

Deep Learning for Sparse Signal Processing

Presented by:

Rubin Jose Peter

SPC Lab

October 20, 2018

Overview

- 1 Deep Learning for Sparse Signal Processing
- 2 Generative Adversarial Networks
- 3 Compressive Sensing Using Generative Models
- 4 New Framework for Sparse Signal Processing
- 5 Coupled Dictionary Learning

DNN for Sparse Signal Recovery

● Learning to optimize

- Signal processing algorithm is approximated by a Deep Neural Network (DNN)
- DNN requires only simple arithmetic operations to approximate the algorithm
- Effectiveness of the proposed approach was demonstrated by implementing WMMSE algorithm using DNN

● Sparse signal recovery using DNN: approaches

- Training a DNN using ground truth(y, x)
- Training a DNN using the input/output of a sparse recovery algorithm
- Approximating each layer of a neural network by the input/output of an iterative sparse recovery algorithm

● Observation

- Performance of the DNN based sparse signal recovery depends on the architecture of the neural network and number of training data
- **Extended Target Detection** problem: DNN based implementation may resolve boundary and block size mismatches

Generative Adversarial Networks (GAN)

Generative Adversarial Network

- Simultaneously learn two models:¹
 - A generative model \mathbf{G} : captures the data distribution
 - A discriminative model \mathbf{D} : estimates the probability that a sample came from training data rather than \mathbf{G}
- Training data: $\mathbf{x} \sim \mathbf{p}_{\text{data}}$
- Generator distribution: $\mathbf{G}(\mathbf{z}) \sim \mathbf{p}_{\mathbf{g}}$
- \mathbf{D} maximizes: the probability of assigning correct label to both training samples and samples from \mathbf{G}
- \mathbf{G} minimizes: $\log(\mathbf{1} - \mathbf{D}(\mathbf{G}(\mathbf{z})))$

$$\min_{\mathbf{G}} \max_{\mathbf{D}} V(\mathbf{D}, \mathbf{G}) \tag{1}$$
$$V(\mathbf{D}, \mathbf{G}) = \mathbf{E}_{\mathbf{x} \sim \mathbf{p}_{\text{data}}} [\log \mathbf{D}(\mathbf{x})] + \mathbf{E}_{\mathbf{z} \sim \mathbf{p}_{\mathbf{z}}(\mathbf{z})} [\log(\mathbf{1} - \mathbf{D}(\mathbf{G}(\mathbf{z})))]$$

¹ Ian J. Goodfellow et al. "Generative Adversarial Networks". In: *CoRR* (2014). arXiv: 1406.2661. URL: <http://arxiv.org/abs/1406.2661>.

Generative Adversarial Network

- For G is fixed, the optimal discriminator D is

$$\mathbf{D}_G^*(\mathbf{x}) = \frac{\mathbf{p}_{\text{data}}(\mathbf{x})}{\mathbf{p}_g(\mathbf{x}) + \mathbf{p}_{\text{data}}(\mathbf{x})} \quad (2)$$

- Proof

$$V(\mathbf{D}, \mathbf{G}) = \int_{\mathbf{x}} \mathbf{p}_{\text{data}}(\mathbf{x}) \log(\mathbf{D}(\mathbf{x})) + \mathbf{p}_g(\mathbf{x}) \log(\mathbf{1} - \mathbf{D}(\mathbf{x})) d\mathbf{x} \quad (3)$$

- Maximum of $a \log(y) + b \log(1 - y)$ is at $\frac{a}{a+b}$ in $y \in [0, 1]$

- \mathbf{D} maximizes $P(y|\mathbf{x})$
- Y indicates whether \mathbf{x} from \mathbf{p}_g or \mathbf{p}_{data}

Generative Adversarial Network

- Cost function during the training of generator

$$\begin{aligned} \mathbf{C}(\mathbf{G}) &= \max_{\mathbf{D}} V(\mathbf{D}, \mathbf{G}) \\ &= \mathbf{E}_{\mathbf{x} \sim \mathbf{p}_{\text{data}}} \log\left(\frac{\mathbf{p}_{\text{data}}(\mathbf{x})}{\mathbf{p}_{\mathbf{g}}(\mathbf{x}) + \mathbf{p}_{\text{data}}(\mathbf{x})}\right) + \mathbf{E}_{\mathbf{x} \sim \mathbf{p}_{\mathbf{g}}} \log\left(\frac{\mathbf{p}_{\text{data}}(\mathbf{x})}{\mathbf{p}_{\mathbf{g}}(\mathbf{x}) + \mathbf{p}_{\text{data}}(\mathbf{x})}\right) \end{aligned} \quad (4)$$

- The global minimum of $\mathbf{C}(\mathbf{G})$ is achieved if and only if $\mathbf{p}_{\mathbf{g}} = \mathbf{p}_{\text{data}}$
- $\mathbf{C}(\mathbf{G}) = -\log 4$
 - Proof

