathbf{z}) \subseteq S^{n+1}} \|A\mathbf{z} - A\mathbf{x}\|_2 = A_{S^{n+1}}^\dagger A\mathbf{x}.$$

Let $T = S \cup S^n \cup S^{n+1}$. Then $|T| \leq 3s$. By analyzing the least squares residual on T , one can show:

$$\|\mathbf{x}^{n+1} - \mathbf{x}\|_2 \leq \sqrt{\frac{3\delta_{3s}^2}{1 - \delta_{3s}^2}} \|\mathbf{x}^n - \mathbf{x}\|_2.$$

If $\delta_{3s} < 1/\sqrt{3}$, the contraction factor $\rho = \sqrt{3\delta_{3s}^2/(1 - \delta_{3s}^2)} < 1$. Since there are only $\binom{N}{s}$ possible supports, and each iteration with a change must strictly reduce the error, the support must eventually stabilize. Once the correct support is found ($S^{n+1} = S$), the LS step yields $\mathbf{x}^{n+1} = \mathbf{x}$ immediately.

Note: 6.21 for reference

□

Theorem 7.6 (CoSaMP)

If $\delta_{4s}(A) < (\sqrt{11/3} - 1)/2 \approx 0.4782$, CoSaMP achieves $\|\mathbf{x} - \mathbf{x}^n\|_2 \leq 2^{-n}\|\mathbf{x}\|_2 + \tau\eta$, where $\eta = \|\mathbf{e}\|_2$.

Proof. The proof proceeds in four steps. Let \mathbf{x} be s -sparse and $\mathbf{v}^n = \mathbf{x} - \mathbf{x}^n$. *Step 1: Signal proxy.* The identified atoms $J = L_{2s}(A^* \mathbf{r}_n)$ contain most of the energy of \mathbf{v}^n . By RIP: $\|\mathbf{v}_J^n\|_2 \leq 0.5\delta_{4s}\|\mathbf{v}^n\|_2 + C\eta$. *Step 2: LS estimation.* On the enlarged support $\Omega = \text{supp}(\mathbf{x}^n) \cup J$, the LS solution \mathbf{u}^{n+1} satisfies: $\|\mathbf{u}^{n+1} - \mathbf{x}\|_2 \leq \frac{1}{\sqrt{1-\delta_{4s}^2}}\|\mathbf{x} - \mathbf{x}_\Omega\|_2 + C'\eta$. *Step 3: Pruning.* $\mathbf{x}^{n+1} = H_s(\mathbf{u}^{n+1})$. By the triangle inequality: $\|\mathbf{x}^{n+1} - \mathbf{x}\|_2 \leq 2\|\mathbf{u}^{n+1} - \mathbf{x}\|_2$. Combining steps: $\|\mathbf{x}^{n+1} - \mathbf{x}\|_2 \leq \rho\|\mathbf{x}^n - \mathbf{x}\|_2 + \tau\eta$. For $\delta_{4s} < 0.4782$, $\rho < 0.5$. Iterating gives geometric convergence.

Note: 6.27 for reference

□

Theorem 7.7 (OMP via RIP, Zhang 2011)

If $\delta_{31s}(A) < 1/3$, OMP recovers every s -sparse \mathbf{x} in $30s$ iterations.

Proof. Standard OMP analysis under coherence requires $\mu(A) < 1/(2s - 1)$. Under RIP, Zhang (2011) used a potential function argument. Let \mathbf{r}_n be the residual after n steps. The energy decrease is $\|\mathbf{r}_n\|_2^2 - \|\mathbf{r}_{n+1}\|_2^2 = |\langle \mathbf{a}_{j_{n+1}}, \mathbf{r}_n \rangle|^2$. By picking the atom with largest correlation, we ensure:

$$|\langle \mathbf{a}_{j_{n+1}}, \mathbf{r}_n \rangle|^2 \geq \frac{1 - \delta_{s+|S_n|}}{s} \|\mathbf{r}_n\|_2^2.$$

However, OMP might pick "wrong" atoms. Zhang showed that for $\delta_{31s} < 1/3$, the number of times OMP can pick an atom outside the true support S is bounded. Suppose n steps have been taken and $S \not\subseteq S_n$. The residual \mathbf{r}_n is still large. Using restricted orthogonality, one can bound the maximum correlation outside S relative to that inside S . A counting argument shows that after at most $30s$ steps, all s true atoms must have been selected. Once $S \subseteq S_n$, the LS step ensures $\mathbf{r}_n = \mathbf{0}$ and exact recovery is achieved. □

Remark 7.8 (The RIP threshold table)

Algorithm	RIP order	Threshold
Basis Pursuit (sharp)	δ_{2s}	$< 4/\sqrt{41} \approx 0.6246$
Basis Pursuit (textbook)	δ_{2s}	$< \sqrt{2} - 1 \approx 0.414$
Iterative Hard Thresholding	δ_{3s}	$< 1/2$
Hard Thresholding Pursuit	δ_{3s}	$< 1/\sqrt{3} \approx 0.577$
CoSaMP	δ_{4s}	< 0.478
Orthogonal Matching Pursuit	δ_{31s}	$< 1/3$

The greedy bound (δ_{31s}) is uniform and conservative; coherence-based analysis gives sharper guarantees in the high-coherence regime.

7.4 Information-theoretic lower bound

Theorem 7.9 (Information-theoretic lower bound on m)

If A satisfies $\delta_{2s}(A) < 1/2$ for some $A \in \mathbb{C}^{m \times N}$ then $m \geq Cs \log(N/s)$ for an absolute constant $C > 0$.

Proof. For each subset $S \subseteq [N]$ with $|S| = s$ and each sign pattern $\sigma \in \{\pm 1\}^s$, form the vector $\mathbf{x}_{S,\sigma} \in \mathbb{R}^N$ with σ on S and 0 elsewhere. The collection

$$\mathcal{P} = \{\mathbf{x}_{S,\sigma} : |S| = s, \sigma \in \{\pm 1\}^s\}$$

has cardinality $|\mathcal{P}| = \binom{N}{s} 2^s$, and any two distinct elements $\mathbf{x}, \mathbf{x}' \in \mathcal{P}$ satisfy $\|\mathbf{x} - \mathbf{x}'\|_0 \leq 2s$ and $\|\mathbf{x} - \mathbf{x}'\|_2 \geq \sqrt{2}$. (If they share a support and differ in k signs, $\|\mathbf{x} - \mathbf{x}'\|_2 = 2\sqrt{k} \geq 2$; if their supports differ, then some coordinate is ± 1 in one and 0 in the other, contributing at least 1, and there are at least two such coordinates by symmetry.)

For any $\mathbf{x}, \mathbf{x}' \in \mathcal{P}$, the difference $\mathbf{v} = \mathbf{x} - \mathbf{x}'$ is $2s$ -sparse. By RIP with $\delta_{2s} < 1/2$,

$$\|A\mathbf{v}\|_2 \geq \sqrt{1 - \delta_{2s}} \|\mathbf{v}\|_2 > \frac{1}{\sqrt{2}} \|\mathbf{v}\|_2 \geq 1.$$

Conversely $\|A\mathbf{x}\|_2 \leq \sqrt{1 + \delta_{2s}} \sqrt{s} < \sqrt{3s/2}$, so all images $\{A\mathbf{x} : \mathbf{x} \in \mathcal{P}\}$ lie in the ball of radius $\sqrt{3s/2}$ in \mathbb{R}^m and are pairwise at distance > 1 . Hence the images form a $\frac{1}{2}$ -packing of that ball.

If K disjoint open balls of radius r fit inside a ball of radius R in \mathbb{R}^m , then $K r^m \leq R^m$, so $K \leq (R/r)^m$. With $r = 1/2$ and $R \leq \sqrt{3s/2} + 1/2 \leq 2\sqrt{s}$ (for $s \geq 1$), we get

$$\binom{N}{s} 2^s = |\mathcal{P}| \leq (4\sqrt{s})^m.$$

Taking logarithms and using $\binom{N}{s} \geq (N/s)^s$,

$$s \log(N/s) + s \log 2 \leq m \log(4\sqrt{s}) = m(\log 4 + \frac{1}{2} \log s).$$

Provided $s \leq N/2$ (otherwise the bound is trivial), the left side is at least $\frac{s}{2} \log(N/s)$ and the right side is $O(m \log s)$; dividing gives $m \geq C s \log(N/s) / \log s$. The sharper bound without the $\log s$ factor follows from a standard refinement: replace \mathcal{P} with a Gilbert–Varshamov subcode of $\{\pm 1, 0\}^N$ of cardinality $\binom{N}{s}/4$ all of whose pairwise differences have weight $\geq s$, recovering $m \geq C s \log(N/s)$. \square

This matches the upper bound proved next, so $m \asymp s \log(N/s)$ is optimal.

Part IV

Random Matrices, Sensing, and Algorithms

8 Random Matrices and Sensing-Matrix Design

8.1 Subgaussian random matrices

Definition 8.1 (Subgaussian random variable)

X is *subgaussian* with parameter $c > 0$ if $\mathbb{E} e^{\theta X} \leq e^{c\theta^2}$ for all $\theta \in \mathbb{R}$. Examples: Gaussian $\mathcal{N}(0, 1)$ ($c = 1/2$), Rademacher (uniform on $\{\pm 1\}$, $c = 1/2$), bounded $|X| \leq K$ ($c = K^2/2$).

Theorem 8.2 (RIP for subgaussian matrices)

Let $A \in \mathbb{R}^{m \times N}$ have i.i.d. mean-zero subgaussian entries with variance $1/m$. There exists a universal constant $C > 0$ depending only on the subgaussian parameter such that, for any $\delta, \epsilon \in (0, 1)$, if

$$m \geq C\delta^{-2}(s \ln(eN/s) + \ln(2/\epsilon)),$$

then with probability $\geq 1 - \epsilon$, $\delta_s(A) \leq \delta$.

Proof. The proof follows a standard concentration-and-covering strategy. *Step 1: Concentration for a fixed vector.* For any $\mathbf{x} \in \mathbb{R}^N$, the random variable $\|\mathbf{Ax}\|_2^2 = \frac{1}{m} \sum_{i=1}^m |(A\sqrt{m}\mathbf{x})_i|^2$ is a sum of independent subgaussian squares. By Bernstein's inequality:

$$\Pr\left(\left|\|\mathbf{Ax}\|_2^2 - \|\mathbf{x}\|_2^2\right| \geq t\|\mathbf{x}\|_2^2\right) \leq 2e^{-ct^2m}.$$

Step 2: Covering the s -sparse sphere. Fix a support S with $|S| = s$. Let $B_S = \{\mathbf{z} \in \mathbb{R}^N : \text{supp}(\mathbf{z}) \subseteq S, \|\mathbf{z}\|_2 \leq 1\}$. A ρ -net \mathcal{U}_S for B_S exists with $|\mathcal{U}_S| \leq (1 + 2/\rho)^s$. Pick $\rho = \delta/4$. By a union bound over \mathcal{U}_S :

$$\Pr\left(\max_{\mathbf{u} \in \mathcal{U}_S} \left|\|\mathbf{Au}\|_2^2 - 1\right| \geq \delta/2\right) \leq 2(1 + 8/\delta)^s e^{-c\delta^2m/4}.$$

Step 3: Extending to the whole sphere. If the bound holds on the net, then for any $\mathbf{x} \in B_S$, let $\mathbf{u} \in \mathcal{U}_S$ be the closest point. One can show $|\langle \mathbf{Ax}, \mathbf{Ax} \rangle - 1| \leq \frac{1}{1-2\rho} \max_{\mathbf{u} \in \mathcal{U}} \left|\|\mathbf{Au}\|_2^2 - 1\right| \leq \delta$. *Step 4: Union bound over supports.* There are $\binom{N}{s} \leq (eN/s)^s$ supports. Total probability of failure:

$$P_{fail} \leq (eN/s)^s \cdot 2(1 + 8/\delta)^s e^{-c\delta^2m/4} \leq 2e^{s(\ln(eN/s) + \ln(9)) - c\delta^2m/4}.$$

Setting this $\leq \epsilon$ and solving for m yields the required bound. □

Theorem 8.3 (Theorem 8.2 \Rightarrow uniform recovery)

Under the same hypothesis, with the same probability A allows uniform recovery of every s -sparse signal via BP, IHT, HTP, or CoSaMP, by combining with the respective RIP thresholds (Section 7).

Proof. Each algorithm in the list has a deterministic recovery guarantee that holds whenever an RIP constant of suitable order is below an explicit threshold:

- BP succeeds for every s -sparse signal whenever $\delta_{2s} < \sqrt{2} - 1$ (Theorem 6.9).
- IHT converges with $\delta_{3s} < 1/2$ (Theorem 6.18).
- HTP converges with $\delta_{3s} < 1/\sqrt{3}$ (Theorem 6.21).
- CoSaMP recovers up to a constant factor of $\sigma_s(\mathbf{x})_1$ when $\delta_{4s} < 0.4782$ (Theorem 6.27).

Fix the algorithm of interest and let (δ^*, k^*) be the corresponding threshold/order pair from this list. Pick any $\delta < \delta^*$.

Apply Theorem 8.2 with sparsity level k^*s and tolerance δ : the bound becomes

$$m \geq C \delta^{-2}(k^*s \ln(eN/(k^*s)) + \ln(2/\varepsilon)),$$

i.e. the same shape with $s \mapsto k^*s$. With probability $\geq 1 - \varepsilon$ the realized matrix A then satisfies $\delta_{k^*s}(A) \leq \delta < \delta^*$, so its restricted-isometry constant lies in the recovery regime of the chosen algorithm. The deterministic recovery theorem applies pointwise to A and yields uniform recovery of every s -sparse signal on the same probability event. \square

8.2 Gaussian matrices: explicit constants

Theorem 8.4 (Gaussian basis pursuit, explicit constants)

Let $A \in \mathbb{R}^{m \times N}$ have i.i.d. $\mathcal{N}(0, 1/m)$ entries. For any $\rho \in (0, 1)$ and $\varepsilon \in (0, 1)$, if

$$\frac{m^2}{m+1} \geq 2s \left[\rho^{-1} + D(s/N) + \sqrt{\frac{\ln(\varepsilon^{-1})}{s \ln(eN/s)}} \right]^2,$$

where $D(\alpha)$ is bounded by 2.05 for $\alpha \in (0, 1)$, then with probability $\geq 1 - \varepsilon$, BP recovers every \mathbf{x} from $\mathbf{y} = A\mathbf{x}$ with $\|\mathbf{x} - \mathbf{x}^\#\|_1 \leq \frac{2(1+\rho)}{1-\rho} \sigma_s(\mathbf{x})_1$. In particular, setting $\rho = 1$, $m \geq 8s \log(eN/s)$ suffices for exact uniform recovery.

Proof. This proof utilizes Gordon’s *Escape-through-the-mesh* theorem. Recovery via BP is exact if $\ker A \cap T(\mathbf{x}) = \{\mathbf{0}\}$, where $T(\mathbf{x})$ is the tangent cone of the ℓ_1 -norm at \mathbf{x} . For Gaussian A , this holds with high probability if $m \geq w(T(\mathbf{x}) \cap S^{N-1})^2 + 1$. The Gaussian width w of the intersection of the cone with

the sphere can be bounded using the distance to the dual cone (the normal cone $N_{\|\cdot\|_1}(\mathbf{x})$):

$$w(T(\mathbf{x}) \cap S^{N-1})^2 \leq \mathbb{E} \text{dist}(\mathbf{g}, N_{\|\cdot\|_1}(\mathbf{x}))^2.$$

For an s -sparse vector \mathbf{x} , the normal cone consists of vectors \mathbf{z} with $z_j = \lambda \text{sgn}(x_j)$ for $j \in S$ and $|z_j| \leq \lambda$ for $j \notin S$. Optimizing over λ and summing independent Gaussian components yields the $w^2 \approx 2s \ln(N/s)$ estimate. The stable NSP extension follows by considering the larger cone $T_{\rho,s} = \{\mathbf{v} : \|\mathbf{v}_S\|_1 \geq \rho \|\mathbf{v}_{\bar{S}}\|_1\}$, whose width is slightly larger but of the same order. Concentration of the singular values of A completes the probability bound. \square

8.3 Singular-value concentration of Gaussian matrices

Theorem 8.5 (Singular-value concentration for Gaussian matrices)

For $A \in \mathbb{R}^{m \times s}$ Gaussian with $m > s$, the singular values $\sigma_{\min}, \sigma_{\max}$ of $\frac{1}{\sqrt{m}}A$ satisfy, for $t > 0$,

$$\Pr(\sigma_{\max} \geq 1 + \sqrt{s/m} + t) \leq e^{-mt^2/2},$$

$$\Pr(\sigma_{\min} \leq 1 - \sqrt{s/m} - t) \leq e^{-mt^2/2}.$$

Proof. Write $G = A$ and view the entries of G as a vector $g \in \mathbb{R}^{ms}$ with $g \sim \mathcal{N}(0, I_{ms})$. The maps

$$F_+(g) = \sigma_{\max}(G), \quad F_-(g) = -\sigma_{\min}(G)$$

are 1-Lipschitz with respect to the Frobenius (Euclidean) norm on g , because for any unit \mathbf{x} and any perturbation $G \mapsto G + H$,

$$| |(G + H)\mathbf{x} \|_2 - \|G\mathbf{x} \|_2 | \leq \|H\mathbf{x} \|_2 \leq \|H\|_F,$$

and singular values are extrema of $\mathbf{x} \mapsto \|G\mathbf{x} \|_2$.

Apply the Gaussian Lipschitz concentration inequality (Borell–TIS): for any 1-Lipschitz F on \mathbb{R}^d with $g \sim \mathcal{N}(0, I_d)$, $\Pr(F(g) \geq \mathbb{E}F + t) \leq e^{-t^2/2}$ for all $t > 0$. This gives

$$\Pr(\sigma_{\max} \geq \mathbb{E}\sigma_{\max} + t) \leq e^{-t^2/2}, \quad \Pr(\sigma_{\min} \leq \mathbb{E}\sigma_{\min} - t) \leq e^{-t^2/2}.$$

It remains to bound the two means. By the Gordon–Slepian comparison theorem, for the Gaussian process $X_{\mathbf{x},\mathbf{y}} = \langle G\mathbf{x}, \mathbf{y} \rangle$ on $S^{s-1} \times S^{m-1}$,

$$\mathbb{E} \sup_{\mathbf{x},\mathbf{y}} X_{\mathbf{x},\mathbf{y}} \leq \mathbb{E} \sup_{\mathbf{x},\mathbf{y}} (\langle \mathbf{g}, \mathbf{x} \rangle + \langle \mathbf{h}, \mathbf{y} \rangle) = \mathbb{E} \|\mathbf{g}\|_2 + \mathbb{E} \|\mathbf{h}\|_2 \leq \sqrt{s} + \sqrt{m},$$

with independent $\mathbf{g} \sim \mathcal{N}(0, I_s)$, $\mathbf{h} \sim \mathcal{N}(0, I_m)$. The left side equals $\mathbb{E} \sigma_{\max}(G)$. A symmetric Gordon argument using the min-max representation $\sigma_{\min} = \min_{\mathbf{x}} \max_{\mathbf{y}} \langle G\mathbf{x}, \mathbf{y} \rangle$ (with the signs of \mathbf{x}, \mathbf{y} exchanged) yields $\mathbb{E} \sigma_{\min}(G) \geq \sqrt{m} - \sqrt{s}$.

Combining the mean bounds with the deviation bounds and rescaling $G \mapsto \frac{1}{\sqrt{m}}G$ (which divides every singular value by \sqrt{m} and the deviation t becomes $t\sqrt{m}$) gives

$$\Pr(\sigma_{\max}(\frac{1}{\sqrt{m}}G) \geq 1 + \sqrt{s/m} + t) \leq e^{-mt^2/2},$$

and analogously for σ_{\min} .

Note: 9.24 for reference

□

8.4 Johnson–Lindenstrauss embeddings and the JL \Leftrightarrow RIP duality

Theorem 8.6 (Johnson–Lindenstrauss lemma)

For any finite set $\{\mathbf{x}_1, \dots, \mathbf{x}_M\} \subset \mathbb{R}^N$ and $\delta \in (0, 1/2)$, if $m \geq C\delta^{-2} \ln M$ (universal C), there exists a linear map $f : \mathbb{R}^N \rightarrow \mathbb{R}^m$ (e.g., a Gaussian random projection) such that

$$(1 - \delta)\|\mathbf{x}_i - \mathbf{x}_j\|_2 \leq \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2 \leq (1 + \delta)\|\mathbf{x}_i - \mathbf{x}_j\|_2 \quad \forall i, j.$$

Proof. Let $A \in \mathbb{R}^{m \times N}$ have i.i.d. $\mathcal{N}(0, 1)$ entries and put $f(\mathbf{x}) = \frac{1}{\sqrt{m}}A\mathbf{x}$. By linearity it suffices to show that for any fixed $\mathbf{u} \in \mathbb{R}^N$,

$$\Pr\left((1 - \delta)\|\mathbf{u}\|_2 \leq \|f(\mathbf{u})\|_2 \leq (1 + \delta)\|\mathbf{u}\|_2\right) \geq 1 - 2e^{-cm\delta^2},$$

for an absolute constant $c > 0$.

By homogeneity assume $\|\mathbf{u}\|_2 = 1$. Then $A\mathbf{u} \sim \mathcal{N}(0, I_m)$, so $\|A\mathbf{u}\|_2^2$ is a χ_m^2 random variable with mean m . Standard χ^2 -concentration (Laurent–Massart) gives, for $t \in (0, 1)$,

$$\Pr\left(\left|\frac{1}{m}\|A\mathbf{u}\|_2^2 - 1\right| \geq t\right) \leq 2\exp(-mt^2/8).$$

Setting $t = \delta(2 - \delta)/(1 + \delta) \geq \delta/2$ converts the squared-norm deviation into the desired norm deviation: if $|m^{-1}\|A\mathbf{u}\|_2^2 - 1| \leq t$ then $1 - \delta \leq m^{-1/2}\|A\mathbf{u}\|_2 \leq 1 + \delta$.

Apply this to each difference vector $\mathbf{u}_{ij} = \mathbf{x}_i - \mathbf{x}_j$. There are $\binom{M}{2} \leq M^2/2$ such pairs, so a union bound gives a failure probability of at most $M^2e^{-cm\delta^2}$. Choosing $m \geq C\delta^{-2} \ln M$ with C large enough makes this < 1 , proving existence. □

Theorem 8.7 (Krahmer–Ward 2011 – RIP \Rightarrow JL)

If A has $\delta_{2s}(A) \leq \eta/4$ for some $s \geq 16 \ln(4M/\epsilon)$, then the randomized matrix AD_ϵ , where $D_\epsilon = \text{diag}(\epsilon_j)$ with $\epsilon_j \in \{\pm 1\}$ random, satisfies the JL bound on E with probability $\geq 1 - \epsilon$.

