



# Deep Sparse Coding and Dictionary Learning

Subhadip Mukherjee and Chandra Sekhar Seelamantula  
 Department of Electrical Engineering  
 Indian Institute of Science, Bangalore 560012, India  
 Email: {subhadip, chandra.sekhar}@ee.iisc.ernet.in



SPECTRUM LAB

## 1. Sparse Coding

► **Problem statement:** Given a signal  $y \in \mathbb{R}^m$  and an overcomplete dictionary  $A \in \mathbb{R}^{m \times n}$ , find an  $s$ -sparse vector  $x \in \mathbb{R}^n$ ,  $s \ll n$ , such that  $y \approx Ax \implies$  **Basis selection:** Express  $y = \sum_{i \in S} x_i a_i$ ,  $|S| = s$

► **Solution:** Seek the minimum  $\ell_0$ - (quasi)-norm solution: NP-hard  

$$\min_x \|x\|_0 \text{ subject to } y = Ax$$

## 2. Iterative Shrinkage-Thresholding Algorithm (ISTA)

► **Relaxation techniques: Solve sparsity-regularized least-squares**

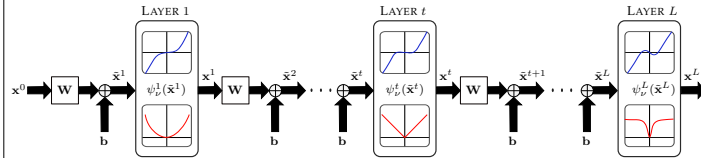
$$\hat{x} = \arg \min_x \underbrace{\frac{1}{2} \|y - Ax\|_2^2}_{f(x)} + \underbrace{\lambda \mathcal{R}(x)}_{\text{promotes sparsity}}$$

► **ISTA update:**  $x^{t+1} = T_{\lambda \eta} (x^t - \eta \nabla f(x^t)) \stackrel{\text{neural net.}}{\equiv} \psi^t(Wx^t + b)$ , where  $W = I - \eta A^T A$  and  $b = \eta A^T y$

► **Our approach:** Learning the activation  $\psi^t$  in a data-driven manner  
 ► Advantage-1: Number of parameters to learn does not grow with  $n$   
 ► Advantage-2: Possible to learn a rich variety of activations and regularizers

## 3. Architecture and Training of LETnet

**LETnet architecture:** Model  $\psi^t(u) = \sum_{k=1}^K c_k^t \phi_k(u)$ ,  $\phi_k(u) = u \exp\left(-\frac{(k-1)u^2}{2\tau^2}\right)$



**LETnet training:**

► **Training cost**  $J(c) = \frac{1}{2} \sum_{q=1}^N \|x_q^L(y_q, c) - x_q\|_2^2$ , where  $c = (c^1; c^2; \dots; c^L)$

► **Second-order Hessian-free optimization**

►  $J_i^{(q)}(c_i + \delta_c) = J(c_i) + \delta_c^T g_i + \frac{1}{2} \delta_c^T H_i \delta_c$   
 ► Compute optimal direction using conjugate-gradient (CG) at every epoch  $i$

$$\delta_c^* = \arg \min_{\delta_c} J_i^{(q)}(c_i + \delta_c) + \gamma \|\delta_c\|_2^2$$

► Update parameters as  $c_{i+1} \leftarrow c_i + \delta_c^*$

► **Two ingredients of CG**

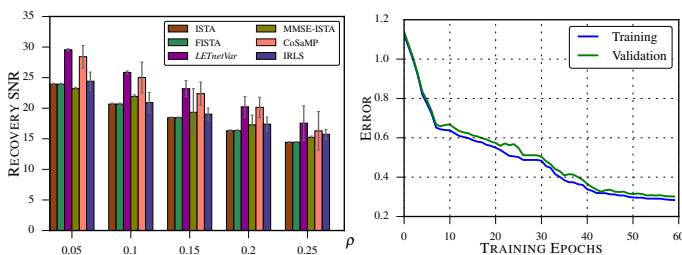
► **Gradient**  $g_i = \nabla J(c_i)|_{c=c_i}$  and the **Hessian-vector product**  $H_i u$ , for any  $u$   
 ► **Computing the hessian-vector product**  
 ►  $\mathcal{R}_u(h(c)) = \lim_{\alpha \rightarrow 0} \frac{h(c+\alpha u) - h(c)}{\alpha} \implies H_i u = \mathcal{R}_u(\nabla J(c_i)|_{c=c_i})$   
 ►  $g_i$  and  $H_i u$  are computed using back-propagation

