



# Sparse representation, dictionary Learning, and deep neural networks: Their connections and new algorithms

**Mostafa Sadeghi**

Electrical Engineering Department  
Sharif University of Technology  
Tehran, Iran.

June 2018

## 1 Proximal algorithms

- Proximal operator
- Forward-backward Splitting

## 2 Sparse representation

- Background
- Iterative Sparsification-Projection
- Iterative Proximal Projection

## 3 Dictionary Learning

- Background
- learning low coherence dictionaries

## 4 Deep Neural Networks

- Background
- Progressive Neural Networks
- Structured Weight Matrices for Neural Networks

## 5 Conclusions

## 1 Proximal algorithms

- Proximal operator
- Forward-backward Splitting

## 2 Sparse representation

## 3 Dictionary Learning

## 4 Deep Neural Networks

## 5 Conclusions

# Proximal Algorithms

- Efficient **first order** algorithms. Suitable for **nonsmooth, constrained, large-scale** problems.

# Proximal Algorithms

- Efficient **first order** algorithms. Suitable for **nonsmooth, constrained, large-scale** problems.

Proximal mapping [Parikh and Boyd, 2014]

$g : \text{dom}_g \rightarrow \mathbb{R} \cup \{+\infty\}$ : proper, lower-semicontinuous

$$\text{prox}_g(\mathbf{u}) = \underset{\mathbf{x} \in \text{dom}_g}{\operatorname{argmin}} g(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{u}\|_2^2$$

# Proximal Algorithms

- Efficient **first order** algorithms. Suitable for **nonsmooth, constrained, large-scale** problems.

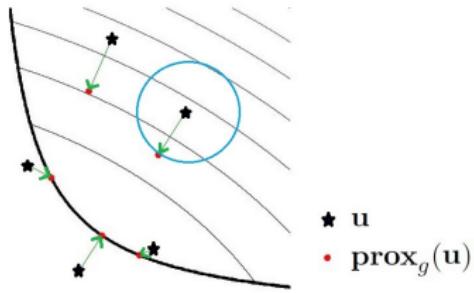
Proximal mapping [Parikh and Boyd, 2014]

$g : \text{dom}_g \rightarrow \mathbb{R} \cup \{+\infty\}$ : proper, lower-semicontinuous

$$\text{prox}_g(\mathbf{u}) = \underset{\mathbf{x} \in \text{dom}_g}{\operatorname{argmin}} g(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{u}\|_2^2$$

$$\underset{\mathbf{x} \in \text{dom}_g}{\operatorname{argmin}} g(\mathbf{x}) \quad \text{s.t.} \quad \|\mathbf{x} - \mathbf{u}\|_2 \leq \tau$$

$$\text{prox}_{\lambda g}(\mathbf{u}) \simeq \mathbf{u} - \lambda \nabla g(\mathbf{u})$$



## 1 Proximal algorithms

- Proximal operator
- Forward-backward Splitting

## 2 Sparse representation

## 3 Dictionary Learning

## 4 Deep Neural Networks

## 5 Conclusions

# Forward-backward Splitting

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + g(\mathbf{x})$$

- $f : \text{dom}_f \rightarrow \mathbb{R}$  smooth (convex/non-convex)
- $g : \text{dom}_g \rightarrow \mathbb{R} \cup \{+\infty\}$  non-smooth (convex/non-convex)

# Forward-backward Splitting

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + g(\mathbf{x})$$

- $f : \text{dom}_f \rightarrow \mathbb{R}$  smooth (convex/non-convex)
- $g : \text{dom}_g \rightarrow \mathbb{R} \cup \{+\infty\}$  non-smooth (convex/non-convex)

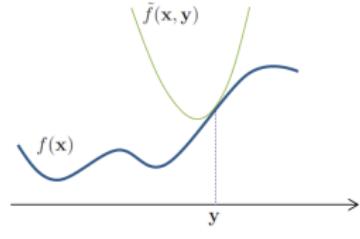
## Descent lemma

$f : \text{dom}_f \rightarrow \mathbb{R}$ , smooth and  $L$ -gradient Lipschitz\*,  $\mu \in (0, 1/L]$

$$\forall \mathbf{x}, \mathbf{y} \in \text{dom} f : f(\mathbf{x}) \leq \tilde{f}(\mathbf{x}, \mathbf{y}) \triangleq f(\mathbf{y}) + \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y}) + \frac{1}{2\mu} \|\mathbf{x} - \mathbf{y}\|_2^2$$

\* $\forall x, y \in \text{dom} f :$

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2$$



# Forward-backward Splitting

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + g(\mathbf{x})$$

$$\boxed{\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \tilde{f}(\mathbf{x}, \mathbf{x}_k) + g(\mathbf{x})}$$

# Forward-backward Splitting

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + g(\mathbf{x})$$

$$\boxed{\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \tilde{f}(\mathbf{x}, \mathbf{x}_k) + g(\mathbf{x})}$$

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - (\mathbf{x}_k - \mu \nabla f(\mathbf{x}_k))\|_2^2 + \mu \cdot g(\mathbf{x})$$

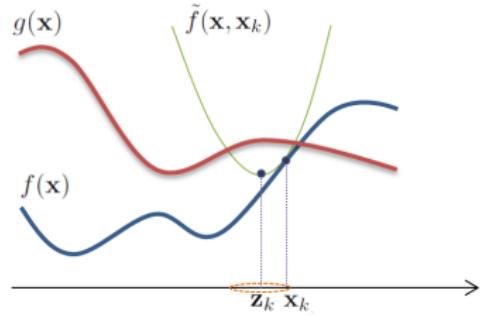
# Forward-backward Splitting

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + g(\mathbf{x})$$

$$\boxed{\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \tilde{f}(\mathbf{x}, \mathbf{x}_k) + g(\mathbf{x})}$$

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - (\mathbf{x}_k - \mu \nabla f(\mathbf{x}_k))\|_2^2 + \mu \cdot g(\mathbf{x})$$

- ➊ Forward step:  $\mathbf{z}_k = \mathbf{x}_k - \mu \nabla f(\mathbf{x}_k)$
- ➋ Backward step:  $\mathbf{x}_{k+1} = \operatorname{prox}_{\mu \cdot g}(\mathbf{z}_k)$



1 Proximal algorithms

2 Sparse representation

- Background
- Iterative Sparsification-Projection
- Iterative Proximal Projection

3 Dictionary Learning

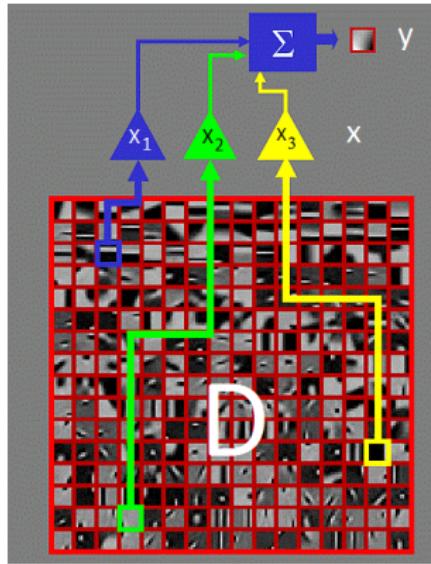
4 Deep Neural Networks

5 Conclusions

# Background

## Sparse representation

$$\mathbf{y} \approx x_1 \mathbf{d}_1 + x_2 \mathbf{d}_2 + \dots + x_n \mathbf{d}_m = \mathbf{Dx} \quad \text{most } x_i \text{'s are zero}$$

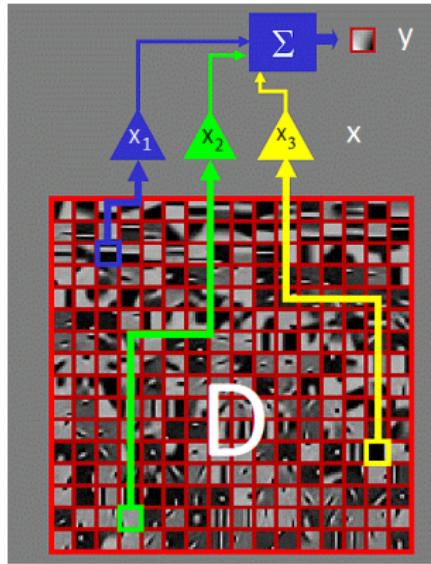


\* Adopted from M. Elad's slides

# Background

## Sparse representation

$$\mathbf{y} \approx x_1 \mathbf{d}_1 + x_2 \mathbf{d}_2 + \dots + x_n \mathbf{d}_m = \mathbf{Dx} \quad \text{most } x_i \text{'s are zero}$$



\* Adopted from M. Elad's slides

- Signal restoration:

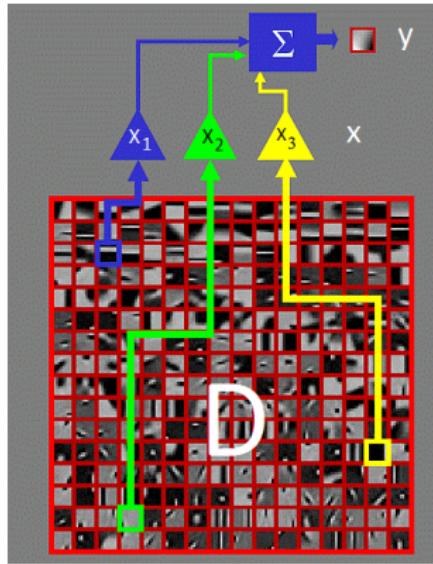
$$\mathbf{y} = \mathbf{Hx} + \mathbf{e}$$

De-noising ( $\mathbf{H}$  = identity), inpainting  
( $\mathbf{H}$  = random rows of identity), de-blurring  
( $\mathbf{H}$  = blurring matrix), super resolution  
( $\mathbf{H}$  = down sampling matrix), ...

# Background

## Sparse representation

$$\mathbf{y} \approx x_1 \mathbf{d}_1 + x_2 \mathbf{d}_2 + \dots + x_n \mathbf{d}_m = \mathbf{Dx} \quad \text{most } x_i \text{'s are zero}$$



\* Adopted from M. Elad's slides

- Signal restoration:

$$\mathbf{y} = \mathbf{Hx} + \mathbf{e}$$

De-noising ( $\mathbf{H}$  = identity), inpainting ( $\mathbf{H}$  = random rows of identity), de-blurring ( $\mathbf{H}$  = blurring matrix), super resolution ( $\mathbf{H}$  = down sampling matrix), ...

$$\mathbf{x} \simeq \mathbf{Da}, \mathbf{a} : \text{sparse}$$

$$\min_{\mathbf{x}, \mathbf{a}} \|\mathbf{y} - \mathbf{Hx}\|_2^2 + \alpha \|\mathbf{x} - \mathbf{Da}\|_2^2 + \beta \|\mathbf{a}\|_1$$

# Sparse recovery algorithms

- Thresholding based algorithms

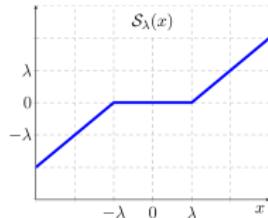
$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2 + \lambda \|\mathbf{x}\|_1$$

# Sparse recovery algorithms

- Thresholding based algorithms

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2 + \lambda \|\mathbf{x}\|_1$$

$$\mathbf{x}^{k+1} = \mathcal{S}_{\mu_k \lambda}(\mathbf{x}^k - \mu_k (\mathbf{D}\mathbf{x}_k - \mathbf{y}))$$



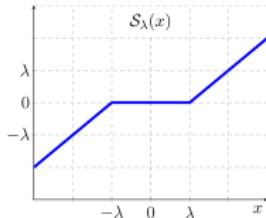
# Sparse recovery algorithms

- Thresholding based algorithms

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2 + \lambda \|\mathbf{x}\|_1$$

$$\mathbf{x}^{k+1} = \mathcal{S}_{\mu_k \lambda}(\mathbf{x}^k - \mu_k (\mathbf{D}\mathbf{x}_k - \mathbf{y}))$$

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2 + \lambda \|\mathbf{x}\|_0$$

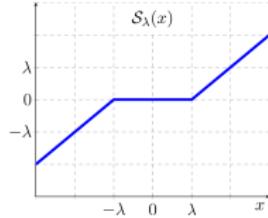


# Sparse recovery algorithms

- Thresholding based algorithms

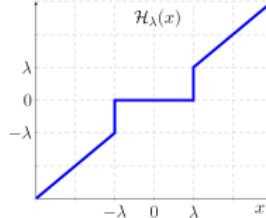
$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2 + \lambda \|\mathbf{x}\|_1$$

$$\mathbf{x}^{k+1} = \mathcal{S}_{\mu_k \lambda}(\mathbf{x}^k - \mu_k (\mathbf{D}\mathbf{x}_k - \mathbf{y}))$$



$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2 + \lambda \|\mathbf{x}\|_0$$

$$\mathbf{x}^{k+1} = \mathcal{H}_{\mu_k \lambda}(\mathbf{x}^k - \mu_k (\mathbf{D}\mathbf{x}_k - \mathbf{y}))$$



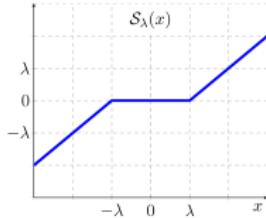
$$\mu_k \in (0, 1/\sigma_{\max}(\mathbf{D}))$$

# Sparse recovery algorithms

- Thresholding based algorithms

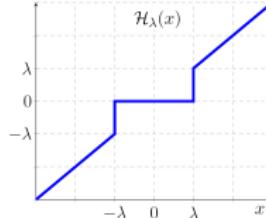
$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2 + \lambda \|\mathbf{x}\|_1$$

$$\mathbf{x}^{k+1} = \mathcal{S}_{\mu_k \lambda}(\mathbf{x}^k - \mu_k (\mathbf{D}\mathbf{x}_k - \mathbf{y}))$$



$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2 + \lambda \|\mathbf{x}\|_0$$

$$\mathbf{x}^{k+1} = \mathcal{H}_{\mu_k \lambda}(\mathbf{x}^k - \mu_k (\mathbf{D}\mathbf{x}_k - \mathbf{y}))$$



$$\mu_k \in (0, 1/\sigma_{\max}(\mathbf{D}))$$

Examples: **IST** [Daubechies et al., 2004], **GPSR** [Figueiredo et al., 2007], **IHT** [Blumensath and Davies, 2009], **AMP** [Donoho et al., 2009], **EMGMAMP** [Vila et al., 2013], **NESTA** [Becker et al., 2009], **SCAD** [Gasso et al., 2009].