Proof. Fix $\mathbf{x} \in E$ with $\|\mathbf{x}\|_2 = 1$. Partition the coordinates of \mathbf{x} into blocks S_0, S_1, \dots of size s , where S_0 indexes the s largest entries by magnitude, S_1 the next s , and so on. Write $\mathbf{x}_k = \mathbf{x}_{S_k}$ and $\mathbf{y}_k = AD_\epsilon \mathbf{x}_k =$

$A\mathbf{x}_k \odot \boldsymbol{\epsilon}_{S_k}$ (with the abusive notation that the sign acts on the columns).

Expand the squared norm

$$\|AD_\epsilon \mathbf{x}\|_2^2 = \sum_k \|A\mathbf{x}_k\|_2^2 + 2 \sum_{k < \ell} \langle AD_\epsilon \mathbf{x}_k, AD_\epsilon \mathbf{x}_\ell \rangle.$$

The diagonal terms are deterministic and well-controlled: each \mathbf{x}_k is s -sparse, so RIP gives $(1 - \delta_{2s})\|\mathbf{x}_k\|_2^2 \leq \|A\mathbf{x}_k\|_2^2 \leq (1 + \delta_{2s})\|\mathbf{x}_k\|_2^2$. Summing over k and using $\sum_k \|\mathbf{x}_k\|_2^2 = \|\mathbf{x}\|_2^2 = 1$,

$$1 - \delta_{2s} \leq \sum_k \|A\mathbf{x}_k\|_2^2 \leq 1 + \delta_{2s}.$$

For the cross-terms, fix $k < \ell$ and write the inner product as a Rademacher chaos in the signs $\boldsymbol{\epsilon}$:

$$\langle AD_\epsilon \mathbf{x}_k, AD_\epsilon \mathbf{x}_\ell \rangle = \sum_{i \in S_k} \sum_{j \in S_\ell} \epsilon_i \epsilon_j x_i x_j \langle \mathbf{a}_i, \mathbf{a}_j \rangle.$$

This is a centered second-order Rademacher chaos $\boldsymbol{\epsilon}^\top M \boldsymbol{\epsilon}$ with off-diagonal matrix $M_{ij} = x_i x_j \langle \mathbf{a}_i, \mathbf{a}_j \rangle$ supported on $S_k \times S_\ell$. The Hanson–Wright inequality bounds its tail by $\Pr(|\boldsymbol{\epsilon}^\top M \boldsymbol{\epsilon}| \geq t) \leq 2 \exp(-c \min(t^2/\|M\|_F^2, t/\|M\|_{op}))$. Both norms are estimable via the RIP of A on $S_k \cup S_\ell$ (size $\leq 2s$): $\|M\|_{op} \leq \delta_{2s} \|\mathbf{x}_k\|_2 \|\mathbf{x}_\ell\|_2$ and $\|M\|_F^2 \leq \delta_{2s}^2 \|\mathbf{x}_k\|_2^2 \|\mathbf{x}_\ell\|_2^2$, since the cross-Gram matrix $A_{S_k}^* A_{S_\ell}$ has operator norm $\leq \delta_{2s}$.

Summing the chaos bounds over the $\binom{N/s}{2}$ pairs (k, ℓ) and using Cauchy–Schwarz to control $\sum_{k < \ell} \|\mathbf{x}_k\|_2 \|\mathbf{x}_\ell\|_2 \leq 1$, the total cross-term contribution is bounded by $C\delta_{2s}\eta$ in absolute value with probability $\geq 1 - 2e^{-cs}$ for $s \geq 16 \ln(4M/\epsilon)$.

Combining diagonal and off-diagonal estimates, $|\|AD_\epsilon \mathbf{x}\|_2^2 - 1| \leq \eta$ with probability $\geq 1 - 2e^{-cs}$. A union bound over the $|E| = M$ points gives failure probability $\leq 2Me^{-cs} \leq \epsilon$ under the stated s -bound, proving the JL guarantee. \square

This says: *RIP and JL are essentially the same property*; one implies the other up to a column-sign randomization.

8.5 Fast and structured constructions

Dense Gaussian matrices cost $O(mN)$ to apply. Faster alternatives:

- **Fast JL Transform (Ailon–Chazelle 2006)**: $A = (\text{sparse Gaussian}) \cdot (\text{Hadamard}) \cdot (\text{random sign flips})$, applied in $O(N \log m)$ time, achieves $m = O(\delta^{-2} \log P)$ JL embedding.
- **Structurally Random Matrices**: $A = DFR$, R a randomizer (sign flips or permutation), F a fast transform (FFT, Walsh–Hadamard, DCT), D a random down-sampler. Cost $O(N \log N)$ with RIP guarantees $m = O(s \log^c N)$.
- **Partial Fourier matrices**: pick m rows uniformly at random from the $N \times N$ DFT. Then RIP with

$m \geq Cs \log^4 N$ holds with high probability (

Note: 12.31 for reference

, restricted-isometry result for bounded orthonormal systems).

- **2D-SME:** separable measurement ensembles for images. $Y = D_r F R X R^T F^T D_c^T$, sampling rows and columns separately.

8.6 Non-uniform recovery

Theorem 8.8 (Gaussian non-uniform recovery via BP)

Fix an s -sparse $\mathbf{x} \in \mathbb{R}^N$. For $\varepsilon \in (0, 1)$, if

$$\frac{m^2}{m+1} \geq 2s \left[\sqrt{\ln(2.34N/s)} + \sqrt{\ln(\varepsilon^{-1})/s} \right]^2,$$

then with probability $\geq 1 - \varepsilon$, BP recovers \mathbf{x} from $\mathbf{y} = A\mathbf{x}$.

Proof. By the tangent-cone characterization (Theorem 4.34), \mathbf{x} is the unique BP solution of $Az = A\mathbf{x}$ iff

$$\ker A \cap T(\mathbf{x}) = \{\mathbf{0}\},$$

where $T(\mathbf{x})$ is the tangent cone of the ℓ_1 -ball at $\mathbf{x}/\|\mathbf{x}\|_1$ (equivalently, the cone of feasible descent directions for the ℓ_1 -norm at \mathbf{x}).

Let $E = T(\mathbf{x}) \cap S^{N-1}$. Gordon's theorem gives, for $A \in \mathbb{R}^{m \times N}$ with i.i.d. $\mathcal{N}(0, 1)$ entries,

$$\mathbb{E} \inf_{\mathbf{v} \in E} \|\mathbf{A}\mathbf{v}\|_2 \geq \lambda_m - w(E), \quad \lambda_m = \mathbb{E} \|\mathbf{g}\|_2 = \sqrt{2} \Gamma(\frac{m+1}{2}) / \Gamma(\frac{m}{2}),$$

with $w(E) = \mathbb{E} \sup_{\mathbf{v} \in E} \langle \mathbf{g}, \mathbf{v} \rangle$ the Gaussian width. The exact identity $\lambda_m^2 = m^2/(m+1)$ (a Gamma-function computation) is what produces the $m^2/(m+1)$ on the LHS of the theorem. Whenever $\lambda_m > w(E)$, Gaussian Lipschitz concentration of $\mathbf{v} \mapsto \inf_{\mathbf{u} \in E} \|\mathbf{A}\mathbf{u}\|_2$ around its mean (a 1-Lipschitz function of A) gives $\Pr(\inf_{\mathbf{v} \in E} \|\mathbf{A}\mathbf{v}\|_2 > 0) \geq 1 - e^{-(\lambda_m - w(E))^2/2}$, hence $\ker A \cap T(\mathbf{x}) = \{\mathbf{0}\}$ with the same probability.

The width of an ℓ_1 tangent cone at an s -sparse point admits the closed-form upper bound $w(T(\mathbf{x}) \cap S^{N-1})^2 \leq 2s \ln(2.34N/s)$. The standard route is duality: $w(E) \leq \mathbb{E} \text{dist}(\mathbf{g}, N(\mathbf{x}))$ where $N(\mathbf{x})$ is the polar cone (the ℓ_1 -subdifferential cone, generated by sign vectors on the support and $[-1, 1]$ -bounded entries off the support). Since the off-support distance contribution is the soft-thresholded \mathbf{g}_τ , $\mathbb{E} \text{dist}(\mathbf{g}, N(\mathbf{x}))^2 \leq 2s \mathbb{E} \eta_\tau(g_1)^2$ with $\tau = \sqrt{2 \ln(N/s)}$, evaluating to $\leq 2s \ln(2.34N/s)$ after optimizing the threshold.

The condition $\lambda_m > w(E) + \sqrt{2 \ln \varepsilon^{-1}}$ guarantees recovery with probability $\geq 1 - \varepsilon$ via the concentration step. Using $\lambda_m^2 \geq m^2/(m+1)$ and the width estimate above produces exactly the $m^2/(m+1) \geq 2s[\sqrt{\ln(2.34N/s)} + \sqrt{\ln(\varepsilon^{-1})/s}]^2$ bound stated in the theorem.

Note: 9.16 for reference

□

Part V

Sparse Representation Classification

9 Sparse Representation-based Classification

9.1 Setup

Stack training images of class $i \in \{1, \dots, C\}$ as columns of $A_i \in \mathbb{R}^{m \times n_i}$, and concatenate

$$A = [A_1 | A_2 | \dots | A_C] \in \mathbb{R}^{m \times n}, \quad n = \sum_i n_i.$$

A test sample \mathbf{y} from class i ideally lies near $\text{span}(A_i)$, so a near-block-sparse representation $\mathbf{y} \approx A\mathbf{x}$ should have \mathbf{x} supported on the i -th block.

9.2 The SRC algorithm (Wright et al. 2009)

Algorithm 9 Sparse Representation-based Classification (SRC)

Require: dictionary $A = [A_1 | \dots | A_C]$, test sample \mathbf{y} , tolerance ε

- 1: Normalize columns of A to unit ℓ_2 .
 - 2: Solve $\hat{\mathbf{x}} = \arg \min \|\mathbf{x}\|_1$ s.t. $\|A\mathbf{x} - \mathbf{y}\|_2 \leq \varepsilon$.
 - 3: **for** $i = 1, \dots, C$ **do**
 - 4: $\delta_i(\hat{\mathbf{x}}) \leftarrow$ zero out coefficients outside class i
 - 5: $r_i(\mathbf{y}) \leftarrow \| \mathbf{y} - A\delta_i(\hat{\mathbf{x}}) \|_2$
 - 6: **end for**
 - 7: **return** $\hat{c} = \arg \min_i r_i(\mathbf{y})$.
-

9.3 Robust SRC under occlusion

For corruption \mathbf{e} that is itself sparse (occlusions, lighting artifacts, salt-and-pepper noise),

$$\mathbf{y} = A\mathbf{x}_0 + \mathbf{e}, \quad \|\mathbf{e}\|_0 \ll m.$$

Augment the dictionary by an identity:

$$B = [A | I_m] \in \mathbb{R}^{m \times (n+m)},$$

and solve

$$\min_{\mathbf{w}} \|\mathbf{w}\|_1 \quad \text{s.t.} \quad B\mathbf{w} = \mathbf{y}, \quad \mathbf{w} = \begin{bmatrix} \mathbf{x} \\ \mathbf{e} \end{bmatrix}.$$

Equivalently $\min_{\mathbf{x}, \mathbf{e}} \|\mathbf{x}\|_1 + \|\mathbf{e}\|_1$ s.t. $\mathbf{y} = A\mathbf{x} + \mathbf{e}$.

9.4 Sparsity Concentration Index (SCI)

A diagnostic quantity:

$$\text{SCI}(\mathbf{x}) := \frac{C \max_i \|\delta_i(\mathbf{x})\|_1 / \|\mathbf{x}\|_1 - 1}{C - 1} \in [0, 1].$$

SCI = 1 if \mathbf{x} is fully concentrated in a single class; 0 if uniform. A threshold τ on SCI lets the classifier reject low-confidence samples (“don’t know” rather than mis-classify).

9.5 Theoretical correctness

Theorem 9.1 (Sparse-representation classification (SRC) correctness, Wright et al. 2009)

If the per-class dictionaries A_i are sufficiently incoherent and the test sample lies in $\text{span}(A_i)$ up to noise ε , then SRC labels \mathbf{y} correctly when ε is below a threshold determined by the inter-class separation.

Proof. Suppose the test sample admits a noisy decomposition $\mathbf{y} = A_i \mathbf{x}_i^* + \mathbf{e}$ with $\|\mathbf{e}\|_2 \leq \varepsilon$ and the ground-truth class is i . Let $\mathbf{x}^* \in \mathbb{R}^N$ be the vector with \mathbf{x}_i^* on the block of class i and 0 elsewhere; then $\|\mathbf{x}^*\|_0 \leq n_i$.

Let $A = [A_1 \mid \dots \mid A_C]$ be the concatenated dictionary. By assumption A satisfies a coherence (or RIP) condition strong enough that BPDN with noise level ε recovers any s -sparse signal stably for all $s \geq n_i$ — concretely, $\mu_1(s) + \mu_1(s - 1) < 1$ or $\delta_{2s}(A) < \sqrt{2} - 1$ depending on which guarantee is invoked. The BPDN solution $\hat{\mathbf{x}} = \arg \min \|\mathbf{x}\|_1$ s.t. $\|A\mathbf{x} - \mathbf{y}\|_2 \leq \varepsilon$ then satisfies

$$\|\hat{\mathbf{x}} - \mathbf{x}^*\|_2 \leq D\varepsilon$$

for an absolute constant D depending on the recovery threshold. Hence $\hat{\mathbf{x}}_i \approx \mathbf{x}_i^*$ and $\hat{\mathbf{x}}_j \approx \mathbf{0}$ for $j \neq i$, both to error $D\varepsilon$.

The class- i residual is

$$r_i(\mathbf{y}) = \|\mathbf{y} - A_i \hat{\mathbf{x}}_i\|_2 \leq \|\mathbf{y} - A_i \mathbf{x}_i^*\|_2 + \|A_i(\hat{\mathbf{x}}_i - \mathbf{x}_i^*)\|_2 \leq \varepsilon + \|A\|_{op} D\varepsilon = (1 + D\|A\|_{op})\varepsilon.$$

For $j \neq i$, using $\hat{\mathbf{x}}_j$ small,

$$r_j(\mathbf{y}) \geq \|\mathbf{y} - A_j \mathbf{x}_j^*\|_2 - \|A_j\|_{op} \|\hat{\mathbf{x}}_j - \mathbf{x}_j^*\|_2 \geq \text{dist}(\mathbf{y}, \text{span}(A_j)) - \|A\|_{op} D\varepsilon.$$

Since \mathbf{y} lies ε -close to $\text{span}(A_i)$ and the per-class subspaces are separated by $\Delta_{ij} := \text{dist}(\text{span}(A_i), \text{span}(A_j))$, we have $\text{dist}(\mathbf{y}, \text{span}(A_j)) \geq \Delta_{ij} - \varepsilon$.

Therefore $r_j(\mathbf{y}) - r_i(\mathbf{y}) \geq \Delta_{ij} - (2 + 2D\|A\|_{op})\varepsilon$, which is positive — i.e. class i wins — whenever

$$\varepsilon < \frac{\min_{j \neq i} \Delta_{ij}}{2 + 2D\|A\|_{op}}.$$

This is the threshold determined by the inter-class separation. □

9.6 Comparison with classical classifiers

- **Eigenfaces (PCA)**: project onto top eigenvectors of the covariance, classify by nearest neighbor. Sensitive to occlusion (corrupted pixels pollute all projections).
- **Fisherfaces (LDA)**: maximize Fisher's criterion $\max_W |W^T S_B W| / |W^T S_W W|$. Better class separation than PCA but still distance-based.
- **k -NN**: classify by majority vote among nearest training points. Simple but high variance, sensitive to dimensionality.
- **SVM**: maximum-margin hyperplane $\min \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum \xi_i$ s.t. $c_i(\mathbf{w}^T \mathbf{y}_i + b) \geq 1 - \xi_i$. Excellent for binary problems, kernel extension for non-linearity.
- **SRC**: leverages CS theory; handles occlusion natively via the identity-augmented dictionary; can certify reliability via SCI.

9.7 Applications

Face recognition under occlusion, hyperspectral image classification, distributed sensor identification, activity recognition from wearable sensors, video surveillance.

Part VI

ℓ_1 -Minimization Algorithms

10 Convex Algorithms for ℓ_1 Minimization

10.1 The toolkit

We solve large-scale instances of

$$\text{BP: } \min \|\mathbf{x}\|_1 \text{ s.t. } \mathbf{Ax} = \mathbf{y}, \quad (7)$$

$$\text{BPDN: } \min \|\mathbf{x}\|_1 \text{ s.t. } \|\mathbf{Ax} - \mathbf{y}\|_2 \leq \eta, \quad (8)$$

$$\text{LASSO: } \min \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1, \quad (9)$$

$$\text{Elastic Net: } \min \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \lambda_1 \|\mathbf{x}\|_1 + \lambda_2 \|\mathbf{x}\|_2^2, \quad (10)$$

$$\text{Dantzig: } \min \|\mathbf{x}\|_1 \text{ s.t. } \|A^*(\mathbf{Ax} - \mathbf{y})\|_\infty \leq \tau. \quad (11)$$

The choice of algorithm depends on problem size, structure of A (dense, fast-multiplied, factored), and noise model.

10.2 Geometry of why ℓ_1 promotes sparsity

The unit ℓ_1 -ball is a polytope with vertices on the coordinate axes; the ℓ_2 -ball is round. As we shrink the ℓ_p -norm subject to $\mathbf{Ax} = \mathbf{y}$, the first contact with the affine set lands at a vertex (sparse) for $p = 1$ but on a generic face (dense) for $p = 2$.

10.3 Soft thresholding and proximal operators

Definition 10.1 (Soft thresholding)

$$\mathcal{S}_\tau(\mathbf{u}) = \text{sign}(\mathbf{u}) \odot \max(|\mathbf{u}| - \tau, 0),$$

component-wise.

Lemma 10.2 (Proximal operator of ℓ_1)

$$\mathcal{S}_\tau(\mathbf{u}) = \arg \min_{\mathbf{x}} \left(\frac{1}{2} \|\mathbf{x} - \mathbf{u}\|_2^2 + \tau \|\mathbf{x}\|_1 \right).$$

Proof. The objective separates across coordinates; for each x_i minimize $\frac{1}{2}(x_i - u_i)^2 + \tau|x_i|$. The KKT condition gives the soft-threshold formula. \square

10.4 ISTA and FISTA

Algorithm 10 ISTA – Iterative Soft-Thresholding for LASSO

Require: step $\alpha \leq 1/\|A\|_{2 \rightarrow 2}^2$, $\lambda > 0$

- 1: $\mathbf{x}^0 \leftarrow \mathbf{0}$
- 2: **for** $k = 0, 1, \dots$ **do**
- 3: $\mathbf{x}^{k+1} \leftarrow \mathcal{S}_{\alpha\lambda}(\mathbf{x}^k - \alpha A^*(A\mathbf{x}^k - \mathbf{y}))$
- 4: **end for**

Algorithm 11 FISTA – Beck–Teboulle 2009

- 1: $\mathbf{x}^0 \leftarrow \mathbf{0}$, $\mathbf{z}^1 \leftarrow \mathbf{x}^0$, $t_1 \leftarrow 1$
- 2: **for** $k = 1, 2, \dots$ **do**
- 3: $\mathbf{x}^k \leftarrow \mathcal{S}_{\alpha\lambda}(\mathbf{z}^k - \alpha A^*(A\mathbf{z}^k - \mathbf{y}))$
- 4: $t_{k+1} \leftarrow \frac{1 + \sqrt{1 + 4t_k^2}}{2}$
- 5: $\mathbf{z}^{k+1} \leftarrow \mathbf{x}^k + \frac{t_k - 1}{t_{k+1}}(\mathbf{x}^k - \mathbf{x}^{k-1})$
- 6: **end for**

Theorem 10.3 (FISTA convergence rate)

For LASSO with $L = \|A\|_{2 \rightarrow 2}^2$ and step $\alpha = 1/L$, ISTA satisfies $F(\mathbf{x}^k) - F^* = O(L\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2/k)$; FISTA improves to $O(L\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2/k^2)$, optimal among first-order methods (Nesterov).

Proof. Write $F = f + g$ with $f(\mathbf{x}) = \frac{1}{2}\|A\mathbf{x} - \mathbf{y}\|_2^2$ (L -smooth, $L = \|A\|_{2 \rightarrow 2}^2$) and $g(\mathbf{x}) = \lambda\|\mathbf{x}\|_1$. The proximal-gradient operator $T_\alpha(\mathbf{z}) := \text{prox}_{\alpha g}(\mathbf{z} - \alpha \nabla f(\mathbf{z}))$ at step $\alpha = 1/L$ satisfies the standard descent lemma: for every \mathbf{x} ,

$$F(T_\alpha \mathbf{z}) - F(\mathbf{x}) \leq \frac{L}{2}(\|\mathbf{z} - \mathbf{x}\|_2^2 - \|T_\alpha \mathbf{z} - \mathbf{x}\|_2^2). \quad (*)$$

This is obtained by combining $f(T_\alpha \mathbf{z}) \leq f(\mathbf{z}) + \langle \nabla f(\mathbf{z}), T_\alpha \mathbf{z} - \mathbf{z} \rangle + \frac{L}{2}\|T_\alpha \mathbf{z} - \mathbf{z}\|_2^2$ (smoothness) with the prox subgradient inequality $g(T_\alpha \mathbf{z}) \leq g(\mathbf{x}) + \langle (\mathbf{z} - T_\alpha \mathbf{z})/\alpha - \nabla f(\mathbf{z}), T_\alpha \mathbf{z} - \mathbf{x} \rangle$ and rearranging.

For ISTA, $\mathbf{x}^{k+1} = T_\alpha \mathbf{x}^k$ with $\mathbf{z} = \mathbf{x}^k$. Apply (*) first with $\mathbf{x} = \mathbf{x}^k$ to get $F(\mathbf{x}^{k+1}) \leq F(\mathbf{x}^k)$, and then with $\mathbf{x} = \mathbf{x}^*$ to get $F(\mathbf{x}^{k+1}) - F^* \leq \frac{L}{2}(\|\mathbf{x}^k - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{k+1} - \mathbf{x}^*\|_2^2)$. Summing the second inequality from 0 to $k - 1$ telescopes; using monotonicity of $F(\mathbf{x}^k)$ yields $F(\mathbf{x}^k) - F^* \leq L\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2/(2k)$, the $O(1/k)$ rate.

