



Minimax Reconstruction in (Convolutional) Sparse Dictionary Learning

Shashank Singh, Barnabás Póczos, and Jian Ma
Carnegie Mellon University

Introduction

- ▶ Sparse dictionary learning (a.k.a. sparse coding) is widely used to denoise data
- ▶ *Convolutional* Sparse Dictionary Learning (CSDL) is popular for data with translation-invariant features (e.g., images, sound, movies, genomics, etc.) [1]
 - ▷ Translation invariant dictionary \Rightarrow smaller dictionary and greater sparsity
- ▶ We study minimax denoising of convolutionally sparse data

Contributions

1. First bounds on minimax reconstruction/denoising risk of CSDL.
2. Most work in compressed sensing assumes mutually independent noise; we show this often-unrealistic assumption is not necessary for CSDL denoising.
3. Prior theory for sparse dictionary learning makes strong assumptions to ensure identifiability of dictionary (e.g. incoherence or restricted isometry properties). We show that, unlike dictionary recovery, sparse dictionary denoising *requires no assumptions whatsoever on dictionary*.

Notation

- ▶ **Multi-convolution:** For two matrices $R \in \mathbb{R}^{(N-n+1) \times K}$ and $D \in \mathbb{R}^{n \times K}$ with equal numbers of columns, we define multi-convolution \otimes by

$$R \otimes D = \sum_{k=1}^K R_k * D_k \in \mathbb{R}^N,$$

where $*$ denotes the usual convolution operator.

- ▶ **Matrix Norms:** For $A \in \mathbb{R}^{n \times m}$ and $p \in [0, \infty]$, we write

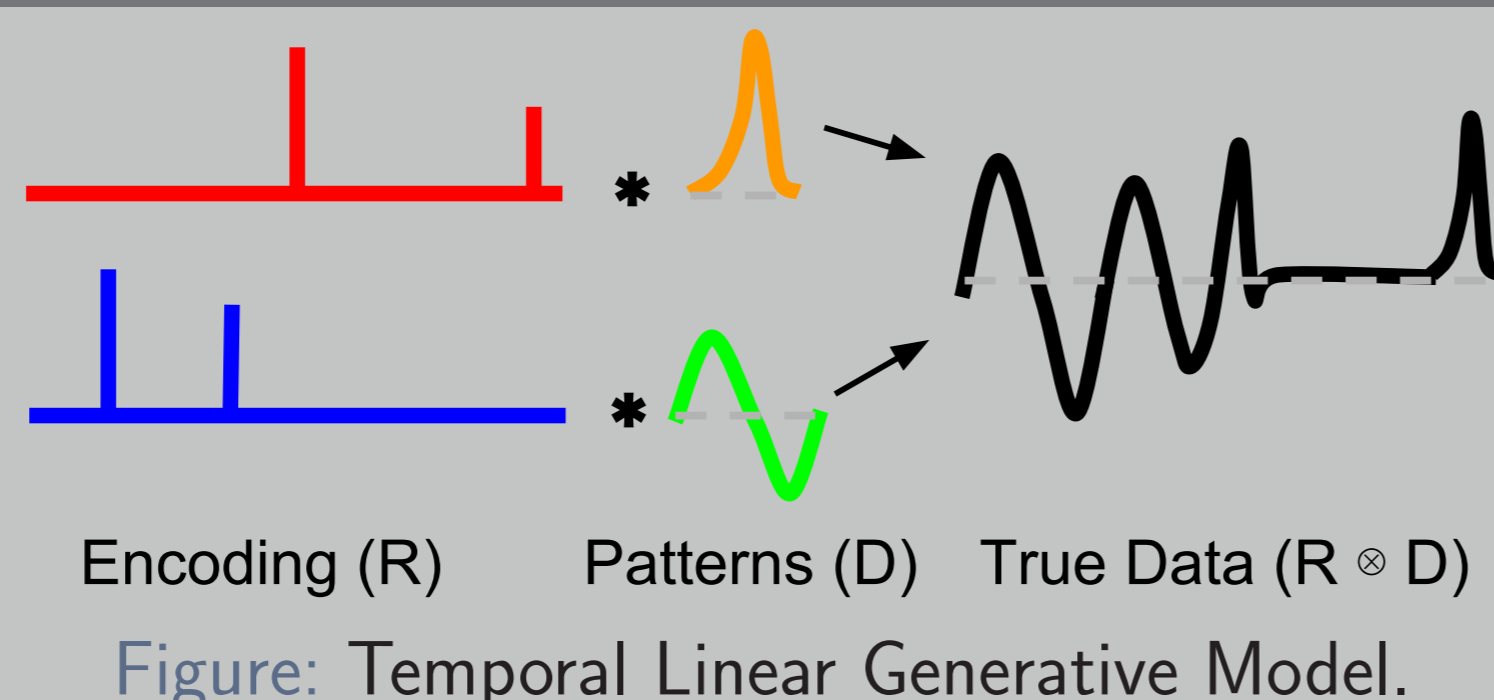
$$\|A\|_{p,q} := \left(\sum_{j=1}^m \left(\sum_{i=1}^n |A_{i,j}|^p \right)^{q/p} \right)^{1/q}.$$

