# HOW STRUCTURE CAN IMPROVE THE THEORY AND PRACTICE IN NEURAL NETWORKS?

**RAJA GIRYES**

**TEL AVIV UNIVERSITY**

ACCESS workshop, KTH, Stockholm, Sweden

May 14, 2017

# DEEP LEARNING IMPACT



| Model | Top-1 (val) | Top-5 (val) | Top-5 (test) |
|---|---|---|---|
| *SIFT + FVs [7]* | — | — | 26.2% |
| 1 CNN | 40.7% | 18.2% | — |
| 5 CNNs | 38.1% | 16.4% | **16.4%** |
| 1 CNN* | 39.0% | 16.6% | — |
| 7 CNNs* | 36.7% | 15.4% | **15.3%** |

Today we get 3.5% by 152 layers



- Imagenet dataset
- 1,400,000 images
- 1000 categories
- 150000 for testing,
- 50000 for validation

# DEEP NEURAL NETWORKS (DNN)

- One layer of a neural net

$$V \in \mathbb{R}^d \longrightarrow \boxed{X} \xrightarrow{VX} \boxed{\psi} \longrightarrow \psi(VX) \in \mathbb{R}^m$$

$X$ is a linear operation

$F$ is a non-linear function

- Concatenation of the layers creates the whole net

$$\Phi(X^1, X^2, \ldots, X^K) = \psi\left(\psi\left(\psi(VX^1)X^2\right) \ldots X^K\right)$$

$$V \in \mathbb{R}^d \longrightarrow \boxed{X^1} \longrightarrow \boxed{\psi} \dashrightarrow \boxed{X^i} \longrightarrow \boxed{\psi} \dashrightarrow \boxed{X^K} \longrightarrow \boxed{\psi} \longrightarrow$$

# CONVOLUTIONAL NEURAL NETWORKS (CNN)

$$V \in \mathbb{R}^d \longrightarrow \boxed{X} \xrightarrow{VX} \boxed{\psi} \longrightarrow \psi(VX) \in \mathbb{R}^m$$

$X$ is a linear operation

$F$ is a non-linear function

- In many cases, $X$ is selected to be a convolution.
- This operator is shift invariant.
- CNN are commonly used with images as they are typically shift invariant.

# THE NON-LINEAR PART

- Usually $\psi = g \circ f$.

$$X \rightarrow \psi \rightarrow$$

- $f$ is the (point-wise) activation function

ReLU
$f(\mathrm{x}) = \max(\mathrm{x}, 0)$

Sigmoid
$f(x) = \dfrac{1}{1 + e^{-x}}$

Hyperbolic tangent
$f(x) = \tanh(x)$

- $g$ is a pooling or an aggregation operator.

$$V_1 \quad V_2 \quad V_3 \quad V_4 \quad \ldots \quad \ldots \quad \ldots \quad \ldots \quad V_r$$

Max pooling
$$\max_i V_i$$

Mean pooling
$$\frac{1}{n} \sum_{i=1}^{n} V_i$$

$l_p$ pooling
$$\sqrt[p]{\sum_{i=1}^{n} V_i^p}$$

# WHY DNN WORK?

What is so special with the DNN structure?

What is the role of the depth of DNN?

What is the capability of DNN?

How many training samples do we need?

What is the role of the activation function?

What is the role of pooling?

What happens to the data throughout the layers?

# SAMPLE OF RELATED EXISTING THEORY

- Universal approximation for any measurable Borel functions [Hornik et. al., 1989, Cybenko 1989]

- Depth of a network provides an exponential complexity compared to the number parameters [Montúfar et al. 2014], invariance to more complex deformations [Bruna & Mallat, 2013] and better modeling of correlations of the input [Cohen et al. 2016]

- Number of training samples scales as the number of parameters [Shalev-Shwartz & Ben-David 2014] or the norm of the weights in the DNN [Neyshabur et al. 2015]

- Pooling stage provides shift invariance [Bruna et al. 2013]

- Relation of pooling and phase retrieval [Bruna et al. 2014]

- Deeper networks have more local minima that are close to the global one and less saddle points [Saxe et al. 2014], [Dauphin et al. 2014], [Choromanska et al. 2015], [Haeffele & Vidal, 2015]

- Relation to dictionary learning [Papayan et al. 2016].