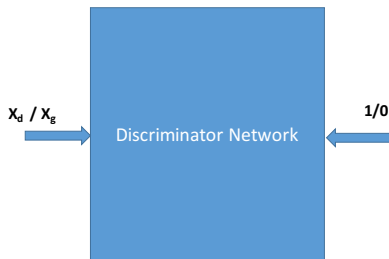
$$\begin{aligned} \mathbf{C}(\mathbf{G}) &= \max_{\mathbf{D}} V(\mathbf{D}, \mathbf{G}) \\ &= -\log 4 + KL(\mathbf{p}_{\text{data}} \parallel \frac{\mathbf{p}_{\mathbf{g}}(\mathbf{x}) + \mathbf{p}_{\text{data}}(\mathbf{x})}{2}) + KL(\mathbf{p}_{\mathbf{g}} \parallel \frac{\mathbf{p}_{\mathbf{g}}(\mathbf{x}) + \mathbf{p}_{\text{data}}(\mathbf{x})}{2}) \end{aligned} \quad (5)$$

- If \mathbf{G} and \mathbf{D} have enough capacity, and at each stage of the training, the discriminator is allowed to reach its optimum given \mathbf{G} , then $\mathbf{p}_{\mathbf{g}}$ converges to \mathbf{p}_{data}

GAN: Discriminator Training Scheme

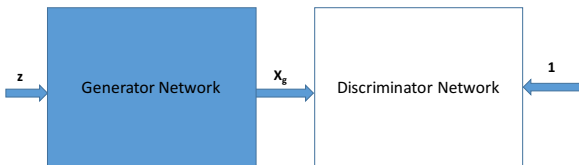
- **Training of Discriminator**

- $\{x_d^{(1)}, x_d^{(2)} \dots x_d^{(m)}\}$: samples from data distribution (labels 1)
- $\{x_g^{(1)}, x_g^{(2)} \dots x_g^{(m)}\}$: samples from generative networks (labels 0)



GAN: Generator Training Scheme

- **Training of Generator**
- Discriminator is frozen
- Generator Network is trained with the desired label at the discriminator output as 1



Generative Adversarial Network

- 1: **for** Number of training iterations **do**
- 2: **for** k steps **do**
- 3: Sample minibatch of m noise samples $\{z^{(1)}, z^{(2)} \dots z^{(m)}\}$ from noise prior $\mathbf{p}_z(\mathbf{z})$
- 4: Sample minibatch of m example $\{x^{(1)}, x^{(2)} \dots x^{(m)}\}$ from data generation distribution $\mathbf{p}_{\text{data}}(\mathbf{x})$
- 5: Update the discriminator by ascending its stochastic gradient

$$\nabla_{\theta_D} \left\{ \frac{1}{m} \sum_{i=1}^m [\log \mathbf{D}(x^i) + \log (1 - \mathbf{D}(\mathbf{G}(z^i)))] \right\} \quad (6)$$

- 6: **end for**
- 7: Sample minibatch of m noise samples $\{z^{(1)}, z^{(2)} \dots z^{(m)}\}$ from noise prior $\mathbf{p}_z(\mathbf{z})$
- 8: Update the generator by descending its stochastic gradient

$$\nabla_{\theta_G} \left\{ \frac{1}{m} \sum_{i=1}^m [\log (1 - \mathbf{D}(\mathbf{G}(z^i)))] \right\} \quad (7)$$