# Sparse recovery algorithms

- $\ell_0$  norm approximation. Approximate  $\ell_0$  norm with a smooth function.

**Smoothed L0 (SL0)** [Mohimani et al., 2009], **SCSA** [Malek-Mohammadi et al., 2016]

# Sparse recovery algorithms

- $\ell_0$  norm approximation. Approximate  $\ell_0$  norm with a smooth function.

**Smoothed L0 (SL0)** [Mohimani et al., 2009], **SCSA** [Malek-Mohammadi et al., 2016]

$$\|\mathbf{x}\|_{\sigma} = n - \sum_{i=1}^n f_{\sigma}(x_i)$$

$$f_{\sigma}(x) = \exp\left(-\frac{x^2}{\sigma^2}\right)$$

# Sparse recovery algorithms

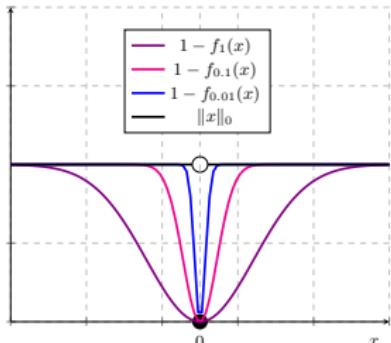
- $\ell_0$  norm approximation. Approximate  $\ell_0$  norm with a smooth function.

**Smoothed L0 (SL0)** [Mohimani et al., 2009], **SCSA** [Malek-Mohammadi et al., 2016]

$$\|\mathbf{x}\|_{\sigma} = n - \sum_{i=1}^n f_{\sigma}(x_i)$$

$$f_{\sigma}(x) = \exp\left(-\frac{x^2}{\sigma^2}\right)$$

When  $\sigma \rightarrow 0$  :  $\|\mathbf{x}\|_{\sigma} \rightarrow \|\mathbf{x}\|_0$



# Sparse recovery algorithms

- $\ell_0$  norm approximation. Approximate  $\ell_0$  norm with a smooth function.

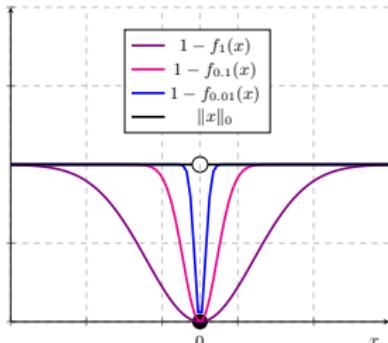
**Smoothed L0 (SL0)** [Mohimani et al., 2009], **SCSA** [Malek-Mohammadi et al., 2016]

$$\|\mathbf{x}\|_{\sigma} = n - \sum_{i=1}^n f_{\sigma}(x_i)$$

$$\min_{\mathbf{x}} \|\mathbf{x}\|_{\sigma} \quad \text{s.t.} \quad \mathbf{y} = \mathbf{D}\mathbf{x}$$

$$f_{\sigma}(x) = \exp\left(-\frac{x^2}{\sigma^2}\right)$$

When  $\sigma \rightarrow 0$  :  $\|\mathbf{x}\|_{\sigma} \rightarrow \|\mathbf{x}\|_0$



# Sparse recovery algorithms

- $\ell_0$  norm approximation. Approximate  $\ell_0$  norm with a smooth function.

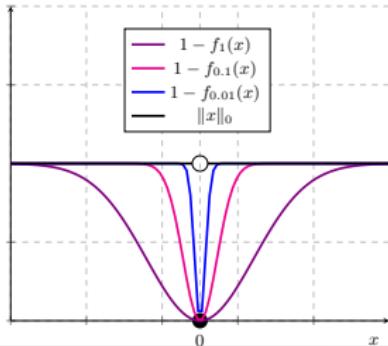
**Smoothed L0 (SL0)** [Mohimani et al., 2009], **SCSA** [Malek-Mohammadi et al., 2016]

$$\|\mathbf{x}\|_{\sigma} = n - \sum_{i=1}^n f_{\sigma}(x_i)$$

$$\min_{\mathbf{x}} \|\mathbf{x}\|_{\sigma} \quad \text{s.t.} \quad \mathbf{y} = \mathbf{D}\mathbf{x}$$

$$f_{\sigma}(x) = \exp\left(-\frac{x^2}{\sigma^2}\right)$$

When  $\sigma \rightarrow 0$  :  $\|\mathbf{x}\|_{\sigma} \rightarrow \|\mathbf{x}\|_0$



```

$$\begin{cases} k = 0 \\ \mathbf{x}^0 = \mathbf{D}^\dagger \mathbf{y} \\ \text{For } i = 1, 2, \dots \\ \quad \text{For } j = 1, 2, \dots \\ \quad \quad \mathbf{x}^{k+1} = \mathbf{x}^k - \mu_{\sigma_i} \nabla \|\mathbf{x}_k\|_{\sigma_i} \\ \quad \quad \mathbf{x}^{k+1} = \mathbf{x}^{k+1} - \mathbf{D}^\dagger (\mathbf{D}\mathbf{x}_{k+1} - \mathbf{y}) \\ \quad \quad k \leftarrow k + 1 \\ \quad \text{End} \\ \sigma_{i+1} = \sigma_i \cdot c \quad (0 < c < 1) \\ \text{End} \end{cases}$$

```

1 Proximal algorithms

2 Sparse representation

- Background
- Iterative Sparsification-Projection
- Iterative Proximal Projection

3 Dictionary Learning

4 Deep Neural Networks

5 Conclusions

# Iterative Sparsification-Projection

## <sup>1</sup> Iterative Sparsification-Projection (ISP)

---

1

M. Sadeghi, M. Babaie-Zadeh, "Iterative Sparsification-Projection: Fast and Robust Sparse Signal Approximation", *IEEE Trans. Sig. Proc.*, vol. 64, no. 21, pp. 5536-5548, November 2016.

# Iterative Sparsification-Projection

## 1 Iterative Sparsification-Projection (ISP)

- Revisiting the SL0 algorithm:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \underbrace{\sum_{i=1}^n \left(1 - \exp\left(-\frac{x_i^2}{\sigma^2}\right)\right)}_{f_\sigma(\mathbf{x})} \quad \text{s.t. } \|\mathbf{y} - \mathbf{Dx}\|_2 \leq \epsilon$$

---

1

M. Sadeghi, M. Babaie-Zadeh, "Iterative Sparsification-Projection: Fast and Robust Sparse Signal Approximation", *IEEE Trans. Sig. Proc.*, vol. 64, no. 21, pp. 5536-5548, November 2016.

# Iterative Sparsification-Projection

## 1 Iterative Sparsification-Projection (ISP)

- Revisiting the SL0 algorithm:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \underbrace{\sum_{i=1}^n \left(1 - \exp\left(-\frac{x_i^2}{\sigma^2}\right)\right)}_{f_\sigma(\mathbf{x})} \quad \text{s.t. } \|\mathbf{y} - \mathbf{Dx}\|_2 \leq \epsilon$$

### Lemma (Lipschitz constant)

The function  $f_\sigma$  (defined above) is gradient Lipschitz with constant  $L = \frac{2}{\sigma^2}$ .

---

1

M. Sadeghi, M. Babaie-Zadeh, "Iterative Sparsification-Projection: Fast and Robust Sparse Signal Approximation", *IEEE Trans. Sig. Proc.*, vol. 64, no. 21, pp. 5536-5548, November 2016.

# Iterative Sparsification-Projection

## 1 Iterative Sparsification-Projection (ISP)

- Revisiting the SL0 algorithm:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \underbrace{\sum_{i=1}^n \left(1 - \exp\left(-\frac{x_i^2}{\sigma^2}\right)\right)}_{f_\sigma(\mathbf{x})} \quad \text{s.t. } \|\mathbf{y} - \mathbf{Dx}\|_2 \leq \epsilon$$

### Lemma (Lipschitz constant)

The function  $f_\sigma$  (defined above) is gradient Lipschitz with constant  $L = \frac{2}{\sigma^2}$ .

$$\min_{\mathbf{x} \in \mathbb{R}^n} \underbrace{\sum_{i=1}^n \left(1 - \exp\left(-\frac{x_i^2}{\sigma^2}\right)\right)}_{f(\mathbf{x})} + \underbrace{\delta_C(\mathbf{x})}_{g(\mathbf{x})} \rightarrow \boxed{\mathbf{x}_{k+1} = \mathbf{prox}_{\delta_C}(\mathbf{x}_k - \mu_\sigma \nabla f(\mathbf{x}_k))}$$

$\delta_C(\mathbf{x}) = 0$  if  $\|\mathbf{y} - \mathbf{Dx}\|_2 \leq \epsilon$  and  $\infty$  otherwise.

<sup>1</sup> M. Sadeghi, M. Babaie-Zadeh, "Iterative Sparsification-Projection: Fast and Robust Sparse Signal Approximation", *IEEE Trans. Sig. Proc.*, vol. 64, no. 21, pp. 5536-5548, November 2016.

# Iterative Sparsification-Projection

## Theorem (SL0 convergence)

Let  $\{\mathbf{x}_k\}$  be the sequence generated by the SL0 algorithm for a fixed  $\sigma$ .

Then:

- This sequence is bounded and convergent, which means that any accumulation point  $\mathbf{x}^*$  of  $\{\mathbf{x}_k\}$  is a critical point, and
- the sequence of objective values, i.e.,  $\{f_\sigma(\mathbf{x}_k) + \delta_{\mathcal{C}}(\mathbf{x}_k)\}_{k \geq 0}$ , is non-increasing.

# Iterative Sparsification-Projection

## Theorem (SL0 convergence)

Let  $\{\mathbf{x}_k\}$  be the sequence generated by the SL0 algorithm for a fixed  $\sigma$ .

Then:

- This sequence is bounded and convergent, which means that any accumulation point  $\mathbf{x}^*$  of  $\{\mathbf{x}_k\}$  is a critical point, and
- the sequence of objective values, i.e.,  $\{f_\sigma(\mathbf{x}_k) + \delta_{\mathcal{C}}(\mathbf{x}_k)\}_{k \geq 0}$ , is non-increasing.

## Corollary

If  $\mu_\sigma \in (0, 2/L]$ , or equivalently, if  $\mu \in (0, 1]$ , then the SL0 algorithm converges.

# Iterative Sparsification-Projection

- ISP

Motivation:

$$\begin{cases} \mathbf{z}_k = \mathbf{x}_k - \mu_\sigma \nabla f(\mathbf{x}_k) \triangleq \mathcal{T}_\sigma^0(\mathbf{x}_k) \\ \mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{C}}(\mathbf{z}_k) \end{cases}$$

# Iterative Sparsification-Projection

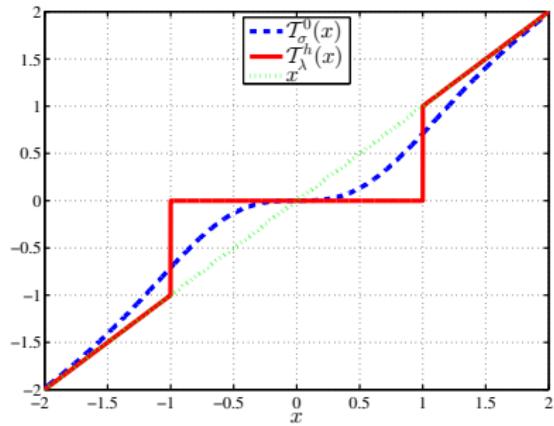
- ISP

Motivation:

$$\begin{cases} \mathbf{z}_k = \mathbf{x}_k - \mu_\sigma \nabla f(\mathbf{x}_k) \triangleq \mathcal{T}_\sigma^0(\mathbf{x}_k) \\ \mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{C}}(\mathbf{z}_k) \end{cases}$$

$$\mathcal{T}_\sigma^0(x) = x \cdot (1 - \exp(-\frac{x^2}{\sigma^2}))$$

$$(\mu_\sigma = \frac{\sigma^2}{2})$$



# Iterative Sparsification-Projection

☞ Gradient descent interpretation of proximal mapping!

