

OVERVIEW

- We design a deep neural network to drastically improve LASSO speed and quality.
- Our network has fewer parameters and is easier to train.
- Our network achieves global linear convergence, better than sublinear and eventual-linear convergence of ISTA/FISTA.

UNFOLD ISTA TO NEURAL NETWORK

Problem: Recover a sparse vector x^* from its noisy measurements:

$$b = Ax^* + \varepsilon,$$

LASSO:

$$\underset{x}{\text{minimize}} \frac{1}{2} \|b - Ax\|_2^2 + \lambda \|x\|_1$$

Iterative shrinkage thresholding algorithm (ISTA) or FPC:

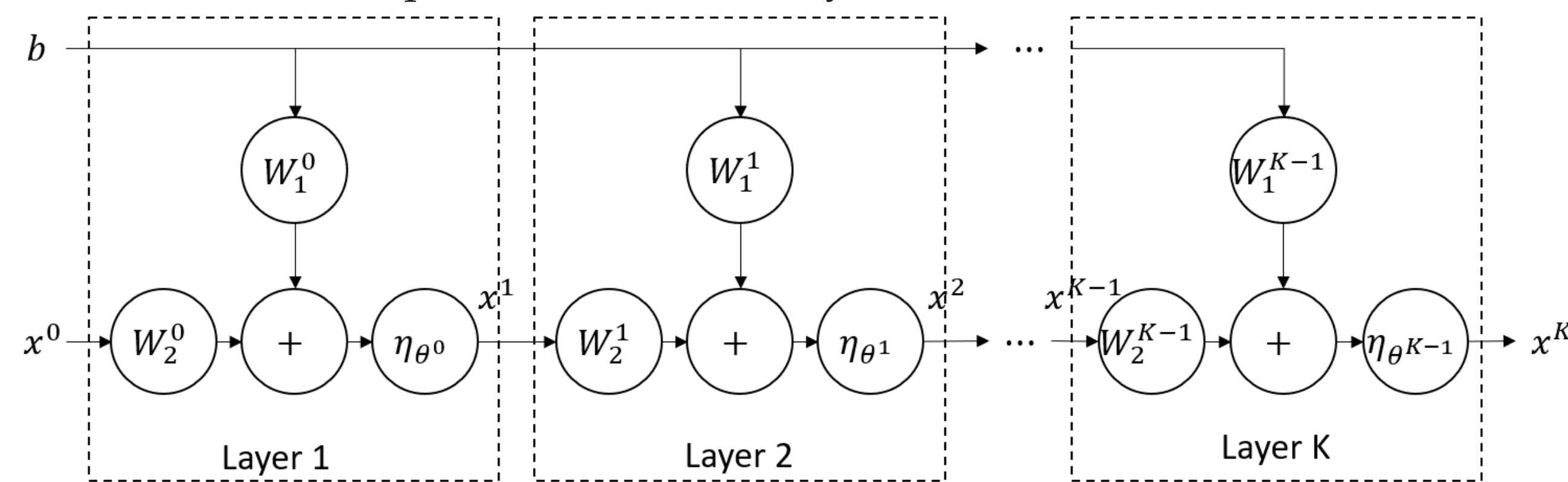
$$x^{k+1} = \eta_{\lambda/L} \left(x^k + \frac{1}{L} A^T (b - Ax^k) \right), \quad k = 0, 1, 2, \dots \quad (\text{ISTA})$$

where η_θ is soft-thresholding, λ and L are selected by hand or cross-validation. ISTA converges **sublinearly** and **eventually-linearly** to a **LASSO solution, not x^*** .

Neural network: unrolls ISTA to a feed-forward neural network, replace A, A^T in ISTA by free matrices, and truncates it to K iterations (known as Learned ISTA or LISTA [1]):

$$x^{k+1} = \eta_{\theta^k} (W_1^k b + W_2^k x^k), \quad k = 0, 1, \dots, K-1, \quad (\text{LISTA})$$

Inputs are x^0 and b . Output x^K is our recovery.