- 9: **end for**

Generation of Handwritten Digits using GAN

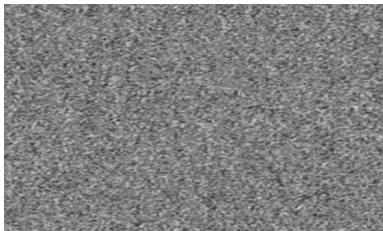


Figure: Iteration # 1

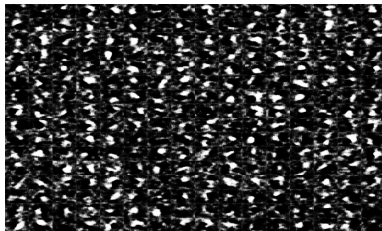


Figure: Iteration # 20



Figure: Iteration # 40



Figure: Iteration # 100

Generation of Sparse Signal Vectors using GAN

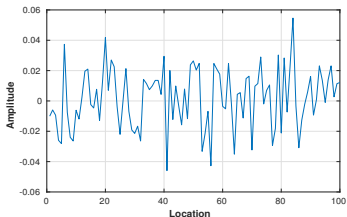


Figure: Iteration # 1

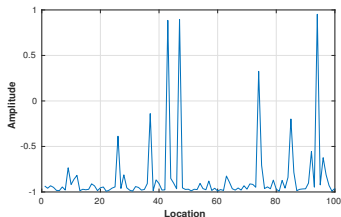


Figure: Iteration # 10

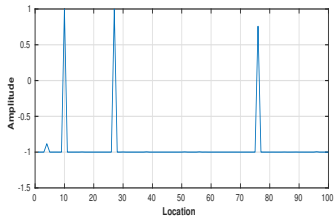


Figure: Iteration # 20

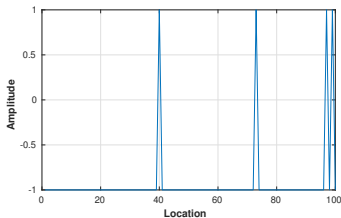


Figure: Iteration # 30

Compressive Sensing Using Generative Models

Compressive Sensing Using GAN

- System model

$$\begin{aligned} \mathbf{y} &= \mathbf{A}\mathbf{x} + \mathbf{n} \quad \mathbf{y} \in R^{m \times 1}, \mathbf{A} \in R^{m \times n}, \mathbf{x} \in R^{n \times 1} \\ \|\mathbf{x}\|_0 &= k \quad k \in \{1, 2, 3, \dots\} \end{aligned} \quad (8)$$

- The generative models learns a mapping from low dimensional representation space $\mathbf{z} \in R^k$ to the high dimensional sample space $\mathbf{G}(\mathbf{z}) \in R^n$
- Proposed algorithm: Find a mapping between observation vectors \mathbf{y} and the vectors in the latent space \mathbf{z}
- Mapping between measurement space and latent space is obtained by minimizing the following loss function²

$$V(\mathbf{z}) = \|\mathbf{A}\mathbf{G}(\mathbf{z}) - \mathbf{y}\|^2 \quad (9)$$

² Ashish Bora et al. "Compressed Sensing using Generative Models". In: *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*. 2017, pp. 537-546. URL: <http://proceedings.mlr.press/v70/bora17a.html>.

Sparse Signal Recovery using GAN

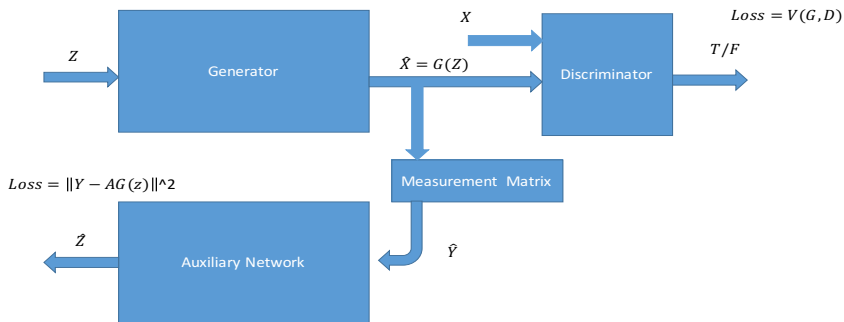
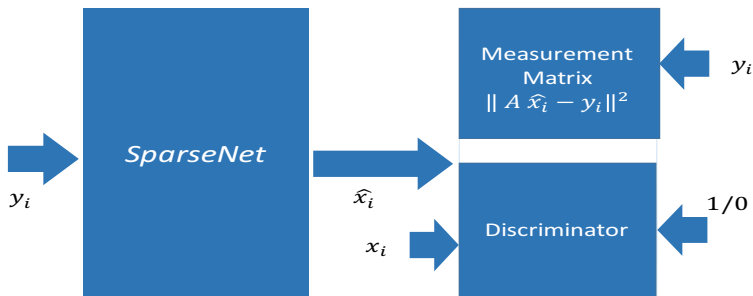


Figure: Compressive Sensing using GAN

New Framework for Sparse Signal Processing

New Framework for Sparse Signal Recovery



- *SparseNet* : DNN for sparse signal recovery
- Discriminator network can ensure sparsity
- May be useful to ensure more general features like block sparsity

Sparse Signal Recovery using New Framework

- Training of *Discriminator*

- Discriminator : Ensures sparsity of \mathbf{x}
- Trained using $\{\text{Data,Label}\} = \{\{\hat{\mathbf{x}}_i, 0\}, \{\mathbf{x}_k, 1\}, \dots\}$

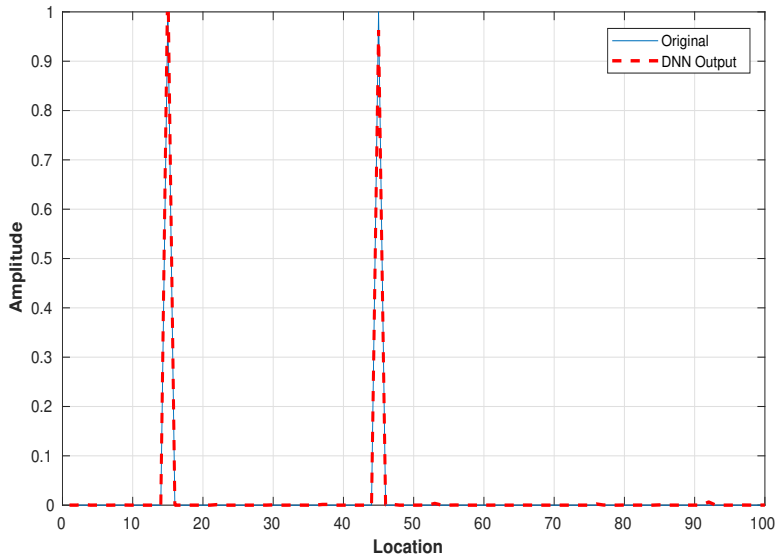
- Training of *SparseNet*

- Discriminator is frozen
- DNN is trained by simultaneously minimizing the loss function