$\|\mathbf{x}\|_\sigma$  smoothed of  $\|\mathbf{x}\|_0 \Rightarrow \mathbf{x} - \mu_\sigma \nabla \|\mathbf{x}\|_\sigma$  behaves like  $\text{prox}_{\mu_\sigma}(\mathbf{x}) = \text{hard-thr}$

# Iterative Sparsification-Projection

☞ Gradient descent interpretation of proximal mapping!

$\|\mathbf{x}\|_\sigma$  smoothed of  $\|\mathbf{x}\|_0 \Rightarrow \mathbf{x} - \mu_\sigma \nabla \|\mathbf{x}\|_\sigma$  behaves like  $\text{prox}_{\mu_\sigma}(\mathbf{x}) = \text{hard-thr}$

---

**Algorithm 1** SL0

```
1: Require:  $\mathbf{y}, \mathbf{D}, \sigma_0, \sigma_{\min}, c > 0, \mu, I$ 
2: Initialization:  $\mathbf{x} = \mathbf{D}^\dagger \mathbf{y}, \sigma = \sigma_0$ 
3: while  $\sigma > \sigma_{\min}$  do
4:   for  $i = 1, 2, \dots, I$  do
5:      $\tilde{\mathbf{x}} = \mathbf{x} - \mu \cdot \sigma^2 \nabla \|\mathbf{x}\|_\sigma$ 
6:      $\mathbf{x} = \tilde{\mathbf{x}} - \mathbf{D}^\dagger (\mathbf{D}\tilde{\mathbf{x}} - \mathbf{y})$ 
7:   end for
8:    $\sigma \leftarrow c \cdot \sigma$ 
9: end while
10: Output:  $\mathbf{x}$ 
```

---

**Algorithm 2** ISP

```
1: Require:  $\mathbf{y}, \mathbf{D}, \tau_0, \tau_{\min}, c > 0, I$ 
2: Initialization:  $\mathbf{x} = \mathbf{D}^\dagger \mathbf{y}, \tau = \tau_0$ 
3: while  $\tau > \tau_{\min}$  do
4:   for  $i = 1, 2, \dots, I$  do
5:      $\tilde{\mathbf{x}} = \mathcal{T}_\tau(\mathbf{x})$ 
6:      $\mathbf{x} = \tilde{\mathbf{x}} - \mathbf{D}^\dagger (\mathbf{D}\tilde{\mathbf{x}} - \mathbf{y})$ 
7:   end for
8:    $\tau \leftarrow c \cdot \tau$ 
9: end while
10: Output:  $\mathbf{x}$ 
```

---

# Iterative Sparsification-Projection

👉 Gradient descent interpretation of proximal mapping!

$\|\mathbf{x}\|_\sigma$  smoothed of  $\|\mathbf{x}\|_0 \Rightarrow \mathbf{x} - \mu_\sigma \nabla \|\mathbf{x}\|_\sigma$  behaves like  $\text{prox}_{\mu_\sigma}(\mathbf{x}) = \text{hard-thr}$

---

**Algorithm 1** SL0

```
1: Require:  $\mathbf{y}, \mathbf{D}, \sigma_0, \sigma_{\min}, c > 0, \mu, I$ 
2: Initialization:  $\mathbf{x} = \mathbf{D}^\dagger \mathbf{y}, \sigma = \sigma_0$ 
3: while  $\sigma > \sigma_{\min}$  do
4:   for  $i = 1, 2, \dots, I$  do
5:      $\tilde{\mathbf{x}} = \mathbf{x} - \mu \cdot \sigma^2 \nabla \|\mathbf{x}\|_\sigma$ 
6:      $\mathbf{x} = \tilde{\mathbf{x}} - \mathbf{D}^\dagger (\mathbf{D}\tilde{\mathbf{x}} - \mathbf{y})$ 
7:   end for
8:    $\sigma \leftarrow c \cdot \sigma$ 
9: end while
10: Output:  $\mathbf{x}$ 
```

---

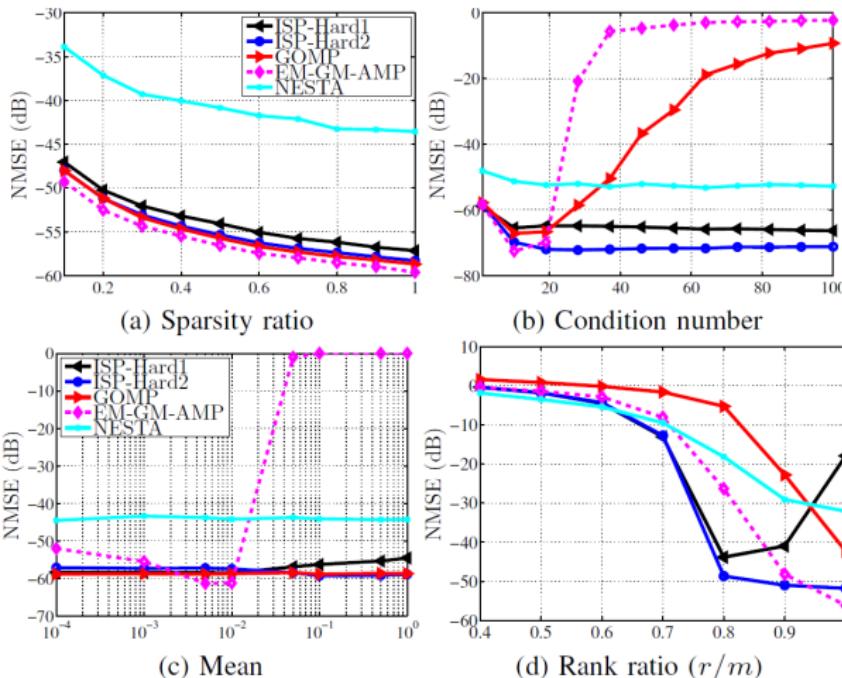
**Algorithm 2** ISP

```
1: Require:  $\mathbf{y}, \mathbf{D}, \tau_0, \tau_{\min}, c > 0, I$ 
2: Initialization:  $\mathbf{x} = \mathbf{D}^\dagger \mathbf{y}, \tau = \tau_0$ 
3: while  $\tau > \tau_{\min}$  do
4:   for  $i = 1, 2, \dots, I$  do
5:      $\tilde{\mathbf{x}} = \mathcal{T}_\tau(\mathbf{x})$ 
6:      $\mathbf{x} = \tilde{\mathbf{x}} - \mathbf{D}^\dagger (\mathbf{D}\tilde{\mathbf{x}} - \mathbf{y})$ 
7:   end for
8:    $\tau \leftarrow c \cdot \tau$ 
9: end while
10: Output:  $\mathbf{x}$ 
```

👉  $\mathcal{T}_\tau$  a sparsifying function =  $\begin{cases} \mathbf{x} - \mu_\sigma \nabla f_\sigma(\mathbf{x}) & (\text{ISP} - \ell_0, \text{ISP} - \ell_1) \\ \text{prox}_{\mu_\sigma}(\mathbf{x}) & (\text{ISP-Hard}, \text{ISP-Soft}) \end{cases}$

# Simulations

Recovery performance. Synthetic data:  $\mathbf{y}_{m \times 1} = \mathbf{D}_{m \times n} \mathbf{x}_{n \times 1} + \mathbf{e}_{m \times 1}$ . Bernoulli-Gaussian sparse signal. Gaussian  $\mathbf{D}$  and noise.  $m = 400$ ,  $n = 1000$ . Different measurement matrices.  $\mathbf{D}$ : sparse, ill-conditioned, non-zero mean, low-rank.



1 Proximal algorithms

2 Sparse representation

- Background
- Iterative Sparsification-Projection
- Iterative Proximal Projection

3 Dictionary Learning

4 Deep Neural Networks

5 Conclusions

# Iterative Proximal Projection

$$\min_{\mathbf{x}} J(\mathbf{x}) \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{Ax}\|_2 \leq \epsilon$$

👉  $J$  : non-smooth, non-convex

---

2

F. Ghayem, M. Sadeghi, M. Babaie-Zadeh, S. Chatterjee, M. Skoglund, and C. Jutten, "Sparse signal recovery via iterative proximal projection", *IEEE Trans. Sig. Proc.*, vol. 66, no. 4, pp. 879–894, 2018.



# Iterative Proximal Projection

$$\min_{\mathbf{x}} J(\mathbf{x}) \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{Ax}\|_2 \leq \epsilon$$

👉  $J$  : non-smooth, non-convex

2

IPP

Main idea

$$\min_{\mathbf{x}, \mathbf{z}} J(\mathbf{z}) + \delta_C(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{z} = \mathbf{x}$$

2

F. Ghayem, M. Sadeghi, M. Babaie-Zadeh, S. Chatterjee, M. Skoglund, and C. Jutten, "Sparse signal recovery via iterative proximal projection", *IEEE Trans. Sig. Proc.*, vol. 66, no. 4, pp. 879–894, 2018.

# Iterative Proximal Projection

$$\boxed{\min_{\mathbf{x}} J(\mathbf{x}) \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{Ax}\|_2 \leq \epsilon}$$

☞  $J$  : non-smooth, non-convex

2

IPP

Main idea

$$\min_{\mathbf{x}, \mathbf{z}} J(\mathbf{z}) + \delta_C(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{z} = \mathbf{x}$$

$$\min_{\mathbf{x}, \mathbf{z}} J(\mathbf{z}) + \delta_C(\mathbf{x}) + \frac{1}{2\alpha} \|\mathbf{x} - \mathbf{z}\|_2^2$$

2

F. Ghayem, M. Sadeghi, M. Babaie-Zadeh, S. Chatterjee, M. Skoglund, and C. Jutten, "Sparse signal recovery via iterative proximal projection", *IEEE Trans. Sig. Proc.*, vol. 66, no. 4, pp. 879–894, 2018.

# Iterative Proximal Projection

$$\boxed{\min_{\mathbf{x}} J(\mathbf{x}) \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{Ax}\|_2 \leq \epsilon}$$

☞  $J$  : non-smooth, non-convex

2

IPP

Main idea

$$\min_{\mathbf{x}, \mathbf{z}} J(\mathbf{z}) + \delta_C(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{z} = \mathbf{x}$$

$$\min_{\mathbf{x}, \mathbf{z}} J(\mathbf{z}) + \delta_C(\mathbf{x}) + \frac{1}{2\alpha} \|\mathbf{x} - \mathbf{z}\|_2^2$$

$$\begin{cases} \mathbf{z}_{k+1} = \operatorname{argmin}_{\mathbf{z}} \alpha J(\mathbf{z}) + \frac{1}{2} \|\mathbf{z} - \mathbf{x}_k\|_2^2 \\ \mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \delta_C(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{z}_{k+1}\|_2^2 \end{cases}$$

2

F. Ghayem, M. Sadeghi, M. Babaie-Zadeh, S. Chatterjee, M. Skoglund, and C. Jutten, "Sparse signal recovery via iterative proximal projection", *IEEE Trans. Sig. Proc.*, vol. 66, no. 4, pp. 879–894, 2018.

# Iterative Proximal Projection

$$\boxed{\min_{\mathbf{x}} J(\mathbf{x}) \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{Ax}\|_2 \leq \epsilon}$$

☞  $J$  : non-smooth, non-convex

2

IPP

Main idea

$$\min_{\mathbf{x}, \mathbf{z}} J(\mathbf{z}) + \delta_C(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{z} = \mathbf{x}$$

$$\min_{\mathbf{x}, \mathbf{z}} J(\mathbf{z}) + \delta_C(\mathbf{x}) + \frac{1}{2\alpha} \|\mathbf{x} - \mathbf{z}\|_2^2$$

$$\begin{cases} \mathbf{z}_{k+1} = \operatorname{argmin}_{\mathbf{z}} \alpha J(\mathbf{z}) + \frac{1}{2} \|\mathbf{z} - \mathbf{x}_k\|_2^2 \\ \mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \delta_C(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{z}_{k+1}\|_2^2 \end{cases} \quad \begin{cases} \tilde{\mathbf{x}}_k = \mathbf{x}_k + w \cdot (\mathbf{x}_k - \mathbf{x}_{k-1}) \\ \mathbf{x}_{k+1} = \mathcal{P}_C \left( \operatorname{prox}_{\alpha \cdot J}(\tilde{\mathbf{x}}_k) \right) \end{cases}$$

2

F. Ghayem, M. Sadeghi, M. Babaie-Zadeh, S. Chatterjee, M. Skoglund, and C. Jutten, "Sparse signal recovery via iterative proximal projection", *IEEE Trans. Sig. Proc.*, vol. 66, no. 4, pp. 879–894, 2018.

# Iterative Proximal Projection

- Difficult to prove convergence. An alternative algorithm with convergence proof!

# Iterative Proximal Projection

- Difficult to prove convergence. An alternative algorithm with convergence proof!

IPP

An approximate solver:

$$\begin{cases} \tilde{\mathbf{x}}_k = \mathbf{x}_k + w \cdot (\mathbf{x}_k - \mathbf{x}_{k-1}) \\ \mathbf{x}_{k+1} = \mathcal{P}_C \left( \text{prox}_{\alpha \cdot J}(\tilde{\mathbf{x}}_k) \right) \end{cases}$$

# Iterative Proximal Projection

- Difficult to prove convergence. An alternative algorithm with convergence proof!