For FISTA, the analysis is more delicate because the prox is taken at the extrapolation point \mathbf{z}^k , not at \mathbf{x}^k . Apply (*) once at $\mathbf{z} = \mathbf{z}^k$, $\mathbf{x} = \mathbf{x}^k$ and once at $\mathbf{z} = \mathbf{z}^k$, $\mathbf{x} = \mathbf{x}^*$, multiply the first by $(t_k - 1)$ and the

second by 1, and add. Using $t_k^2 = t_{k+1}^2 - t_{k+1}$ (from the recurrence $t_{k+1} = (1 + \sqrt{1 + 4t_k^2})/2$) and the extrapolation identity $t_{k+1}z^{k+1} = x^{k+1} + (t_k - 1)(x^{k+1} - x^k)$, one obtains the Lyapunov sequence

$$u_k := t_k^2(F(x^k) - F^*) + \frac{L}{2} \|t_k z^{k+1} - (t_k - 1)x^k - x^*\|_2^2,$$

satisfying $u_{k+1} \leq u_k$. Therefore $u_k \leq u_0 = \frac{L}{2} \|x^0 - x^*\|_2^2$. Since the recurrence forces $t_k \geq (k + 1)/2$ by induction (base $t_1 = 1$, $t_{k+1}^2 - t_{k+1} = t_k^2 \geq (k + 1)^2/4$ implies $t_{k+1} \geq (k + 2)/2$),

$$F(x^k) - F^* \leq \frac{u_k}{t_k^2} \leq \frac{2L \|x^0 - x^*\|_2^2}{(k + 1)^2}.$$

The matching lower bound — that no first-order oracle method achieves better than $\Omega(L \|x^0 - x^*\|_2^2/k^2)$ in the worst case — is Nesterov’s classical construction using a particular tridiagonal quadratic. \square

10.5 ADMM for BP

Algorithm 12 ADMM for $\min \|x\|_1$ s.t. $Ax = y$

- 1: $x^0, z^0, u^0 \leftarrow \mathbf{0}; \rho > 0$
 - 2: **for** $t = 0, 1, \dots$ **do**
 - 3: $x^{t+1} \leftarrow \mathcal{S}_{1/\rho}(z^t - u^t)$ \triangleright soft threshold
 - 4: $z^{t+1} \leftarrow x^{t+1} + u^t + A^*(AA^*)^{-1}(y - A(x^{t+1} + u^t))$
 - 5: $u^{t+1} \leftarrow u^t + x^{t+1} - z^{t+1}$
 - 6: **end for**
-

The projection step uses a single Cholesky factor of AA^* amortized over all iterations. Theoretical rate: $O(1/k)$ ergodic; linear under additional strong-convexity hypotheses.

10.6 Primal–dual splitting (Chambolle–Pock 2011)

For the saddle-point problem $\min_x \max_\xi \langle Ax, \xi \rangle + G(x) - F^*(\xi)$:

Algorithm 13 Chambolle–Pock primal-dual algorithm

Require: $\theta \in [0, 1]$, $\sigma, \tau > 0$ with $\sigma\tau\|A\|_{2 \rightarrow 2}^2 < 1$

- 1: $\mathbf{x}^0 \in \mathbb{R}^N$, $\boldsymbol{\xi}^0 \in \mathbb{R}^m$, $\bar{\mathbf{x}}^0 \leftarrow \mathbf{x}^0$
- 2: **for** $n = 0, 1, \dots$ **do**
- 3: $\boldsymbol{\xi}^{n+1} \leftarrow \text{prox}_{\sigma F^*}(\boldsymbol{\xi}^n + \sigma A\bar{\mathbf{x}}^n)$
- 4: $\mathbf{x}^{n+1} \leftarrow \text{prox}_{\tau G}(\mathbf{x}^n - \tau A^*\boldsymbol{\xi}^{n+1})$
- 5: $\bar{\mathbf{x}}^{n+1} \leftarrow \mathbf{x}^{n+1} + \theta(\mathbf{x}^{n+1} - \mathbf{x}^n)$
- 6: **end for**

Theorem 10.4 (ℓ_1 -instance optimality)

Under the step-size condition, with $\theta = 1$, the iterates converge to a saddle point. The ergodic gap decays at rate $O(1/k)$ in the convex case, $O(1/k^2)$ if either F^* or G is strongly convex, and linearly if both are.

Proof. Let $(\mathbf{x}^*, \boldsymbol{\xi}^*)$ be a saddle point. The proximal updates produce the variational inequalities (subgradient definitions of the prox)

$$\frac{1}{\tau}(\mathbf{x}^n - \mathbf{x}^{n+1}) - A^*\boldsymbol{\xi}^{n+1} \in \partial G(\mathbf{x}^{n+1}), \quad \frac{1}{\sigma}(\boldsymbol{\xi}^n - \boldsymbol{\xi}^{n+1}) + A\bar{\mathbf{x}}^n \in \partial F^*(\boldsymbol{\xi}^{n+1}).$$

Pairing each with the corresponding gap functional and using monotonicity of the subdifferentials gives

$$\mathcal{L}(\mathbf{x}^{n+1}, \boldsymbol{\xi}^*) - \mathcal{L}(\mathbf{x}^*, \boldsymbol{\xi}^{n+1}) \leq \frac{1}{2\tau}\|\mathbf{x}^n - \mathbf{x}^*\|_2^2 - \frac{1}{2\tau}\|\mathbf{x}^{n+1} - \mathbf{x}^*\|_2^2 + \frac{1}{2\sigma}\|\boldsymbol{\xi}^n - \boldsymbol{\xi}^*\|_2^2 - \frac{1}{2\sigma}\|\boldsymbol{\xi}^{n+1} - \boldsymbol{\xi}^*\|_2^2 - R_n,$$

where $\mathcal{L}(\mathbf{x}, \boldsymbol{\xi}) = \langle A\mathbf{x}, \boldsymbol{\xi} \rangle + G(\mathbf{x}) - F^*(\boldsymbol{\xi})$ and $R_n = \frac{1}{2\tau}\|\mathbf{x}^{n+1} - \mathbf{x}^n\|_2^2 + \frac{1}{2\sigma}\|\boldsymbol{\xi}^{n+1} - \boldsymbol{\xi}^n\|_2^2 - \langle A(\mathbf{x}^{n+1} - \mathbf{x}^n), \boldsymbol{\xi}^{n+1} - \boldsymbol{\xi}^n \rangle$. Cauchy–Schwarz on the cross-term and the step-size condition $\sigma\tau\|A\|_{2 \rightarrow 2}^2 < 1$ imply $R_n \geq 0$. So the Lyapunov function $V_n = \frac{1}{2\tau}\|\mathbf{x}^n - \mathbf{x}^*\|_2^2 + \frac{1}{2\sigma}\|\boldsymbol{\xi}^n - \boldsymbol{\xi}^*\|_2^2$ is non-increasing.

Summing the inequality from $n = 0$ to $N - 1$ and using convexity of \mathcal{L} in \mathbf{x} and concavity in $\boldsymbol{\xi}$ on the ergodic averages $\bar{\mathbf{x}}^N = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^n$, $\bar{\boldsymbol{\xi}}^N = \frac{1}{N} \sum_{n=1}^N \boldsymbol{\xi}^n$,

$$\mathcal{L}(\bar{\mathbf{x}}^N, \boldsymbol{\xi}^*) - \mathcal{L}(\mathbf{x}^*, \bar{\boldsymbol{\xi}}^N) \leq \frac{V_0}{N}.$$

This is the $O(1/N)$ rate on the primal-dual gap.

If G is γ -strongly convex (the case $\gamma > 0$), the prox-step inequality picks up an extra $\frac{\gamma}{2}\|\mathbf{x}^{n+1} - \mathbf{x}^*\|_2^2$ term. Choosing time-varying $\tau_n, \sigma_n, \theta_n$ that exploit this acceleration via Chambolle–Pock’s adaptive scheme yields $V_n = O(1/n^2)$. If F^* is also strongly convex, the Lyapunov bound contracts geometrically, giving linear convergence with explicit rate $1/(1 + 2\sqrt{\gamma_F\gamma_G\sigma\tau})$. \square

10.7 Iteratively Reweighted Least Squares (IRLS)

The identity $|t| = t^2/|t|$ inspires solving a sequence of weighted least-squares problems:

Algorithm 14 IRLS, Daubechies–DeVore–Fornasier–Güntürk

- 1: $\mathbf{w}^0 \leftarrow \mathbf{1}, \varepsilon_0 \leftarrow 1$
- 2: **for** $n = 0, 1, \dots$ **do**
- 3: $\mathbf{x}^{n+1} \leftarrow \arg \min_{\mathbf{x}} \sum_j w_j^n |x_j|^2$ s.t. $\mathbf{A}\mathbf{x} = \mathbf{y}$
- 4: *closed form:* $\mathbf{x}^{n+1} = D^{-1} \mathbf{A}^* (\mathbf{A} D^{-1} \mathbf{A}^*)^{-1} \mathbf{y}, D = \text{diag}(w_j^n)$
- 5: $\varepsilon_{n+1} \leftarrow \min\{\varepsilon_n, \gamma (x^{n+1})_{s+1}^*\}$ $\triangleright (x)_{s+1}^* = (s+1)$ -st largest entry
- 6: $w_j^{n+1} \leftarrow (|x_j^{n+1}|^2 + \varepsilon_{n+1}^2)^{-1/2}$
- 7: **end for**

Theorem 10.5 (Gelfand width lower bound)

If \mathbf{A} satisfies the stable NSP of order s with $\rho < 1$, IRLS with $\gamma = 1/N$ converges. If $\varepsilon_n \rightarrow 0$, $\mathbf{x}^n \rightarrow \mathbf{x}^\#$, the BP solution; for sparse signals, convergence is super-linear with rate $\mu = \rho(1 + \rho)/(1 - \kappa)(1 + 1/(s + 1 - \tilde{s})) < 1$.

Proof. Define the auxiliary functional

$$\mathcal{J}(\mathbf{x}, \mathbf{w}, \varepsilon) = \frac{1}{2} \sum_j w_j x_j^2 + \frac{1}{2} \sum_j (\varepsilon^2 w_j + 1/w_j).$$

Two key inequalities hold. First, the weighted-LS step minimizes $\mathbf{x} \mapsto \mathcal{J}(\mathbf{x}, \mathbf{w}^n, \varepsilon_n)$ subject to $\mathbf{A}\mathbf{x} = \mathbf{y}$, so $\mathcal{J}(\mathbf{x}^{n+1}, \mathbf{w}^n, \varepsilon_n) \leq \mathcal{J}(\mathbf{x}^n, \mathbf{w}^n, \varepsilon_n)$. Second, fixing \mathbf{x}^{n+1} and $\varepsilon_{n+1} \leq \varepsilon_n$, the choice $w_j = (|x_j^{n+1}|^2 + \varepsilon_{n+1}^2)^{-1/2}$ minimizes $\mathbf{w} \mapsto \mathcal{J}$, giving $\mathcal{J}(\mathbf{x}^{n+1}, \mathbf{w}^{n+1}, \varepsilon_{n+1}) \leq \mathcal{J}(\mathbf{x}^{n+1}, \mathbf{w}^n, \varepsilon_{n+1})$. Combined with $\mathcal{J}(\cdot, \cdot, \varepsilon_{n+1}) \leq \mathcal{J}(\cdot, \cdot, \varepsilon_n)$ (since $\varepsilon_{n+1} \leq \varepsilon_n$), the sequence $\mathcal{J}_n := \mathcal{J}(\mathbf{x}^n, \mathbf{w}^n, \varepsilon_n)$ is monotone non-increasing.

At the optimum of \mathbf{w} , $\mathcal{J}(\mathbf{x}, \mathbf{w}^*, \varepsilon) = \sum_j \sqrt{x_j^2 + \varepsilon^2}$, a smoothed ℓ_1 -norm. Stable NSP of order s with $\rho < 1$ provides a coercivity estimate: any feasible \mathbf{x} with $\sum_j \sqrt{x_j^2 + \varepsilon^2}$ bounded has $\|\mathbf{x}\|_2$ bounded uniformly in ε . Hence $\{\mathbf{x}^n\}$ is bounded; pass to a convergent subsequence $\mathbf{x}^{n_k} \rightarrow \mathbf{x}^\#$.

By construction $\mathbf{A}\mathbf{x}^n = \mathbf{y}$ throughout, and the first-order conditions at the optimum of the weighted LS read $\mathbf{A}^* \boldsymbol{\lambda}^n = \text{diag}(w^n) \mathbf{x}^n$. Passing to the limit and using $\varepsilon_n \rightarrow 0$ on the support, the sign vector $\boldsymbol{\lambda}^\#$ becomes a $\partial \|\cdot\|_1$ -subgradient at $\mathbf{x}^\#$, so $\mathbf{x}^\#$ satisfies the BP optimality conditions and equals $\mathbf{x}_{BP}^\#$.

For the super-linear rate, write $\mathbf{x}^n = \mathbf{x}^\# + \mathbf{e}^n$ and linearize the weight update. The map $\mathbf{x}^n \mapsto \mathbf{x}^{n+1}$ becomes a contraction on \mathbf{e}^n near a sparse fixed point; the stable NSP gives a contraction factor $\mu = \rho(1 + \rho)/(1 - \kappa)(1 + 1/(s + 1 - \tilde{s})) < 1$, and the coupling between ε_n and the $(s + 1)$ -st largest entry of \mathbf{x}^{n+1} drives the error super-linearly to zero on the sparse part. \square

10.8 LARS / Homotopy method

Algorithm 15 LARS-Homotopy for the LASSO path

- 1: $\mathbf{x}^0 \leftarrow \mathbf{0}$, $\lambda^0 \leftarrow \|A^*\mathbf{y}\|_\infty$, active set $T_0 \leftarrow \emptyset$
 - 2: $j_1 \leftarrow \arg \max_j |(A^*\mathbf{y})_j|$; $T_1 \leftarrow \{j_1\}$
 - 3: **for** $n = 1, 2, \dots$ **do**
 - 4: Compute direction \mathbf{d}^n on T_n : $\mathbf{d}_{T_n}^n = (A_{T_n}^* A_{T_n})^{-1} \text{sign}(A^*(\mathbf{y} - A\mathbf{x}^{n-1}))_{T_n}$
 - 5: Determine the next breakpoint γ^n either:
 - an inactive index hits the boundary $|(A^*(\mathbf{y} - A\mathbf{x}^{n-1} - \gamma A\mathbf{d}^n))_\ell| = \lambda^{n-1} - \gamma$, or
 - an active coefficient changes sign.
 - 6: Update $\mathbf{x}^n = \mathbf{x}^{n-1} + \gamma^n \mathbf{d}^n$, $\lambda^n = \lambda^{n-1} - \gamma^n$.
 - 7: Update T_{n+1} accordingly (add or remove one index).
 - 8: **end for**
-

Theorem 10.6 (Quotient property and BP stability)

If the BP solution is unique and breakpoints are non-degenerate, LARS terminates in finitely many steps and outputs the full LASSO regularization path. At $\lambda = 0$, the BP solution is recovered.

Proof. The LASSO KKT conditions at parameter λ are: there exists $\mathbf{s} \in \partial\|\mathbf{x}\|_1$ (so $s_j = \text{sign}(x_j)$ if $x_j \neq 0$ and $|s_j| \leq 1$ otherwise) such that $A^*(A\mathbf{x} - \mathbf{y}) + \lambda\mathbf{s} = \mathbf{0}$. Equivalently, the residual correlation $\mathbf{c}(\mathbf{x}) = A^*(\mathbf{y} - A\mathbf{x})$ satisfies $c_j = \lambda \text{sign}(x_j)$ for $j \in T$ (active) and $|c_j| \leq \lambda$ for $j \notin T$ (inactive).

Fix a sign pattern $\boldsymbol{\eta} \in \{\pm 1\}^T$ on the active set T . Within this pattern the KKT system becomes the linear equation $A_T^* A_T \mathbf{x}_T = A_T^* \mathbf{y} - \lambda \boldsymbol{\eta}$, whose solution is $\mathbf{x}_T(\lambda) = (A_T^* A_T)^{-1} (A_T^* \mathbf{y} - \lambda \boldsymbol{\eta})$. Hence $\mathbf{x}^*(\lambda)$ is affine in λ on this regime, and the residual $\mathbf{r}(\lambda) = \mathbf{y} - A_T \mathbf{x}_T(\lambda)$ is also affine, with derivative $A_T (A_T^* A_T)^{-1} \boldsymbol{\eta}$ — the LARS direction \mathbf{d}^n .

Two events end the current regime as λ decreases:

- An inactive coordinate $\ell \notin T$ has $|c_\ell(\lambda)| = \lambda$ — the dual feasibility becomes tight, forcing ℓ into the active set.
- An active coordinate $j \in T$ has $x_j(\lambda) = 0$ — its sign is no longer determined by KKT, forcing it out of the active set.

Both events are linear conditions on λ , so the next breakpoint γ^n is computable in closed form. Under the non-degeneracy assumption exactly one event occurs per breakpoint, so T_{n+1} differs from T_n by a single element.

There are at most $\binom{N}{|T|} \cdot 2^{|T|}$ distinct $(T, \boldsymbol{\eta})$ regimes, and each is visited at most once because λ^n is strictly decreasing along the path. Hence LARS terminates in finitely many steps.

At $\lambda = 0$ the inactive constraints reduce to $A^*(A\mathbf{x} - \mathbf{y}) = \mathbf{0}$ on the active set with no penalty, while the dual constraint $|c_\ell| \leq 0$ on inactive coordinates forces $\mathbf{c} = \mathbf{0}$ when those coordinates have not been reached. Equivalently, the LARS endpoint is feasible for P_1 and the residual sign vector \mathbf{s} certifies optimality, so $\mathbf{x}^*(0)$ is the BP solution. \square

Per-breakpoint cost: $O(s^2 + sN)$; total: $O(s^3 + s^2N)$ for s -sparse outputs.

10.9 Iteratively reweighted ℓ_1

A different reweighting (Candès–Wakin–Boyd 2008): solve

$$\mathbf{x}^{(t+1)} = \arg \min_{\mathbf{x}} \sum_i w_i^{(t)} |x_i| \text{ s.t. } A\mathbf{x} = \mathbf{y}, \quad w_i^{(t+1)} = \frac{1}{|x_i^{(t)}| + \epsilon}.$$

Penalizes small coefficients more heavily on each pass; empirically beats plain BP.

10.10 Robust PCA and outliers

For $\mathbf{y} = A\mathbf{x} + \mathbf{e}$ with sparse \mathbf{e} , $\min \|\mathbf{x}\|_1 + \|\mathbf{e}\|_1$ s.t. $\mathbf{y} = A\mathbf{x} + \mathbf{e}$. The matrix analog: *Robust PCA*,

$$\min \|L\|_* + \lambda \|S\|_1 \quad \text{s.t.} \quad M = L + S,$$

recovers low-rank L from sparse corruption S (Candès, Li, Ma, Wright 2011).

10.11 Software

spgl1, l1_ls, yall1, tfocs, gpsr, slep, cvx/cvxpy/scs for prototyping. For very large problems use FISTA with continuation in λ , or ADMM with warm starts; for adaptive sparsity use IRLS.

Appendices

A Probability Tools

A.1 Markov, Chebyshev, Chernoff

Theorem A.1 (Markov)

For a non-negative r.v. X and $t > 0$, $\Pr(X \geq t) \leq \mathbb{E}X/t$.

Proof. Pointwise we have $X \geq t \mathbf{1}_{\{X \geq t\}}$: on the event $\{X \geq t\}$ both sides exceed t , and on $\{X < t\}$ the right side is 0 while the left is non-negative. Taking expectations of both sides (which preserves the inequality by linearity of \mathbb{E}),

$$\mathbb{E}X \geq t \mathbb{E} \mathbf{1}_{\{X \geq t\}} = t \Pr(X \geq t).$$

Dividing by $t > 0$ gives the result. □

Theorem A.2 (Chebyshev)

For any r.v. X with finite second moment, $\Pr(|X - \mathbb{E}X| \geq t) \leq \text{Var}(X)/t^2$.

Proof. Set $Y = (X - \mathbb{E}X)^2 \geq 0$. The event $\{|X - \mathbb{E}X| \geq t\}$ is the same as $\{Y \geq t^2\}$. Markov's inequality applied to Y with threshold t^2 gives $\Pr(Y \geq t^2) \leq \mathbb{E}Y/t^2 = \text{Var}(X)/t^2$. □

Theorem A.3 (Chernoff)

For any r.v. X , $\Pr(X \geq t) \leq \inf_{\theta > 0} e^{-\theta t} \mathbb{E} e^{\theta X}$.

Proof. Fix $\theta > 0$. The map $u \mapsto e^{\theta u}$ is monotone increasing, so $\{X \geq t\} = \{e^{\theta X} \geq e^{\theta t}\}$. Markov's inequality applied to the non-negative random variable $e^{\theta X}$ with threshold $e^{\theta t}$ yields

$$\Pr(X \geq t) = \Pr(e^{\theta X} \geq e^{\theta t}) \leq e^{-\theta t} \mathbb{E} e^{\theta X}.$$

Since this holds for every $\theta > 0$, take the infimum over the right side. □

A.2 Cramér and Hoeffding

Theorem A.4 (Cramér)

For independent r.v.s X_1, \dots, X_M with cumulant generating functions $C_\ell(\theta) = \ln \mathbb{E} e^{\theta X_\ell}$,

$$\Pr(\sum_\ell X_\ell \geq t) \leq \exp(\inf_{\theta > 0} \{-\theta t + \sum_\ell C_\ell(\theta)\}).$$

Proof. Set $S = \sum_{\ell=1}^M X_\ell$ and apply Chernoff: for any $\theta > 0$, $\Pr(S \geq t) \leq e^{-\theta t} \mathbb{E} e^{\theta S}$. Independence factors the joint MGF:

$$\mathbb{E} e^{\theta S} = \mathbb{E} \prod_{\ell=1}^M e^{\theta X_\ell} = \prod_{\ell=1}^M \mathbb{E} e^{\theta X_\ell} = \prod_{\ell} e^{C_\ell(\theta)} = \exp\left(\sum_{\ell} C_\ell(\theta)\right),$$

the second equality using the product rule for expectations of independent random variables. Hence $\Pr(S \geq t) \leq \exp(-\theta t + \sum_{\ell} C_\ell(\theta))$. Taking the infimum over $\theta > 0$ gives the claimed bound. \square

Theorem A.5 (Hoeffding)

If X_ℓ are independent with $\mathbb{E}X_\ell = 0$, $a_\ell \leq X_\ell \leq b_\ell$, then for $t > 0$,

$$\Pr\left(\sum_{\ell} X_\ell \geq t\right) \leq \exp\left(-\frac{2t^2}{\sum_{\ell} (b_\ell - a_\ell)^2}\right).$$

Proof. We first establish the auxiliary bound (*Hoeffding's lemma*): if $\mathbb{E}X = 0$ and $X \in [a, b]$ a.s., then $\mathbb{E} e^{\theta X} \leq \exp(\theta^2(b - a)^2/8)$ for every $\theta \in \mathbb{R}$. Indeed, for $x \in [a, b]$ write $x = \lambda b + (1 - \lambda)a$ with $\lambda = (x - a)/(b - a) \in [0, 1]$. Convexity of $u \mapsto e^{\theta u}$ gives $e^{\theta x} \leq \lambda e^{\theta b} + (1 - \lambda)e^{\theta a}$. Taking expectation and using $\mathbb{E}X = 0$ (so $\mathbb{E}\lambda = -a/(b - a)$),

$$\mathbb{E} e^{\theta X} \leq \frac{-a}{b - a} e^{\theta b} + \frac{b}{b - a} e^{\theta a}.$$

Set $u = \theta(b - a)$ and $p = -a/(b - a) \in [0, 1]$, so the right side equals $e^{-pu + \ln(1 - p + pe^u)}$. Define $\varphi(u) = -pu + \ln(1 - p + pe^u)$. Then $\varphi(0) = \varphi'(0) = 0$, and $\varphi''(u) = \frac{p(1-p)e^u}{(1-p+pe^u)^2} \leq 1/4$ for all u (the maximum of $q(1 - q)$ is $1/4$). Taylor's theorem gives $\varphi(u) \leq u^2/8$, so $\mathbb{E} e^{\theta X} \leq e^{\theta^2(b-a)^2/8}$.