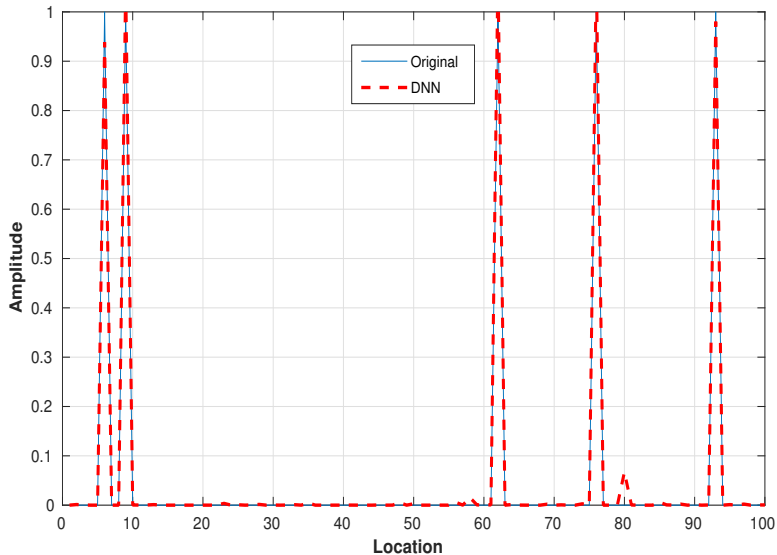
$$\min_{\mathbf{G}} \lambda_1 V(\mathbf{D}, \mathbf{G}) + \lambda_2 \mathbf{E}_{\mathbf{y} \sim \mathbf{p}_y} \|\mathbf{y} - \mathbf{A}\mathbf{G}(\mathbf{y})\|^2 \quad (10)$$
$$V(\mathbf{D}, \mathbf{G}) = \mathbf{E}_{\mathbf{y} \sim \mathbf{p}_y} [\log(1 - \mathbf{D}(\mathbf{G}(\mathbf{y})))]$$