- ▶ **Problem Domain:** For $N, K, n \in \mathbb{N}$ and $\lambda \geq 0$, we write

$$\mathcal{S}_\lambda := \left\{ (R, D) \in \mathbb{R}^{(N-n+1) \times K} \times \mathbb{R}^{n \times K} : \|D\|_{2,\infty} \leq 1, \|R\|_{1,1} \leq \lambda \right\}.$$

Modeling Assumptions

- ▶ Notation:
 - ▷ Single Observation $Y \in \mathbb{R}^N$
 - ▷ True signal $X \in \mathbb{R}^N$
 - ▷ Noise $\epsilon \in \mathbb{R}^N$
 - ▷ Dictionary $D \in \mathbb{R}^{n \times K}$
 - ▷ Encoding $R \in \mathbb{R}^{(N-n+1) \times K}$



- ▶ **Temporal Linear Generative Model (TLGM)** [2]:

$$Y = X + \epsilon, \quad \text{where } X = R \otimes D, \quad \text{for some } (R, D) \in \mathcal{S}_\lambda.$$

- ▶ We consider several possible noise assumptions:

1. ϵ is called **componentwise σ^2 -sub-Gaussian** ($\epsilon \in \text{CSG}(\sigma^2)$) if

$$\max_{i \in \{1, \dots, N\}} \mathbb{E} [e^{t\epsilon_i}] \leq e^{t^2\sigma^2/2}, \quad \text{for all } t \in \mathbb{R}.$$

2. ϵ is called **jointly σ^2 -sub-Gaussian** ($\epsilon \in \text{JSG}(\sigma^2)$) if

$$\mathbb{E} [e^{\langle t, \epsilon \rangle}] \leq e^{\|t\|_2^2 \sigma^2 / 2}, \quad \text{for all } t \in \mathbb{R}^N.$$

Note: In general, $\text{JSG}(\sigma^2) \subseteq \text{CSG}(\sigma^2)$ and $\text{CSG}(\sigma^2) \subseteq \text{JSG}(N\sigma^2)$.

If the entries of ϵ are mutually independent, then $\text{CSG}(\sigma^2) \subseteq \text{JSG}(\sigma^2)$.

3. Paper also has bounds under weaker bounded-moment assumptions.

CSDL Estimator

$$\hat{X}_\lambda = \hat{R}_\lambda \otimes \hat{D}_\lambda \quad \text{where} \quad (\hat{R}_\lambda, \hat{D}_\lambda) := \underset{(R,D) \in \mathcal{S}_\lambda}{\text{argmin}} \|Y - R \otimes D\|_2^2.$$

- ▶ Computable by alternating (between R and D) projected gradient descent
- ▶ Paper also has similar results for $\|R\|_{1,1}$ -penalized version

References

- [1] Vardan Pappyan, Jeremias Sulam, and Michael Elad. Working locally thinking globally-part I: Theoretical guarantees for convolutional sparse coding. *arXiv preprint arXiv:1607.02005*, 2016.
- [2] Bruno A. Olshausen. Probabilistic Models of the Brain. page 257, 2002.
- [3] Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery over the sphere. In *Sampling Theory and Applications (SampTA), 2015 International Conference on*, pages 407–410. IEEE, 2015.

Theoretical Results: Upper Bounds

- ▶ **Lemma 1 (Oracle Inequality):** If $Y = X + \epsilon$, then

$$\|X - \hat{X}_\lambda\|_2^2 \leq \underbrace{\inf_{(R,D) \in \mathcal{S}_\lambda} \|X - R \otimes D\|_2^2}_{\text{model misspecification}} + \underbrace{2\langle \epsilon, \hat{X}_\lambda - R \otimes D \rangle}_{\text{statistical error}}.$$

$\Rightarrow \hat{X}_\lambda$ robust to violation of TLGM assumption.