## 4. Validation on Synthetic Signals

► **Data generation:**

►  $n = 256$ ,  $m = \lceil 0.7n \rceil$   
 ►  $A_{ij} \sim \mathcal{N}(0, 1/m)$   
 ►  $x = x_{\text{supp}} \odot x_{\text{mag}}$ , where  $x_{\text{supp}} \sim \text{Bernoulli}(\rho)$  and  $x_{\text{mag}} \sim \mathcal{N}(0, 1)$   
 ►  $0 < \rho < 1$ : smaller the value of  $\rho$ , sparser the vector  $x$   
 ►  $\lambda$  chosen optimally using cross-validation  
 ► 100 examples used for training and testing  
 ► Performance averaged over 10 independent trials

► **Reconstruction SNR and training error:**



(a)  $SNR_{\text{input}} = 20 \text{ dB}$

(b) Training and validation error

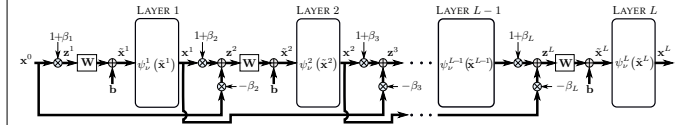
► Both training and validation errors decrease; no overfitting

## 5. fLETnet: A Deep Architecture Motivated by FISTA

► **FISTA iterations unfolded**

1.  $z^t = (1 + \beta_t)x^{t-1} - \beta_t x^{t-2}$   
 2.  $x^t = T_v(z^t)$ , where  $\bar{x}^t = Wz^t + b$  for  $t = 1, 2, \dots, L$

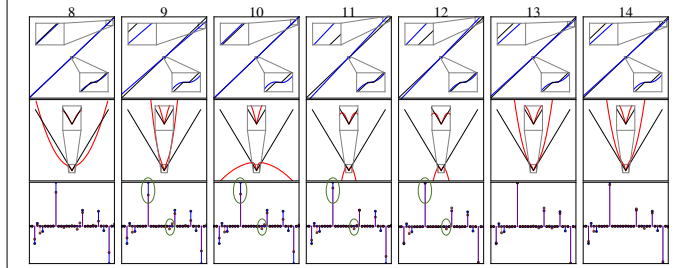
► **fLETnet architecture:**



► **Key Features of fLETnet:**

► Direct links from two previous layers (second-order memory)  
 ► Circumvents the issue of vanishing/exploding gradient  
 ► The resulting architecture **fLETnet** is essentially a deep residual network  
 ► Equivalent performance as the **LETnet** with half as many layers

► **What regularizers does the fLETnet learn?**



## 6. Comparison of Testing Run-times

Algorithm	per iteration/layer run-time (in milliseconds)	# layers/ iterations	total time (in milliseconds)
ISTA	0.0331	1000	33.10
FISTA	0.0394	1000	39.40
LETnet	0.0895	100	8.95
fLETnet	0.1088	50	5.44
MMSE-ISTA	0.6184	1000	618.40
CoSaMP	11.7672	50	588.36
IRLS	5.2784	50	263.92

## 7. Deep Dictionary Learning

► **Problem statement:** Given a set of signals  $\{y_j\}_{j=1}^N \in \mathbb{R}^m$ , learn an overcomplete dictionary  $A \in \mathbb{R}^{m \times n}$  and  $s$ -sparse vectors  $\{x_j\}_{j=1}^N \in \mathbb{R}^n$ ,  $s \ll n$ , such that  $y_j \approx Ax_j$  for all  $j$

► **Proposed approach:**  $\hat{A} = \min_A \sum_{j=1}^N \|A \text{net}_{y_j}(A) - y_j\|_2^2$