- λ_1 & λ_2 : Loss weights can be specified during the training phase

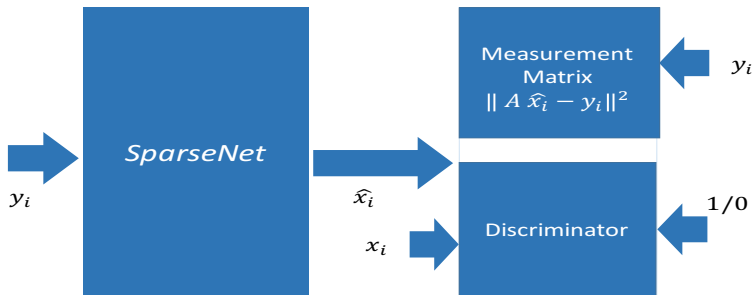
Sparse Signal Recovery



Sparse Signal Recovery

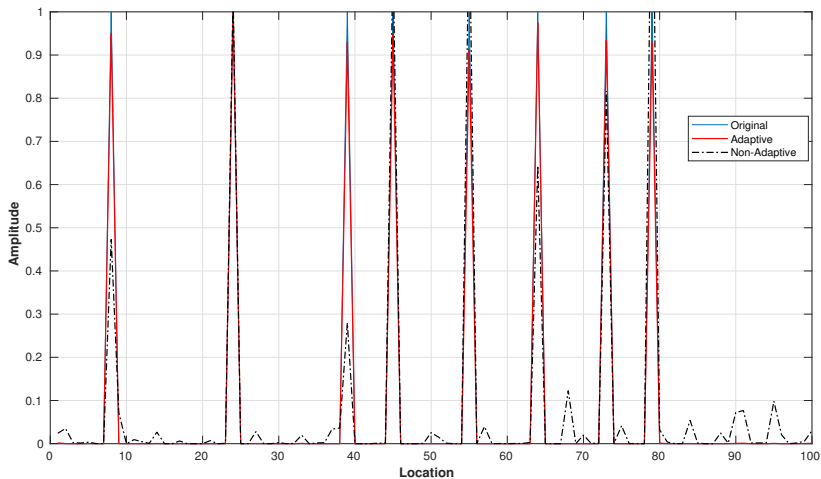


Adaptive Signal Recovery

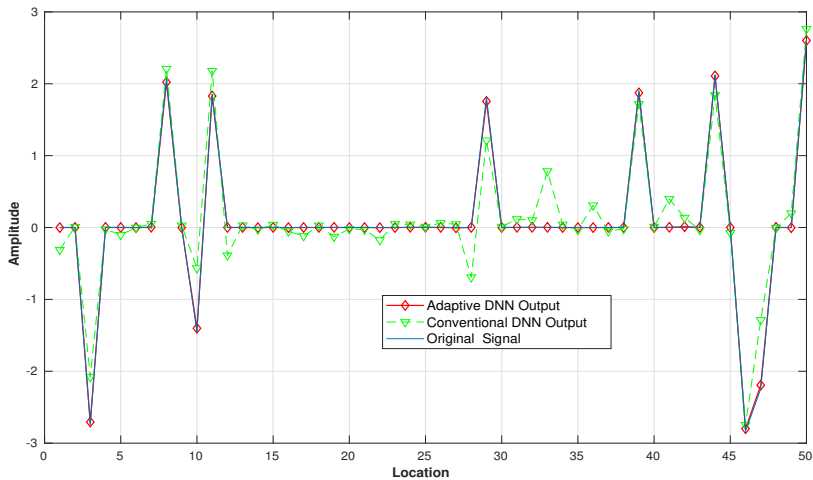


- Training of **G** does not require unknown sparse vector **x**
- Update the weights and biases of the DNN during signal recovery phase