IPP

An approximate solver:

$$\begin{cases} \tilde{\mathbf{x}}_k = \mathbf{x}_k + w \cdot (\mathbf{x}_k - \mathbf{x}_{k-1}) \\ \mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{C}} \left( \text{prox}_{\alpha \cdot J}(\tilde{\mathbf{x}}_k) \right) \end{cases} \rightarrow \begin{cases} \tilde{\mathbf{x}}_k = \mathbf{x}_k + w \cdot (\mathbf{x}_k - \mathbf{x}_{k-1}) \\ \mathbf{z}_{k+1} = \text{prox}_{\mu_z \cdot \alpha J}(\mathbf{z}_k + \mu_z(\tilde{\mathbf{x}}_k - \mathbf{z}_k)) \\ \mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{C}}(\mathbf{x}_k + \mu_x(\mathbf{z}_{k+1} - \mathbf{x}_k)) \end{cases}$$

☞  $0 < \mu_x, \mu_z < 1$ . If  $\mu_x, \mu_z \rightarrow 1$  the two algorithms coincide!

# Iterative Proximal Projection

- Difficult to prove convergence. An alternative algorithm with convergence proof!

IPP

An approximate solver:

$$\begin{cases} \tilde{\mathbf{x}}_k = \mathbf{x}_k + w \cdot (\mathbf{x}_k - \mathbf{x}_{k-1}) \\ \mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{C}} \left( \text{prox}_{\alpha \cdot J}(\tilde{\mathbf{x}}_k) \right) \end{cases} \rightarrow \begin{cases} \tilde{\mathbf{x}}_k = \mathbf{x}_k + w \cdot (\mathbf{x}_k - \mathbf{x}_{k-1}) \\ \mathbf{z}_{k+1} = \text{prox}_{\mu_z \cdot \alpha J}(\mathbf{z}_k + \mu_z(\tilde{\mathbf{x}}_k - \mathbf{z}_k)) \\ \mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{C}}(\mathbf{x}_k + \mu_x(\mathbf{z}_{k+1} - \mathbf{x}_k)) \end{cases}$$

☞  $0 < \mu_x, \mu_z < 1$ . If  $\mu_x, \mu_z \rightarrow 1$  the two algorithms coincide!

---

**Algorithm** Iterative Proximal Projection (IPP)

---

1: **Inputs:**  $\mathbf{y}, \mathbf{A}, \epsilon, \alpha_i, \alpha_f, \tau, 0 < c < 1, w, 0 < \mu_x, \mu_z < 1$   
2: **Initialization:**  $k = 0, \mathbf{x}_0 = \mathbf{z}_0 = \mathbf{A}^\dagger \mathbf{y}, \alpha = \alpha_i$   
3: **while**  $\alpha > \alpha_f$  **do**  
4:     **while**  $\|\mathbf{x}_k - \mathbf{x}_{k-1}\|_2 > \tau$  **do**  
5:          $\tilde{\mathbf{x}}_k = \mathbf{x}_k + w \cdot (\mathbf{x}_k - \mathbf{x}_{k-1})$   
6:          $\mathbf{z}_{k+1} = \text{prox}_{\mu_z \cdot \alpha J}(\mathbf{z}_k + \mu_z(\tilde{\mathbf{x}}_k - \mathbf{z}_k))$   
7:          $\mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{A}_c}(\mathbf{x}_k + \mu_x(\mathbf{z}_{k+1} - \mathbf{x}_k))$   
8:          $k \rightarrow k + 1$   
9:     **end while**  
10:     $\alpha \leftarrow c \cdot \alpha$   
11: **end while**  
12: **Output:**  $\mathbf{x}_k$

---

# Iterative Proximal Projection

## Theorem (IPP convergence)

In the IPP algorithm, assume that  $0 \leq w < \frac{1}{\max(\mu_x, \mu_z)} - 1$ . The sequence  $\left\{ \mathbf{u}_k \triangleq (\mathbf{x}_k, \mathbf{z}_k) \right\}_{k=0}^{\infty}$  generated by IPP for each value of  $\alpha$  (the inner-loop iterations) converges to a critical point,  $\mathbf{u}^*$ , of the cost function. Furthermore, if the cost function satisfies the Kurdyka-Łojasiewicz (KL) property [Bolt et al., 2014] with  $\psi(t) = c \cdot t^{1-\theta}$  for some  $t > 0$  and  $\theta \in [0, 1)$ , then:

- If  $\theta = 0$  then the sequence  $\{\mathbf{u}_k\}_{k \geq 0}$  converges in a finite number of steps.
- If  $\theta \in (0, 1/2]$  then there exist  $d > 0$  and  $\tau \in [0, 1)$  such that  $\|\mathbf{u}_k - \mathbf{u}^*\|_2 \leq d \cdot \tau^k$ .
- If  $\theta \in (1/2, 1)$  then there exist  $d > 0$  such that  $\|\mathbf{u}_k - \mathbf{u}^*\|_2 \leq d \cdot k^{\frac{\theta-1}{2\theta-1}}$ .

# Simulations

Block based compressed image recovery: 50% overlapping blocks, Gaussian measurement,  $\delta$  = sampling rate.

	$\delta = 0.1$			$\delta = 0.2$		
	House	Barbara	Monarch	House	Barbara	Monarch
$\ell_q$	25.13	22.99	19.87	28.56	25.30	22.71
GOMP	25.00	22.38	19.03	26.68	23.94	21.23
SCSA	25.11	22.96	19.88	<b>28.62</b>	25.26	22.64
EMGMAMP	25.02	22.94	19.69	27.96	25.06	22.22
IPP ( $w = 0$ )	<b>25.31</b>	<b>23.08</b>	<b>20.46</b>	27.72	<b>25.36</b>	<b>22.95</b>
IPP ( $w = 0.85$ )	<b>25.57</b>	<b>23.38</b>	<b>20.67</b>	<b>28.65</b>	<b>25.63</b>	<b>23.45</b>

	$\delta = 0.3$			$\delta = 0.4$		
	House	Barbara	Monarch	House	Barbara	Monarch
$\ell_q$	31.17	27.38	24.97	34.15	30.17	27.65
GOMP	30.85	27.14	24.43	33.59	29.69	27.10
SCSA	<b>31.35</b>	27.37	<b>24.99</b>	<b>34.49</b>	30.23	<b>27.79</b>
EMGMAMP	30.91	27.15	24.65	33.64	29.72	27.06
IPP ( $w = 0$ )	31.11	<b>27.61</b>	24.73	33.56	<b>30.41</b>	27.70
IPP ( $w = 0.85$ )	<b>32.50</b>	<b>28.17</b>	<b>25.55</b>	<b>35.01</b>	<b>30.92</b>	<b>28.47</b>

1 Proximal algorithms

2 Sparse representation

3 Dictionary Learning

- Background
- learning low coherence dictionaries

4 Deep Neural Networks

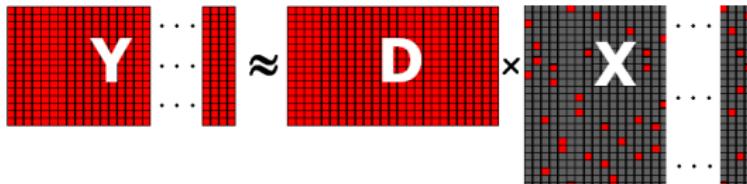
5 Conclusions

# Dictionary Learning

- Learn a sparsifying dictionary from training data:  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_L]$ .

## Dictionary Learning Problem

$$\min_{\mathbf{D}, \mathbf{X}} \sum_{i=1}^L \frac{1}{2} \|\mathbf{y}_i - \mathbf{D}\mathbf{x}_i\|_2 = \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 \quad \text{s.t.} \quad \mathbf{D} \in \mathcal{D}, \mathbf{X} \in \mathcal{X}$$

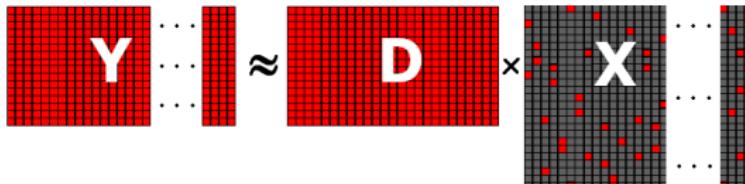


# Dictionary Learning

- Learn a sparsifying dictionary from training data:  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_L]$ .

## Dictionary Learning Problem

$$\min_{\mathbf{D}, \mathbf{X}} \sum_{i=1}^L \frac{1}{2} \|\mathbf{y}_i - \mathbf{D}\mathbf{x}_i\|_2 = \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 \quad \text{s.t. } \mathbf{D} \in \mathcal{D}, \mathbf{X} \in \mathcal{X}$$



## Alternating minimization

- Start with  $(\mathbf{D}^{(0)}, \mathbf{X}^{(0)})$ . Alternate between:

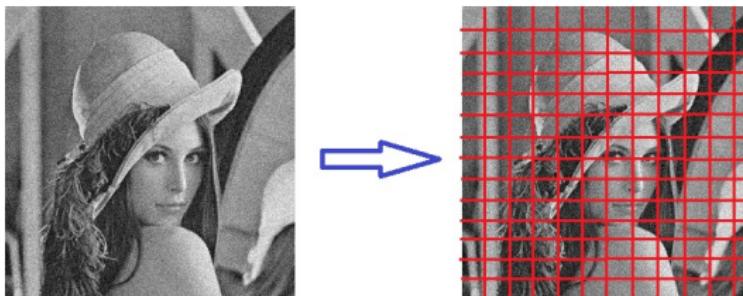
- Sparse representation:**  $\mathbf{X}^{(k+1)} = \operatorname{argmin}_{\mathbf{X} \in \mathcal{X}} \frac{1}{2} \|\mathbf{Y} - \mathbf{D}^{(k)} \mathbf{X}\|_F^2$ 
  - ☞ OMP, IST, SL0, ...
- Dictionary update:**  $\mathbf{D}^{(k+1)} = \operatorname{argmin}_{\mathbf{D} \in \mathcal{D}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}^{(k+1)}\|_F^2$ 
  - ☞ MOD, KSVD, ...

# Applications of dictionary learning

- Signal/image restoration and enhancement
  - Image denoising [Elad et al., 2006]:

$$\mathbf{y} = \mathbf{x} + \mathbf{e}$$

Learn the dictionary from the noisy image itself!



$$\hat{\mathbf{x}} = \underset{\mathbf{x}, \mathbf{D}, \{\alpha_{ij}\}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \sum_{i,j} \|\mathbf{R}_{ij} \mathbf{x} - \mathbf{D} \alpha_{ij}\|_2^2 \quad \text{s.t.} \quad \|\alpha_{ij}\|_0 \leq \tau$$

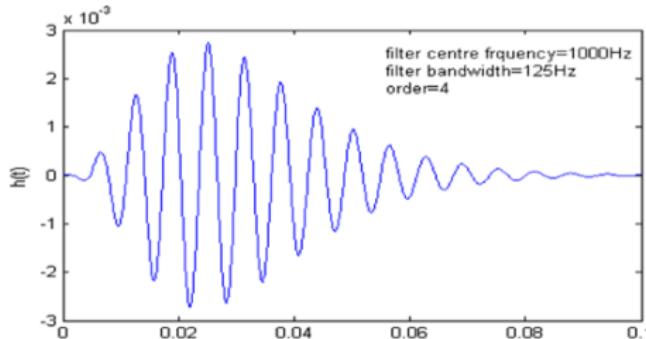
# Applications of dictionary learning

- Speech denoising [Sigg et al., 2012; Jafari et al., 2011]
- Parameter dictionary learning [Yaghoobi et al., 2009; Ataee et al., 2010].