- DNN with Random Gaussian weights are good for classifying the average points in the data and an important goal of training is to classify the boundary points between the different classes in the data. [Giryes et al. 2016]

# OUTLINE

DNN may solve optimization problems

Generalization error depends on the DNN input margin

# Generalization Error

DNN may solve optimization problems

Generalization error depends on the DNN input margin

# NETWORK STRUCTURE



**Two Classes** $\in Y$

$X^1 \rightarrow \psi \rightarrow X^i \rightarrow \psi \rightarrow X^K \rightarrow \psi$

general non-linearity (ReLU, pooling,…)

softmax/ linear classifier

$w$

$w^T \Phi(X^1, X^2, \ldots, X^K) = 0$

Input Space

Feature Space

Class 1
Class 2

# GENERALIZATION ERROR (GE)

- In training, we reduce the classification error $\ell_{\text{training}}$ of the training data as the number of training examples $L$ increases.

- However, we are interested to reduce the error $\ell_{\text{test}}$ of the (unknown) testing data as $L$ increases.

- The difference between the two is the generalization error

$$\text{GE} = \ell_{\text{training}} - \ell_{\text{test}}$$

➡️ It is important to understand the GE of DNN

# REGULARIZATION TECHNIQUES

- Weight decay – penalizing DNN weights [Krogh & Hertz, 1992].
- Dropout - randomly drop units (along with their connections) from the neural network during training [Hinton et al., 2012], [Baldi & Sadowski, 2013], Srivastava et al., 2014].
- DropConnect – dropout extension [Wan et al., 2013]
- Batch normalization [Ioffe & Szegedy, 2015].
- Stochastic gradient descent (SGD) [Hardt, Recht & Singer, 2016].
- Path-SGD [Neyshabur et al., 2015].
- And more [Rifai et al., 2011], [Salimans & Kingma, 2016], [Sun et al, 2016].

# A SAMPLE OF GE BOUNDS

- Using the VC dimension it can be shown that

$$\text{GE} \leq O\left(\sqrt{\text{DNN params} \cdot \frac{\log(L)}{L}}\right)$$

L is the number of training samples

[Shalev-Shwartz and Ben-David, 2014].

- The GE was bounded also by the DNN weights

$$\text{GE} \leq \frac{1}{\sqrt{L}} 2^K \|w\|_2 \prod_i \|X^i\|_{2,2}$$

[Neyshabur et al., 2015].

# A SAMPLE OF GE BOUNDS

- Using the VC dimension it can be shown that

$$\text{GE} \leq O\left(\sqrt{\textcolor{red}{\text{DNN params}} \cdot \frac{\log(L)}{L}}\right)$$

[Shalev-Shwartz and Ben-David, 2014].

- The GE was bounded also by the DNN weights

$$\text{GE} \leq \frac{1}{\sqrt{L}} \textcolor{red}{2^{K-1}} \|w\|_2 \prod_i \left\|X^i\right\|_F$$

[Neyshabur et al., 2015].

- Note that in both cases the GE grows with the depth

# RETHINKING GENERALIZATION

- Networks with the same architecture may generalize well with structured data but overfit if the data is given with random labels [Zhang et al., 2017].

- This phenomena is affected by explicit regularization.

- This shows that taking into account only the network structure for bouding the generalization error is misleading

- We need to seek an alternative to the Rademacher Complexity and VC-dimension based bounds

# DNN INPUT MARGIN

- Theorem 6: If for every input margin $\gamma_{in}(V^i) > \gamma$

  then $$GE \leq \sqrt{N_{\gamma/2}(\Upsilon)}/\sqrt{L}$$ [Sokolic, Giryes, Sapiro, Rodrigues, 2017]

- $N_{\gamma/2}(\Upsilon)$ is the covering number of the data $\Upsilon$.

- $N_{\gamma/2}(\Upsilon)$ gets smaller as $\gamma$ gets larger.

- Bound is independent of depth.

- Our theory relies on the robustness framework [Xu & Mannor, 2012].