Now apply Cramér with $C_\ell(\theta) \leq \theta^2(b_\ell - a_\ell)^2/8$:

$$\Pr\left(\sum_{\ell} X_\ell \geq t\right) \leq \exp\left(-\theta t + \frac{\theta^2}{8} \sum_{\ell} (b_\ell - a_\ell)^2\right).$$

Minimize the right side over $\theta > 0$; the optimal choice $\theta^* = 4t / \sum (b_\ell - a_\ell)^2$ produces the stated bound $\exp(-2t^2 / \sum (b_\ell - a_\ell)^2)$. \square

A.3 Subgaussian random variables

Definition A.6 (Subgaussian)

X is subgaussian if $\Pr(|X| \geq t) \leq \beta e^{-\kappa t^2}$ for all $t > 0$, or equivalently, $\mathbb{E} e^{\theta X} \leq e^{c\theta^2}$ for all $\theta \in \mathbb{R}$ (some $c > 0$).

Proposition A.7 (Tail decay)

If $\mathbb{E} e^{\theta X} \leq e^{c\theta^2}$, then $\Pr(|X| \geq t) \leq 2e^{-t^2/(4c)}$.

Proof. By Chernoff, for any $\theta > 0$, $\Pr(X \geq t) \leq e^{-\theta t} \mathbb{E} e^{\theta X} \leq \exp(-\theta t + c\theta^2)$. The exponent $-\theta t + c\theta^2$ is a quadratic in θ minimized at $\theta^* = t/(2c)$, where its value is $-t^2/(4c)$, hence $\Pr(X \geq t) \leq \exp(-t^2/(4c))$. The same argument applied to $-X$ (whose MGF satisfies the same bound by hypothesis allowing both signs of θ) gives $\Pr(X \leq -t) \leq \exp(-t^2/(4c))$. Adding the two tail probabilities yields $\Pr(|X| \geq t) \leq 2 \exp(-t^2/(4c))$. \square

Theorem A.8 (Subgaussian sum)

For independent zero-mean subgaussian X_ℓ with parameter c and weights $\mathbf{a} \in \mathbb{R}^M$, $\Pr(|\sum a_\ell X_\ell| \geq t) \leq 2e^{-t^2/(4c\|\mathbf{a}\|_2^2)}$.

Proof. Each $a_\ell X_\ell$ is itself subgaussian: substituting $\theta \mapsto a_\ell \theta$ in $\mathbb{E} e^{\theta X_\ell} \leq e^{c\theta^2}$ yields $\mathbb{E} e^{\theta a_\ell X_\ell} \leq e^{c a_\ell^2 \theta^2}$. By independence the joint MGF factors:

$$\mathbb{E} e^{\theta \sum_\ell a_\ell X_\ell} = \prod_{\ell=1}^M \mathbb{E} e^{\theta a_\ell X_\ell} \leq \prod_{\ell} e^{c a_\ell^2 \theta^2} = \exp\left(c\theta^2 \sum_\ell a_\ell^2\right) = \exp(c\|\mathbf{a}\|_2^2 \theta^2).$$

So $\sum_\ell a_\ell X_\ell$ is subgaussian with effective parameter $c\|\mathbf{a}\|_2^2$. The previous proposition then gives the tail bound. \square

A.4 Bernstein inequalities

Theorem A.9 (Bernstein)

Let X_1, \dots, X_M be independent zero-mean with $|X_\ell| \leq K$ and $\mathbb{E} X_\ell^2 \leq \sigma_\ell^2$. Set $\sigma^2 = \sum_\ell \sigma_\ell^2$. Then for $t > 0$,

$$\Pr(|\sum_\ell X_\ell| \geq t) \leq 2 \exp\left(-\frac{t^2/2}{\sigma^2 + Kt/3}\right).$$

Proof. We bound each MGF $\mathbb{E} e^{\theta X_\ell}$ for $\theta \in (0, 3/K)$. Taylor-expand and use $\mathbb{E} X_\ell = 0$:

$$\mathbb{E} e^{\theta X_\ell} = 1 + \sum_{n \geq 2} \frac{\theta^n}{n!} \mathbb{E} X_\ell^n.$$

The boundedness $|X_\ell| \leq K$ together with $\mathbb{E} X_\ell^2 \leq \sigma_\ell^2$ gives $\mathbb{E} |X_\ell|^n \leq K^{n-2} \mathbb{E} X_\ell^2 \leq K^{n-2} \sigma_\ell^2$ for $n \geq 2$. Therefore

$$\sum_{n \geq 2} \frac{\theta^n}{n!} \mathbb{E} |X_\ell|^n \leq \sigma_\ell^2 \sum_{n \geq 2} \frac{\theta^n K^{n-2}}{n!} \leq \sigma_\ell^2 \cdot \frac{\theta^2}{2} \sum_{n \geq 0} \frac{(K\theta)^n}{3^n} = \frac{\theta^2 \sigma_\ell^2 / 2}{1 - K\theta/3},$$

where the inequality $n!/2 \geq 3^{n-2}$ for $n \geq 2$ was used. Combining with $1 + u \leq e^u$ yields

$$\mathbb{E} e^{\theta X_\ell} \leq \exp\left(\frac{\theta^2 \sigma_\ell^2 / 2}{1 - K\theta/3}\right).$$

Apply Cramér's theorem with $C_\ell(\theta) \leq \theta^2 \sigma_\ell^2 / (2(1 - K\theta/3))$:

$$\Pr(\sum_\ell X_\ell \geq t) \leq \exp\left(-\theta t + \frac{\theta^2 \sigma^2 / 2}{1 - K\theta/3}\right).$$

Choose $\theta = t/(\sigma^2 + Kt/3) \in (0, 3/K)$. Then $1 - K\theta/3 = \sigma^2/(\sigma^2 + Kt/3)$, and substituting gives

$$-\theta t + \frac{\theta^2 \sigma^2 / 2}{1 - K\theta/3} = -\theta t + \frac{\theta^2 (\sigma^2 + Kt/3)}{2} = -\frac{t^2 / 2}{\sigma^2 + Kt/3}.$$

This bounds $\Pr(\sum X_\ell \geq t)$. The lower tail follows from the same argument applied to $-X_\ell$ (which satisfies the same hypotheses), and adding the two tails produces the factor of 2 in front. \square

Theorem A.10 (Noncommutative Bernstein, Tropp 2012)

For independent zero-mean self-adjoint random matrices $X_\ell \in \mathbb{C}^{d \times d}$ with $\lambda_{\max}(X_\ell) \leq K$ a.s. and $\sigma^2 := \|\sum_\ell \mathbb{E} X_\ell^2\|_{2 \rightarrow 2}$,

$$\Pr\left(\left\|\sum_\ell X_\ell\right\|_{2 \rightarrow 2} \geq t\right) \leq d \exp\left(-\frac{t^2 / 2}{\sigma^2 + Kt/3}\right).$$

Proof. Markov's inequality lifts to self-adjoint matrices via the trace: for any self-adjoint S ,

$$\Pr(\lambda_{\max}(S) \geq t) \leq \Pr(\text{tr } e^{\theta S} \geq e^{\theta t}) \leq e^{-\theta t} \mathbb{E} \text{tr } e^{\theta S}, \quad \theta > 0,$$

using $\text{tr } e^{\theta S} \geq e^{\theta \lambda_{\max}(S)}$. Apply this to $S = \sum_\ell X_\ell$. The independence step requires Lieb's concavity theorem (Appendix B): for any fixed self-adjoint H , the map $A \mapsto \text{tr} \exp(H + \ln A)$ is concave on positive-definite A . By Jensen, $\mathbb{E} \text{tr} \exp(H + \ln Y) \leq \text{tr} \exp(H + \ln \mathbb{E} Y)$ for any random positive-definite Y . Iterating across the M independent terms (with $H = \theta \sum_{\ell < j} X_\ell + \theta \sum_{\ell > j} X_\ell$ at the j -th step, treating $Y = e^{\theta X_j}$),

$$\mathbb{E} \text{tr } e^{\theta \sum_\ell X_\ell} \leq \text{tr} \exp\left(\sum_\ell \ln \mathbb{E} e^{\theta X_\ell}\right).$$

We bound each matrix MGF $\mathbb{E} e^{\theta X_\ell}$ in the semidefinite order. The same Taylor-series-with-moments argument as in the scalar Bernstein proof goes through: $|X_\ell^n| \preceq K^{n-2} X_\ell^2$ for $n \geq 2$, giving $\mathbb{E} X_\ell^n \preceq K^{n-2} \mathbb{E} X_\ell^2$, hence

$$\mathbb{E} e^{\theta X_\ell} \preceq I + \frac{\theta^2 / 2}{1 - K\theta/3} \mathbb{E} X_\ell^2 \preceq \exp\left(\frac{\theta^2 / 2}{1 - K\theta/3} \mathbb{E} X_\ell^2\right),$$

for $0 < \theta < 3/K$, using $I + A \preceq e^A$ for $A \succeq 0$. Taking matrix log and summing over ℓ , $\sum_\ell \ln \mathbb{E} e^{\theta X_\ell} \preceq \frac{\theta^2 / 2}{1 - K\theta/3} \sum_\ell \mathbb{E} X_\ell^2$. Hence

$$\text{tr} \exp\left(\sum_\ell \ln \mathbb{E} e^{\theta X_\ell}\right) \leq d \exp\left(\lambda_{\max}\left(\sum_\ell \ln \mathbb{E} e^{\theta X_\ell}\right)\right) \leq d \exp\left(\frac{\theta^2 / 2}{1 - K\theta/3} \sigma^2\right),$$

the first inequality being $\text{tr } e^A \leq d e^{\lambda_{\max}(A)}$. Combining,

$$\Pr(\|\sum_\ell X_\ell\|_{2 \rightarrow 2} \geq t) \leq d \exp\left(-\theta t + \frac{\theta^2 \sigma^2 / 2}{1 - K\theta/3}\right).$$

Choosing $\theta = t/(\sigma^2 + Kt/3)$ as in the scalar Bernstein proof gives the stated bound. The two-sided result (operator norm rather than λ_{\max}) follows by a union bound on $\pm S$. \square

A.5 Concentration of measure

Theorem A.11 (Gaussian concentration)

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be L -Lipschitz and $\mathbf{g} \sim \mathcal{N}(0, I_n)$. Then for $t > 0$,

$$\Pr(|f(\mathbf{g}) - \mathbb{E}f(\mathbf{g})| \geq t) \leq 2 \exp(-t^2/(2L^2)).$$

Proof. By rescaling assume $L = 1$. We show the subgaussian MGF bound $\mathbb{E} e^{\lambda(f(\mathbf{g}) - \mathbb{E}f)} \leq e^{\lambda^2 \pi^2/8}$ for all $\lambda \in \mathbb{R}$; the tail bound then follows from Proposition (Tail decay) and a constant adjustment.

Let \mathbf{g}' be an independent copy of \mathbf{g} and rotate: $\mathbf{u}(\theta) = \cos \theta \cdot \mathbf{g} + \sin \theta \cdot \mathbf{g}'$, $\mathbf{u}'(\theta) = -\sin \theta \cdot \mathbf{g} + \cos \theta \cdot \mathbf{g}'$ for $\theta \in [0, \pi/2]$. By rotational invariance of the standard Gaussian, $(\mathbf{u}(\theta), \mathbf{u}'(\theta))$ is jointly distributed as $(\mathbf{g}, \mathbf{g}')$ for every θ . Note $(\mathbf{u}(0), \mathbf{u}'(0)) = (\mathbf{g}, \mathbf{g}')$ and $(\mathbf{u}(\pi/2), \mathbf{u}'(\pi/2)) = (\mathbf{g}', -\mathbf{g})$.

Fix $\lambda \in \mathbb{R}$. Differentiating $f(\mathbf{u}(\theta))$ in θ and using the chain rule plus Cauchy-Schwarz with $|\nabla f| \leq 1$,

$$f(\mathbf{g}') - f(\mathbf{g}) = \int_0^{\pi/2} \langle \nabla f(\mathbf{u}(\theta)), \mathbf{u}'(\theta) \rangle d\theta.$$

Take $\mathbb{E} e^{\lambda \cdot}$ of both sides and apply Jensen's inequality (after rescaling the θ -integral to a probability measure on $[0, \pi/2]$):

$$\mathbb{E} e^{\lambda(f(\mathbf{g}') - f(\mathbf{g}))} \leq \frac{2}{\pi} \int_0^{\pi/2} \mathbb{E} e^{\frac{\lambda \pi}{2} \langle \nabla f(\mathbf{u}(\theta)), \mathbf{u}'(\theta) \rangle} d\theta.$$

Conditioning on $\mathbf{u}(\theta)$ (which makes $\nabla f(\mathbf{u}(\theta))$ a fixed vector) and using that $\mathbf{u}'(\theta) \sim \mathcal{N}(0, I_n)$ independently, the inner expectation equals $\exp(\frac{\lambda^2 \pi^2}{8} |\nabla f(\mathbf{u}(\theta))|^2) \leq e^{\lambda^2 \pi^2/8}$. By symmetry $\mathbb{E} e^{\lambda(f(\mathbf{g}) - \mathbb{E}f)} \leq \sqrt{\mathbb{E} e^{\lambda(f(\mathbf{g}') - f(\mathbf{g}))}} \leq e^{\lambda^2 \pi^2/16}$. Optimizing the resulting Chernoff bound and absorbing the $\pi^2/16$ constant yields a Gaussian-tail bound with constant $\pi^2/16$; the sharper textbook constant $1/2$ follows from a refinement using Borell's isoperimetric inequality, but the argument above already gives the qualitative form $\Pr(|f(\mathbf{g}) - \mathbb{E}f| \geq t) \leq 2e^{-ct^2/L^2}$. \square

A.6 Khintchine, Slepian, Gordon, Dudley

Lemma A.12 (Symmetrization)

For independent random vectors ξ_ℓ in a normed space and a convex F , $\mathbb{E}F(\sum(\xi_\ell - \mathbb{E}\xi_\ell)) \leq \mathbb{E}F(2\sum \varepsilon_\ell \xi_\ell)$, where ε_ℓ are independent Rademacher signs.

Proof. Let $\{\xi'_\ell\}$ be an independent copy of $\{\xi_\ell\}$, so $\mathbb{E}\xi'_\ell = \mathbb{E}\xi_\ell$. By Jensen's inequality applied to the

conditional expectation given ξ_1, \dots, ξ_M ,

$$F\left(\sum_{\ell}(\xi_{\ell} - \mathbb{E}\xi_{\ell})\right) = F\left(\mathbb{E}_{\xi'}\left[\sum_{\ell}(\xi_{\ell} - \xi'_{\ell})\right]\right) \leq \mathbb{E}_{\xi'}F\left(\sum_{\ell}(\xi_{\ell} - \xi'_{\ell})\right).$$

Taking expectation over ξ as well yields $\mathbb{E}F(\sum_{\ell}(\xi_{\ell} - \mathbb{E}\xi_{\ell})) \leq \mathbb{E}F(\sum_{\ell}(\xi_{\ell} - \xi'_{\ell}))$. The vector $\xi_{\ell} - \xi'_{\ell}$ is symmetric (its distribution is invariant under multiplication by -1), and the M -tuple of these differences is jointly symmetric in each coordinate independently. So multiplying each by an independent Rademacher ε_{ℓ} preserves the joint law: $\sum_{\ell}(\xi_{\ell} - \xi'_{\ell})$ has the same distribution as $\sum_{\ell}\varepsilon_{\ell}(\xi_{\ell} - \xi'_{\ell})$. Now use the triangle inequality (and convexity/symmetry of F):

$$F\left(\sum_{\ell}\varepsilon_{\ell}(\xi_{\ell} - \xi'_{\ell})\right) \leq \frac{1}{2}F\left(2\sum_{\ell}\varepsilon_{\ell}\xi_{\ell}\right) + \frac{1}{2}F\left(-2\sum_{\ell}\varepsilon_{\ell}\xi'_{\ell}\right).$$

Taking expectations and using that $\sum_{\ell}\varepsilon_{\ell}\xi_{\ell}$ and $\sum_{\ell}\varepsilon_{\ell}\xi'_{\ell}$ are identically distributed, the right-hand side reduces to $\mathbb{E}F(2\sum_{\ell}\varepsilon_{\ell}\xi_{\ell})$. □

Theorem A.13 (Khintchine)

For $\mathbf{a} \in \mathbb{C}^M$ and Rademacher $\varepsilon_1, \dots, \varepsilon_M$,

$$\left(\mathbb{E}\left|\sum_{\ell}\varepsilon_{\ell}a_{\ell}\right|^p\right)^{1/p} \asymp \sqrt{p}\|\mathbf{a}\|_2, \quad p \geq 1.$$

Consequence: $\Pr\left(\left|\sum_{\ell}\varepsilon_{\ell}a_{\ell}\right| \geq u\|\mathbf{a}\|_2\right) \leq 2e^{-u^2/2}$.

Proof. By splitting real and imaginary parts, assume $\mathbf{a} \in \mathbb{R}^M$. Write $S = \sum_{\ell}\varepsilon_{\ell}a_{\ell}$. For the MGF, factor by independence:

$$\mathbb{E}e^{\theta S} = \prod_{\ell}\mathbb{E}e^{\theta\varepsilon_{\ell}a_{\ell}} = \prod_{\ell}\cosh(\theta a_{\ell}) \leq \prod_{\ell}e^{\theta^2 a_{\ell}^2/2} = e^{\theta^2\|\mathbf{a}\|_2^2/2},$$

using $\cosh u \leq e^{u^2/2}$ (which follows from comparing Taylor series: $\cosh u = \sum_{n \geq 0} u^{2n}/(2n)! \leq \sum_{n \geq 0} (u^2/2)^n/n! = e^{u^2/2}$ since $(2n)! \geq 2^n n!$). Thus S is subgaussian with parameter $\|\mathbf{a}\|_2^2/2$, and the tail bound $\Pr(|S| \geq u\|\mathbf{a}\|_2) \leq 2e^{-u^2/2}$ follows from Proposition (Tail decay).

For the moment statement, the upper bound $(\mathbb{E}|S|^p)^{1/p} \leq c_1\sqrt{p}\|\mathbf{a}\|_2$ follows from the subgaussian MGF: integrating $\Pr(|S| \geq u\|\mathbf{a}\|_2) \leq 2e^{-u^2/2}$ via $\mathbb{E}|S|^p = p \int_0^{\infty} u^{p-1} \Pr(|S| \geq u) du$ gives the bound after a change of variables. The lower bound $(\mathbb{E}|S|^p)^{1/p} \geq c_2\sqrt{p}\|\mathbf{a}\|_2$ is sharper and uses the explicit formula for even moments $\mathbb{E}S^{2n} = \sum_{k_1+\dots+k_M=2n} \binom{2n}{2k_1, \dots, 2k_M} \prod_{\ell} a_{\ell}^{2k_{\ell}}$, which by AM-GM is at least $(2n)!/(2^n n!) \cdot \|\mathbf{a}\|_2^{2n}$. Stirling gives $(2n)!/(2^n n!) \sim \sqrt{2}(2n/e)^n$, so $(\mathbb{E}S^{2n})^{1/(2n)} \geq c\sqrt{n}\|\mathbf{a}\|_2$. Hölder interpolates between $p = 2n$ and $p = 2(n+1)$. □

Theorem A.14 (Slepian's lemma)

Let \mathbf{X}, \mathbf{Y} be mean-zero Gaussian vectors with $\mathbb{E}|X_i - X_j|^2 \leq \mathbb{E}|Y_i - Y_j|^2$ for all i, j . Then $\mathbb{E}\max_i X_i \leq \mathbb{E}\max_i Y_i$.

Proof. We may assume \mathbf{X}, \mathbf{Y} are independent. Define the interpolation $\mathbf{Z}(t) = \sqrt{t}\mathbf{Y} + \sqrt{1-t}\mathbf{X}$ for $t \in [0, 1]$, so $\mathbf{Z}(0) = \mathbf{X}$, $\mathbf{Z}(1) = \mathbf{Y}$, and $\mathbf{Z}(t)$ is a centered Gaussian vector with covariance $K(t) = tK_Y + (1-t)K_X$.

The function $z \mapsto \max_i z_i$ is non-smooth, so we work with the smoothed maximum $\Phi_\beta(\mathbf{z}) = \frac{1}{\beta} \log \sum_i e^{\beta z_i}$, which satisfies $\max_i z_i \leq \Phi_\beta(\mathbf{z}) \leq \max_i z_i + \frac{\log n}{\beta}$, and is smooth and convex with $\partial_i \Phi_\beta(\mathbf{z}) = \frac{e^{\beta z_i}}{\sum_k e^{\beta z_k}}$.

Compute $\frac{d}{dt} \mathbb{E} \Phi_\beta(\mathbf{Z}(t))$. By the multivariate Gaussian integration-by-parts identity (Stein's lemma applied componentwise),

$$\frac{d}{dt} \mathbb{E} \Phi_\beta(\mathbf{Z}(t)) = \frac{1}{2} \sum_{i,j} ((K_Y)_{ij} - (K_X)_{ij}) \mathbb{E} \partial_i \partial_j \Phi_\beta(\mathbf{Z}(t)).$$

Direct computation gives $\partial_i \partial_j \Phi_\beta = \beta(p_i \delta_{ij} - p_i p_j)$ where $p_i = e^{\beta z_i} / \sum_k e^{\beta z_k} \geq 0$. The hypothesis $\mathbb{E}|X_i - X_j|^2 \leq \mathbb{E}|Y_i - Y_j|^2$ rewrites as $(K_Y)_{ii} + (K_Y)_{jj} - 2(K_Y)_{ij} \geq (K_X)_{ii} + (K_X)_{jj} - 2(K_X)_{ij}$. A direct (if tedious) algebraic check shows the derivative above is non-negative under this hypothesis, so $t \mapsto \mathbb{E} \Phi_\beta(\mathbf{Z}(t))$ is non-decreasing. Hence $\mathbb{E} \Phi_\beta(\mathbf{X}) \leq \mathbb{E} \Phi_\beta(\mathbf{Y})$, and letting $\beta \rightarrow \infty$ gives $\mathbb{E} \max_i X_i \leq \mathbb{E} \max_i Y_i$. \square

Theorem A.15 (Gordon's min-max lemma)

Under doubly-indexed Gaussian comparison hypotheses, $\mathbb{E} \min_j \max_i X_{i,j} \geq \mathbb{E} \min_j \max_i Y_{i,j}$.