## Learning structured atoms:

- Example. Gammatone filters have shown similarities with the human auditory system:

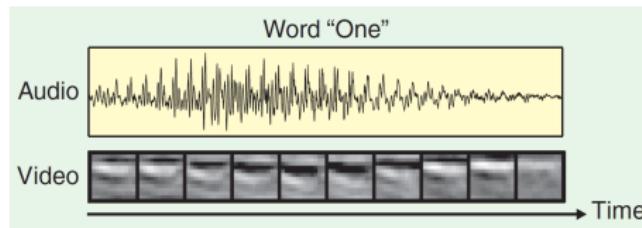
$$g(t) = at^{n-1}e^{-2\pi bBt} \cos(2\pi f_c t)$$



[https://www.researchgate.net/publication/264837246\\_Voice\\_biometric\\_feature\\_using\\_Gammatone\\_filterbank\\_and\\_ICA/figures?lo=1](https://www.researchgate.net/publication/264837246_Voice_biometric_feature_using_Gammatone_filterbank_and_ICA/figures?lo=1)

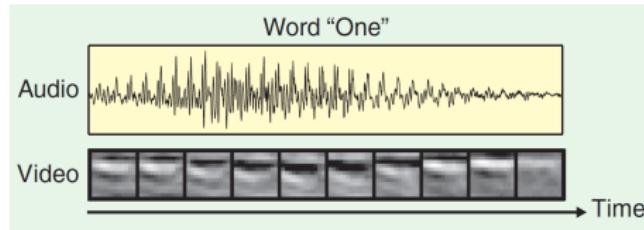
# Applications of dictionary learning

- Multi-modal dictionary learning [Monaci et al., 2007; Zhuang et al., 2013]
  - Learning multi-modal atoms to describe underlying generating cause, speaker localization, and so on



# Applications of dictionary learning

- Multi-modal dictionary learning [Monaci et al., 2007; Zhuang et al., 2013]
  - Learning multi-modal atoms to describe underlying generating cause, speaker localization, and so on



- Stereo image representation [Tosic et al., 2011]
  - Efficient image representation to perform different vision task like camera pose estimation

The diagram illustrates the decomposition of two stereo images into sparse linear combinations of learned atoms. The top row shows the left image ( $L$ ) being represented as:

$$L = a_1 \cdot \text{atom}_1 + a_2 \cdot \text{atom}_2 + a_3 \cdot \text{atom}_3 + \eta_L$$

The bottom row shows the right image ( $R$ ) being represented as:

$$R = b_1 \cdot \text{atom}_1 + b_2 \cdot \text{atom}_2 + b_3 \cdot \text{atom}_3 + \eta_R$$

# Applications of dictionary learning

- Supervised dictionary learning [Mairal et al., 2010; Zhang et al., 2010]

$$\min_{\mathbf{D}, \mathbf{X}} \underbrace{\frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2}_{\text{representation power}} + \underbrace{\lambda \|\mathbf{T} - \mathbf{WX}\|_F^2}_{\text{discrimination power}} \quad \text{s.t. } \|\mathbf{X}\|_0 \leq \tau$$

- $\mathbf{T}$ : label matrix
- $\mathbf{W}$ : linear classifier

# Learning Low Coherence Dictionaries

## Mutual Coherence

For a dictionary  $\mathbf{D} \in \mathbb{R}^{n \times m}$ , its mutual coherence is defined as

$$\mu(\mathbf{D}) \triangleq \max_{i \neq j} \frac{|\langle \mathbf{d}_i, \mathbf{d}_j \rangle|}{\|\mathbf{d}_i\|_2 \cdot \|\mathbf{d}_j\|_2}$$

The Welch bound:

$$\mu \geq \sqrt{\frac{m-n}{n(m-1)}}$$

# Learning Low Coherence Dictionaries

## Mutual Coherence

For a dictionary  $\mathbf{D} \in \mathbb{R}^{n \times m}$ , its mutual coherence is defined as

$$\mu(\mathbf{D}) \triangleq \max_{i \neq j} \frac{|\langle \mathbf{d}_i, \mathbf{d}_j \rangle|}{\|\mathbf{d}_i\|_2 \cdot \|\mathbf{d}_j\|_2}$$

The Welch bound:

$$\mu \geq \sqrt{\frac{m-n}{n(m-1)}}$$

## Low MC dictionaries

A dictionary with low MC is desired for sparse recovery algorithms.

 Guaranteed sparse recovery as long as  $\|\mathbf{x}\|_0 \leq \frac{1}{2}(1 + \frac{1}{\mu(\mathbf{D})})$

# Learning Low Coherence Dictionaries

- Regularized:

$$\min_{\mathbf{D} \in \mathcal{D}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 + \lambda \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_F^2$$

$$\|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_F^2 = \sum_{i \neq j} |\langle \mathbf{d}_i, \mathbf{d}_j \rangle|^2 + \sum_i (\langle \mathbf{d}_i, \mathbf{d}_i \rangle - 1)^2$$

Bounded Self Coherence (BSC) [Sigg et. al, 2012], Gradient Projection (GP) [Bao et. al, 2015]

# Learning Low Coherence Dictionaries

- Regularized:

$$\min_{\mathbf{D} \in \mathcal{D}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 + \lambda \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_F^2$$

$$\|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_F^2 = \sum_{i \neq j} |\langle \mathbf{d}_i, \mathbf{d}_j \rangle|^2 + \sum_i (\langle \mathbf{d}_i, \mathbf{d}_i \rangle - 1)^2$$

Bounded Self Coherence (BSC) [Sigg et. al, 2012], Gradient Projection (GP) [Bao et. al, 2015]

- Constrained:

$$\min_{\mathbf{D} \in \mathcal{D}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 \quad \text{s.t.} \quad \mu(\mathbf{D}) \leq \mu_0$$

Iterative Projection Rotation (IPR) [Barchiesi and Plumbley, 2013]

# Learning Low Coherence Dictionaries

- Regularized:

$$\min_{\mathbf{D} \in \mathcal{D}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 + \lambda \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_F^2$$

$$\|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_F^2 = \sum_{i \neq j} |\langle \mathbf{d}_i, \mathbf{d}_j \rangle|^2 + \sum_i (\langle \mathbf{d}_i, \mathbf{d}_i \rangle - 1)^2$$

Bounded Self Coherence (BSC) [Sigg et. al, 2012], Gradient Projection (GP) [Bao et. al, 2015]

- Constrained:

$$\min_{\mathbf{D} \in \mathcal{D}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 \quad \text{s.t.} \quad \mu(\mathbf{D}) \leq \mu_0$$

Iterative Projection Rotation (IPR) [Barchiesi and Plumbley, 2013]



Uses a two-step procedure: decorrelation + rotation

1 Proximal algorithms

2 Sparse representation

3 Dictionary Learning

- Background
- learning low coherence dictionaries

4 Deep Neural Networks

5 Conclusions

# Mutual coherence-based DL

## Low coherence DL

$$\mu(\mathbf{D}) = \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_\infty \quad \|\mathbf{A}\|_\infty \triangleq \max_{i,j} |a_{ij}|$$

Proposed problems<sup>a</sup>:

$$\begin{cases} \min_{\mathbf{D} \in \mathcal{D}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 + \lambda \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_\infty \\ \min_{\mathbf{D} \in \mathcal{D}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 \quad \text{s.t.} \quad \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_\infty \leq \mu_0 \end{cases}$$

<sup>a</sup>

M. Sadeghi and M. Babaie-Zadeh, "Learning low-coherence dictionaries for sparse representation", *Sig. Proc.*, 2018 (submitted).

# Mutual coherence-based DL

## Low coherence DL

$$\mu(\mathbf{D}) = \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_\infty \quad \|\mathbf{A}\|_\infty \triangleq \max_{i,j} |a_{ij}|$$

Proposed problems<sup>a</sup>:

$$\begin{cases} \min_{\mathbf{D} \in \mathcal{D}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 + \lambda \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_\infty \\ \min_{\mathbf{D} \in \mathcal{D}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 \quad \text{s.t.} \quad \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_\infty \leq \mu_0 \end{cases}$$

<sup>a</sup> M. Sadeghi and M. Babaie-Zadeh, "Learning low-coherence dictionaries for sparse representation", *Sig. Proc.*, 2018 (submitted).

## Main idea

- Introduce  $\mathbf{G} = \mathbf{D}^T \mathbf{D} - \mathbf{I}$ . Use penalty method + proximal algorithm!

$$\min_{\mathbf{D} \in \mathcal{D}, \mathbf{G}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 + \frac{1}{2\alpha} \|\mathbf{G} - \mathbf{D}^T \mathbf{D} + \mathbf{I}\|_F^2 + \lambda g(\mathbf{G})$$

# Mutual coherence-based DL

## Low coherence DL

$$\mu(\mathbf{D}) = \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_\infty \quad \|\mathbf{A}\|_\infty \triangleq \max_{i,j} |a_{ij}|$$

Proposed problems<sup>a</sup>:

$$\begin{cases} \min_{\mathbf{D} \in \mathcal{D}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 + \lambda \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_\infty \\ \min_{\mathbf{D} \in \mathcal{D}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 \quad \text{s.t.} \quad \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_\infty \leq \mu_0 \end{cases}$$

<sup>a</sup> M. Sadeghi and M. Babaie-Zadeh, "Learning low-coherence dictionaries for sparse representation", *Sig. Proc.*, 2018 (submitted).

## Main idea

- Introduce  $\mathbf{G} = \mathbf{D}^T \mathbf{D} - \mathbf{I}$ . Use penalty method + proximal algorithm!

$$\min_{\mathbf{D} \in \mathcal{D}, \mathbf{G}} \frac{1}{2} \|\mathbf{Y} - \mathbf{DX}\|_F^2 + \frac{1}{2\alpha} \|\mathbf{G} - \mathbf{D}^T \mathbf{D} + \mathbf{I}\|_F^2 + \lambda g(\mathbf{G})$$

$$g(\mathbf{G}) = \begin{cases} \|\mathbf{G}\|_\infty & \text{(regularized)} \\ \delta_C(\mathbf{G}), & C \triangleq \{\mathbf{G} \mid \|\mathbf{G}\|_\infty \leq \mu_0\} \quad \text{(constrained)} \end{cases}$$

## Updating G

Let  $g$  denote the function  $\eta \|\cdot\|_\infty : \mathbb{R}^{N \times N} \rightarrow \mathbb{R}$ . The proximal mapping of  $g$  is given by

$$\text{prox}_g(\mathbf{U}) = \mathbf{U} - P_{\frac{\eta}{1}}(\mathbf{U})$$

where,  $P_{\frac{\eta}{1}}(\cdot) : \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^{N \times N}$  is the projection onto the  $\ell_1$  norm-ball of radius  $\eta$ .

# Mutual coherence-based DL

## Updating G

Let  $g$  denote the function  $\eta \|\cdot\|_\infty : \mathbb{R}^{N \times N} \rightarrow \mathbb{R}$ . The proximal mapping of  $g$  is given by

$$\text{prox}_g(\mathbf{U}) = \mathbf{U} - P_{\frac{1}{\eta}}(\mathbf{U})$$

where,  $P_{\frac{1}{\eta}}(\cdot) : \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^{N \times N}$  is the projection onto the  $\ell_1$  norm-ball of radius  $\eta$ .

## Updating D

The gradient of  $\frac{1}{2}\|\mathbf{Y} - \mathbf{DX}\|_F^2 + \frac{1}{2\alpha}\|\mathbf{G} - \mathbf{D}^T \mathbf{D} + \mathbf{I}\|_F^2$  with respect to  $\mathbf{D}$  is Lipschitz continuous over  $\mathcal{D}$  with constant

$$L = \|\mathbf{X}\mathbf{X}^T\| + \frac{6N + 2\|\mathbf{G}\|_F}{2\alpha}$$

# Mutual coherence-based DL

## Regularized Incoherent DL (RINC-DL) and Constrained Incoherent DL (CINC-DL)

---

**Algorithm 1** RINC-DL

---

```
1: Require:  $\mathbf{Y}, \mathbf{D}_0, \tau, \lambda, c, L_2, \epsilon, I, J$ 
2: Initialization:  $\mathbf{D} = \mathbf{D}_0, \mathbf{G} = 0$ 
3: while stopping criterion for DL not met do
4:   1. Sparse approximation:  $\mathbf{X} = \text{SD}(\mathbf{Y}, \mathbf{D}, \tau)$ 
5:   2. Dictionary update:
6:      $L_1 = \|\mathbf{X}^T \mathbf{X}\|$ 
7:      $\alpha = 3 \cdot \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_\infty$ 
8:      $i = 1$ 
9:     while  $i \leq I$  and  $\|\mathbf{G} - \mathbf{D}^T \mathbf{D}\|_F > \epsilon$  do
10:       $\mu_d = 1/(L_1 + \alpha^{-1} L_2)$ 
11:      for  $j = 1, 2, \dots, J$  do
12:         $\mathbf{G} = \mathbf{D}^T \mathbf{D} - P_{\mathcal{B}_1^{\lambda, \alpha}}(\mathbf{D}^T \mathbf{D} - \mathbf{I})$ 
           
$$\boxed{\mathbf{G} = \mathbf{I} + P_{\mathcal{B}_{\infty}^{c\alpha}}(\mathbf{D}^T \mathbf{D} - \mathbf{I})}$$
 CINC-DL
13:        
$$|\quad \mathbf{D} = P_{\mathcal{D}}(\mathbf{D} - \mu_d \nabla f(\mathbf{D}))$$

14:      end for
15:       $\alpha \leftarrow c \cdot \alpha$ 
16:       $i \leftarrow i + 1$ 
17:    end while
18:  end while
19: Output:  $\mathbf{D}, \mathbf{X}$ 
```

---

# Mutual coherence-based DL

## Regularized Incoherent DL (RINC-DL) and Constrained Incoherent DL (CINC-DL)

---

**Algorithm 1** RINC-DL

---

```
1: Require:  $\mathbf{Y}, \mathbf{D}_0, \tau, \lambda, c, L_2, \epsilon, I, J$ 
2: Initialization:  $\mathbf{D} = \mathbf{D}_0, \mathbf{G} = 0$ 
3: while stopping criterion for DL not met do
4:   1. Sparse approximation:  $\mathbf{X} = \text{SD}(\mathbf{Y}, \mathbf{D}, \tau)$ 
5:   2. Dictionary update:
6:      $L_1 = \|\mathbf{X}^T \mathbf{X}\|$ 
7:      $\alpha = 3 \cdot \|\mathbf{D}^T \mathbf{D} - \mathbf{I}\|_\infty$ 
8:      $i = 1$ 
9:     while  $i \leq I$  and  $\|\mathbf{G} - \mathbf{D}^T \mathbf{D}\|_F > \epsilon$  do
10:       $\mu_d = 1/(L_1 + \alpha^{-1} L_2)$ 
11:      for  $j = 1, 2, \dots, J$  do
12:         $\mathbf{G} = \mathbf{D}^T \mathbf{D} - P_{\mathcal{B}_1^{\lambda, \alpha}}(\mathbf{D}^T \mathbf{D} - \mathbf{I})$ 
13:         $\mathbf{G} = \mathbf{I} + P_{\mathcal{B}_{\infty}^{c\alpha}}(\mathbf{D}^T \mathbf{D} - \mathbf{I})$  CINC-DL
14:         $\mathbf{D} = P_{\mathcal{D}}(\mathbf{D} - \mu_d \nabla f(\mathbf{D}))$ 
15:      end for
16:       $\alpha \leftarrow c \cdot \alpha$ 
17:       $i \leftarrow i + 1$ 
18:    end while
19:  end while
20: Output:  $\mathbf{D}, \mathbf{X}$ 
```