►  $\text{net}_{y_j}(A)$  is the output of ISTA corresponding to the signal  $y_j$  and dictionary  $A$   
 ► Gradient descent:  $A \leftarrow A - \mu \nabla J(A)$ ,  $\nabla J(A)$  requires only matrix-vector products

► **Advantages over conventional dictionary learning algorithms:**

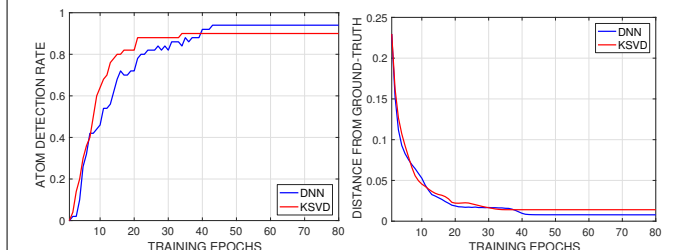
► Online and distributed implementations  
 ► Certain desirable properties on the dictionary, such as incoherence, can be promoted by adding a penalty and appropriately modifying the gradient

► **Performance metrics:**

► **Atom detection rate:**  $\frac{\#\{\text{recovered atoms}\}}{n}$   
 ► **Distance of  $\hat{A}$  from  $A^*$ :**  $\kappa = \frac{1}{n} \sum_{i=1}^n \min_{1 \leq j \leq n} (1 - |a_i^T a_j^*|)$

► **Experimental validation:**

►  $n = 50$ ,  $m = 20$ , sparsity  $s = 3$ , # examples  $N = 2000$ ,  $SNR_{\text{input}} = 30 \text{ dB}$   
 ► # layers  $L = 200$ , # training epochs  $N_{\text{epoch}} = 80$





## Outline

- What is sparse coding?
- Iterative shrinkage-thresholding algorithms (ISTA)
  - ▶ Iterative unfolding and connections to **deep neural networks** (DNNs)
  - ▶ Building a learnable model for sparse coding
- Our contributions
  - ▶ Modeling the non-linearity using a **linear expansion of thresholds** (LETs)
  - ▶ Parametric flexibility for designing regularizers
  - ▶ Efficient second-order (**Hessian-free**) optimization for learning
  - ▶ Reducing the number of layers and link with **deep residual networks**
  - ▶ Building a deep architecture for **dictionary learning**
- Simulation results and insights
- Summary and future works

## What is sparse coding?

- **Problem statement:** Given a signal  $\mathbf{y} \in \mathbb{R}^m$  and an overcomplete dictionary  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , find an  $s$ -sparse vector  $\mathbf{x} \in \mathbb{R}^n$ ,  $s \ll n$ , such that  $\mathbf{y} = \mathbf{A}\mathbf{x}$

$$\mathbf{y} \in \mathbb{R}^m = \mathbf{A} \in \mathbb{R}^{m \times n} \mathbf{x} \in \mathbb{R}^n$$

- **Basis selection:** Express  $\mathbf{y} = \sum_{i \in S} x_i \mathbf{a}_i$ ,  $|S| = s$ 
  - Estimating  $\mathbf{A}$  from given data is the problem of **dictionary learning**
- **Solution:** Seek the minimum  $\ell_0$ -(quasi)-norm solution: Combinatorially hard

$$\min_{\mathbf{x}} \|\mathbf{x}\|_0 \text{ subject to } \mathbf{y} = \mathbf{A}\mathbf{x}$$

## Iterative shrinkage-thresholding algorithm (ISTA) meets neural network

- Relaxation techniques: Solve sparsity-regularized least-squares

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \underbrace{\frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2}_{f(\mathbf{x})} + \underbrace{\lambda \mathcal{R}(\mathbf{x})}_{\text{promotes sparsity}}$$

- ISTA update rule (Daubechies et al., 2004)
  - ▶  $\mathbf{x}^{t+1} = T_{\lambda\eta}(\mathbf{x}^t - \eta \nabla f(\mathbf{x}^t))$ , where  $T_\nu$  denotes soft-thresholding with threshold  $\nu$
- Unfolding of ISTA iterations

$$\mathbf{x}^{t+1} = \underbrace{\psi^t}_{\text{non lin.}} \left( \underbrace{\mathbf{W}\mathbf{x}^t + \mathbf{b}}_{\text{affine}} \right), \text{ where } \mathbf{W} = \mathbf{I} - \eta \mathbf{A}^\top \mathbf{A} \text{ and } \mathbf{b} = \eta \mathbf{A}^\top \mathbf{y}$$