Proof. Specifically, suppose $\{X_{i,j}\}, \{Y_{i,j}\}$ are centered Gaussian families satisfying $\mathbb{E}|X_{i,j} - X_{i,k}|^2 \leq \mathbb{E}|Y_{i,j} - Y_{i,k}|^2$ for all $i, j \neq k$ (within the same row, X -increments are smaller), and $\mathbb{E}|X_{i,j} - X_{\ell,k}|^2 \geq \mathbb{E}|Y_{i,j} - Y_{\ell,k}|^2$ for all $i \neq \ell, j, k$ (across rows, X -increments are larger). Then $\mathbb{E} \min_j \max_i X_{i,j} \geq \mathbb{E} \min_j \max_i Y_{i,j}$.

Run the same Gaussian interpolation $\mathbf{Z}(t) = \sqrt{t}\mathbf{X} + \sqrt{1-t}\mathbf{Y}$ as in Slepian's lemma. Replace $\Phi_\beta(\mathbf{z}) = \beta^{-1} \log \sum_i e^{\beta z_i}$ by the doubly-smoothed soft-min-max $\Psi_{\alpha,\beta}(\mathbf{z}) = -\alpha^{-1} \log \sum_j \exp(-\alpha \beta^{-1} \log \sum_i e^{\beta z_{i,j}})$, which converges to $\min_j \max_i z_{i,j}$ as $\alpha, \beta \rightarrow \infty$ and is jointly smooth.

Differentiating $\mathbb{E} \Psi_{\alpha,\beta}(\mathbf{Z}(t))$ in t via Gaussian integration by parts yields a sum of covariance differences $((K_X)_{(i,j),(\ell,k)} - (K_Y)_{(i,j),(\ell,k)})$ multiplied by mixed second partials of $\Psi_{\alpha,\beta}$. The mixed partials separate cleanly into "same-row" (positive coefficient) and "cross-row" (negative coefficient) contributions; the two sign-flipped hypotheses ensure the entire sum is non-positive, so $t \mapsto \mathbb{E} \Psi_{\alpha,\beta}(\mathbf{Z}(t))$ is non-increasing. Therefore $\mathbb{E} \Psi_{\alpha,\beta}(\mathbf{X}) \geq \mathbb{E} \Psi_{\alpha,\beta}(\mathbf{Y})$, and letting $\alpha, \beta \rightarrow \infty$ gives the result. \square

Theorem A.16 (Dudley's entropy integral)

For a centered subgaussian process $\{X_t\}_{t \in T}$ with pseudo-metric $d(s, t) = \|X_s - X_t\|_{L^2}$ and covering numbers $N(T, d, u)$,

$$\mathbb{E} \sup_{t \in T} X_t \leq 12 \int_0^{\Delta(T)/2} \sqrt{\ln N(T, d, u)} du, \quad \Delta(T) = \sup_{s,t} d(s, t).$$

Proof. We use the *chaining* method. Set $\Delta = \Delta(T)$ and $\varepsilon_n = \Delta 2^{-n}$ for $n \geq 0$. By definition there exists a ε_n -net $T_n \subseteq T$ with $|T_n| \leq N(T, d, \varepsilon_n)$; choose T_0 to be a single point (so $|T_0| = 1$). Let $\pi_n : T \rightarrow T_n$ be a

nearest-neighbor projection in the metric d , so $d(t, \pi_n(t)) \leq \varepsilon_n$ and $d(\pi_n(t), \pi_{n-1}(t)) \leq \varepsilon_n + \varepsilon_{n-1} = 3\varepsilon_n$.

By continuity of the process (or by truncation) we have the telescoping representation

$$X_t - X_{\pi_0(t)} = \sum_{n \geq 1} (X_{\pi_n(t)} - X_{\pi_{n-1}(t)}).$$

Each increment $X_{\pi_n(t)} - X_{\pi_{n-1}(t)}$ is subgaussian with parameter $d(\pi_n(t), \pi_{n-1}(t))^2 \leq 9\varepsilon_n^2$. The number of such pairs across all $t \in T$ is at most $|T_n| \cdot |T_{n-1}| \leq N(T, d, \varepsilon_n)^2$. Applying the standard maximal-subgaussian bound $\mathbb{E} \max_{j \leq M} |Z_j| \leq c \sigma \sqrt{\log M}$ (for subgaussian Z_j with parameter σ^2) gives

$$\mathbb{E} \sup_{t \in T} |X_{\pi_n(t)} - X_{\pi_{n-1}(t)}| \leq c \cdot 3\varepsilon_n \sqrt{2 \log N(T, d, \varepsilon_n)}.$$

Summing the telescoped chain and using $\varepsilon_n = 2(\varepsilon_n - \varepsilon_{n+1})$ to convert the geometric series into the Riemann-integral form,

$$\mathbb{E} \sup_t X_t \leq \sum_{n \geq 1} c' \varepsilon_n \sqrt{\log N(T, d, \varepsilon_n)} \leq 12 \int_0^{\Delta/2} \sqrt{\log N(T, d, u)} du,$$

with the constant 12 absorbing the chaining overhead and the conversion factor. □

B Convex Analysis

B.1 Convex sets and functions

Definition B.1 (Convex set and convex hull)

$K \subseteq \mathbb{R}^N$ is *convex* if $tx + (1 - t)y \in K$ for all $\mathbf{x}, \mathbf{y} \in K, t \in [0, 1]$. The *convex hull* of $T \subset \mathbb{R}^N$ is $\text{conv}(T) = \{\sum_j t_j \mathbf{x}_j : t_j \geq 0, \sum t_j = 1, \mathbf{x}_j \in T\}$.

Definition B.2 (Convex and strongly convex function)

$F : \mathbb{R}^N \rightarrow (-\infty, \infty]$ is *convex* if for all $\mathbf{x}, \mathbf{y}, t \in [0, 1]$, $F(t\mathbf{x} + (1 - t)\mathbf{y}) \leq tF(\mathbf{x}) + (1 - t)F(\mathbf{y})$. *Strongly convex* with parameter $\gamma > 0$ adds the term $-\frac{\gamma}{2}t(1 - t)\|\mathbf{x} - \mathbf{y}\|_2^2$ on the right side.

Proposition B.3 (First/second order characterizations)

For differentiable F :

- Convex iff $F(\mathbf{x}) \geq F(\mathbf{y}) + \langle \nabla F(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$.
- Strongly convex with γ iff above with $+\frac{\gamma}{2}\|\mathbf{x} - \mathbf{y}\|_2^2$.
- Twice differentiable: convex iff $\nabla^2 F \succeq 0$.

Theorem B.4 (Hyperplane separation)

Disjoint convex sets K_1, K_2 (with non-empty interiors disjoint) admit a separating hyperplane $\{\mathbf{x} : \langle \mathbf{x}, \mathbf{w} \rangle = \lambda\}$ such that $K_1 \subseteq \{\leq \lambda\}, K_2 \subseteq \{\geq \lambda\}$.

Proof. Assume K_1, K_2 are closed (otherwise replace by closures and the argument carries through after a limiting argument). The Minkowski difference $C := K_1 - K_2 = \{\mathbf{a} - \mathbf{b} : \mathbf{a} \in K_1, \mathbf{b} \in K_2\}$ is convex (sum of convex sets) and non-empty. Disjointness $K_1 \cap K_2 = \emptyset$ is equivalent to $\mathbf{0} \notin C$. Assume further that C is closed (true if at least one of K_1, K_2 is compact); the proof in the general case proceeds by separating $\mathbf{0}$ from successively-shrunk closed approximations and taking a weak limit.

Let $\mathbf{w}^* = \arg \min_{\mathbf{c} \in C} \|\mathbf{c}\|_2$, which exists by closedness and convexity (the projection onto a closed convex set is unique). Since $\mathbf{0} \notin C$, $\mathbf{w}^* \neq \mathbf{0}$. The first-order optimality of the projection gives, for every $\mathbf{c} \in C$, $\langle \mathbf{w}^*, \mathbf{c} - \mathbf{w}^* \rangle \geq 0$, i.e. $\langle \mathbf{w}^*, \mathbf{c} \rangle \geq \|\mathbf{w}^*\|_2^2 > 0$. Substituting $\mathbf{c} = \mathbf{a} - \mathbf{b}$, $\langle \mathbf{w}^*, \mathbf{a} \rangle - \langle \mathbf{w}^*, \mathbf{b} \rangle \geq \|\mathbf{w}^*\|_2^2$ for every $\mathbf{a} \in K_1, \mathbf{b} \in K_2$. Set $\alpha = \inf_{\mathbf{a} \in K_1} \langle \mathbf{w}^*, \mathbf{a} \rangle$ and $\beta = \sup_{\mathbf{b} \in K_2} \langle \mathbf{w}^*, \mathbf{b} \rangle$. The displayed inequality gives $\alpha \geq \beta + \|\mathbf{w}^*\|_2^2 > \beta$. Pick any $\lambda \in [\beta, \alpha]$: then $K_1 \subseteq \{\langle \mathbf{w}^*, \cdot \rangle \geq \lambda\}$ and $K_2 \subseteq \{\langle \mathbf{w}^*, \cdot \rangle \leq \lambda\}$. Replacing \mathbf{w}^* by $-\mathbf{w}^*$ if needed gives the stated form. \square

B.2 Convex conjugate

Definition B.5 (Fenchel conjugate)

$$F^*(\mathbf{y}) = \sup_{\mathbf{x}} \{\langle \mathbf{x}, \mathbf{y} \rangle - F(\mathbf{x})\}.$$

Always convex (as a sup of affine functions). *Fenchel–Young inequality:* $\langle \mathbf{x}, \mathbf{y} \rangle \leq F(\mathbf{x}) + F^*(\mathbf{y})$. For closed proper convex F , $F^{**} = F$.

Example B.6 (Norm and dual norm)

$$(\|\cdot\|)^*(\mathbf{y}) = \iota_{\|\mathbf{y}\|_* \leq 1}(\mathbf{y}). \text{ In particular, } (\|\cdot\|_1)^* = \iota_{\|\cdot\|_\infty \leq 1}.$$

Example B.7 (Quadratic)

$$(\frac{1}{2}\|\cdot\|_2^2)^*(\mathbf{y}) = \frac{1}{2}\|\mathbf{y}\|_2^2.$$

B.3 Subdifferential

Definition B.8 (Subdifferential)

$$\partial F(\mathbf{x}) = \{\mathbf{g} : F(\mathbf{z}) \geq F(\mathbf{x}) + \langle \mathbf{g}, \mathbf{z} - \mathbf{x} \rangle \forall \mathbf{z}\}.$$

If F is differentiable at \mathbf{x} , $\partial F(\mathbf{x}) = \{\nabla F(\mathbf{x})\}$.

Example B.9 (Subdifferential of $\|\cdot\|_1$)

$$\partial\|x\|_1 = \{g : g_i = \text{sign}(x_i) \text{ if } x_i \neq 0, |g_i| \leq 1 \text{ if } x_i = 0\}.$$

Theorem B.10 (Subgradient optimality condition)

x^* minimizes F iff $0 \in \partial F(x^*)$.

Proof. (\Leftarrow) Suppose $0 \in \partial F(x^*)$. The defining subgradient inequality with subgradient $g = 0$ reads $F(z) \geq F(x^*) + \langle 0, z - x^* \rangle = F(x^*)$ for every z , so x^* is a global minimizer.

(\Rightarrow) Conversely, if x^* is a global minimizer then $F(z) - F(x^*) \geq 0 = \langle 0, z - x^* \rangle$ for every z , which is exactly the statement that 0 is a subgradient at x^* . \square

Theorem B.11 (Conjugate-subgradient inversion)

$y \in \partial F(x)$ iff $F(x) + F^*(y) = \langle x, y \rangle$, iff $x \in \partial F^*(y)$ (when F is closed proper convex).

Proof. Recall $F^*(y) = \sup_z \{\langle z, y \rangle - F(z)\}$, hence the Fenchel–Young inequality $F(x) + F^*(y) \geq \langle x, y \rangle$ always holds.

($y \in \partial F(x) \Rightarrow F(x) + F^*(y) = \langle x, y \rangle$): the subgradient inequality $F(z) \geq F(x) + \langle y, z - x \rangle$ rearranges to $\langle z, y \rangle - F(z) \leq \langle x, y \rangle - F(x)$ for all z . Taking the supremum on the left gives $F^*(y) \leq \langle x, y \rangle - F(x)$. Combined with Fenchel–Young this is equality.

($F(x) + F^*(y) = \langle x, y \rangle \Rightarrow y \in \partial F(x)$): for any z , $\langle z, y \rangle - F(z) \leq F^*(y) = \langle x, y \rangle - F(x)$, rearranging to $F(z) \geq F(x) + \langle y, z - x \rangle$, which is the subgradient inequality.

Finally, by the biconjugate theorem $F^{**} = F$ for closed proper convex F , the equality $F(x) + F^*(y) = \langle x, y \rangle$ is symmetric in (x, y) if we replace F by $F^{**} = F$ on the left. Reading the same derivation with the roles of x and y swapped (and $F \leftrightarrow F^*$) shows the equality is equivalent to $x \in \partial F^*(y)$. \square

B.4 Lagrangian duality and the BP dual certificate

For

$$\min_x F_0(x) \text{ s.t. } Ax = y, F_j(x) \leq b_j,$$

the Lagrangian is $L(x, \xi, \nu) = F_0 + \xi^T(Ax - y) + \sum_j \nu_j(F_j - b_j)$, $\nu_j \geq 0$. The dual function $H(\xi, \nu) = \inf_x L$ is concave; *weak duality* $H(\xi^*, \nu^*) \leq F_0(x^*)$ always holds. *Strong duality* (equality) holds for convex F_j under Slater’s condition (interior point feasible).

Theorem B.12 (BP primal–dual pair)

Primal: $\min \|x\|_1$ s.t. $Ax = y$. Dual:

$$\max_{\xi} -\langle \xi, y \rangle \text{ s.t. } \|A^* \xi\|_{\infty} \leq 1.$$

Theorem B.13 (Dual certificate \Rightarrow unique BP solution)

Suppose x^* is feasible with support S , A_S is injective, and there exists ν^* with

(i) $(A^* \nu^*)_j = \text{sign}(x_j^*)$ for $j \in S$,

(ii) $|(A^* \nu^*)_{\ell}| < 1$ for $\ell \notin S$.

Then x^* is the unique BP minimizer.

Proof. Conditions (i)–(ii) say $u = A^* \nu^* \in \partial \|x^*\|_1$. For any feasible x , $x - x^* \in \ker A$, so

$$\|x\|_1 \geq \|x^*\|_1 + \langle u, x - x^* \rangle = \|x^*\|_1 + \langle \nu^*, A(x - x^*) \rangle = \|x^*\|_1.$$

The strict inequality on \bar{S} implies any $x \neq x^*$ with the same support has cost strictly larger; combined with injectivity of A_S , uniqueness follows. \square

C Matrix Analysis

C.1 Norms and the SVD

Theorem C.1 (Singular Value Decomposition)

For any $A \in \mathbb{C}^{m \times N}$ there exist unitaries $U \in \mathbb{C}^{m \times m}$, $V \in \mathbb{C}^{N \times N}$ and unique $\sigma_1 \geq \dots \geq \sigma_r > 0$ ($r = \text{rank } A$) with

$$A = U \Sigma V^*, \quad \Sigma = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots).$$

Proof. The matrix $A^*A \in \mathbb{C}^{N \times N}$ is Hermitian and positive semidefinite, so by the spectral theorem there exist a unitary $V \in \mathbb{C}^{N \times N}$ and diagonal $D = \text{diag}(\lambda_1, \dots, \lambda_N)$ with $\lambda_1 \geq \dots \geq \lambda_N \geq 0$ such that $A^*A = V D V^*$. Let r be the number of non-zero λ_i ; this equals $\text{rank}(A^*A) = \text{rank}(A)$. Define $\sigma_i := \sqrt{\lambda_i}$ for $i = 1, \dots, r$.

Let v_1, \dots, v_N be the columns of V . For $i \leq r$ define $u_i := \sigma_i^{-1} A v_i \in \mathbb{C}^m$. These are orthonormal: for $i, j \leq r$,

$$u_i^* u_j = \frac{1}{\sigma_i \sigma_j} v_i^* A^* A v_j = \frac{1}{\sigma_i \sigma_j} v_i^* (V D V^*) v_j = \frac{\lambda_j}{\sigma_i \sigma_j} \delta_{ij} = \delta_{ij},$$

using $V^* v_j = e_j$ and $D e_j = \lambda_j e_j$. Extend $\{u_1, \dots, u_r\}$ to an orthonormal basis $\{u_1, \dots, u_m\}$ of \mathbb{C}^m and form $U = [u_1 \dots u_m]$.

For $i > r$ we have $\lambda_i = 0$, so $\|A v_i\|_2^2 = v_i^* A^* A v_i = 0$, i.e. $A v_i = \mathbf{0}$. Hence $A = A V V^* = (\sum_{i \leq r} \sigma_i u_i v_i^*)$ by direct computation. Writing Σ as the $m \times N$ matrix with σ_i on positions (i, i) for $i \leq r$ and zeros

elsewhere, this is exactly $A = U\Sigma V^*$. Uniqueness of σ_i follows from uniqueness of the eigenvalues of A^*A . \square

Theorem C.2 (Eckart–Young)

The best rank- r approximation of A in operator (or Frobenius) norm is $A_r = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^*$, with $\|A - A_r\|_{2 \rightarrow 2} = \sigma_{r+1}$ and $\|A - A_r\|_F = \sqrt{\sum_{i>r} \sigma_i^2}$.

Proof. Let $A = U\Sigma V^*$ be the SVD. The truncation $A_r = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^*$ has rank r and $A - A_r = \sum_{i>r} \sigma_i \mathbf{u}_i \mathbf{v}_i^*$, an SVD of $A - A_r$ with non-zero singular values $\sigma_{r+1}, \sigma_{r+2}, \dots$. Therefore $\|A - A_r\|_{2 \rightarrow 2} = \sigma_{r+1}$ and $\|A - A_r\|_F^2 = \sum_{i>r} \sigma_i^2$.

We now show optimality. Let B be any matrix of rank at most r . The kernel of B has dimension at least $N - r$. The subspace $W = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_{r+1}\}$ has dimension $r + 1$, so by dimension counting in the ambient space \mathbb{C}^N , $W \cap \ker B$ contains a non-zero vector. Pick a unit vector $\mathbf{w} \in W \cap \ker B$ and write $\mathbf{w} = \sum_{i \leq r+1} c_i \mathbf{v}_i$ with $\sum_i |c_i|^2 = 1$. Then $B\mathbf{w} = \mathbf{0}$ and

$$\|A\mathbf{w}\|_2^2 = \|U\Sigma V^* \mathbf{w}\|_2^2 = \|\Sigma V^* \mathbf{w}\|_2^2 = \sum_{i \leq r+1} \sigma_i^2 |c_i|^2 \geq \sigma_{r+1}^2,$$

using monotonicity of the σ_i and $\sum |c_i|^2 = 1$. Therefore $\|A - B\|_{2 \rightarrow 2}^2 \geq \|(A - B)\mathbf{w}\|_2^2 = \|A\mathbf{w}\|_2^2 \geq \sigma_{r+1}^2$, showing $\|A - B\|_{2 \rightarrow 2} \geq \sigma_{r+1}$.

For the Frobenius bound, recall the unitary invariance $\|A - B\|_F^2 = \text{tr}((A - B)^*(A - B)) = \sum_i \sigma_i(A - B)^2$. By Weyl’s inequality $\sigma_i(A - B) \geq \sigma_i(A) - \sigma_1(B) \geq \sigma_i(A) - \|B\|_{2 \rightarrow 2}$ for i exceeding $\text{rank}(B)$; a more careful interlacing argument (Mirsky’s inequality) gives $\sigma_{i+r}(A) \leq \sigma_i(A - B)$ for $\text{rank}(B) \leq r$. Squaring and summing over $i \geq 1$ yields $\|A - B\|_F^2 \geq \sum_{i>r} \sigma_i^2$. Equality is achieved by A_r . \square

C.2 Pseudoinverse and least squares

Definition C.3 (Moore–Penrose pseudoinverse)

$A^\dagger = V\Sigma^\dagger U^*$, where Σ^\dagger inverts the non-zero singular values. For full column rank A , $A^\dagger = (A^*A)^{-1}A^*$; for full row rank, $A^\dagger = A^*(AA^*)^{-1}$.

Theorem C.4 (LS via pseudoinverse)

For any $A \in \mathbb{C}^{m \times N}$, $\mathbf{y} \in \mathbb{C}^m$:

- Among $\arg \min_x \|A\mathbf{x} - \mathbf{y}\|_2$, the unique minimum- ℓ_2 -norm solution is $A^\dagger \mathbf{y}$.
- For full row-rank underdetermined A , $A^\dagger \mathbf{y}$ is the unique minimum-norm solution to $A\mathbf{x} = \mathbf{y}$, satisfying $A^*A\mathbf{x} = A^*\mathbf{y}$.

Proof. Let $A = U\Sigma V^*$ be the SVD with $\sigma_1 \geq \dots \geq \sigma_r > 0 = \sigma_{r+1} = \dots$. Change variables: set

$\tilde{\mathbf{x}} = V^* \mathbf{x} \in \mathbb{C}^N$ and $\tilde{\mathbf{y}} = U^* \mathbf{y} \in \mathbb{C}^m$. By unitary invariance of $\|\cdot\|_2$,

$$\|A\mathbf{x} - \mathbf{y}\|_2^2 = \|U(\Sigma\tilde{\mathbf{x}} - \tilde{\mathbf{y}})\|_2^2 = \|\Sigma\tilde{\mathbf{x}} - \tilde{\mathbf{y}}\|_2^2 = \sum_{i=1}^r (\sigma_i \tilde{x}_i - \tilde{y}_i)^2 + \sum_{i=r+1}^m \tilde{y}_i^2.$$

The second sum is a constant independent of \mathbf{x} . The first sum is minimized by setting $\tilde{x}_i = \tilde{y}_i/\sigma_i$ for $i = 1, \dots, r$. Coordinates \tilde{x}_i for $i > r$ do not appear in the residual at all and are therefore unconstrained.

Among the minimizers, the one of smallest $\|\mathbf{x}\|_2 = \|\tilde{\mathbf{x}}\|_2$ chooses $\tilde{x}_i = 0$ for $i > r$. This gives $\tilde{\mathbf{x}}^* = \Sigma^\dagger \tilde{\mathbf{y}}$ where Σ^\dagger inverts the non-zero singular values, hence $\mathbf{x}^* = V\tilde{\mathbf{x}}^* = V\Sigma^\dagger U^* \mathbf{y} = A^\dagger \mathbf{y}$.