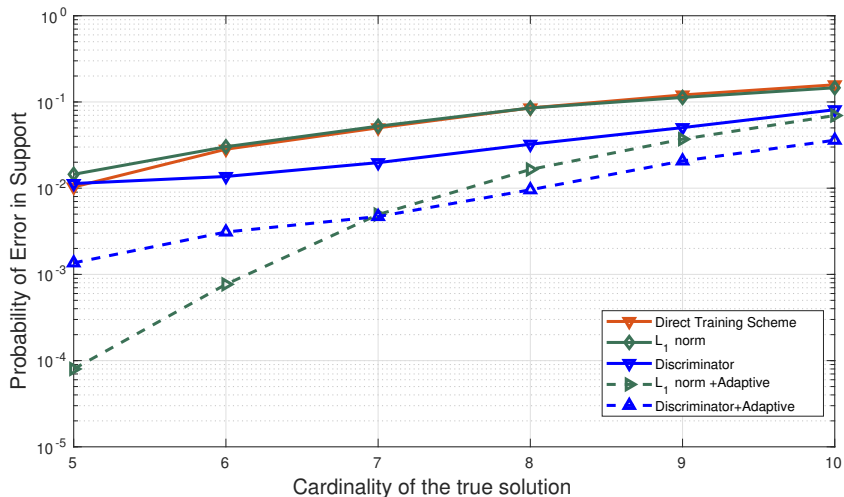
Comparison of Adaptive Vs Non-Adaptive



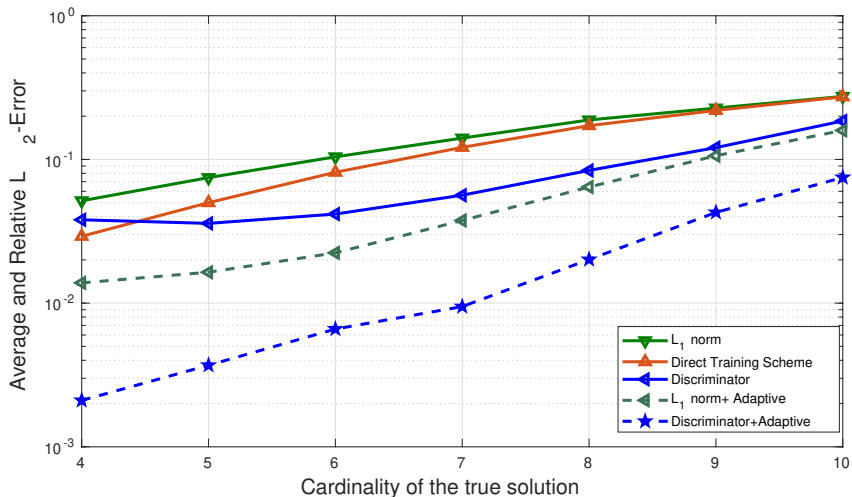
Comparison of Adaptive Vs Non-Adaptive



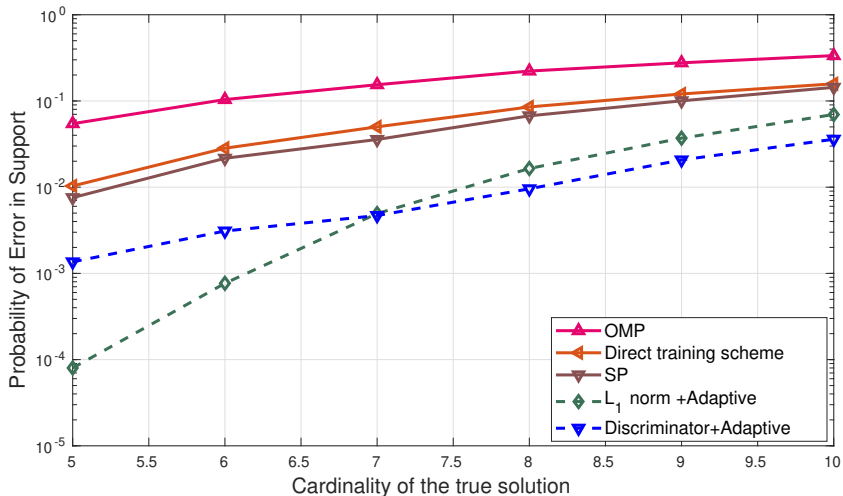
Comparison of Adaptive Vs Non-Adaptive



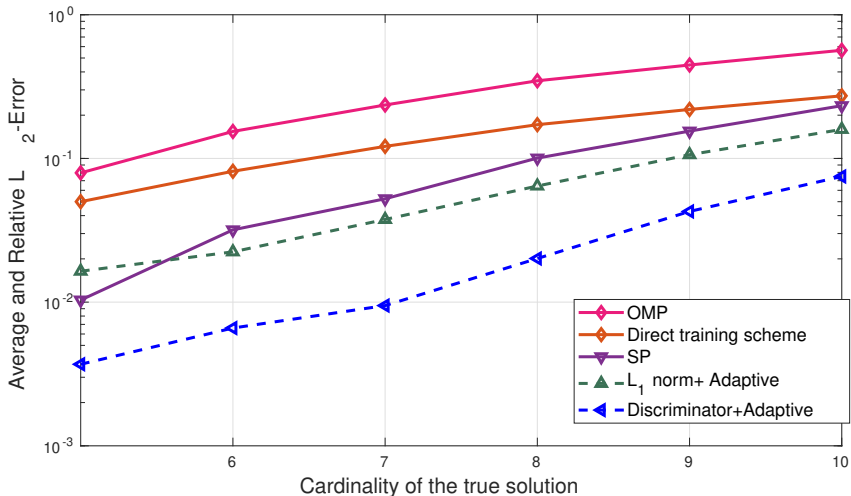
Comparison of Adaptive Vs Non-Adaptive



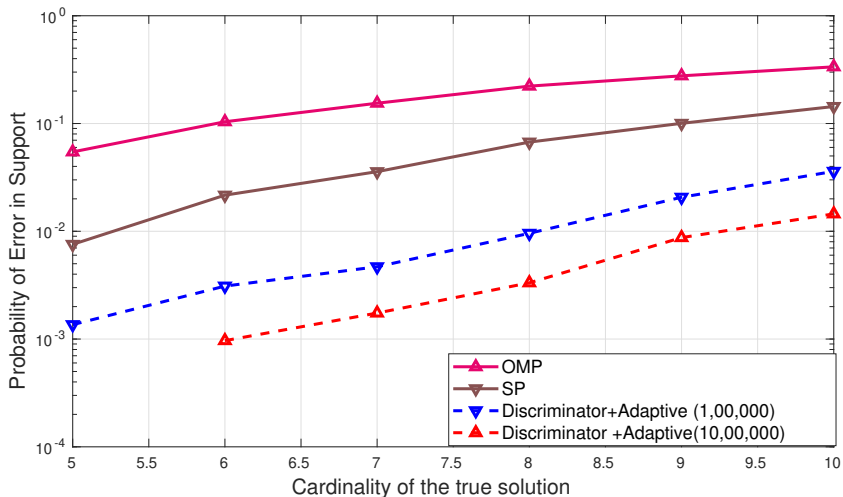
Comparison with Other Algorithms



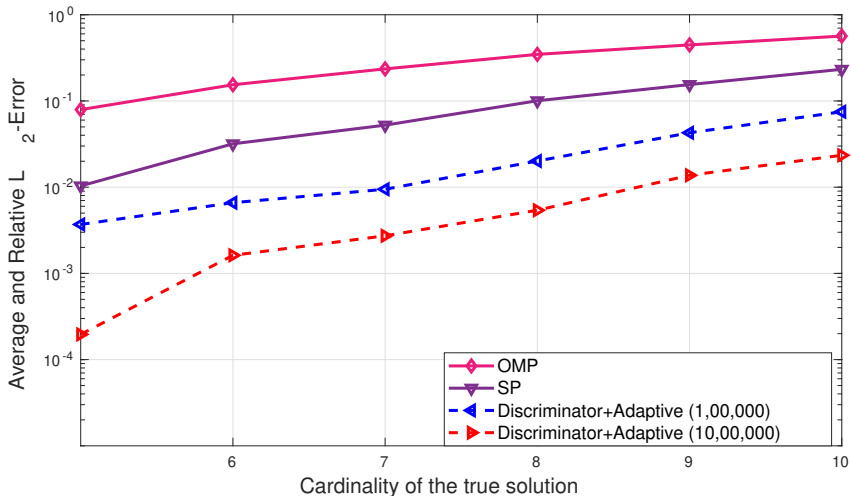
Comparison with Other Algorithms



Training Sets with Different Cardinality



Training Sets with Different Cardinality



Intuitive Explanation Under Bayesian Framework

- **Signal Model:**

$$\begin{aligned} \mathbf{y} &= \mathbf{A}\mathbf{x} + \mathbf{n} \quad \mathbf{y} \in \mathbb{R}^{m \times 1}, \mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{x} \in \mathbb{R}^{n \times 1}, \\ \|\mathbf{x}\|_0 &\leq K \quad \mathbf{n} \sim \mathbf{N}(\mathbf{0}, \frac{\mathbf{I}}{\lambda}) \end{aligned} \quad (11)$$

- **Likelihood** term is given by,

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}, \lambda) &= \left(\frac{\lambda}{2\pi}\right)^{\frac{m}{2}} \exp\left(-\frac{\lambda}{2}\|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2\right) \\ \log(p(\mathbf{y}|\mathbf{x}, \lambda)) &= -\frac{\lambda}{2}\|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + f(\lambda) \end{aligned} \quad (12)$$

- **Maximum Likelihood Estimation** of \mathbf{x} with sparsity constraint is

$$\begin{aligned} \hat{\mathbf{x}} &= \underset{\mathbf{x} \in \mathbf{S}}{\operatorname{argmax}} \quad p(\mathbf{y}|\mathbf{x}, \lambda) \\ \mathbf{S} &= \{\mathbf{x} : \|\mathbf{x}\|_0 \leq K\} \end{aligned} \quad (13)$$