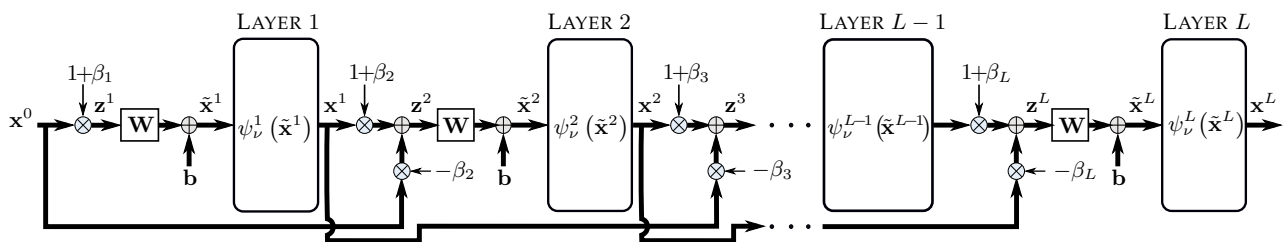
- Building learnable models
  - ▶ Learn the linear transformation parameters  $\mathbf{W}$  and  $\mathbf{b}$  from the data – data-intensive
  - ▶ Learn the nonlinear shrinkage function
- LETnet: Modeling the activation nonlinearity using LETs

- ▶ LET modeling:  $\psi^t(u) = \sum_{k=1}^K c_k^t \phi_k(u)$ , where  $\phi_k(u) = u \exp\left(-\frac{(k-1)u^2}{2\tau^2}\right)$

## Unfolding of FISTA (Nesterov, 1980s) and the *fLETnet*

- FISTA iterations unfolded

- $\mathbf{z}^t = (1 + \beta_t)\mathbf{x}^{t-1} - \beta_t\mathbf{x}^{t-2}$
- $\tilde{\mathbf{x}}^t = \mathbf{W}\mathbf{z}^t + \mathbf{b}$
- $\mathbf{x}^t = T_\nu(\tilde{\mathbf{x}}^t)$ , for  $t = 1, 2, \dots, L$



- fLETnet*: Network architecture motivated by FISTA

- Replace  $T_\nu$  with a parametric activation  $\psi^t$
- Direct links from two previous layers (second-order memory)
- 1 link: identity mapping; thus no new parameters to learn
- Circumvents the issue of vanishing/exploding gradient
- The *fLETnet* architecture is essentially a deep residual network (He et al., 2015)
- Results in equivalent performance as the *LETnet* with half as many layers

## Training the *LETnet* and the *fLETnet*

- Training dataset  $\mathcal{D}$  consists of  $N$  examples  $\{(\mathbf{y}_q, \mathbf{x}_q)\}_{q=1}^N$ , where  $\mathbf{y}_q = \mathbf{A}\mathbf{x}_q + \xi_q$ , for a given  $\mathbf{A}$

- Training cost  $J(\mathbf{c}) = \frac{1}{2} \sum_{q=1}^N \|\mathbf{x}_q^L(\mathbf{y}_q, \mathbf{c}) - \mathbf{x}_q\|_2^2$ , where  $\mathbf{c} = (\mathbf{c}^1; \mathbf{c}^2; \dots; \mathbf{c}^L)$

- Second-order optimization

- ▶ Quadratic approximation in the  $i^{\text{th}}$  training epoch

$$J_i^{(q)}(\mathbf{c}_i + \delta_{\mathbf{c}}) = J(\mathbf{c}_i) + \delta_{\mathbf{c}}^T \mathbf{g}_i + \frac{1}{2} \delta_{\mathbf{c}}^T \mathbf{H}_i \delta_{\mathbf{c}}$$

- ▶ Compute optimal direction using conjugate-gradient (CG) at every epoch  $i$

$$\delta_{\mathbf{c}}^* = \arg \min_{\delta_{\mathbf{c}}} J_i^{(q)}(\mathbf{c}_i + \delta_{\mathbf{c}}) + \gamma \|\delta_{\mathbf{c}}\|_2^2$$