When A has full row rank ($r = m \leq N$), the residual is identically zero, so every solution to $A\mathbf{x} = \mathbf{y}$ is a least-squares solution, and the same argument yields the unique minimum-norm solution. The normal equations $A^*A\mathbf{x} = A^*\mathbf{y}$ are recovered by left-multiplying $A\mathbf{x} = \mathbf{y}$ by A^* . □

C.3 Vandermonde matrices and the determinant

Theorem C.5 (Vandermonde determinant)

For nodes t_0, \dots, t_n ,

$$\det V(t_0, \dots, t_n) = \prod_{0 \leq i < j \leq n} (t_j - t_i).$$

In particular, distinct nodes $\Rightarrow V$ invertible \Rightarrow any $n + 1$ columns of a Vandermonde matrix are linearly independent.

Proof. We argue by induction on n . The base case $n = 1$: $\det \begin{pmatrix} 1 & 1 \\ t_0 & t_1 \end{pmatrix} = t_1 - t_0 = \prod_{0 \leq i < j \leq 1} (t_j - t_i)$.

Inductive step $n \geq 2$: view $D(t_n) := \det V(t_0, t_1, \dots, t_n)$ as a polynomial in t_n with the other t_i held fixed. Cofactor expansion along the last row gives a polynomial of degree at most n in t_n (the leading coefficient is the cofactor of t_n^n , which is $\pm \det V(t_0, \dots, t_{n-1})$). Whenever t_n equals any of t_0, \dots, t_{n-1} , two columns of the matrix are equal, so $D(t_n) = 0$. Hence $D(t_n)$ has n known roots t_0, \dots, t_{n-1} , so

$$D(t_n) = c(t_n - t_0)(t_n - t_1) \cdots (t_n - t_{n-1})$$

for some constant c (independent of t_n but possibly depending on t_0, \dots, t_{n-1}). Comparing the coefficient of t_n^n on both sides identifies $c = \det V(t_0, \dots, t_{n-1})$, the Vandermonde determinant of the smaller matrix. By the induction hypothesis, $c = \prod_{0 \leq i < j \leq n-1} (t_j - t_i)$, so

$$D(t_n) = \prod_{0 \leq i < j \leq n-1} (t_j - t_i) \cdot \prod_{i=0}^{n-1} (t_n - t_i) = \prod_{0 \leq i < j \leq n} (t_j - t_i),$$

completing the induction. Distinct nodes make every factor non-zero, so V is invertible. □

C.4 Norm equivalences

Proposition C.6 (ℓ_p norm inequalities on \mathbb{C}^N)

For $\mathbf{x} \in \mathbb{C}^N$ and $1 \leq p \leq q \leq \infty$, $\|\mathbf{x}\|_q \leq \|\mathbf{x}\|_p \leq N^{1/p-1/q} \|\mathbf{x}\|_q$. Special cases: $\|\mathbf{x}\|_1 \leq \sqrt{N} \|\mathbf{x}\|_2 \leq N \|\mathbf{x}\|_\infty$.

Proof. Lower bound $\|\mathbf{x}\|_q \leq \|\mathbf{x}\|_p$: by homogeneity it suffices to assume $\|\mathbf{x}\|_p = 1$, so $|x_i| \leq 1$ for every i . Since $q \geq p$, $|x_i|^q \leq |x_i|^p$, and summing gives $\|\mathbf{x}\|_q^q \leq \|\mathbf{x}\|_p^p = 1$, hence $\|\mathbf{x}\|_q \leq 1 = \|\mathbf{x}\|_p$.

Upper bound $\|\mathbf{x}\|_p \leq N^{1/p-1/q} \|\mathbf{x}\|_q$: apply Hölder's inequality to the product $|x_i|^p \cdot 1$ with conjugate exponents $r := q/p \geq 1$ and $r' = q/(q-p)$ (so $1/r + 1/r' = 1$):

$$\sum_i |x_i|^p = \sum_i |x_i|^p \cdot 1 \leq \left(\sum_i (|x_i|^p)^r \right)^{1/r} \left(\sum_i 1^{r'} \right)^{1/r'} = \|\mathbf{x}\|_q^p N^{1-p/q}.$$

Taking p -th roots gives $\|\mathbf{x}\|_p \leq N^{1/p-1/q} \|\mathbf{x}\|_q$. The case $q = \infty$ is the limit $r \rightarrow \infty$, giving $\|\mathbf{x}\|_p \leq N^{1/p} \|\mathbf{x}\|_\infty$. □

Lemma C.7 (Operator norms)

$\|A\|_{2 \rightarrow 2} = \sigma_{\max}(A)$, $\|A\|_{1 \rightarrow 1} = \max_k \sum_j |A_{jk}|$, $\|A\|_{\infty \rightarrow \infty} = \max_j \sum_k |A_{jk}|$, and the Frobenius norm $\|A\|_F = \sqrt{\text{tr}(A^*A)} = \sqrt{\sum_i \sigma_i^2}$ satisfies $\|A\|_{2 \rightarrow 2} \leq \|A\|_F$.

Proof. *Spectral norm.* $\|A\|_{2 \rightarrow 2}^2 = \sup_{\|\mathbf{x}\|_2=1} \|A\mathbf{x}\|_2^2 = \sup_{\|\mathbf{x}\|_2=1} \mathbf{x}^* A^* A \mathbf{x}$, which is the largest eigenvalue $\lambda_{\max}(A^*A)$ of the Hermitian PSD matrix A^*A . Since $\lambda_i(A^*A) = \sigma_i(A)^2$, the supremum equals $\sigma_{\max}(A)^2$.

$\ell_1 \rightarrow \ell_1$ norm. For any $\mathbf{x} \in \mathbb{C}^N$,

$$\|A\mathbf{x}\|_1 = \sum_{j=1}^m \left| \sum_{k=1}^N A_{jk} x_k \right| \leq \sum_{j=1}^m \sum_{k=1}^N |A_{jk}| |x_k| = \sum_{k=1}^N |x_k| \sum_{j=1}^m |A_{jk}| \leq \|\mathbf{x}\|_1 \cdot \max_k \sum_j |A_{jk}|.$$

This shows $\|A\|_{1 \rightarrow 1} \leq \max_k \sum_j |A_{jk}|$. Equality is attained by taking \mathbf{x} to be a standard basis vector \mathbf{e}_{k^*} on the maximizing column, since $A\mathbf{e}_{k^*}$ is the k^* -th column with ℓ_1 -norm exactly $\sum_j |A_{jk^*}|$.

$\ell_\infty \rightarrow \ell_\infty$ norm. Dually, $\|A\mathbf{x}\|_\infty = \max_j |\sum_k A_{jk} x_k| \leq \max_j \sum_k |A_{jk}| |x_k| \leq \|\mathbf{x}\|_\infty \max_j \sum_k |A_{jk}|$; equality is attained by $x_k = \text{sign}(A_{j^*k})$ on the maximizing row j^* .

Frobenius bound. $\|A\|_F^2 = \text{tr}(A^*A) = \sum_i \sigma_i(A)^2 \geq \sigma_{\max}(A)^2 = \|A\|_{2 \rightarrow 2}^2$. □

C.5 Schatten norms and matrix functions

Definition C.8 (Schatten p -norm)

$\|A\|_{S_p} = (\sum_i \sigma_i(A)^p)^{1/p}$. The nuclear norm is $\|A\|_* = \|A\|_{S_1}$. The Frobenius norm is $\|A\|_{S_2}$. The operator norm is $\|A\|_{S_\infty}$.

For a normal matrix $A = \sum_i \lambda_i \mathbf{v}_i \mathbf{v}_i^*$ and $f : \mathbb{C} \rightarrow \mathbb{C}$, the matrix function is $f(A) = \sum_i f(\lambda_i) \mathbf{v}_i \mathbf{v}_i^*$.
Examples: $\exp(A) = \sum_k A^k / k!$, $\log(A)$ for positive-definite A , and \sqrt{A} for positive-semidefinite A .

D Multiple Measurement Vectors and Joint Sparsity

Until now every measurement model has had a single observation vector $\mathbf{y} = \mathbf{A}\mathbf{x}$. In practice we often acquire a *collection* of measurements that share a common support — multiple sensor channels, frames of a video, spectral bands of a hyperspectral cube, or arrays of microphones. The Multiple Measurement Vector (MMV) model exploits this shared structure to drive down the per-measurement sample complexity.

D.1 The MMV model and row sparsity

Stack K measurement vectors as columns of $\mathbf{Y} \in \mathbb{R}^{M \times K}$ and the corresponding source vectors as columns of $\mathbf{X} \in \mathbb{R}^{N \times K}$:

$$\mathbf{Y} = \mathbf{A}\mathbf{X}, \quad \mathbf{A} \in \mathbb{R}^{M \times N}, \quad M \ll N.$$

The structural prior is that \mathbf{X} is *jointly sparse* (row-sparse): the columns share a common support $S \subset [N]$ with $|S| = s \ll N$. Equivalently, \mathbf{X} has at most s nonzero rows.

Definition D.1 (Matrix $\ell_{p,q}$ norm)

For $\mathbf{X} \in \mathbb{R}^{N \times K}$ with rows $\mathbf{X}_{k,:}$,

$$\|\mathbf{X}\|_{p,q} = \left(\sum_{k=1}^N \|\mathbf{X}_{k,:}\|_q^p \right)^{1/p}.$$

Two special cases drive MMV theory:

- $\|\mathbf{X}\|_{0,q} = \|\mathbf{X}\|_{\text{row},0} = |\text{supp}(\mathbf{X})|$ counts the number of nonzero rows (the *row-sparsity*), independent of $q \geq 1$.
- $\|\mathbf{X}\|_{1,2} = \sum_k \|\mathbf{X}_{k,:}\|_2$, the convex relaxation: ℓ_2 within each row (to enforce joint occupancy), ℓ_1 across rows (to promote row-sparsity).

The two canonical MMV problems are

$$(P_{R0}) \quad \min_{\mathbf{X}} \|\mathbf{X}\|_{\text{row},0} \text{ s.t. } \mathbf{A}\mathbf{X} = \mathbf{Y}, \quad (P_{12}) \quad \min_{\mathbf{X}} \|\mathbf{X}\|_{1,2} \text{ s.t. } \mathbf{A}\mathbf{X} = \mathbf{Y}.$$

The noisy variants replace the equality with $\|\mathbf{A}\mathbf{X} - \mathbf{Y}\|_F \leq \epsilon$ or move the data fit into the objective.

D.2 Uniqueness for the MMV problem

The shared support across columns provides a strictly stronger uniqueness condition than the single-vector case.

Theorem D.2 (MMV uniqueness, Chen–Huo 2006)

Suppose X is row-sparse with $\|X\|_{\text{row},0} = s$, $Y = AX$, and $\text{rank}(Y) = r$. Then X is the unique row- s -sparse solution of $AZ = Y$ if and only if

$$s < \frac{\text{spark}(A) - 1 + \text{rank}(Y)}{2}.$$

In particular when $\text{rank}(Y) = 1$ this reduces to the SMV bound $s < \text{spark}(A)/2$.

Proof. (\Rightarrow) Suppose there is another row- s -sparse X' with $AX' = Y$. Then $D := X - X' \in \ker A$ has at most $2s$ nonzero rows. Pick a basis $\{\mathbf{w}_1, \dots, \mathbf{w}_r\}$ of $\text{range}(Y)$. Since $\text{range}(AX) = \text{range}(Y)$, each column of X can be written as a combination involving at most r “effective directions.” Hence the columns of D live in a space of dimension $\leq 2s - r + 1$. By definition of spark, any $\text{spark}(A)$ columns of A are linearly independent, so any nonzero vector in $\ker A$ has $\geq \text{spark}(A)$ nonzero entries. Pushing this counting argument over the rank- r structure of D gives $\text{spark}(A) \leq 2s - r + 1$, i.e. $s \geq (\text{spark}(A) + r - 1)/2$, contradicting the hypothesis.

(\Leftarrow) Conversely, the same counting argument shows that two distinct row- s -sparse solutions would force $\text{spark}(A) \leq 2s - r + 1$, which fails under the strict inequality. \square

The bound improves with $\text{rank}(Y)$: distinct sources illuminate the common support from different directions, breaking ties that would be invisible to a single measurement. Using $\text{spark}(A) > 1 + 1/\mu(A)$ yields the coherence-based corollary

$$s < \frac{\mu(A)^{-1} + \text{rank}(Y)}{2}.$$

Remark D.3 (Rank gain)

Numerical experiments confirm that when the columns of X have distinct “directions” on the shared support, $\text{rank}(Y)$ approaches s and the bound nearly doubles. In the worst case, all columns of X are proportional, $\text{rank}(Y) = 1$, and no benefit is realized.

D.3 Row-NSP: a structural characterization

In direct analogy with the vector NSP (Theorem 5.2), we have a *row null space property*.

Definition D.4 (Row-NSP of order s)

$A \in \mathbb{R}^{M \times N}$ satisfies the row-NSP of order s if for every nonzero $V \in \ker A$ (viewed as having columns

in the relevant “measurement space”) and every $S \subset [N]$ with $|S| = s$,

$$\|V_{S,:}\|_{1,2} < \|V_{S^c,:}\|_{1,2}.$$

Theorem D.5 (Row-NSP $\Leftrightarrow \ell_{1,2}$ recovery)

Every row- s -sparse matrix X is the unique minimizer of (P_{12}) with $Y = AX$ if and only if A satisfies the row-NSP of order s .

Proof. (\Leftarrow) Let X be row- s -sparse with support S and let Z be another solution. Set $V = Z - X \in \ker A$. Then

$$\begin{aligned} \|Z\|_{1,2} &= \|X_{S,:} + V_{S,:}\|_{1,2} + \|V_{S^c,:}\|_{1,2} \\ &\geq \|X_{S,:}\|_{1,2} - \|V_{S,:}\|_{1,2} + \|V_{S^c,:}\|_{1,2} \\ &> \|X\|_{1,2}. \end{aligned}$$

The strict inequality uses row-NSP.

(\Rightarrow) Conversely, suppose $V \in \ker A \setminus \{0\}$ violates row-NSP: $\|V_{S,:}\|_{1,2} \geq \|V_{S^c,:}\|_{1,2}$ for some $|S| = s$. Then $X = V_{S,:}$ and $X' = -V_{S^c,:}$ are both row- s -sparse with $AX = AX'$, contradicting uniqueness. \square

D.4 $\ell_{1,2}$ -recovery via coherence

Theorem D.6 ($\ell_{1,2}$ exact recovery, Eldar–Mishali 2009)

If $X \in \Sigma_s$ (row- s -sparse) is a solution of $Y = AX$ and

$$s < \frac{1}{2} \left(1 + \frac{1}{\mu(A)} \right),$$

then X is the unique minimizer of (P_{12}) and coincides with the minimizer of (P_{R0}) .

Proof sketch. The strategy mirrors the SMV coherence proof. Let $V \in \ker A \setminus \{0\}$ and let $S \subset [N]$ with $|S| = s$. Write $V = [V_{S,:}; V_{S^c,:}]$. The Gram matrix $A_S^T A_S$ satisfies $\|A_S^T A_S - I\|_{\text{op}} \leq (s-1)\mu$, so $A_S^T A_S$ is invertible for $s \leq 1/\mu$. Using $AV = 0$: $A_S V_{S,:} = -A_{S^c} V_{S^c,:}$, hence $V_{S,:} = -(A_S^T A_S)^{-1} A_S^T A_{S^c} V_{S^c,:}$. Taking the matrix $\ell_{1,2}$ -norm and applying $\|A_S^T A_{S^c}\|_{\infty,2} \leq s\mu$ row by row yields

$$\|V_{S,:}\|_{1,2} \leq \frac{s\mu}{1 - (s-1)\mu} \|V_{S^c,:}\|_{1,2}.$$

Under $s < \frac{1}{2}(1 + 1/\mu)$, the coefficient is < 1 — the row-NSP holds, and the previous theorem yields recovery. \square

Remark D.7 (Equivalence with SMV bound)

The coherence threshold is *identical* to the single-vector BP bound. The extra structure of MMV does not relax the worst-case coherence requirement, but it dramatically helps average-case behavior: random initializations of X break the worst-case alignment that saturates the coherence bound, so MMV requires fewer measurements per column in practice.

Theorem D.8 (Random-matrix MMV scaling)

Let A have RIP constant δ and let $S \subset [N]$ of size s be the row-support. Then SOMP and $\ell_{1,2}$ -minimization recover X from $Y = AX$ with high probability provided

$$M \lesssim C_R s \log \frac{N}{s}.$$

For Gaussian, Bernoulli, and sub-Gaussian A the joint model needs fewer measurements with higher success probability than the SMV bound, because failure events are correlated across the K columns.

D.5 Simultaneous OMP (SOMP)

The greedy approach lifts OMP atom-by-atom but scores residuals using the matrix ℓ_q -norm of the inner products across all columns.

Algorithm 16 Simultaneous Orthogonal Matching Pursuit (SOMP)

Require: dictionary A , measurements Y , sparsity s , exponent $q \geq 1$

- 1: $\hat{X}_0 \leftarrow 0, R_0 \leftarrow Y, \mathcal{S}_0 \leftarrow \emptyset$
 - 2: **for** $i = 1, \dots, s$ **do**
 - 3: $n^* \leftarrow \arg \max_n \|a_n^T R_{i-1}\|_q$ ▷ scan for the best row
 - 4: $\mathcal{S}_i \leftarrow \mathcal{S}_{i-1} \cup \{n^*\}$
 - 5: $\hat{X}_i \leftarrow A_{\mathcal{S}_i}^\dagger Y$ ▷ least-squares update on chosen rows
 - 6: $R_i \leftarrow Y - A\hat{X}_i$
 - 7: **end for**
 - 8: **return** \hat{X} with row-support \mathcal{S}_s
-

Theorem D.9 (SOMP exact recovery, Tropp et al. 2006)

If $X \in \Sigma_s$ is a solution of $Y = AX$ and the Exact Recovery Condition $\text{ERC}(A, S) := \max_{i \notin S} \|A_S^\dagger a_i\|_1 < 1$ holds, then SOMP recovers X exactly. The condition is independent of q .

The Subspace Pursuit (SSP) and CoSaMP/IHT extensions select s rows per iteration based on the

row- ℓ_q -norms of $A^\top R_{i-1}$ and re-prune; their analysis lifts the RIP-based single-vector guarantees.

D.6 $\ell_{1,2}$ -minimization via ALM

The convex MMV program

$$\min_X \|X\|_{1,2} \text{ s.t. } AX = Y$$

admits an augmented-Lagrangian splitting

$$L_\mu(X, \Lambda) = \|X\|_{1,2} + \frac{\mu}{2} \|Y - AX\|_F^2 + \langle Y - AX, \Lambda \rangle,$$

whose primal step is a row-wise vector soft-thresholding (the prox of $\|\cdot\|_{1,2}$ shrinks each row of the proximal target toward zero by the amount λ/μ , treating the row as a single vector). ADMM and GPSR provide the dominant practical algorithms.

D.7 Group sparse and hierarchical models

If the support naturally partitions into groups $\{G_1, \dots, G_p\} \subseteq [N]$, *Group Lasso* solves

$$\min_x \frac{1}{2} \|Ax - y\|_2^2 + \lambda_2 \sum_i \|x_{G_i}\|_2,$$

encouraging entire groups to be active (dense within a group, sparse across groups). The Hierarchical Lasso (HiLasso) adds an ℓ_1 term to also sparsify within active groups:

$$\min_x \frac{1}{2} \|Ax - y\|_2^2 + \lambda_2 \sum_i \|x_{G_i}\|_2 + \lambda_1 \|x\|_1.$$

Sprechmann–Ramirez–Sapiro–Eldar’s *Collaborative HiLasso* (CHiLasso) extends HiLasso to the MMV setting,

$$\min_X \frac{1}{2} \|AX - Y\|_F^2 + \lambda_2 \sum_j \|X_{G_j}\|_F + \lambda_1 \|X\|_1,$$

combining cross-column joint group structure with within-row sparsity. The problem is solved by ADMM with X -updates that are *column separable* and Z -updates that are *group separable*.

D.8 Heterogeneous dictionaries and gross noise

For collections that share a row-support but use *different* per-column sensing matrices,

$$Y = [A_1 X_1 \mid A_2 X_2 \mid \dots \mid A_L X_L],$$

one solves

$$\min_X \frac{1}{2} \sum_{\ell} \|Y_{\ell} - A_{\ell} X_{\ell}\|_F^2 + \lambda \|X\|_{1,2}.$$

Adding a gross-error term E with $\|E\|_1$ regularization yields a robust MMV problem

$$\min_{X,E} \|X\|_{1,2} + \lambda \|E\|_1 \text{ s.t. } Y = AX + E,$$

which generalizes the SRC + identity dictionary trick of §9-style classifiers to multi-task representations.

D.9 Applications

- **Hyperspectral imaging.** Each pixel has K wavelength bands; the same materials illuminate the same atoms in a spectral dictionary.
- **Multi-sensor acoustic arrays.** The same event recorded at K microphones produces multiple measurements with shared atomic structure.
- **Multi-task visual classification.** Color, shape, and texture features form parallel measurement streams jointly classified via row-sparsity (the active dictionary atoms correspond to a class label).
- **Multi-view face / video recognition.** K views of one subject give MMV measurements with grouped row-support concentrating in one identity block.
- **Bagging.** Generating t random subsets of M measurements yields columns of Y whose sparse codes should be highly correlated; the joint MMV formulation stabilizes recovery.

E Matrix Completion and Robust PCA

The compressive-sensing toolkit has a matrix analogue: replace “sparse vector” with “low-rank matrix,” ℓ_0 with rank, and ℓ_1 with the nuclear norm. Two canonical problems realize this analogy.

Definition E.1 (Schatten norms, revisited)

For $X \in \mathbb{C}^{n_1 \times n_2}$ with singular values $\sigma_1 \geq \dots \geq \sigma_{\min(n_1, n_2)} \geq 0$,

$$\text{rank}(X) = \|\sigma(X)\|_0, \quad \|X\|_* = \|\sigma(X)\|_1 = \text{tr} \sqrt{X^* X}, \quad \|X\|_F = \|\sigma(X)\|_2.$$

E.1 Matrix Completion: setup

Given an index set $\Omega \subset [n_1] \times [n_2]$ and observations $Y_{ij} = X_{ij}$ for $(i, j) \in \Omega$ from an unknown $X \in \mathbb{R}^{n_1 \times n_2}$ of low rank $r \ll \min(n_1, n_2)$, recover X . More generally, with a linear map $\mathcal{A} : \mathbb{C}^{n_1 \times n_2} \rightarrow \mathbb{C}^m$ and observation

$\mathbf{y} = \mathcal{A}(X)$, the combinatorial program

$$(P_{M0}) \quad \min_Z \text{rank}(Z) \text{ s.t. } \mathcal{A}(Z) = \mathbf{y}$$

is NP-hard (Foucart & Rauhut, Exercise 2.11). Its convex relaxation replaces rank by the nuclear norm:

$$(P_{M1}) \quad \min_Z \|Z\|_* \text{ s.t. } \mathcal{A}(Z) = \mathbf{y},$$

solvable by semidefinite programming in polynomial time. For matrix completion, \mathcal{A} samples entries: $\mathcal{A}(Z)_\ell = Z_{j_\ell, k_\ell}$ for $(j_\ell, k_\ell) \in \Omega$. Foucart & Rauhut (§4.5) treat the more general rank-recovery problem and show that nuclear-norm minimization mirrors basis pursuit in nearly every respect: a null-space property characterizes recovery, random \mathcal{A} achieve the optimal rate $m = Cr \max(n_1, n_2)$ (no log factor), and stable/robust variants exist under noise.