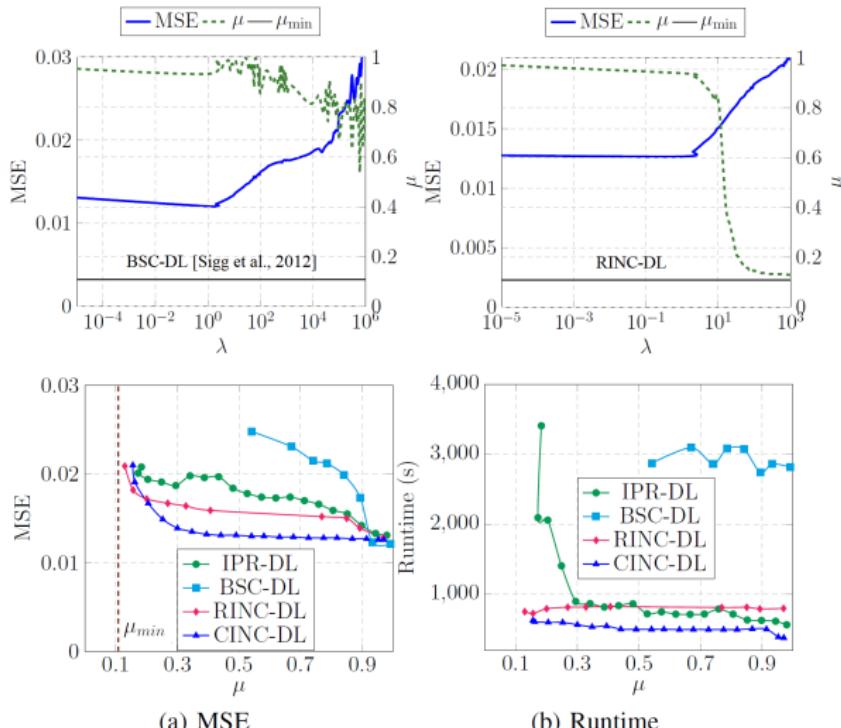
---

Algorithm	Complexity
IPR-DL	$\mathcal{O}(nNM + MN^2 + 3n^2N + 2n^3 + 2N^3)$
CINC-DL	$\mathcal{O}(nNM + nN^2)$
RINC-DL	$\mathcal{O}(nNM + nN^2)$

- $n$  : signal dimension
- $N$  : number of atoms
- $M$  : number of training signals

# Simulations

Low-coherence DL for  $8 \times 8$  image blocks:  $M = 50,000$ ,  $\mathbf{D}_0 = \text{DCT}_{64 \times 256}$   
( $\mu_{\min} = 0.1085$ ),  $s = 10$ .



1 Proximal algorithms

2 Sparse representation

3 Dictionary Learning

4 Deep Neural Networks

- Background

- Progressive Neural Networks

- Structured Weight Matrices for Neural Networks

5 Conclusions

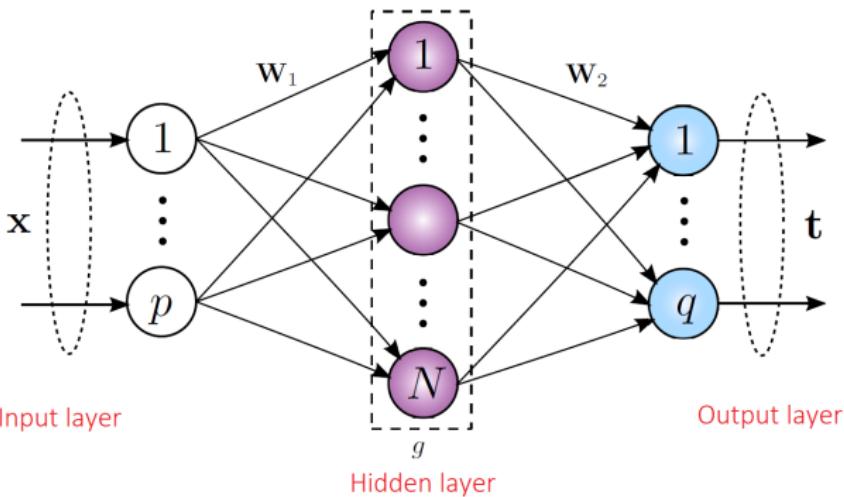
# Artificial Neural Networks

Estimate the mapping function  $t = f(x)$  given some training data  $\{(x_i, t_i)\}_i$

Single layer NN:

$$\hat{t} = \mathbf{W}_2 \cdot g(\mathbf{W}_1 \cdot \mathbf{x})$$

$g$ : non-linear activation function



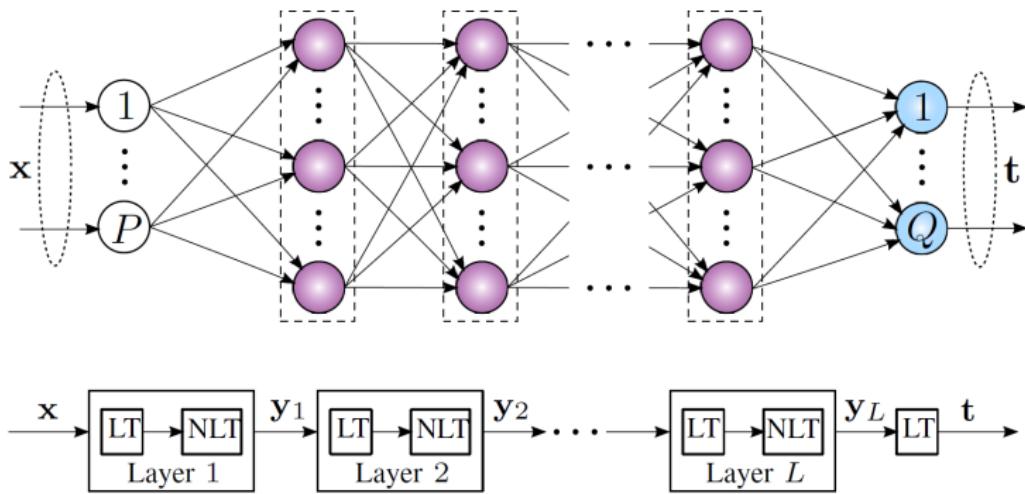
It can approximate any continuous function under mild conditions (**universal approximation**)

# Artificial Neural Networks

- To approximate complex functions, increase the number of hidden nodes
- Leads to very wide networks!

**Solution:** Deepen the network instead of widening it!

Multilayer NN:



# Artificial Neural Networks

Training algorithms:

$$\min_{\mathbf{W}_1, \mathbf{W}_2} \frac{1}{M} \sum_{i=1}^M \|\mathbf{t}_i - \mathbf{W}_2 \cdot g(\mathbf{W}_1 \cdot \mathbf{x}_i)\|_2^2$$

- Stochastic Gradient Descent (SGD) – Backpropagation
- Gradient-free training using ADMM [Taylor et al. 2016]
- Proximal backpropagation [Frerix et al. 2017]
- Extreme learning machine (ELM) [Huang et al. 2006]:

☞ Set  $\mathbf{W}_1$  randomly

☞ Find  $\mathbf{W}_2$  using least squares:  $\min_{\mathbf{W}_2} \frac{1}{M} \sum_{i=1}^M \|\mathbf{t}_i - \mathbf{W}_2 \cdot g(\mathbf{W}_1 \cdot \mathbf{x}_i)\|_2^2$

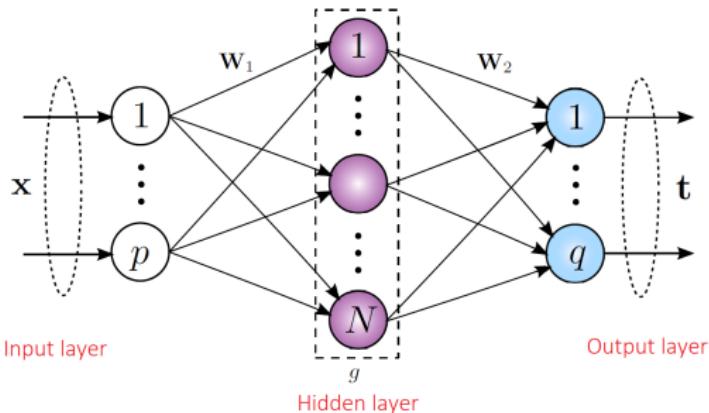
# Dictionary Learning and Neural Networks

Dictionary learning has some similarities with single-layer NN:

- DL:  $\min_{\mathbf{D}, \mathbf{X}} \|\mathbf{Y} - \mathbf{DX}\|_F^2 \text{ s.t. } \mathbf{X} \text{ is sparse}$

- SNN::  $\min_{\mathbf{W}_1, \mathbf{W}_2} \|\mathbf{T} - \mathbf{W}_2 \cdot g(\mathbf{W}_1 \cdot \mathbf{X})\|_F^2$

- ☞ Both  $\mathbf{X}$  and  $g(\mathbf{W}_1 \cdot \mathbf{X})$  with  $g$  being ReLU:  $g(x) = \max(x, 0)$  are non-linear and sparse representation of data



# Dictionary Learning and Neural Networks

More explicit connections between DL (dictionary learning) and DL (deep learning):

- V. Petyan, Y. Romano, J. Sulam, and M. Elad, “**Theoretical Foundations of Deep Learning via Sparse Representations**,” to appear in *IEEE Signal Processing Magazine*.
- J. Sulam, V. Petyan, Y. Romano, and M. Elad, “**Multi-Layer Convolutional Sparse Modeling: Pursuit and Dictionary Learning**,” to appear in *IEEE Trans. on Signal Processing*.

From multi-layer convolutional sparse coding (CSC) to convolutional neural networks (CNNs)

# Artificial Neural Networks

Questions to address in a multilayer NN:

- How to choose number of layers in a network?
- How to choose number of nodes in each layer?
- How to guarantee that increase in size results in better (non-increasing) optimized cost for training data?
- How to design with appropriate regularization of network parameters to avoid over-fitting to training data?
- Can we use random weight matrices to keep the number of parameters to learn in balance?

1 Proximal algorithms

2 Sparse representation

3 Dictionary Learning

4 Deep Neural Networks

- Background
- **Progressive Neural Networks**
- Structured Weight Matrices for Neural Networks

5 Conclusions

# Progression Learning Network

## Definition (Progression Property)

A non-linear  $g(\cdot)$  function holds the progression property (PP) if there are two known linear transformations  $\mathbf{V} \in \mathbb{R}^{M \times N}$  and  $\mathbf{U} \in \mathbb{R}^{N \times M}$  such that  $\forall \gamma \in \mathbb{R}^N$ :

$$\mathbf{U}g(\mathbf{V}\gamma) = \gamma$$

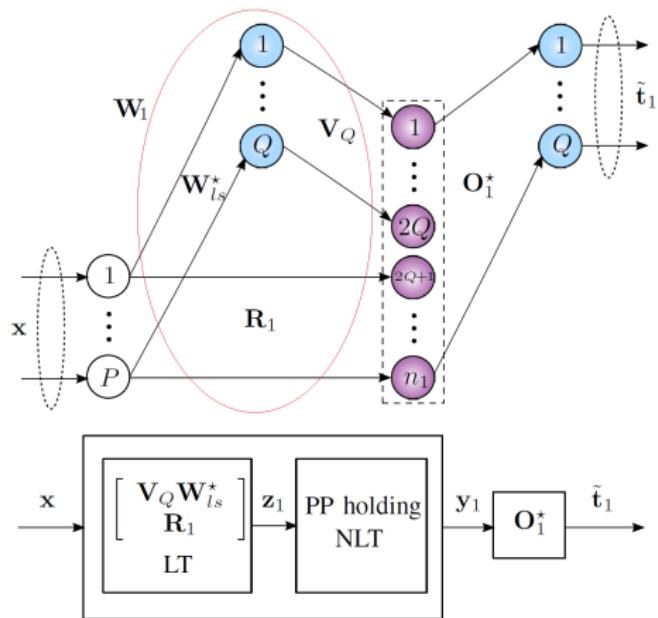
<sup>a</sup> S. Chatterjee, A. M. Javid, M. Sadeghi, P. P. Mitra and M. Skoglund, "Progressive learning for systematic design of large neural networks", *IEEE Trans. Neural Networks and Learning Systems*, 2017 (submitted).