Training (deciding θ^k, W_1^k, W_2^k) For fixed A , we want to obtain parameters $\Theta^K = \{(W_1^k, W_2^k, \theta^k)\}_{k=0}^{K-1}$ such that x^K is close to x^* (**the ground truth**) for input $b = Ax^* + \varepsilon$ for almost all x^*, ε following certain distribution. In another word, given the distributions of x^* and ε , we

$$\underset{\Theta}{\text{minimize}} \frac{1}{2} \mathbb{E}_{x^*, \varepsilon} \|x^K(\Theta^K, b, x^0) - x^*\|_2^2.$$

Stochastic gradient descent (SGD) can be applied to solve the above minimization problem. The gradient of x^K on Θ^K can be obtained by the chain rule (back propagation).

Issues: largely many free parameters, training is slow

REFERENCES

- [1] K. Gregor and Y. LeCun, "Learning fast approximations of sparse coding," in *ICML*, 2010.
- [2] X. Chen, J. Liu, Z. Wang, and W. Yin, "Theoretical linear convergence of unfolded ista and its practical weights and thresholds," in *NIPS*, 2018.

LINKS

arXiv preprint:



Source codes:



IMPROVE BY WEIGHT COUPLING (CP)

New idea: exploit certain dependencies among W_1^k, W_2^k, θ^k to simplify the network and improve the recovery result.

Theorem 1 (Necessary Condition) Suppose $K = \infty$ and there is no noise $\varepsilon = 0$. Let $\{x^k\}_{k=1}^\infty$ be generated by (LISTA). If $x^k(\Theta^k, b, x^0) \rightarrow x^*$ as $k \rightarrow \infty$ uniformly for all sparse x^* , then the parameters $\{W_1^k, W_2^k, \theta^k\}_{k=0}^\infty$ are **not independent to each other** but must satisfy

$$W_2^k - (I - W_1^k A) \rightarrow 0, \quad \theta^k \rightarrow 0, \quad \text{as } k \rightarrow \infty. \quad (1)$$

Weight simplification: Couple $W_2^k = I - W_1^k A$ and simplify LISTA to:

$$x^{k+1} = \eta_{\theta^k} \left(x^k + W_1^k (b - Ax^k) \right), \quad k = 0, 1, \dots, K-1. \quad (\text{LISTA-CP})$$

Now, only $\bar{\Theta}^K = \{W_1^k, \theta^k\}_{k=0}^{K-1}$ need to be trained, yet recovery is still fast.

Theorem 2 (LISTA-CP trainability) Suppose $K = \infty$ and let $\{x^k\}_{k=1}^\infty$ be generated by (LISTA-CP). There exists a sequence of parameters $\{W_1^k, \theta^k\}$ such that

$$\|x^k(\bar{\Theta}^k, b, x^0) - x^*\|_2 \leq C_1 \exp(-ck) + C_2 \sigma, \quad \forall k = 1, 2, \dots,$$

holds for all (x^*, ε) satisfying some assumptions (see [2]), where $c, C_1, C_2 > 0$ are constants that depend only on A and the distribution of x^* , and σ is the noise level.

If $\sigma = 0$ (noiseless case), the k th layer output x^k converges to x^* linearly:

$$\|x^k - x^*\|_2 \leq C_1 e^{-ck}.$$

IMPROVE BY SUPPORT SELECTION (SS)

Before applying soft thresholding in each layer, trust a percentages of largest entries as "true support" to bypass thresholding.

$$x^{k+1} = \eta_{\text{ss}\theta^k} \left(x^k + W_1^k (b - Ax^k) \right), \quad k = 0, 1, \dots, K-1. \quad (\text{LISTA-CPSS})$$

Theorem 3 (Convergence of LISTA-CPSS) Suppose $K = \infty$ and let $\{x^k\}_{k=1}^\infty$ be generated by (LISTA-CPSS). There exists a sequence of parameters $\{W_1^k, \theta^k\}$ such that