- ▶ Update parameters as  $\mathbf{c}_{i+1} \leftarrow \mathbf{c}_i + \delta_{\mathbf{c}}^*$

- Two ingredients of CG

- ▶ Gradient  $\mathbf{g}_i = \nabla J(\mathbf{c})|_{\mathbf{c}=\mathbf{c}_i}$

- ▶ Hessian-vector product  $\mathbf{H}_i \mathbf{u}$  for any vector  $\mathbf{u}$ , where  $\mathbf{H}_i = \nabla^2 J(\mathbf{c})|_{\mathbf{c}=\mathbf{c}_i}$

- $\mathcal{R}_{\mathbf{u}}(\mathbf{h}(\mathbf{c})) = \lim_{\alpha \rightarrow 0} \frac{\mathbf{h}(\mathbf{c} + \alpha \mathbf{u}) - \mathbf{h}(\mathbf{c})}{\alpha} \implies \mathbf{H}_i \mathbf{u} = \mathcal{R}_{\mathbf{u}}(\nabla_{\mathbf{c}} J(\mathbf{c}))|_{\mathbf{c}=\mathbf{c}_i}$

## Experimental details and parameter settings

- Data generation

- ▶  $n = 256, m = \lceil 0.7n \rceil$
- ▶  $A_{i,j} \sim \mathcal{N}(0, 1/m)$
- ▶  $\mathbf{x} = \mathbf{x}_{\text{supp}} \odot \mathbf{x}_{\text{mag}}$ , where  $\mathbf{x}_{\text{supp}} \sim \text{Bernoulli}(\rho)$  and  $\mathbf{x}_{\text{mag}} \sim \mathcal{N}(0, 1)$
- ▶  $0 < \rho < 1$ : smaller the value of  $\rho$ , sparser the vector  $\mathbf{x}$
- ▶  $\lambda$  chosen optimally using cross-validation
- ▶  $N_{\text{train}} = 100$  examples used for training
- ▶  $N_{\text{test}} = 100$  examples used for testing
- ▶ Performance averaged over 10 independent trials



## Performance assessment of *LETnetVar*

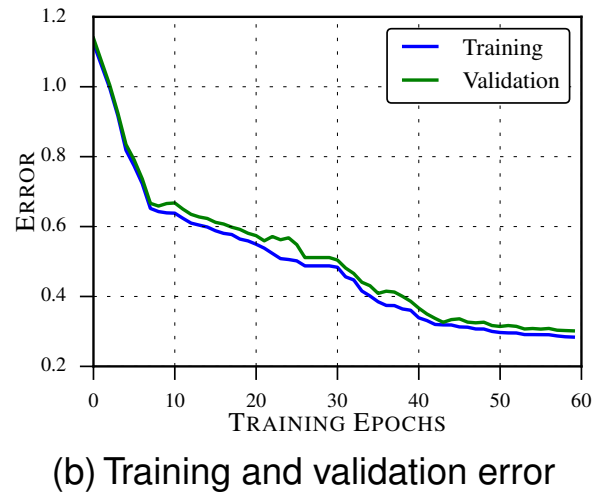
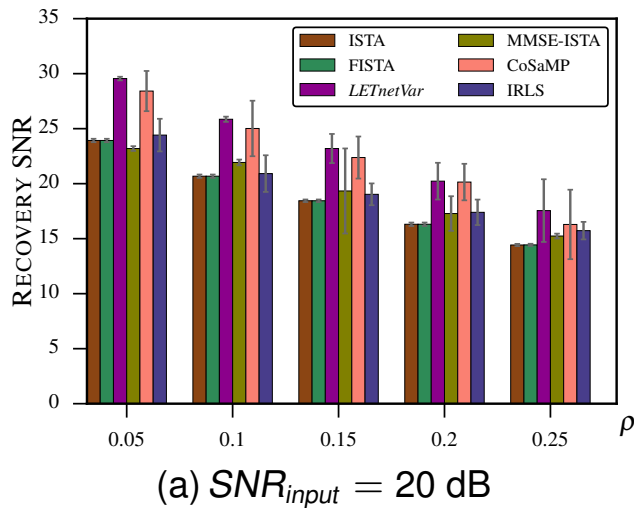


Figure: Comparison of ensemble-averaged reconstruction SNR and its standard deviation.