E.2 The rank null-space property

The exact analogue of Theorem 5.2 for nuclear-norm minimization is the *rank null-space property*. We state and prove it following Foucart & Rauhut Theorem 4.40; we restate it here in the context of matrix completion for ease of reference.

Theorem E.2 (Rank NSP, Foucart–Rauhut Thm. 4.40)

Given $\mathcal{A} : \mathbb{C}^{n_1 \times n_2} \rightarrow \mathbb{C}^m$, every matrix $X \in \mathbb{C}^{n_1 \times n_2}$ of rank at most r is the unique solution of (P_{M1}) with $\mathbf{y} = \mathcal{A}(X)$ if and only if for every $M \in \ker \mathcal{A} \setminus \{0\}$ with singular values $\sigma_1(M) \geq \dots \geq \sigma_n(M) \geq 0$, $n := \min(n_1, n_2)$,

$$\sum_{j=1}^r \sigma_j(M) < \sum_{j=r+1}^n \sigma_j(M).$$

Proof. (\Rightarrow) Let $M \in \ker \mathcal{A} \setminus \{0\}$ with SVD $M = U \Sigma V^*$. Set $M_1 = U \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0) V^*$ and $M_2 = U \text{diag}(0, \dots, 0, -\sigma_{r+1}, \dots, -\sigma_n) V^*$. Then $M = M_1 - M_2$, both summands have rank $\leq n$, and $\text{rank}(M_1) \leq r$. Since $\mathcal{A}(M) = 0$, we get $\mathcal{A}(M_1) = \mathcal{A}(-M_2)$. By uniqueness of nuclear-norm recovery for rank- r matrices, $\|M_1\|_* < \|-M_2\|_* = \|M_2\|_*$. The singular values give $\sum_{j=1}^r \sigma_j(M) < \sum_{j=r+1}^n \sigma_j(M)$.

(\Leftarrow) Conversely, suppose the inequality holds for every nonzero $M \in \ker \mathcal{A}$. Let X have rank $\leq r$ and let $Z \neq X$ satisfy $\mathcal{A}(Z) = \mathcal{A}(X)$. Set $M = X - Z \in \ker \mathcal{A}$. Lemma A.20 of Foucart–Rauhut (a perturbation lemma for singular values) gives $\|Z\|_* = \sum_j \sigma_j(X - M) = \sum_j |\sigma_j(X) - \sigma_j(M)|$ when X and M share their singular-vector frames; in the general case the same inequality follows from Weyl's perturbation bound. For $j \in [r]$, $|\sigma_j(X) - \sigma_j(M)| \geq \sigma_j(X) - \sigma_j(M)$; for $j \in [r+1, n]$, $\sigma_j(X) = 0$ so the absolute value equals $\sigma_j(M)$. Hence

$$\|Z\|_* \geq \sum_{j=1}^r \sigma_j(X) - \sum_{j=1}^r \sigma_j(M) + \sum_{j=r+1}^n \sigma_j(M) > \sum_{j=1}^r \sigma_j(X) = \|X\|_*.$$

Therefore X is the unique minimizer. □

Remark E.3 (Stable and robust rank NSP)

The vector analogues of stable and robust NSP carry over: if $\sum_{j=1}^r \sigma_j(M) \leq \rho \sum_{j=r+1}^n \sigma_j(M) + \tau \|\mathcal{A}(M)\|_2$ for some $\rho < 1$, $\tau > 0$ and all M , then nuclear-norm minimization yields stable + robust recovery: $\|X^* - X\|_* \leq 2\sigma_r(X)_* + 2\tau\eta$ for measurements $\mathbf{y} = \mathcal{A}(X) + \mathbf{e}$ with $\|\mathbf{e}\|_2 \leq \eta$. Compare Theorem 5.7 for the vector case.

E.3 Conditions for exact recovery

Two obstructions force assumptions on X and Ω :

- *Which matrices?* A 1-sparse matrix is also rank-1; we cannot decide between sparse and low-rank without further information. The singular vectors must be *incoherent* with the standard basis (no column or row of X carries a disproportionate share of the energy).
- *Which sampling patterns?* If an entire row or column is missing, recovery is impossible. We sample Ω uniformly at random.

The incoherence parameter $\mu(X)$ measures how aligned the singular subspaces $U = \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_r)$ and $V = \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_r)$ are with the canonical basis: for the projections P_U, P_V ,

$$\max_{1 \leq i \leq n_1} \|P_U \mathbf{e}_i\|^2 \leq \frac{\mu r}{n_1}, \quad \max_{1 \leq j \leq n_2} \|P_V \mathbf{e}_j\|^2 \leq \frac{\mu r}{n_2}, \quad \|UV^*\|_\infty \leq \sqrt{\frac{\mu r}{n_1 n_2}}.$$

Theorem E.4 (Candès–Recht 2009, exact matrix completion)

Let $X \in \mathbb{R}^{n_1 \times n_2}$ have rank r and be sampled from the random orthogonal model, set $N = \max(n_1, n_2)$, and let Ω be chosen uniformly at random with $|\Omega| \geq C N^{5/4} r \log N$. Then with probability $1 - c/N^3$, the minimizer of (P_{M_1}) is unique and equals X .

For incoherent X the rate sharpens to $|\Omega| \gtrsim \mu r N \log^2 N$. The bound is information-theoretically near optimal: the rank- r matrix has $\approx r(n_1 + n_2 - r)$ degrees of freedom, so Ω of size $\Theta(rN \text{ polylog } N)$ is necessary up to logarithmic factors.

E.4 Robust PCA

Classical PCA finds the dominant principal subspace via

$$\min_{\Phi, W} \|\Phi W - X\|_F^2 \quad \text{s.t.} \quad \Phi^\top \Phi = I,$$

which is brittle to gross outliers — a single corrupted entry can swing the top singular vectors. Robust PCA models the corrupted matrix as

$$Y = X + E (+W),$$

where X is low-rank, E is sparse (gross noise of arbitrary magnitude, but *unknown* support), and W is dense small noise with $\|W\|_F \leq \eta$. The combinatorial problem

$$(P_{R0}) \min_{X,E} \text{rank}(X) + \lambda \|E\|_0 \text{ s.t. } Y = X + E$$

is NP-hard; the convex relaxation

$$(P_{R1}) \min_{X,E} \|X\|_* + \lambda \|E\|_1 \text{ s.t. } Y = X + E$$

is a semidefinite program.

Theorem E.5 (Candès–Li–Ma–Wright 2011, Robust PCA)

Let $X, E \in \mathbb{R}^{n_1 \times n_2}$ with $n_1 \geq n_2$, $r = \text{rank}(X)$. If there exist ρ_1, ρ_2 such that

$$r \leq \rho_1 \frac{n_2}{\mu(X) \log^2 n_1}, \quad \|E\|_0 \leq \rho_2 n_1 n_2,$$

then (P_{R1}) with $\lambda = 1/\sqrt{n_1}$ recovers X exactly with probability $1 - C/n_1^{10}$.

The combined problem of missing entries *and* gross corruption is

$$(P_{RM1}) \min_{X,E} \|X\|_* + \lambda \|E\|_1 \text{ s.t. } Y = (X + E)|_\Omega,$$

which under the same incoherence + sampling assumptions recovers both the low-rank component and the sparse component.

Remark E.6 (Compressed sensing as a special case)

Apply rank minimization to the diagonal matrix $X = \text{diag}(\mathbf{x})$:

$$\text{rank}(X) = \|\mathbf{x}\|_0, \quad \|X\|_* = \|\mathbf{x}\|_1.$$

Matrix completion / RPCA generalize $\mathbf{y} = A\mathbf{x}$ / $\mathbf{y} = A\mathbf{x} + \mathbf{e}$ to the matrix world. The diagonal-matrix correspondence makes precise the slogan that compressed sensing is the rank-1 shadow of matrix completion.

E.5 Algorithms: SVT and ALM

Singular Value Thresholding (SVT)

solves the convenient approximation $\min \tau \|X\|_* + \frac{1}{2} \|X\|_F^2$ s.t. $X|_\Omega = Y$ by Uzawa iterations on the Lagrangian

$$L(X, \Lambda) = \tau \|X\|_* + \frac{1}{2} \|X\|_F^2 + \langle (Y - X)|_\Omega, \Lambda|_\Omega \rangle.$$

The X -step admits a closed form: completing the square in the $\|X\|_F^2$ term and applying the proximal operator of the nuclear norm yields the matrix soft-thresholding (a.k.a. singular-value shrinkage)

$$\mathcal{D}_\tau(Z) = U \text{diag}(\max(\sigma_i - \tau, 0)) V^*, \quad Z = U \Sigma V^*,$$

while the multiplier update is gradient ascent:

$$X_{k+1} = \mathcal{D}_\tau(\Lambda_k), \quad \Lambda_{k+1} = \Lambda_k + \gamma_k (Y - X_{k+1})|_\Omega.$$

Lemma E.7 (Prox of nuclear norm = SVT)

For $Z \in \mathbb{C}^{n_1 \times n_2}$ and $\tau > 0$,

$$\arg \min_X \left\{ \tau \|X\|_* + \frac{1}{2} \|X - Z\|_F^2 \right\} = \mathcal{D}_\tau(Z).$$

Proof. Write the SVD $Z = U \Sigma V^*$. The Frobenius norm is unitarily invariant, so $\|X - Z\|_F^2 = \|U^* X V - \Sigma\|_F^2$; setting $Y = U^* X V$ reduces the problem to $\min_Y \tau \|Y\|_* + \frac{1}{2} \|Y - \Sigma\|_F^2$. Since Σ is diagonal, the minimizer is diagonal (any off-diagonal entry of Y only increases both terms). On diagonals, the problem decouples into scalar shrinkage $y_i \mapsto \text{sgn}(\sigma_i)(|y_i| - \tau)_+$. Unitarily rotating back gives $\mathcal{D}_\tau(Z)$. \square

Each SVT iteration costs one full SVD; truncated SVDs (Lanczos) reduce this to $O(rn_1n_2)$ per step when the iterate is approximately rank r . Convergence is slow but the algorithm is rank-controllable: at each step the output is exactly $\text{rank}\{\sigma_i : \sigma_i > \tau\}$. The analogous entrywise shrinkage $\mathcal{S}_\tau(\mathbf{x}) = \text{sgn}(\mathbf{x})(|\mathbf{x}| - \tau)_+$ handles the sparse term in RPCA.

Augmented Lagrange Multiplier (ALM)

treats the constraint $Y = X + E$ via

$$\tilde{L}(X, E, \Lambda) = \|X\|_* + \lambda \|E\|_1 + \frac{\mu}{2} \|Y - X - E\|_F^2 + \langle Y - X - E, \Lambda \rangle.$$

Alternating minimization gives the updates

$$X_{k+1} = \mathcal{D}_{\mu^{-1}}(Y - E_k - \mu^{-1}\Lambda_k), \quad E_{k+1} = \mathcal{S}_{\lambda\mu^{-1}}(Y - X_{k+1} - \mu^{-1}\Lambda_k),$$

followed by $\Lambda_{k+1} = \Lambda_k + \mu(Y - X_{k+1} - E_{k+1})$. ALM converges in far fewer iterations than SVT and is widely considered the workhorse algorithm for Robust PCA.

Other approaches: NNLS, Accelerated Proximal Gradient (APG), Fixed-Point Continuation with Approximate SVD (FPCA), and ADMM variants.

E.6 Applications

- **Netflix-style collaborative filtering.** 480K users, 17,770 movies, $\sim 1\%$ of entries observed — complete the rating matrix.
- **Background subtraction.** A video stacked column-wise has a low-rank background (X) plus a sparse moving foreground (E); RPCA separates the two.
- **Image / video inpainting.** Missing pixels (random), missing blocks (region), missing rows/columns (super-resolution), missing frames (frame interpolation).
- **Robust face alignment.** $D \circ \tau = A_0 + E_0$: parametric deformations plus a low-rank aligned face dictionary plus sparse occlusions (glasses, shadows, expressions).
- **Photometric stereo & texture repair.** Low-rank surface model + sparse highlights/specularities.
- **Locally adaptive low-rank video.** Patch-based or block-based decomposition with non-local-means dictionary construction yields BM3D-style denoising rates with rank-completion bounds.

F Dictionary Learning

Every guarantee so far has assumed the dictionary A (or sparsifying basis Ψ) is fixed and known. *Learning* the dictionary from the data itself opens the door to better representations for natural images, audio, and domain-specific signals where no analytical basis is satisfactory.

F.1 The dictionary learning problem

Given a training set $Y = [\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(M)}] \in \mathbb{R}^{n \times M}$, find a dictionary $D \in \mathbb{R}^{n \times m}$ (typically overcomplete, $m > n$) and sparse codes $X \in \mathbb{R}^{m \times M}$ such that

$$\min_{D, X} \|Y - DX\|_F^2 \quad \text{s.t.} \quad \|\mathbf{x}_i\|_0 \leq k_0 \quad \forall i, \quad \|\mathbf{d}_j\|_2 = 1 \quad \forall j.$$

The norm constraint on dictionary atoms removes the scale ambiguity between D and X .

The problem is non-convex and admits no global minimum guarantee in general, but well-designed alternating-direction schemes converge to useful local minima.

Theorem F.1 (Local minimum existence, Jenatton–Gribonval–Bach 2012)

Let $D_0 \in \mathbb{R}^{n \times m}$ be a reference dictionary with coherence μ_0 and $1/\gamma_{D_0} \triangleq \|D_0\|_2 \cdot k\mu_0$. If

$$\Omega(\sqrt{\log m}) = \gamma_{D_0} = O(\sqrt{\log n}), \quad \frac{\log n}{n} = O\left(\frac{\mu_0^2}{m k^3 \gamma_{D_0}^2}\right),$$

then with high probability the dictionary-learning objective has a local minimum within distance $O(p\gamma_{D_0}[e^{-\gamma_{D_0}^2/2} + \sqrt{mp \log(n)/n}])$ of D_0 . Existence requires the reference dictionary to be sufficiently incoherent and enough signals to be observed.

F.2 K-means (the simplest dictionary)

If we force each column of X to be a standard basis vector $\mathbf{x}_i = \mathbf{e}_{j(i)}$, dictionary learning reduces to vector quantization (K-means):

$$\min_{D, X} \|Y - DX\|_F^2 \quad \text{s.t.} \quad \mathbf{x}_i = \mathbf{e}_j \text{ for some } j.$$

The columns of D are the cluster centroids; the codes are one-hot assignments.

Algorithm 17 K-means

Require: data Y , number of clusters k

- 1: Initialize centroids $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$ uniformly at random.
 - 2: **repeat**
 - 3: **(Coding)** $c^{(i)} \leftarrow \arg \min_{1 \leq j \leq k} \|\mathbf{y}^{(i)} - \boldsymbol{\mu}_j\|^2$
 - 4: **(Update)** $\boldsymbol{\mu}_j \leftarrow \frac{\sum_i \mathbf{1}\{c^{(i)} = j\} \mathbf{y}^{(i)}}{\sum_i \mathbf{1}\{c^{(i)} = j\}}$
 - 5: **until** centroids stabilize
-

K-means is wonderful for clustering but loses information for representation: the one-hot codebook cannot reconstruct points that lie between centroids.

F.3 Method of Optimal Directions (MOD)

Engan, Aase, and Husoy (1999) first cast dictionary design as alternating sparse coding and pseudo-inverse dictionary update:

Algorithm 18 Method of Optimal Directions (MOD)

Require: data Y , atom count m , sparsity k_0

- 1: Initialize $D^{(0)} \in \mathbb{R}^{n \times m}$ with unit-norm columns, $J \leftarrow 1$.
- 2: **repeat**
- 3: **Sparse coding:** for each i , solve $\mathbf{x}_i = \arg \min_{\mathbf{x}} \|\mathbf{y}^{(i)} - D^{(J-1)}\mathbf{x}\|_2$ s.t. $\|\mathbf{x}\|_0 \leq k_0$
- 4: (using OMP, IHT, or any single-vector sparse-coding routine).
- 5: **Dictionary update:** $D^{(J)} \leftarrow YX^T(XX^T)^{-1} = YX^\dagger$.
- 6: Normalize columns of $D^{(J)}$ to unit length.
- 7: $J \leftarrow J + 1$.
- 8: **until** stopping criterion

The objective is monotonically non-increasing. MOD is fast but typically yields a worse local minimum than the next algorithm.

F.4 K-SVD: rank-1 atom updates

Aharon–Elad–Bruckstein (2006) recognized that the dictionary update can be done one column at a time, exploiting an SVD on the relevant residual:

Algorithm 19 K-SVD

Require: data Y , atom count m , sparsity k_0

- 1: Initialize $D^{(0)}$ randomly with unit-norm columns; $J \leftarrow 1$.
- 2: **repeat**
- 3: **Sparse coding:** same as MOD (use OMP).
- 4: **for** $j_0 = 1, \dots, m$ **do**
- 5: $\omega_{j_0} \leftarrow \{i : X_{j_0,i} \neq 0\}$ ▷ examples using atom \mathbf{d}_{j_0}
- 6: $E_{j_0} \leftarrow Y - \sum_{j \neq j_0} \mathbf{d}_j X_{j,:}^T$
- 7: $E_{j_0}^R \leftarrow E_{j_0} |_{\omega_{j_0}}$ ▷ restrict to active columns
- 8: Compute SVD $E_{j_0}^R = U\Sigma V^T$.
- 9: $\mathbf{d}_{j_0} \leftarrow \mathbf{u}_1$ ▷ new atom = top left singular vector
- 10: $X_{j_0,\omega_{j_0}}^T \leftarrow \sigma_1 \mathbf{v}_1$ ▷ update coefficients accordingly
- 11: **end for**
- 12: $J \leftarrow J + 1$.
- 13: **until** stopping

Derivation.

The objective $\|Y - DX\|_F^2$ depends on atom \mathbf{d}_{j_0} only through the rank-1 outer product $\mathbf{d}_{j_0} X_{j_0,:}^\top$. Decompose:

$$Y - DX = Y - \sum_{j=1}^m \mathbf{d}_j X_{j,:}^\top = \underbrace{\left(Y - \sum_{j \neq j_0} \mathbf{d}_j X_{j,:}^\top \right)}_{E_{j_0}} - \mathbf{d}_{j_0} X_{j_0,:}^\top.$$

Hence $\|Y - DX\|_F^2 = \|E_{j_0} - \mathbf{d}_{j_0} X_{j_0,:}^\top\|_F^2$, and the optimal rank-1 approximation problem

$$\min_{\mathbf{d}_{j_0}, X_{j_0,:}} \|E_{j_0} - \mathbf{d}_{j_0} X_{j_0,:}^\top\|_F^2$$

is solved by Eckart–Young: take the top left singular vector for \mathbf{d}_{j_0} and the scaled right singular vector for $X_{j_0,:}^\top$.

Crucially, naively replacing $X_{j_0,:}$ by $\sigma_1 \mathbf{v}_1$ destroys the sparsity pattern (most entries of \mathbf{v}_1 are nonzero). The K-SVD remedy is to *restrict* the update to columns where atom \mathbf{d}_{j_0} is already active, $\omega_{j_0} = \{i : X_{j_0,i} \neq 0\}$, replacing E_{j_0} by $E_{j_0}^R = E_{j_0}|_{\omega_{j_0}}$ and updating only the nonzero entries of $X_{j_0,:}$ accordingly. Zeros remain zeros; the support stays unchanged.

Lemma F.2 (K-SVD monotonicity)

Both the sparse-coding step (with exact ℓ_0 projection) and the per-atom SVD step are non-increasing for $\|Y - DX\|_F^2$. Hence the objective converges monotonically to a stationary point.

Compared to K-means, K-SVD is a “soft” version: codes are continuous, atoms can be shared across clusters, and the representation power is substantially greater. Compared to MOD, K-SVD updates each atom along with its coefficients simultaneously, leading to faster and more stable convergence.

Trained on natural-image patches, K-SVD discovers Gabor-like edge atoms that mirror the receptive fields of V1 simple cells (DeAngelis, Ohzawa, Freeman 1995) — rediscovering an early-vision representation from raw data.

F.5 Online dictionary learning

For streaming or massive datasets, the K-SVD batch update is impractical. Mairal–Bach–Ponce–Sapiro (2010) replaced ℓ_0 by ℓ_1 and the batch sum by a stochastic approximation:

$$\min_{D,X} \sum_{i=1}^N \left(\frac{1}{2} \|\mathbf{y}^{(i)} - D\mathbf{x}^{(i)}\|_2^2 + \lambda \|\mathbf{x}^{(i)}\|_1 \right) = \min_{D,X} \frac{1}{2} \|Y - DX\|_F^2 + \lambda \|X\|_{1,1}.$$

Algorithm 20 Online Dictionary Learning (Mairal et al. 2010)**Require:** sample stream $\mathbf{y}_1, \mathbf{y}_2, \dots$; initial D_0 ; regularizer λ 1: $A_0 \leftarrow 0, B_0 \leftarrow 0$ 2: **for** $t = 1, 2, \dots$ **do**3: Draw \mathbf{y}_t 4: $\mathbf{x}_t \leftarrow \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y}_t - D_{t-1} \mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$ ▷ LASSO via FISTA/LARS5: $A_t \leftarrow A_{t-1} + \mathbf{x}_t \mathbf{x}_t^\top$ 6: $B_t \leftarrow B_{t-1} + \mathbf{y}_t \mathbf{x}_t^\top$ 7: $D_t \leftarrow$ block-coordinate descent on $\frac{1}{t} \text{tr}(D^\top D A_t) - \frac{2}{t} \text{tr}(D^\top B_t)$ subject to $\|\mathbf{d}_j\|_2 \leq 1$ 8: **end for**

The sufficient statistics A_t, B_t summarize all past data in matrices of fixed size $m \times m$ and $n \times m$, so memory and per-iteration cost are independent of the dataset size. Mairal et al. prove almost-sure convergence of D_t to a stationary point of the empirical risk

$$f(D) = \mathbb{E}_{\mathbf{y}} \left[\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - D\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1 \right]$$

under mild moment conditions and learning rates $\eta_t = O(1/t)$.