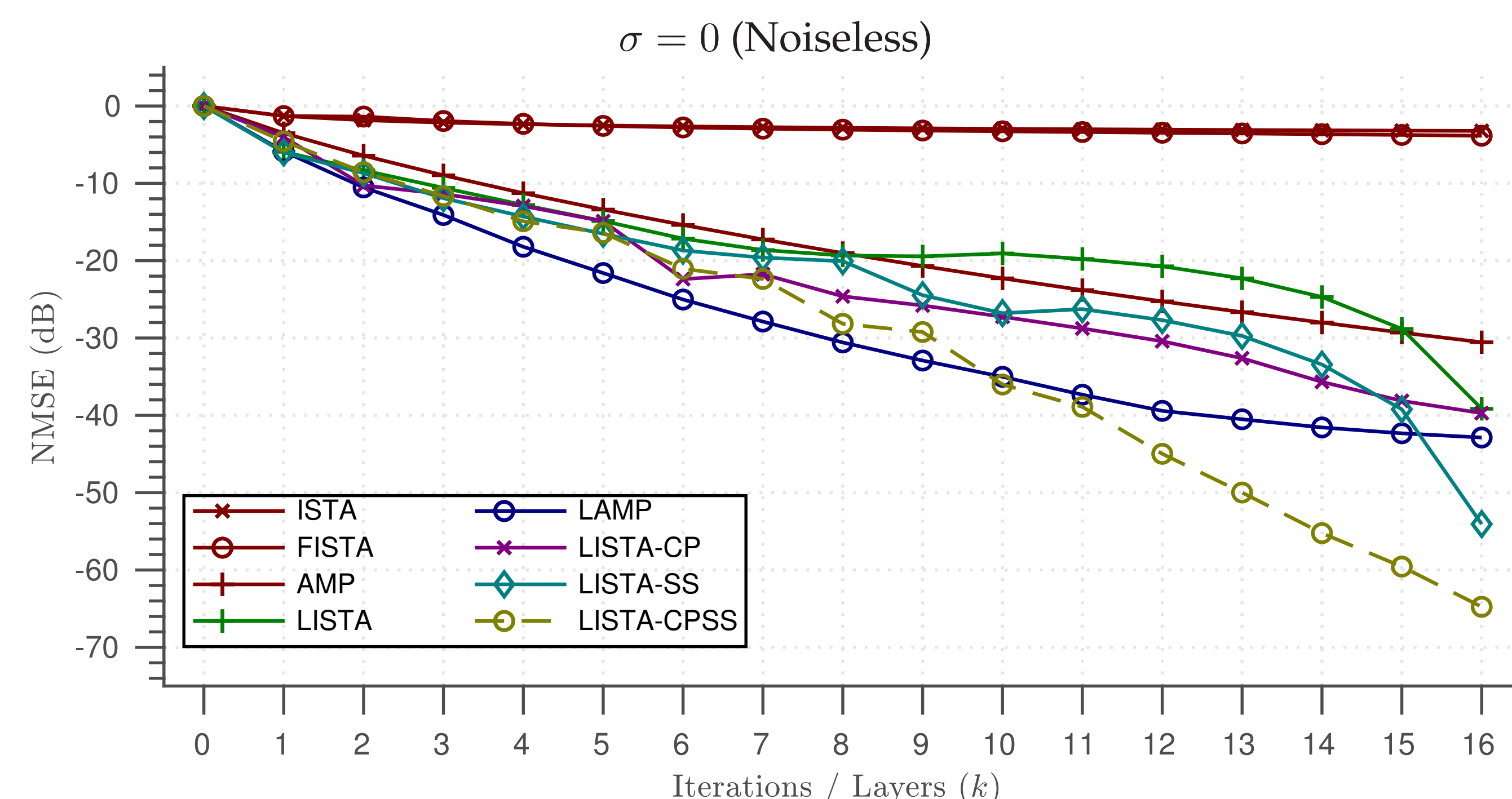
$$\|x^k(\bar{\Theta}^k, b, x^0) - x^*\|_2 \leq C_1 \exp\left(-\sum_{t=0}^{k-1} \tilde{c}_{\text{ss}}^t\right) + \tilde{C}_{\text{ss}} \sigma, \quad \forall k = 1, 2, \dots,$$

holds for all (x^*, ε) satisfying some assumptions. The convergence rate is better: $\tilde{c}_{\text{ss}}^k > c$ for large enough k . The recovery error is better: $\tilde{C}_{\text{ss}} < C_2$.

NUMERICAL VALIDATION — SUPPORT SELECTION

Validation of support selection:

LISTA with support selection (LISTA-SS) achieves linear convergence and better final performance. Coupled LISTA with support selection (LISTA-CPSS) yields the best performance.