- ▶ **Theorem 2:** Under TLGM with $\epsilon \in \text{CSG}(\sigma^2)$,

$$\frac{1}{N} \mathbb{E} \left[\|\hat{X}_\lambda - X\|_2^2 \right] \leq \frac{4\lambda\sigma\sqrt{2n\log(2N)}}{N}.$$

- ▶ **Theorem 3:** Under TLGM with $\epsilon \in \text{JSG}(\sigma^2)$,

$$\frac{1}{N} \mathbb{E} \left[\|\hat{X}_\lambda - X\|_2^2 \right] \leq \frac{4\lambda\sigma\sqrt{2\log(2(N-n+1))}}{N}.$$

- ▶ **Note:** Under TLGM, we always have $\frac{1}{N} \mathbb{E} \left[\|\hat{X}_0 - X\|_2^2 \right] \leq \frac{\lambda^2}{N} \Rightarrow$ Under extreme sparsity/noise ($\lambda \ll \sigma\sqrt{n\log N}$), trivial estimate $\hat{X} = \mathbf{0}$ is better.

Theoretical Results: Lower Bounds

- ▶ **Minimax Error:** For $\lambda \in [0, \infty]$, $N > n \in \mathbb{N}$, and a class \mathcal{E} of \mathbb{R}^N -valued random variables,

$$M(\lambda, N, n, \mathcal{E}) := \inf_{\hat{X}: \mathbb{R}^N \rightarrow \mathbb{R}^N} \sup_{(R,D) \in \mathcal{S}_\lambda, \epsilon \in \mathcal{E}} \frac{1}{N} \mathbb{E} \left[\|\hat{X}(Y) - X\|_2^2 \right].$$

- ▶ **Theorem 4 (Componentwise Sub-Gaussian Noise):**

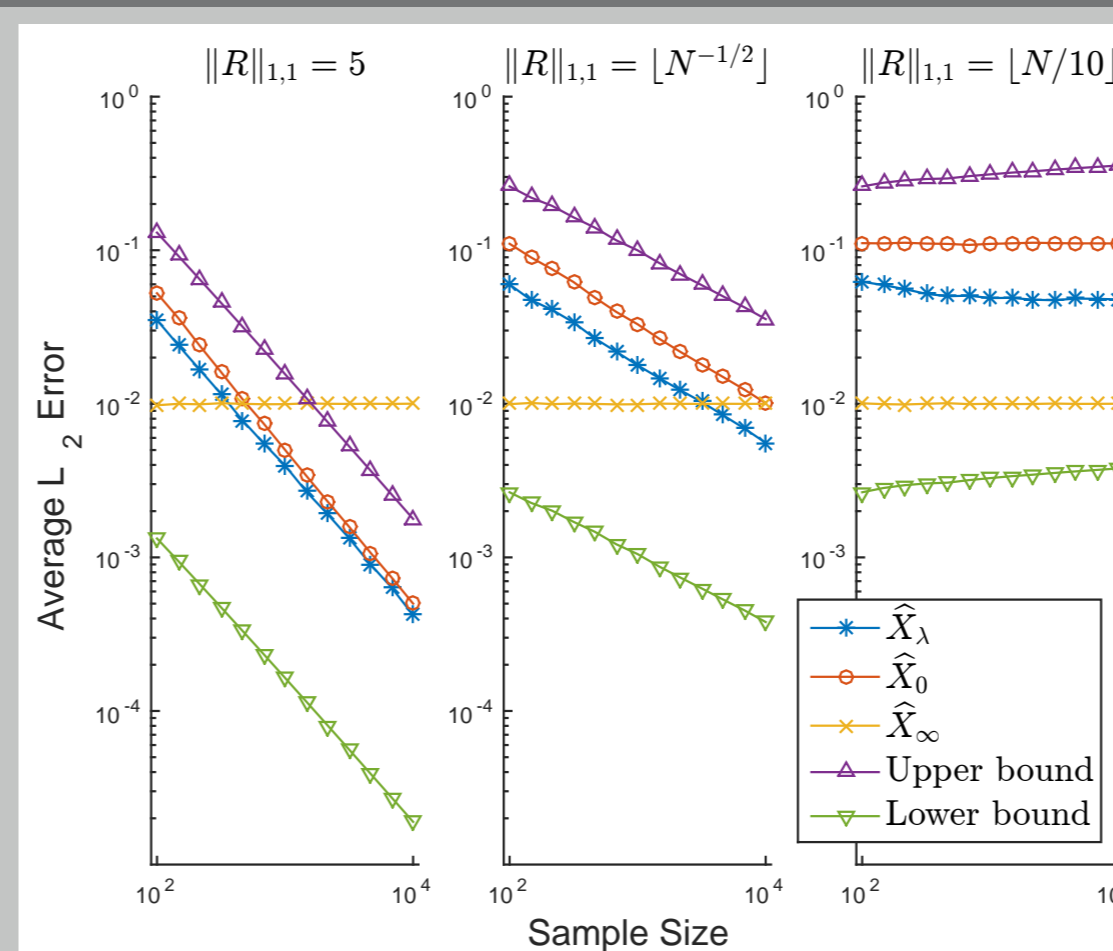
$$M(\lambda, N, n, \text{CSG}(\sigma^2)) \geq \frac{\lambda}{8N} \min \left\{ \lambda, \sigma\sqrt{n\log(N-n+1)} \right\}.$$