- **Example:** The rectified linear unit (ReLU) function [Glorot et al. 2011]

$$g(\gamma) = \max(\gamma, 0) = \begin{cases} \gamma, & \text{if } \gamma \geq 0 \\ 0, & \text{if } \gamma < 0. \end{cases}$$

If  $\mathbf{V} \triangleq \mathbf{V}_N = [\mathbf{I}_N, -\mathbf{I}_N]^T \in \mathbb{R}^{2N \times N}$  and  $\mathbf{U} \triangleq \mathbf{U}_N = [\mathbf{I}_N \ -\mathbf{I}_N] \in \mathbb{R}^{N \times 2N}$  then ReLU holds PP. Here  $\mathbf{I}_N$  denotes identity matrix of size  $N$  ( $M = 2N$ ).

# Single layer PLN



- $\tilde{t} = O_1 g(W_1 x)$
- $R_1$  is a random matrix

# Single layer PLN

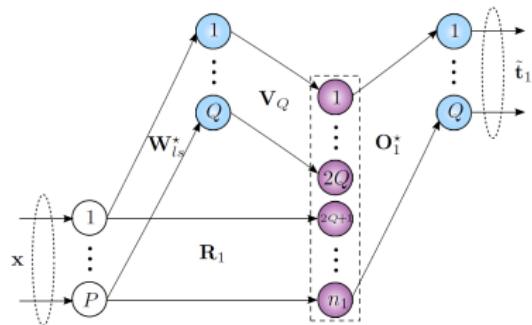
- $\mathbf{W}_{ls}^* = \arg \min_{\mathbf{W}_{ls}} \sum_j \|\mathbf{t}^{(j)} - \mathbf{W}_{ls} \mathbf{x}^{(j)}\|_p^p$  s.t.  $\|\mathbf{W}_{ls}\|_q^q \leq \epsilon$
- $\mathbf{O}_1^* = \arg \min_{\mathbf{O}_1} \sum_j \|\mathbf{t}^{(j)} - \mathbf{O}_1 \mathbf{y}_1^{(j)}\|_p^p$  such that  $\|\mathbf{O}_1\|_q^q \leq \alpha \|\mathbf{U}_Q\|_q^q$ ,

$$C_{ls}^* \triangleq C(\mathbf{W}_{ls}^*) = \sum_j \|\mathbf{t}^{(j)} - \mathbf{W}_{ls}^* \mathbf{x}^{(j)}\|_p^p$$

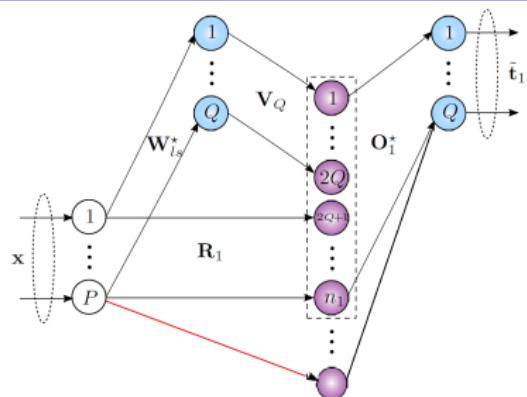
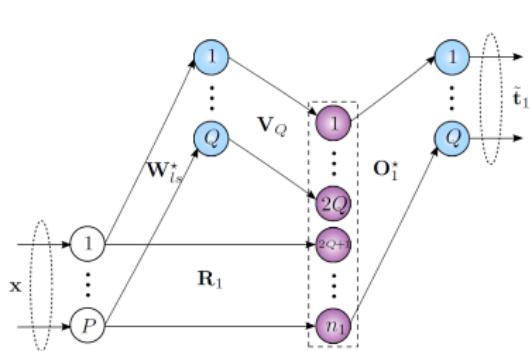
$$C_1^* = C(\mathbf{O}_1^*) = \sum_j \|\mathbf{t}^{(j)} - \mathbf{O}_1^* \mathbf{y}_1^{(j)}\|_p^p = \sum_j \|\mathbf{t}^{(j)} - \mathbf{O}_1^* \mathbf{g}(\mathbf{W}_1 \mathbf{x}^{(j)})\|_p^p$$

- ☞ Relation between optimal linear system and single layer PLN costs:  $C_1^* \leq C_{ls}^*$ . At the equality condition:  $\mathbf{O}_1^* = [\mathbf{U}_Q \ \mathbf{0}]$
- ☞ Adding nodes to the layer:  $C_1^*(n_1 + \Delta) \leq C_1^*(n_1)$

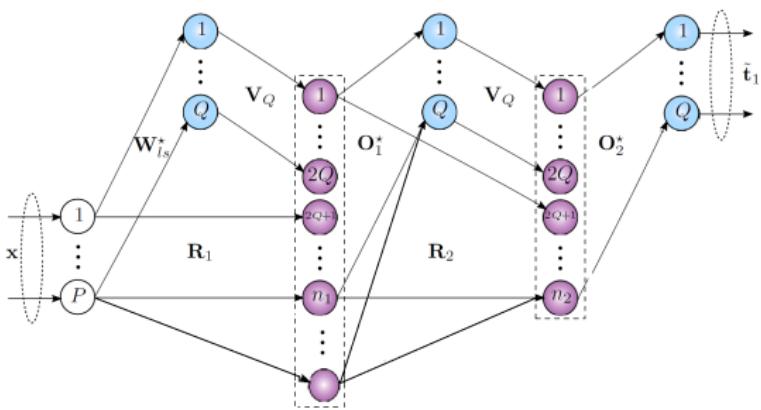
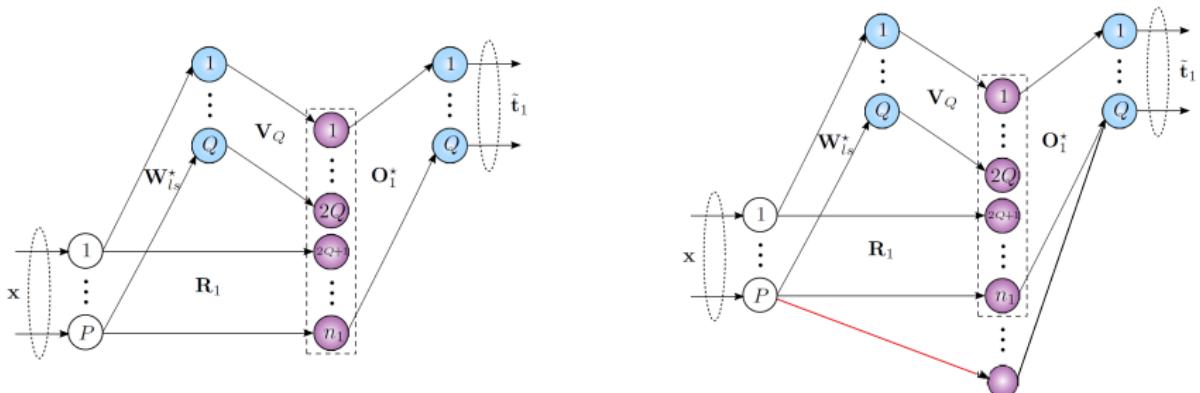
# Multilayer layer (deep) PLN



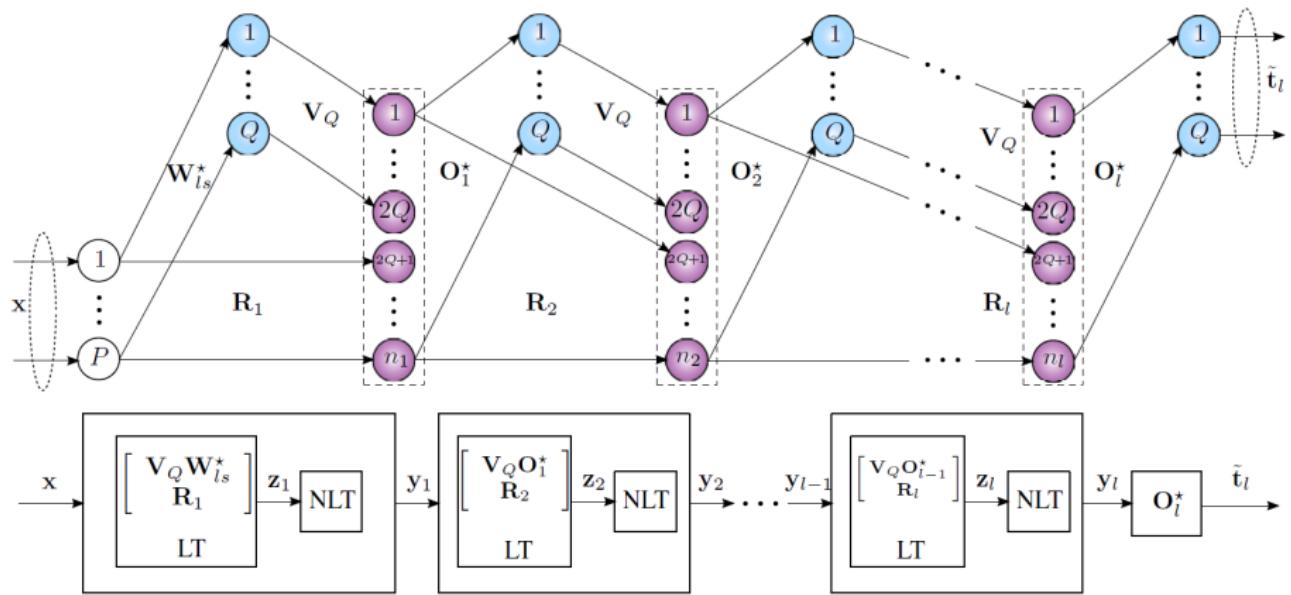
## Multilayer layer (deep) PLN



## Multilayer layer (deep) PLN



## Multilayer layer (deep) PLN



# Multilayer layer (deep) PLN

- $\mathbf{O}_l^* = \arg \min_{\mathbf{O}_l} \sum_j \|\mathbf{t}^{(j)} - \mathbf{O}_l \mathbf{y}_l^{(j)}\|_p^p$  such that  $\|\mathbf{O}_l\|_q^q \leq \alpha \|\mathbf{U}_Q\|_q^q$ ,
- $C_l^* = C(\mathbf{O}_l^*) = \sum_j \|\mathbf{t}^{(j)} - \mathbf{O}_l^* \mathbf{y}_l^{(j)}\|_p^p$

## Proposition (Small approximation error)

Using PP and under the technical condition  $\forall l, \mathbf{O}_l^* \neq [\mathbf{U}_Q \ \mathbf{0}]$  where  $\mathbf{0}$  denotes a zero matrix of size  $Q \times (n_l - 2Q)$ , the optimized cost is monotonically decreasing with increase in number of layers, that is  $C_l^* < C_{l-1}^*$ . For a large number of layers, that means when  $l \rightarrow \infty$ , we have  $C_l^* \leq \kappa$  where  $\kappa$  is an arbitrarily small non-negative real scalar.

- ☞ If we increase  $\Delta$  nodes (random nodes) in the  $l$ 'th layer then we have  
 $C_l^*(n_l + \Delta) \leq C_l^*(n_l)$

# Simulation results

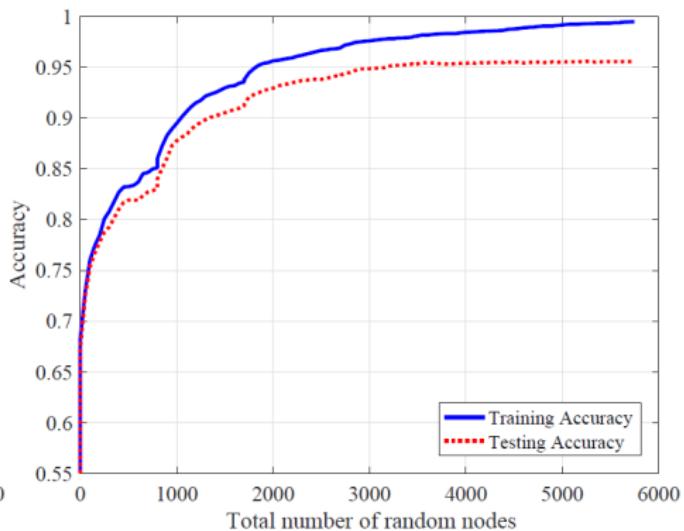
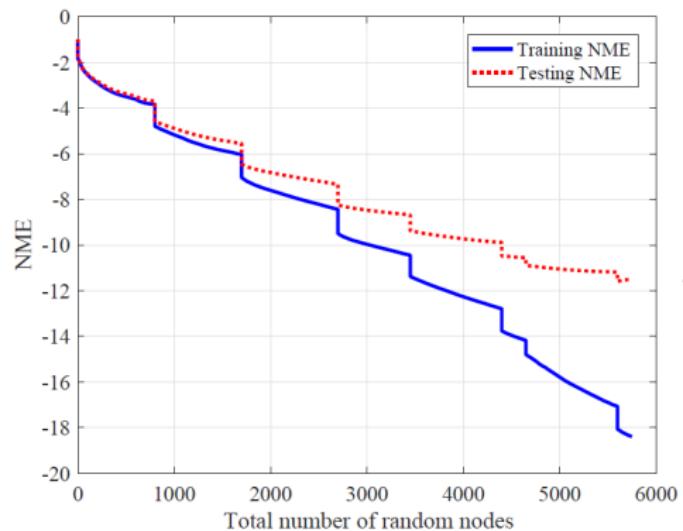
## Classification:

Dataset	Regularized LS				Regularized ELM				PLN			
	Training NME	Testing NME	Test Accuracy	Training Time(s)	Training NME	Testing NME	Test Accuracy	Training Time(s)	Training NME	Testing NME	Test Accuracy	Training Time(s)
Vowel	-1.06	-0.81	28.1 $\pm$ 0.0	0.0035	-6.083	-1.49	53.8 $\pm$ 1.7	0.0549	<b>-72.54</b>	<b>-2.21</b>	<b>60.2</b> $\pm$ 2.4	1.2049
Extended YaleB	-7.51	-4.34	96.9 $\pm$ 0.6	0.0194	-12.75	-6.39	<b>97.8</b> $\pm$ 0.5	0.3908	<b>-49.97</b>	<b>-12.0</b>	97.7 $\pm$ 0.5	2.5776
AR	-3.82	-1.82	96.1 $\pm$ 0.6	0.0297	-9.019	-2.10	97.2 $\pm$ 0.7	0.5150	<b>-35.53</b>	<b>-7.69</b>	<b>97.6</b> $\pm$ 0.6	4.0691
Satimage	-2.82	-2.73	68.1 $\pm$ 0.0	0.0173	-7.614	-5.22	84.6 $\pm$ 0.5	0.8291	<b>-11.73</b>	<b>-7.92</b>	<b>89.9</b> $\pm$ 0.5	1.4825
Scene15	-8.68	-5.03	99.1 $\pm$ 0.2	0.6409	-7.821	-5.78	97.6 $\pm$ 0.3	2.7224	<b>-42.94</b>	<b>-14.7</b>	<b>99.1</b> $\pm$ 0.3	4.1209
Caltech101	-3.22	-1.29	66.3 $\pm$ 0.6	1.1756	-4.784	-1.21	63.4 $\pm$ 0.8	8.1560	<b>-14.66</b>	<b>-4.13</b>	<b>76.1</b> $\pm$ 0.8	5.3712
Letter	-1.00	-0.99	55.0 $\pm$ 0.8	0.0518	-9.217	-6.29	95.7 $\pm$ 0.2	20.987	<b>-18.60</b>	<b>-11.5</b>	<b>95.7</b> $\pm$ 0.2	12.926
NORB	-2.47	-1.54	80.4 $\pm$ 0.0	1.7879	<b>-15.97</b>	-6.77	<b>89.8</b> $\pm$ 0.5	23.207	-13.39	<b>-6.90</b>	86.1 $\pm$ 0.2	10.507
Shuttle	-6.17	-6.31	89.2 $\pm$ 0.0	0.1332	-18.31	-12.2	99.6 $\pm$ 0.1	1.8940	<b>-26.26</b>	<b>-25.0</b>	<b>99.8</b> $\pm$ 0.1	4.6345
MNIST	-4.07	-4.04	85.3 $\pm$ 0.0	0.8122	-9.092	-8.46	<b>96.9</b> $\pm$ 0.1	27.298	<b>-11.42</b>	<b>-10.9</b>	95.7 $\pm$ 0.1	14.181
CIFAR-10	-1.33	-1.33	50.3 $\pm$ 0.0	10.753	-2.004	-2.01	60.3 $\pm$ 0.3	53.842				
CIFAR-100	-0.20	-0.13	14.9 $\pm$ 0.0	12.883								

\* The vowel database is for vowel recognition task (a speech recognition application) and all other databases are for image classification (computer vision applications).

# Simulation results

NME and accuracy versus number of nodes for the “Letter” dataset with 8-layer PLN:



1 Proximal algorithms

2 Sparse representation

3 Dictionary Learning

4 Deep Neural Networks

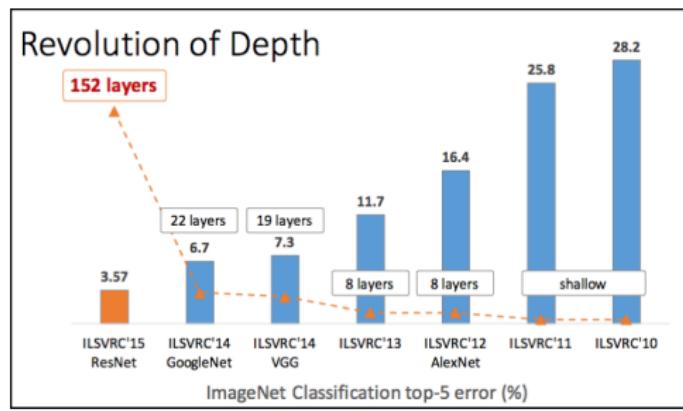
- Background
- Progressive Neural Networks
- Structured Weight Matrices for Neural Networks

5 Conclusions

# Background

Deep networks create deep trouble!

- Deep neural networks have many layers and nodes, with hundreds of millions of parameters → hundreds of megabytes for storage
- Needs more storage and computational resources
- Limiting their application in real-time tasks, and smart phones/ wearable devices



# Background

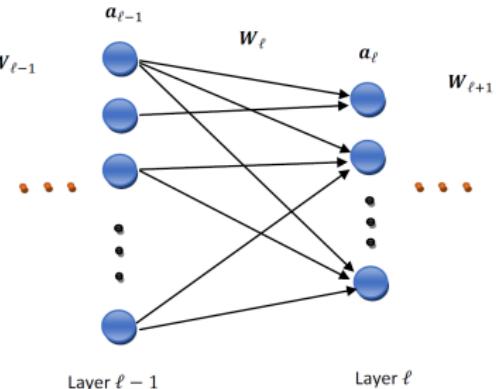
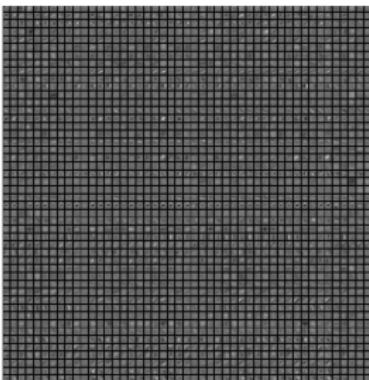
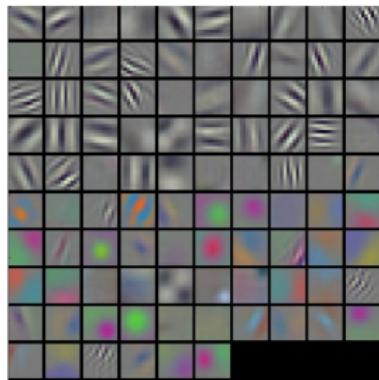
Solutions [see Cheng et al., 2018]:

- **Parameter pruning and sharing:** Reducing redundant parameters that are not sensitive to the performance, using e.g. pruning weak connections, network quantization
- **Low-rank factorization:** Using matrix/tensor decomposition to estimate the informative parameters
- **Transferred/compact convolutional filters:** Designing special structural convolutional filters to save parameters

# Using sparse representation

## Motivation:

Weight vectors and filters corresponding to each node in a neural network exhibit some structure. So, they can be written as sparse linear combinations of e.g. DCT atoms.



Typical-looking filters on the first CONV layer (left), and the 2nd CONV layer (right) of a trained AlexNet.

<http://cs231n.github.io/understanding-cnn/>

Let  $\mathbf{w}_\ell$  be a row of  $\mathbf{W}_\ell$ . Then,

$$\mathbf{w}_\ell = \mathbf{s}_\ell \Phi = \sum_i s_\ell^i \cdot \phi_i$$

and  $\mathbf{s}_\ell$  is sparse.

# Using sparse representation

$$\mathbf{W}_\ell = \mathbf{S}_\ell \Phi$$

- $\Phi$  is a complete (i.e., square matrix) like DCT for images or Gabor dictionary for speech data
- $\mathbf{S}_\ell$  is a matrix with sparse rows

## Advantages:

- ① Low memory consumption, as  $\mathbf{S}_\ell$ 's are sparse and the basis  $\Phi$  is shared among all the layers
- ② Low computational complexity due to the sparseness of  $\mathbf{S}_\ell$ 's and that the multiplications with  $\Phi$  can be done very efficiently for particular transforms like DCT and Fourier
- ③ Preventing overfitting
- ④ Efficient training: SGD + projection

# Using sparse representation

Another sparse structure:

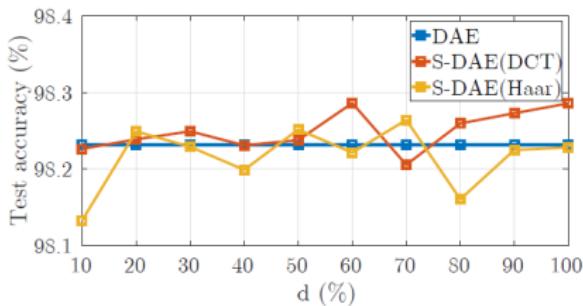
$$\mathbf{W}_\ell = \mathbf{S}_\ell^1 \cdot \mathbf{S}_\ell^2$$

- Both factors  $\mathbf{S}_\ell^1$  and  $\mathbf{S}_\ell^2$  are sparse matrices
- Sparsity is global not row-wise

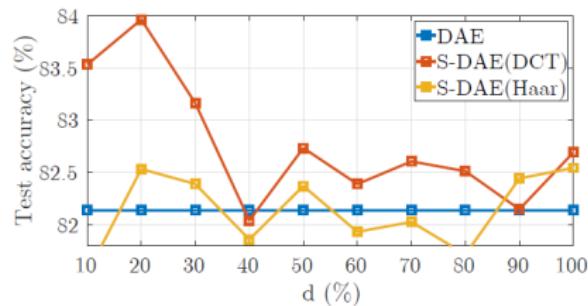
Again, training is easy: SGD (via backpropagation) + projection

# Simulation results

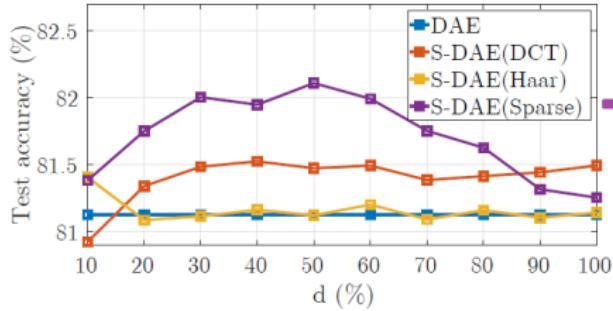
Test accuracy vs the percentage of non-zeros in each row of  $\mathbf{S}$  for a 3-layer structure:  $\mathbf{W}_\ell = \mathbf{S}_\ell \Phi$



(a) MNIST dataset



(b) NORB dataset

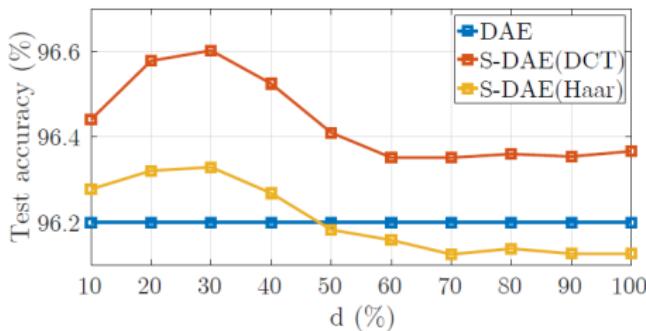


(c) SVHN dataset

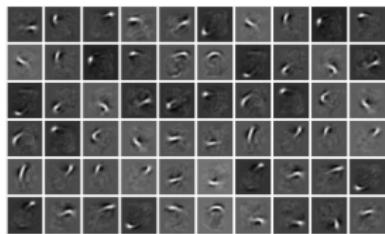
$$\mathbf{W}_\ell = \mathbf{S}_\ell^1 \cdot \mathbf{S}_\ell^2$$

# Simulation results

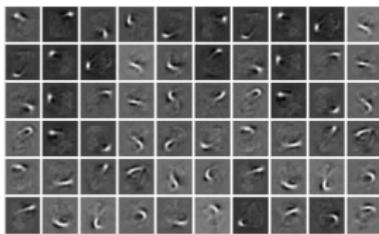
Overfitting effect: MNIST( $50,000 \rightarrow 10,000$  training samples)



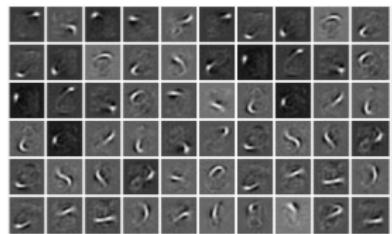
Receptive fields:



(a) KLD method



(b) S-KLD (Haar) method



(c) S-KLD (DCT) method

- 1 Proximal algorithms
- 2 Sparse representation
- 3 Dictionary Learning
- 4 Deep Neural Networks
- 5 Conclusions

# Conclusions

- New algorithms were proposed for sparse recovery and dictionary learning problems based on penalty and ADMM methods combined with proximal algorithms
- The proposed sparse recovery algorithms gave new insights into some previous algorithms like SL0
- Inspired by a progression property, we develop progressive neural networks to learn architecture of neural networks
- Structured weight matrices were proposed using sparse representation to save memory and computation in deep networks

THANK YOU!