### ● Observations

- ▶ *LETnetVar* with 100 layers performs 3 to 4 dB better than ISTA, with optimally chosen hyper-parameter ( $\lambda$ ).
- ▶ When measurements are more noisy, the improvement in recovery SNR of *LETnet* was found to be higher.
- ▶ Training and validation errors reduce monotonically with training epochs

What regularizers did the *fLETnet* learn?

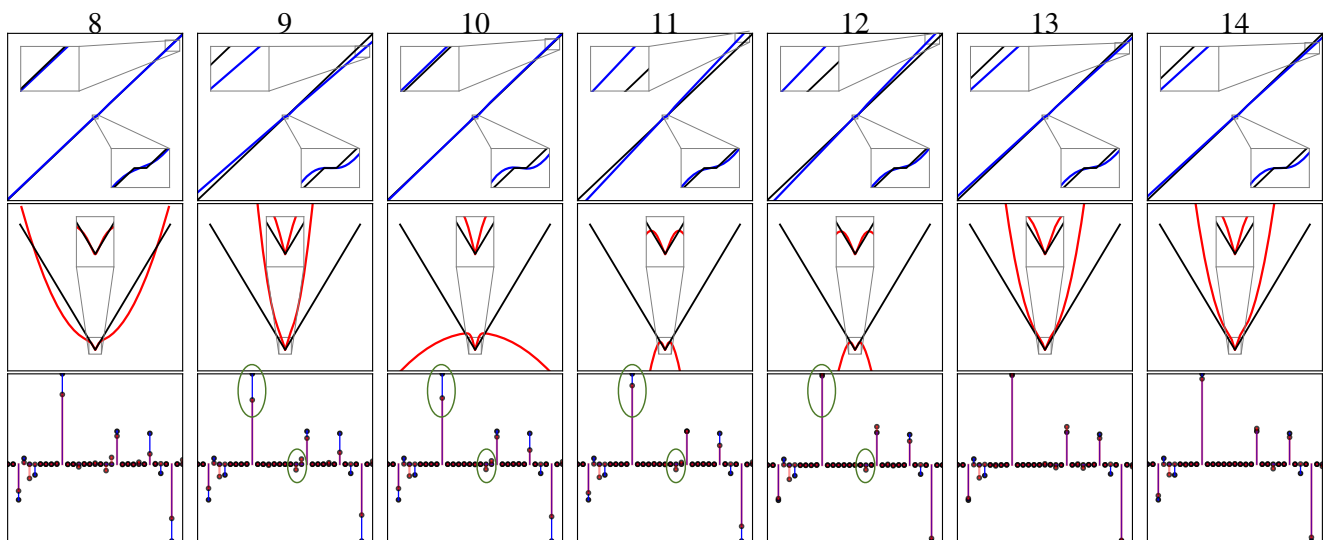


Figure: The learnt LET functions (blue) of *fLETnet* and the corresponding induced regularizers (red).

- Key observations

- ▶ Can learn a wide variety of regularizers
- ▶ Signal component missed in a layer can be recovered subsequently
- ▶ Balance between signal preservation and noise cancellation

## A comparison of testing run-times

Algorithm	per iteration/layer run-time (in milliseconds)	number of layers/ iterations	total time (in milliseconds)
ISTA	0.0331	1000	33.10
FISTA	0.0394	1000	39.40
<i>LETnet</i>	0.0895	100	8.95
<i>fLETnet</i>	0.1088	50	5.44
<i>MMSE-ISTA</i>	0.6184	1000	618.40
CoSaMP	11.7672	50	588.36
IRLS	5.2784	50	263.92

## Deep dictionary learning

- **Problem statement:** Given a set of signals  $\{\mathbf{y}_j\}_{j=1}^N \in \mathbb{R}^m$ , learn an overcomplete dictionary  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $s$ -sparse vectors  $\{\mathbf{x}_j\}_{j=1}^N \in \mathbb{R}^n$ ,  $s \ll n$ , such that  $\mathbf{y}_j \approx \mathbf{A}\mathbf{x}_j$
- **Proposed approach:**

$$\hat{\mathbf{A}} = \min_{\mathbf{A}} \sum_{j=1}^N \|\mathbf{A} \operatorname{net}_{\mathbf{y}_j}(\mathbf{A}) - \mathbf{y}_j\|_2^2$$