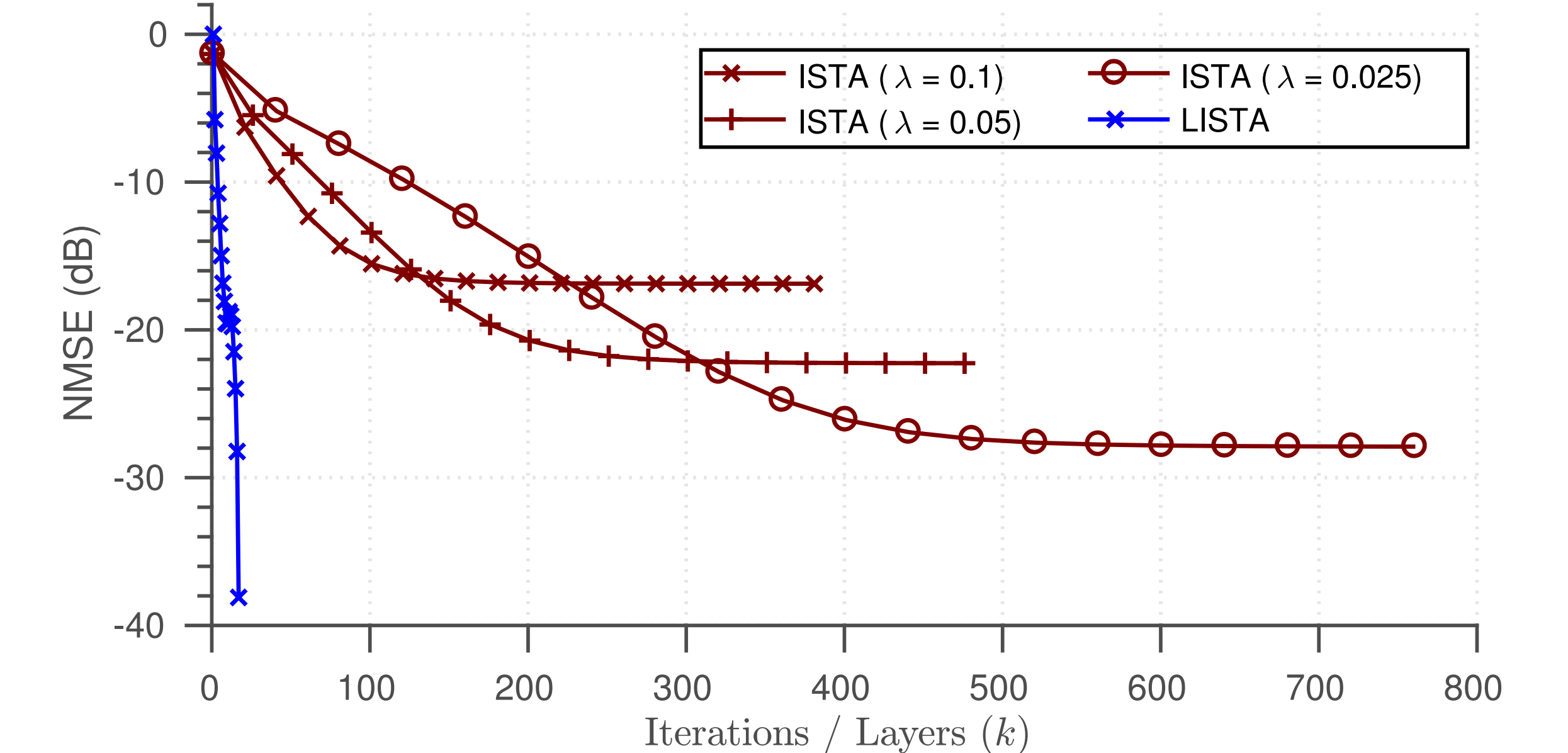


NUMERICAL VALIDATION

Data: fix $A \in \mathbb{R}^{250 \times 500}$, $A_{ij} \sim N(0, 1)$. Columns of A are normalized. Sample x^* with 10% nonzeros, each from the normal distribution. All plots below used the same 1000 samples.

Recovery speeds: ISTA \ll LISTA[1] $<$ LISTA-CP[2] $<$ LISTA-CPSS[2]

Baseline LISTA vs ISTA:



Validation of coupled LISTA ($\sigma=0$): Coupled LISTA (LISTA-CP) achieves linear convergence and stabilize intermediate steps compared to LISTA.