Intuitive Explanation Under Bayesian Framework

- For G is fixed, the optimal discriminator D is

$$\mathbf{D}_G^*(\mathbf{x}) = \frac{\mathbf{p}_{\text{data}}(\mathbf{x})}{\mathbf{p}_g(\mathbf{x}) + \mathbf{p}_{\text{data}}(\mathbf{x})} \quad (14)$$

- The global minimum of $\mathbf{C}(\mathbf{G})$ is achieved if and only if $\mathbf{p}_g = \mathbf{p}_{\text{data}}$.

$$\begin{aligned} \mathbf{C}(\mathbf{G}) &= \max_D V(\mathbf{D}, \mathbf{G}) \\ &= -\log 4 + KL(\mathbf{p}_{\text{data}} \parallel \frac{\mathbf{p}_g(\mathbf{x}) + \mathbf{p}_{\text{data}}(\mathbf{x})}{2}) + KL(\mathbf{p}_g \parallel \frac{\mathbf{p}_g(\mathbf{x}) + \mathbf{p}_{\text{data}}(\mathbf{x})}{2}) \end{aligned} \quad (15)$$

- Discriminator ensures $\mathbf{p}_g = \mathbf{p}_{\text{data}} \implies \hat{x} = G(y) \in S$

Intuitive Explanation Under Bayesian Framework

- Optimization problem during the training phase becomes

$$\min_{\mathbf{G}: \mathbf{G}(\mathbf{y}) \in S} (-\log 4)\lambda_1 + \lambda_2 \frac{1}{m} \sum_{i=1}^m \|\mathbf{y}_i - \mathbf{A}\mathbf{G}(\mathbf{y}_i)\|^2 \quad (16)$$

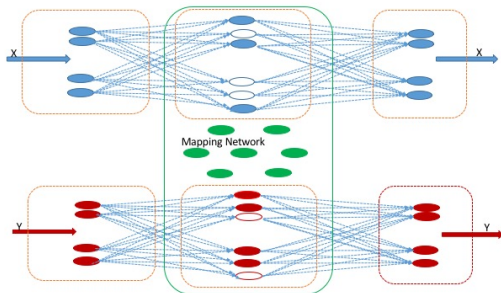
m : Number of samples in a minibatch

- Above cost function is proportional to the log likelihood of $\{\mathbf{y}_1, \mathbf{y}_2 \dots \mathbf{y}_m\}$
- $\mathbf{P}_{\mathbf{X}}(x)$: uniform prior over S
- The new framework tries to give a MAP estimate of \mathbf{x} with prior distribution $\mathbf{P}_{\mathbf{X}}(x)$ or ML estimate on the set of k sparse vectors