- ▶  $\operatorname{net}_{\mathbf{y}_j}(\mathbf{A})$  is the output of ISTA corresponding to the signal  $\mathbf{y}_j$  and dictionary  $\mathbf{A}$
  - ▶ Gradient descent:  $\mathbf{A} \leftarrow \mathbf{A} - \mu \nabla J(\mathbf{A})$
  - ▶ Computing  $\nabla J(\mathbf{A})$  requires only matrix-vector products
- **Advantages over conventional dictionary learning algorithms:**
    - ▶ Online implementation
    - ▶ Distributed implementation for a large training dataset
    - ▶ Certain desirable properties on the dictionary, such as incoherence, can be promoted by adding a penalty and appropriately modifying the gradient

## Deep dictionary learning: Performance validation

- Data generation

- ▶  $n = 50, m = 20$
- ▶ Number of examples  $N = 2000$
- ▶ Sparsity  $s = 3$
- ▶ Number of layers  $L = 200$
- ▶ Add noise to the training dataset such that  $\text{SNR}_{\text{input}} = 30 \text{ dB}$

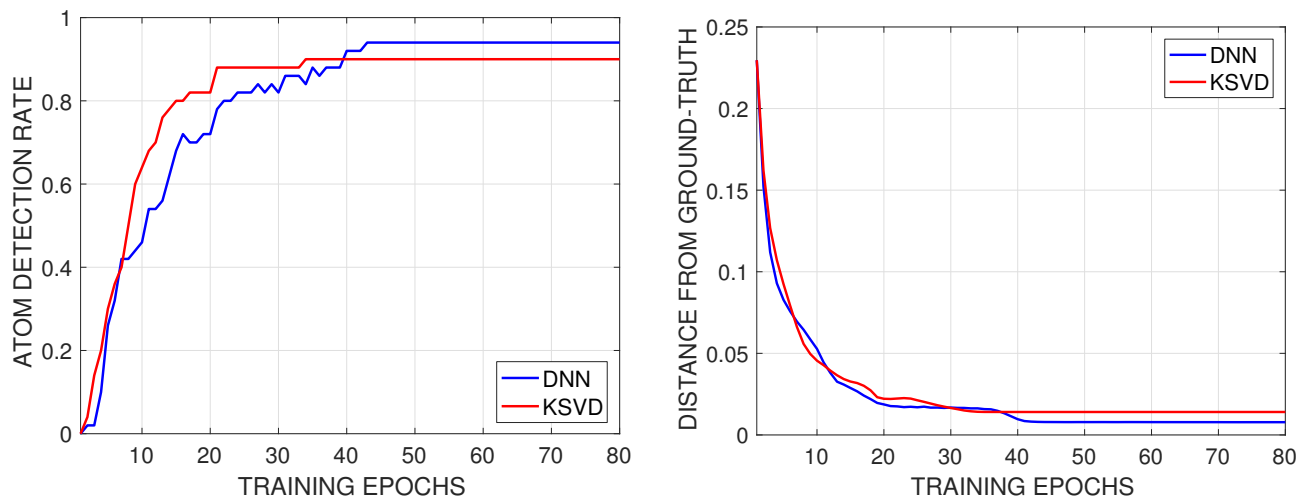


Figure: Comparison of DNN-based dictionary learning with K-SVD

## Summary and future works

### Summary

- Sparse coding as a function approximation problem
- Link between ISTA and a DNN
- Learning data-driven parameterized nonlinearities
- Efficient Hessian-free second-order optimization for training the network
- Link between FISTA and deep residual network
- The induced regularizers are nonstandard! Go beyond  $\ell_1$  and  $\ell_0$
- Extension to dictionary learning

### Future works

- Restricted to a fixed dictionary. How about adaptive/slowly varying dictionaries?
- A generic DNN solution to any iterative algorithm for inverse problems?