- ▶ **Theorem 5 (Jointly Sub-Gaussian Noise):**

$$M(\lambda, N, n, \text{JSG}(\sigma^2)) \geq \frac{\lambda}{8N} \min \left\{ \lambda, \sigma\sqrt{\log(N-n+1)} \right\}.$$

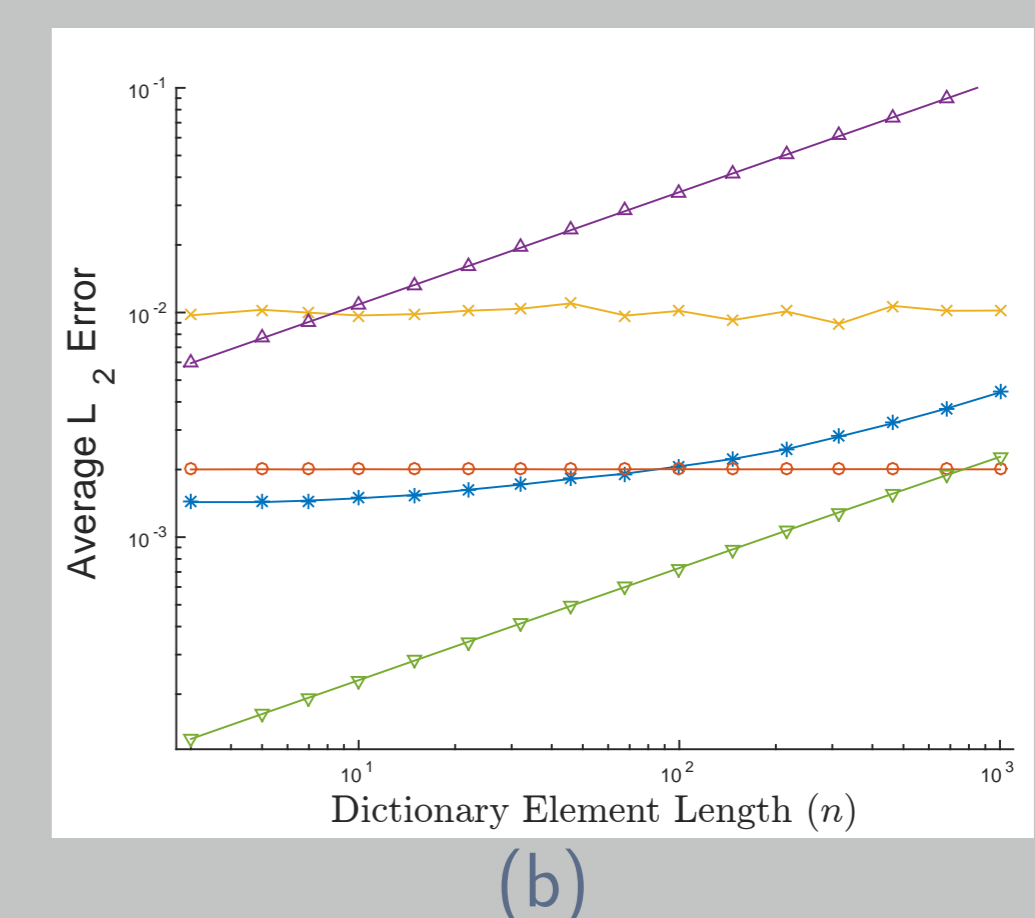
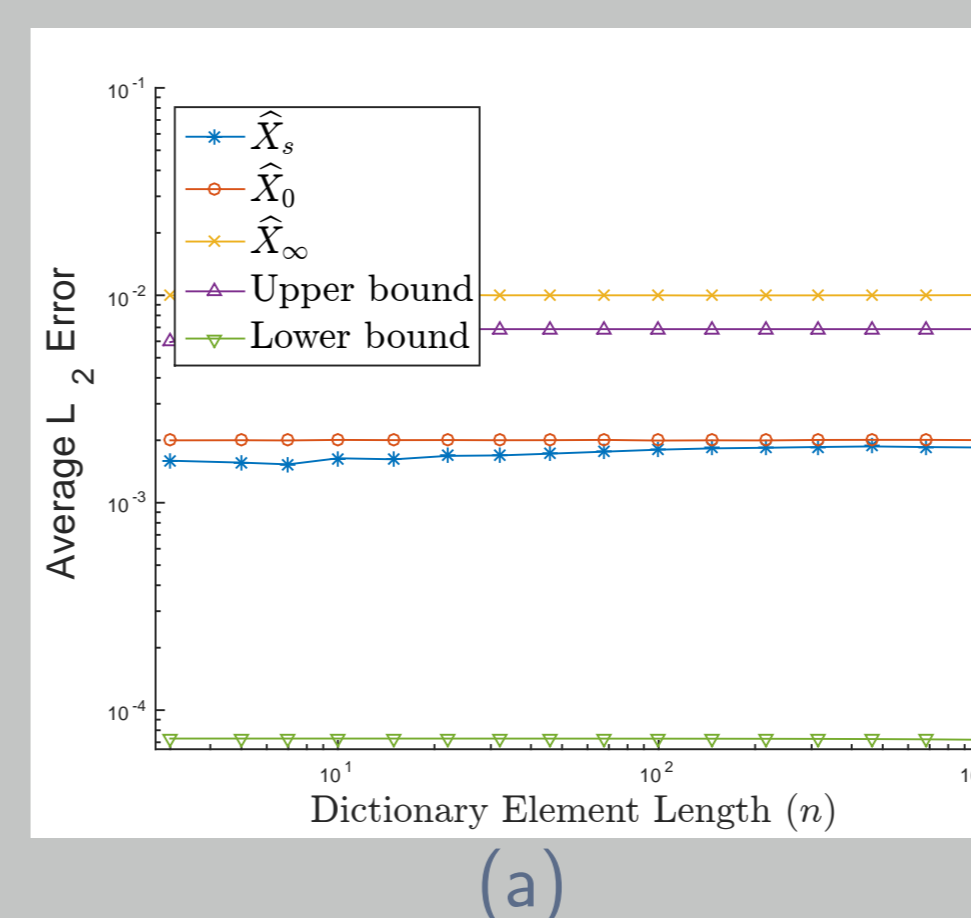
- ▶ **Note:** Lower bounds hold even when D is known in advance; estimating X is about as hard as estimating R (convolutional sparse recovery).