Coupled Dictionary Learning

Coupled Dictionary Learning using DNN

- Dictionary Learning of \mathbf{x} and \mathbf{y}
 - \mathbf{z}_x : Sparse representation of \mathbf{x}
 - \mathbf{z}_y : Sparse representation of \mathbf{y}
- Coupled Dictionary Learning
 - Train a mapping network between \mathbf{z}_x and \mathbf{z}_y
 - \mathbf{z}_x : Sparse representation of \mathbf{x} and \mathbf{y} with respect to coupled dictionary
 - Dictionary for \mathbf{x} :Decoding network of \mathbf{x}
 - Dictionary for \mathbf{y} :Decoding network of \mathbf{y} and mapping network



Dictionary Learning using K-SVD

$$\% \text{ of Recovered Atoms} = \frac{R}{n} 100, \quad R = \sum_{i=1}^n \mathbb{1}_{(x < .01)}, \quad x = 1 - \max(\hat{\mathbf{d}}_i^T \mathbf{d}_j)$$

$$\text{Average representation error, } E = \frac{1}{N} \sum_{i=1}^N \frac{\|\mathbf{Y}_i - \hat{\mathbf{D}}\hat{\mathbf{Z}}_i\|^2}{m} \quad (17)$$

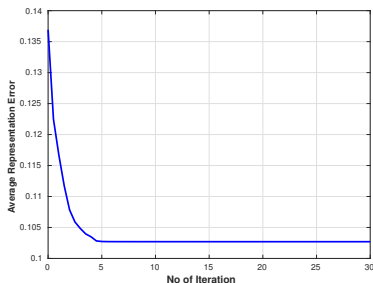
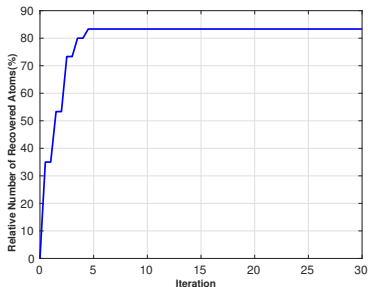


Figure: K-SVD

Dictionary Learning using DNN

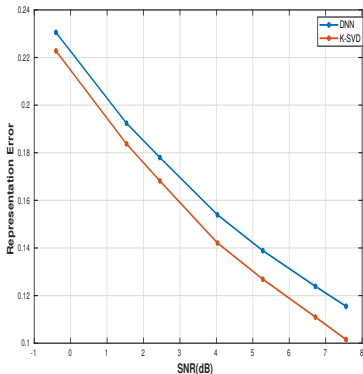
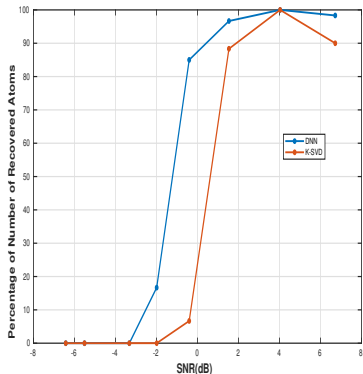
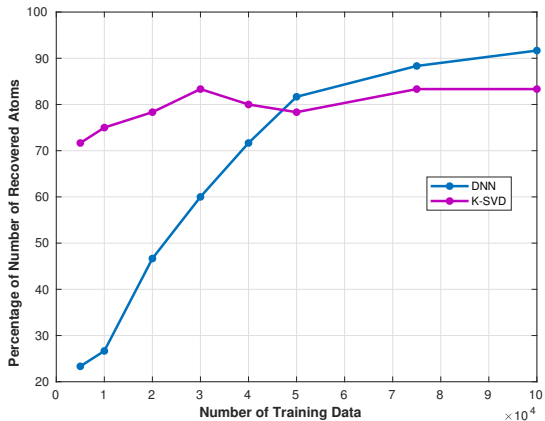
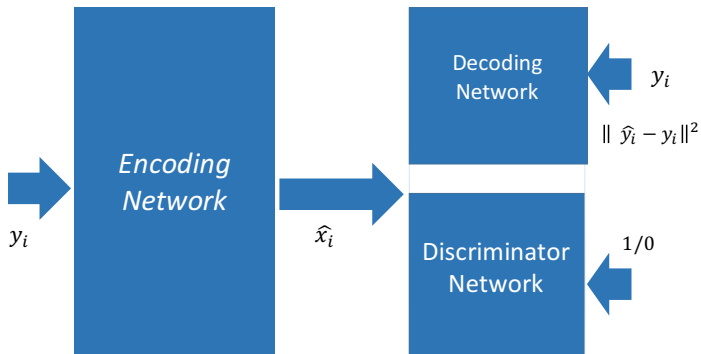


Figure: K-SVD / DNN

Dictionary Learning using DNN



Dictionary Learning under new Framework



- Decoding network is a single layer MLP with linear activation functions

Conclusion

- The new framework allows to update the inverse function during testing phase (Adaptive Signal Recovery)
- The proposed discriminator based scheme can be extended for arbitrary prior distribution
- More general features like block sparsity may be ensured using adversarial training
- The new framework may be useful for other sparse signal processing applications like dictionary learning, coupled dictionary learning etc.