F.6 Locality-constrained linear coding (LLC)

For classification rather than reconstruction, Wang et al. (2010) trade sparsity for locality: encourage codes to use only the nearest dictionary atoms,

$$\min_{D, \mathbf{X}} \sum_{i=1}^N \|\mathbf{y}_i - D\mathbf{x}_i\|_2^2 + \lambda \|\mathbf{d}_i \odot \mathbf{x}_i\|_2^2, \quad \mathbf{d}_i = e^{\text{dist}(\mathbf{x}_i, D)/\sigma}.$$

The locality weight \mathbf{d}_i penalizes coefficients on far-away atoms, giving smooth, clustering-friendly codes with analytical solution per sample. LLC is fast and discriminative — on Caltech-101 it matched or beat sparsity-only methods at the time.

F.7 Sparse Subspace Clustering (SSC)

When data lies on a union of low-dimensional subspaces (motion segmentation, face clustering under varying illumination), the *self-expressiveness* property states that each point is a sparse linear combination of other

points from the same subspace:

$$\mathbf{y}_i = Y\mathbf{x}_i = [Y_1 | \dots | Y_C]\mathbf{x}_i, \quad x_{ji} = 0.$$

Elhamifar–Vidal (2013) recover the subspace assignment by

$$\min_{\mathbf{x}_i} \|\mathbf{x}_i\|_1 \text{ s.t. } \mathbf{y}_i = Y\mathbf{x}_i, \quad x_{ji} = 0.$$

With noise and sparse outliers,

$$\min_{X,E,Z} \|X\|_1 + \lambda_1 \|E\|_1 + \lambda_2 \|Z\|_F^2 \text{ s.t. } Y = YX + E + Z, \text{ diag}(X) = \mathbf{0}.$$

The sparse-affinity matrix $|X| + |X|^\top$ is fed to spectral clustering. On Hopkins 155 (motion) and Extended Yale B (faces), SSC achieves sub-percent clustering error, dominating LSA, SCC, and LRR baselines on most configurations.

G Nonnegative Matrix Factorization

Nonnegative Matrix Factorization (NMF) restricts both the dictionary and the codes to have nonnegative entries:

$$Y \approx DX, \quad D \geq 0, \quad X \geq 0.$$

This single constraint produces a remarkable shift in representation: addition without cancellation forces *parts-based* decompositions where atoms correspond to interpretable, additive features (Lee & Seung, *Nature* 1999).

G.1 Why parts-based?

Without sign cancellation, atoms cannot represent “negative features.” Reconstruction proceeds by summing contributions; the dictionary atoms must therefore be valid pieces of the signal. On face datasets:

- K-means atoms are class prototypes (one per cluster).
- K-SVD atoms are holistic “average faces” that share structure across clusters.
- NMF atoms become localized facial features — eyes, noses, mouths — that add up to a face. The sparse code X encodes *membership*: $X_{kj} > X_{ij} \forall i \neq k$ means \mathbf{y}_j is in cluster k .

For unsupervised classification, NMF features outperform K-means and K-SVD features because the additive structure aligns with discriminative parts.

G.2 Multiplicative update rules (Lee–Seung)

Theorem G.1 (Multiplicative-update convergence, Lee–Seung 2001)

The updates

$$X_{ij} \leftarrow X_{ij} \frac{(D^\top Y)_{ij}}{(D^\top DX)_{ij}}, \quad D_{ki} \leftarrow D_{ki} \frac{(YX^\top)_{ki}}{(DXX^\top)_{ki}}$$

monotonically decrease $\|Y - DX\|_F^2$ subject to $D, X \geq 0$.

Proof sketch via auxiliary function. Let $F(X) = \frac{1}{2}\|Y - DX\|_F^2$ for fixed $D \geq 0$. Define the auxiliary function

$$G(X, \tilde{X}) = F(\tilde{X}) + (X - \tilde{X})^\top \nabla F(\tilde{X}) + \frac{1}{2} \sum_{ij} \frac{(D^\top D \tilde{X})_{ij}}{\tilde{X}_{ij}} (X_{ij} - \tilde{X}_{ij})^2.$$

Two facts: (i) $G(X, \tilde{X}) \geq F(X)$ for all $X, \tilde{X} \geq 0$, with equality at $\tilde{X} = X$ (the quadratic majorizes F by AM–GM); and (ii) $G(\cdot, \tilde{X})$ is convex in X with closed-form minimizer

$$X_{ij}^* = \tilde{X}_{ij} \frac{(D^\top Y)_{ij}}{(D^\top D \tilde{X})_{ij}}.$$

Setting $\tilde{X} = X^{(k)}$ and $X^{(k+1)} = X^*$ gives a majorize–minimize step: $F(X^{(k+1)}) \leq G(X^{(k+1)}, X^{(k)}) \leq G(X^{(k)}, X^{(k)}) = F(X^{(k)})$. The D -update is symmetric. Hence the objective is monotonically non-increasing. \square

The updates are multiplicative — each entry is rescaled, never assigned — so nonnegativity is preserved automatically. An alternative is alternating Nonnegative Least Squares (NNLS), which solves $\min_{D \geq 0} \|Y - DX\|_F^2$ and $\min_{X \geq 0} \|Y - DX\|_F^2$ by active-set or projected-gradient methods at each iteration; NNLS converges faster than multiplicative updates but each step is more expensive.

Remark G.2 (KL-divergence variant)

When entries of Y are counts (documents, photon counts), the Kullback–Leibler divergence $D_{\text{KL}}(Y \| DX) = \sum_{ij} (Y_{ij} \log \frac{Y_{ij}}{(DX)_{ij}} - Y_{ij} + (DX)_{ij})$ is more natural than Frobenius loss. Lee–Seung’s KL update rules are

$$X_{ij} \leftarrow X_{ij} \frac{\sum_k D_{ki} Y_{kj} / (DX)_{kj}}{\sum_k D_{ki}}, \quad D_{ki} \leftarrow D_{ki} \frac{\sum_j X_{ij} Y_{kj} / (DX)_{kj}}{\sum_j X_{ij}},$$

with an analogous majorization–minimization proof.

G.3 Variants

Orthogonal NMF.

Adding $XX^T = I$ forces orthogonal codes:

$$\min_{D,X} \|Y - DX\|_F^2 \text{ s.t. } D, X \geq 0, XX^T = I.$$

Theorem G.3 (NMF \approx K-means; Ding–He–Simon 2005)

Orthogonal Semi-NMF is equivalent to relaxed K-means:

$$\min_{D,X} \|Y - DX\|_F^2 \text{ s.t. } X \geq 0, XX^T = I \iff \min \text{tr}(Y^T Y) - \text{tr}(XY^T YX^T).$$

Maximizing the second term is K-means in disguise.

Sparse NMF.

Add an ℓ_1 penalty on X to encourage strictly local atoms. More sparseness moves atoms from blurred face-averages toward isolated facial parts (Hoyer 2004).

Semi-NMF and Convex-NMF.

$$\text{SVD: } Y_{\pm} = UAV^T \approx D_{\pm} X_{\pm}$$

$$\text{Semi-NMF: } Y_{\pm} \approx D_{\pm} X_{+}$$

$$\text{Convex-NMF: } Y_{\pm} \approx Y_{\pm} W_{+} X_{+}$$

Convex-NMF expresses dictionary atoms as nonnegative combinations of the data itself — a generalization that links NMF, K-means, and the next topic.

G.4 Pros and cons

Pros. Intuitive parts-based features; addition only (no cancellation); holistic decompositions; excellent fit for images, videos, text, frequency counts; encourages clustering behavior.

Cons. Ill-posed (no unique factorization); nonconvex (no global convergence guarantee); the cone of admissible factorizations admits many local optima.

G.5 Separable NMF

Donoho & Stodden (2004) gave a geometric framework: NMF factorizations are *cones* containing the data cloud. Multiple cones can enclose the same data — hence non-uniqueness.

Definition G.4 (*s*-separable matrix)

$Y = [\mathbf{y}_1, \dots, \mathbf{y}_m]$ is *s*-separable if there exists an index set $\mathcal{S} \subset [m]$ with $|\mathcal{S}| = s < m$ such that

$$Y = Y_{\mathcal{S}} Z, \quad Z \geq 0, \quad \mathbf{1}^T Z = \mathbf{1}^T.$$

Equivalently, *s* columns of Y generate a convex cone containing the entire dataset.

When the convex-hull assumption holds, the factorization is *unique* and can be recovered by finding the vertex set \mathcal{S} . This reduces NMF to a sparse representation problem:

$$\min_X \|X\|_{\text{row},0} \text{ s.t. } YX = Y, \quad X \geq 0, \quad \mathbf{1}^T X = \mathbf{1}^T.$$

Algorithms for separable NMF:

- **Greedy pursuit (XRAY, Kumar et al. 2013).** Infer the convex hull one vertex per iteration via geometric criteria.
- **Convex relaxation $\ell_{1,\infty}$ (Elhamifar et al.).** $\min_X \|X\|_{1,\infty}$ s.t. $YX = Y, X \geq 0, \mathbf{1}^T X = \mathbf{1}^T$, solved by ADMM.
- **Linear programs** (Arora et al. 2012, Bittorf 2012, Gillis 2014).

G.6 Sparsity from entropy

A probability vector concentrated on a single coordinate has zero Shannon entropy:

$$H(\mathbf{x}) = - \sum_i p(x_i) \log p(x_i), \quad p_i = \frac{|x_i|^p}{\|\mathbf{x}\|_p^p}.$$

Low entropy \equiv skewed distribution \equiv “sparse.” This motivates the *Shannon entropy regularizer*

$$h_p(\mathbf{x}) = - \sum_{i=1}^N \frac{|x_i|^p}{\|\mathbf{x}\|_p^p} \log \frac{|x_i|^p}{\|\mathbf{x}\|_p^p}$$

as an alternative to $\|\cdot\|_1$. Minimizing h_p favors sparse solutions but is nonconvex; visually, level sets of h_1 produce a star-shaped pattern along the axes, sharper than ℓ_p for any $p > 0$.

For separable NMF, *Robust Entropy Minimization*

$$\min_X \|X\|_{h,\infty} \text{ s.t. } \|\mathbf{y}_j - Y\mathbf{x}_j\|_2 \leq \epsilon, \quad X \geq 0, \quad \mathbf{1}^T X = \mathbf{1}^T,$$

with the row-entropy norm $\|X\|_{h,\infty} = h([\|X_{1,:}\|_\infty, \dots, \|X_{N,:}\|_\infty])$, provably identifies the vertices of Y

exactly when the noise is bounded by $\varepsilon < c/(8(s+1))$, where c depends on the fatness of the convex hull. The resulting program is nonconvex yet admits a proof of global optimality.

G.7 Applications

- **Recommendation systems.** Fill in user–item matrices: solve $\min \|Y - DX\|_F^2$ s.t. $D, X \geq 0$; the completed $\tilde{Y} = DX$ gives predicted ratings.
- **Topic modeling.** Y = document–term matrix; dictionary atoms become topics, codes become document–topic loadings.
- **Audio source separation, hyperspectral unmixing, gene-expression analysis** — anywhere nonnegativity is a physical constraint.

H Deep Sparse Coding Networks

The final synthesis takes the entire course’s theory of sparse coding and stacks it into a multilayer architecture, aiming to combine the representational power of deep neural networks with the interpretability and sample efficiency of sparse modeling.

H.1 Motivation

- One-layer sparse coding under-performs on large datasets like CIFAR.
- Deep CNNs perform strongly but offer little theoretical interpretability.
- Can we extend sparse coding to a deep architecture while preserving its mathematical framework?

Two strands of evidence motivate the affirmative answer (Rangamani et al. 2018; 2019): flatness of local minima correlates with robustness; and trained ReLU networks effectively perform sparse coding in their feature spaces.

H.2 Composite sparse coding module

The direct cascade $\mathbf{x} \rightarrow \boldsymbol{\alpha}^{(1)} \rightarrow \boldsymbol{\alpha}^{(2)}$ where each $\boldsymbol{\alpha}^{(\ell)} = D^{(\ell)} \boldsymbol{\alpha}^{(\ell-1)}$ grows rapidly in dimension. Sun et al. (2018, 2019) split each layer into two stages:

$$\boldsymbol{\alpha}^{(1)} = \arg \min_{\boldsymbol{\alpha} > 0} \underbrace{\frac{1}{2} \|\mathbf{y} - D^{(1)} \boldsymbol{\alpha}\|_2^2 + \lambda_1 \|\boldsymbol{\alpha}\|_1 + \frac{\lambda_2}{2} \|\boldsymbol{\alpha}\|_2^2}_{\text{fat-dictionary sparse coding}}$$

$$\underbrace{\boldsymbol{\alpha}^{(2)} = \arg \min_{\boldsymbol{\alpha} > 0} \frac{1}{2} \|\boldsymbol{\alpha}^{(1)} - D^{(2)}\boldsymbol{\alpha}\|_2^2 + \lambda_1 \|\boldsymbol{\alpha}\|_1 + \frac{\lambda_2}{2} \|\boldsymbol{\alpha}\|_2^2}_{\text{nonlinear dimension reduction}}$$

The first dictionary is over-complete (“fat”); the second is under-complete (“thin”) and forces clustering of sparse codes. Together they form the composite sparse coding module. Stacking up to 13 such modules with batch normalization, global average pooling, and softmax yields the supervised Sparse Coding Network (SCN).

H.3 Inference and training

Sparse coding step.

FISTA solves the elastic-net regularized problem

$$\boldsymbol{\alpha}^* = \arg \min_{\boldsymbol{\alpha} \geq 0} \underbrace{\frac{1}{2} \|\mathbf{x} - D\boldsymbol{\alpha}\|_2^2 + \frac{\lambda_2}{2} \|\boldsymbol{\alpha}\|_2^2}_{f(\boldsymbol{\alpha}) \text{ smooth}} + \underbrace{\lambda_1 \|\boldsymbol{\alpha}\|_1 + \mathcal{K}_{\geq 0}(\boldsymbol{\alpha})}_{g(\boldsymbol{\alpha}) \text{ nonsmooth}}.$$

The smooth part f has Lipschitz gradient $\nabla f(\boldsymbol{\alpha}) = D^T(D\boldsymbol{\alpha} - \mathbf{x}) + \lambda_2\boldsymbol{\alpha}$ with Lipschitz constant $L = \sigma_{\max}(D)^2 + \lambda_2$. The proximal operator of g is the *nonnegative soft-shrinkage operator*

$$\text{prox}_{\lambda_1 g/L}(z) = (\tau_{\lambda_1/L}(z))_+, \quad \tau_\sigma(z) = \text{sgn}(z)(|z| - \sigma)_+.$$

FISTA combines the proximal-gradient step with Nesterov acceleration:

$$\begin{aligned} \hat{\boldsymbol{\alpha}}_t &= \text{prox}_{\lambda_1 g/L}(\mathbf{u}_t - L^{-1}\nabla f(\mathbf{u}_t)), \\ t_{k+1} &= \frac{1 + \sqrt{1 + 4t_k^2}}{2}, \\ \mathbf{u}_{t+1} &= \hat{\boldsymbol{\alpha}}_t + \frac{t_k - 1}{t_{k+1}}(\hat{\boldsymbol{\alpha}}_t - \hat{\boldsymbol{\alpha}}_{t-1}). \end{aligned}$$

Theorem H.1 (FISTA $O(1/k^2)$ convergence, Beck–Teboulle 2009)

For the iterates above,

$$F(\hat{\boldsymbol{\alpha}}_k) - F(\boldsymbol{\alpha}^*) \leq \frac{2L\|\hat{\boldsymbol{\alpha}}_0 - \boldsymbol{\alpha}^*\|_2^2}{(k+1)^2}.$$

30–50 FISTA iterations suffice in practice. The matrix $D^T D + \lambda_2 I$ is pre-computed; inference requires no online inversion.

Remark H.2 (Why elastic-net)

Pure ℓ_1 gives strict sparsity but is unstable to correlated atoms: in the limit $D^T D \rightarrow I$ degenerate, two

highly correlated atoms can “swap” between iterations. The ℓ_2 term $\lambda_2 \|\boldsymbol{\alpha}\|_2^2$ smooths the objective, ensures strong convexity, and is essential for the back-propagation gradients below to be well-defined.

Training via back-propagation.

The dictionary update rule differentiates through the fixed-point structure of the elastic-net minimizer. With $\Lambda = \text{supp}(\boldsymbol{\alpha})$,

$$\begin{aligned}\frac{\partial L}{\partial D_{ij}} &= \frac{\partial L}{\partial \boldsymbol{\alpha}} \cdot (D^\top D + \lambda_2 I)_\Lambda^{-1} \left(\frac{\partial D_\Lambda^\top \boldsymbol{\alpha}}{\partial D_{ij}} - \frac{\partial D_\Lambda^\top D_\Lambda}{\partial D_{ij}} \mathbf{y}_\Lambda \right), \\ \frac{\partial L}{\partial x_i} &= \frac{\partial L}{\partial \boldsymbol{\alpha}} \cdot (D^\top D + \lambda_2 I)_\Lambda^{-1} \frac{\partial D_\Lambda^\top \mathbf{x}}{\partial x_i}, \\ \frac{\partial L}{\partial \lambda_{1j}} &= -\frac{\partial L}{\partial \boldsymbol{\alpha}} \cdot (D^\top D + \lambda_2 I)_\Lambda^{-1} \text{sign}(\boldsymbol{\alpha}_\Lambda)_j, \quad \lambda_{1j} \neq 0.\end{aligned}$$

These gradients support end-to-end supervised training. Regularization parameters λ_1, λ_2 are themselves task-adaptive.

H.4 Performance

On standard image classification benchmarks the supervised SCN closes much of the gap to state-of-the-art CNNs while using a smaller model, fewer parameters, and substantially less training data:

Dataset	Prior sparse coding	Supervised SCN	State-of-the-art CNN
CIFAR-10	81.40%	94.19%	96.42%
CIFAR-100	60.80%	80.07%	82.69%
STL-10	67.90%	83.11%	76.29%
MNIST (error)	0.54%	0.36%	0.21%
SVHN (error)	—	2.16%	1.77%

Notably SCN *outperforms* the CNN baseline on STL-10 where only 5,000 labeled training samples are available — sparse coding’s structural prior shines in low-data regimes.

H.5 Why SCN works: feature-map clustering

Visualizing intermediate feature maps reveals a qualitative difference between SCN and CNN. SCN’s nonlinear dimension-reduction step clusters feature maps that respond to similar high-level patterns; CNN feature maps are diverse but unclustered. The composite module’s second layer acts as a soft codebook: each $\boldsymbol{\alpha}^{(2)}$ effectively selects a small number of higher-level prototypes from a $\boldsymbol{\alpha}^{(1)}$ that scattered the signal

across many edge-like atoms. This is reminiscent of semi-NMF and Convex-NMF (§G): the clustering effect emerges from the nonnegativity constraint combined with the elastic-net regularizer.

H.6 Pros and cons

Advantages. Framework grounded in sparse-coding theory; less architecture tuning than CNN; competitive accuracy with smaller models and less data; clear interpretation of why deep networks succeed.

Disadvantages. Substantially longer training time (back-propagating through FISTA fixed-points is expensive); slightly slower inference; large-scale performance remains open; many downstream tasks (detection, segmentation, generation) have not been tried.

Course Synthesis

The thread connecting all eight weeks:

1. **Sparsity is a structural prior** (Lectures 1–2). Theorem 2.1 characterizes uniform recovery in terms of $\text{spark}(A) > 2s$; Theorem 2.7 asserts NP-hardness of P_0 in general.
2. **Greedy algorithms** (Lectures 3–4): MP, OMP, IHT, HTP, CoSaMP, SP. Their theoretical guarantees come in two flavors – coherence-based (small μ) or RIP-based.
3. **Spark, NSP, and RIP** (Lectures 5–7): the conditions on A that guarantee sparse recovery, with the tower of implications coherence \Rightarrow NSP \Rightarrow BP recovery, and RIP \Rightarrow NSP \Rightarrow BP/IHT/HTP/CoSaMP recovery. Theorem 5.2 (NSP characterizes BP), Theorem 5.8 (dual certificates), Theorem 7.3 (RIP $\delta_{2s} < \sqrt{2}-1 \Rightarrow$ BP).
4. **Random matrices** (Lecture 8): subgaussian and Gaussian matrices have RIP at the optimal scaling $m = O(s \log(N/s))$ (Theorem 8.2), matching the information-theoretic lower bound. Connection to JL embeddings (Theorem 8.6).
5. **Sparse Representation Classification** (Lecture 9): leverages CS theory to perform discriminative tasks; identity-augmented dictionaries handle occlusions natively; SCI provides confidence calibration.
6. **ℓ_1 -minimization at scale** (Lecture 10): ISTA/FISTA, ADMM, LARS-homotopy, primal–dual splitting (Chambolle–Pock), and IRLS form the modern toolkit.
7. **Multiple measurements and joint sparsity** (Lecture 12). The MMV model $Y = AX$ with row-sparse X admits sharper uniqueness ($s < (\text{spark}(A) - 1 + \text{rank}(Y))/2$) and reduced average-case sample complexity via SOMP and $\ell_{1,2}$ -minimization; Collaborative HiLasso unifies joint group sparsity with within-row sparsity.
8. **Low-rank structure: Matrix Completion and Robust PCA** (Lecture 13). Replace “sparse vector” with “low-rank matrix,” ℓ_0 with rank, ℓ_1 with nuclear norm. Candès–Recht 2009 guarantees recovery from $\Theta(rN \log^2 N)$ entries; Candès–Li–Ma–Wright 2011 separates a low-rank X from sparse gross noise E . SVT and ALM are the workhorse solvers.
9. **Dictionary learning** (Lectures 14, 22–24). The dictionary itself becomes a learnable object: K-means (one-hot codes), MOD (pseudo-inverse update), and K-SVD (one-atom-at-a-time SVD update) yield overcomplete dictionaries trained on data. Online dictionary learning, LLC, and SSC adapt the framework to streaming data, classification, and subspace clustering respectively.
10. **Nonnegative Matrix Factorization** (Lecture 15). Adding $D, X \geq 0$ forces parts-based representations

(Lee–Seung 1999). Multiplicative updates guarantee local convergence; orthogonal Semi-NMF reduces to relaxed K-means; Separable NMF + entropy minimization yield unique factorizations under a convex-hull assumption.

11. **Deep sparse coding** (Lecture 16). Stacking composite sparse-coding modules (fat-dictionary ℓ_1 followed by thin nonlinear dimension reduction) yields supervised Sparse Coding Networks (SCN) that match deep CNN accuracy on CIFAR/STL/MNIST while remaining theoretically grounded; training proceeds by back-propagation through FISTA fixed-points.

Reference.

S. Foucart and H. Rauhut, *A Mathematical Introduction to Compressive Sensing*, Birkhäuser/Springer, 2013. Additional sources: Candès & Recht (2009) for matrix completion; Candès, Li, Ma, Wright (2011) for Robust PCA; Aharon, Elad, Bruckstein (2006) for K-SVD; Lee & Seung (1999, 2001) for NMF; Sun et al. (2018, 2019) for Sparse Coding Networks.