Simulation Results



Experiment 1: Average \mathcal{L}_2 -error as a function of signal length N , with sparsity scaling as:

1. $\|R\|_{1,1} = 5$ (Panel 1)
2. $\|R\|_{1,1} = \lfloor \sqrt{N} \rfloor$ (Panel 2)
3. $\|R\|_{1,1} = \lfloor N/10 \rfloor$ (Panel 3)



Experiment 2: Average \mathcal{L}_2 -error as a function of dictionary element length n , when entries of noise ϵ are (a) IID and (b) perfectly correlated.

Conclusions

- ▶ (Convolutional) sparse dictionary denoising is essentially **assumption-free**.
- ▶ For fixed n , CSDL is worst-case consistent (in reconstruction risk) if and only if

$$\frac{\lambda\sigma\sqrt{\log(N)}}{N} \rightarrow 0.$$

- ▶ When noise is independent, error is independent of dictionary atom length n .
- ▶ Similar results hold for classical dictionary learning (replace $n \rightarrow d, N \rightarrow \frac{N}{d}$).
- ▶ Alternating minimization appears minimax-rate optimal, consistent with recent results suggesting that local optima in SDL are often global optima [3]

Acknowledgments

This material is based upon work supported by the NSF Graduate Research Fellowship DGE-1252522 (to S.S.). The work was also partly supported by NSF IIS-1563887, Darpa D3M program, and AFRL (to B.P.), and NSF IIS-1717205 and NIH HG007352 (to J.M.).