# INPUT MARGIN BOUND

- Maximizing the input margin directly is hard
- Our strategy: relate the input margin to the output margin $\gamma_{out}(V^i)$ and other DNN properties
- Theorem 7:

$$\gamma_{in}(V^i) \geq \frac{\gamma_{out}(V^i)}{\sup\limits_{V \in \Upsilon}\left\|\frac{V}{\|V\|_2}J(V)\right\|_2}$$

$$\geq \frac{\gamma_{out}(V^i)}{\prod_{1 \leq i \leq K}\|X^i\|_2}$$

$$\geq \frac{\gamma_{out}(V^i)}{\prod_{1 \leq i \leq K}\|X^i\|_F}$$

[Sokolic, Giryes, Sapiro, Rodrigues, 2017]



Class 1
Class 2

$\gamma_{out}(V^i)$

$\Phi(V^i)$

# OUTPUT MARGIN

- **Theorem 7:** $\gamma_{in}(V^i) \geq \dfrac{\textcolor{red}{\gamma_{out}(V^i)}}{\sup\limits_{V \in \Upsilon}\left\|\frac{V}{\|V\|_2}J(V)\right\|_2}$

$\geq \dfrac{\textcolor{red}{\gamma_{out}(V^i)}}{\prod_{1 \leq i \leq K}\|X^i\|_2} \geq \dfrac{\textcolor{red}{\gamma_{out}(V^i)}}{\prod_{1 \leq i \leq K}\|X^i\|_F}$

- Output margin is easier to maximize – SVM problem

- Maximized by many cost functions, e.g., hinge loss.



Class 1
Class 2

$\gamma_{out}(V^i)$

$\Phi(V^i)$

# GE AND WEIGHT DECAY

- Theorem 7:   $\gamma_{in}(V^i) \geq \dfrac{\gamma_{out}(V^i)}{\sup\limits_{V \in \Upsilon}\left\|\frac{V}{\|V\|_2}J(V)\right\|_2} \geq \dfrac{\gamma_{out}(V^i)}{\prod_{1 \leq i \leq K}\|X^i\|_2}$

$\geq \dfrac{\gamma_{out}(V^i)}{\prod_{1 \leq i \leq K}\|X^i\|_F}$



Class 1
Class 2

$\gamma_{out}(V^i)$
$\Phi(V^i)$

- Bounding the weights increases the input margin

- Weight decay regularization decreases the GE

- Related to regularization used by [Haeffele & Vidal, 2015]

# JACOBIAN BASED REGULARIZATION

- Theorem 7: $\gamma_{in}(V^i) \geq \dfrac{\gamma_{out}(V^i)}{\sup\limits_{V \in \Upsilon}\left\|\frac{V}{\|V\|_2}{\color{red}J(V)}\right\|_2} \geq \dfrac{\gamma_{out}(V^i)}{\prod_{1 \leq i \leq K}\|X^i\|_2}$

$\geq \dfrac{\gamma_{out}(V^i)}{\prod_{1 \leq i \leq K}\|X^i\|_F}$

- $J(V)$ is the Jacobian of the DNN at point $V$.

- $J(\cdot)$ is piecewise constant.

- Using the Jacobian of the DNN leads to a better bound.

➡ New regularization technique.

# RESULTS

- Better performance with less training samples

MNIST
Dataset

| loss | # layers | 256 samples | | | 512 samples | | | 1024 samples | | |
|------|----------|-------------|------|------|-------------|------|------|--------------|------|------|
| | | no reg. | WD | LM | no reg. | WD | LM | no reg. | WD | LM |
| hinge | 2 | 88.37 | 89.88 | **93.83** | 93.99 | 94.62 | 95.49 | 95.79 | 96.57 | 97.45 |
| hinge | 3 | 87.22 | 89.31 | 93.22 | 93.41 | 93.97 | **95.76** | 95.46 | 96.45 | **97.60** |
| CCE | 2 | 88.45 | 88.45 | 92.77 | 92.29 | 93.14 | 95.25 | 95.38 | 95.79 | 96.89 |
| CCE | 3 | 89.05 | 89.05 | 93.10 | 91.81 | 93.02 | 95.32 | 95.11 | 95.86 | 97.14 |

- CCE: the categorical cross entropy.

[Sokolic, Giryes, Sapiro, Rodrigues, 2017]

- WD: weight decay regularization.

- LM: Jacobian based regularization for large margin.

- Note that hinge loss generalizes better than CCE and that LM is better than WD as predicted by our theory.

# INVARIANCE

- Our theory extends also to study of the relation between invariance in the data and invariance in the network

- We have proposed also a new strategy to enforce invariance in the network [Sokolic, Giryes, Sapiro, Rodrigues, 2017]

# Minimization by DNN

**DNN may solve optimization problems**

**Generalization error depends on the DNN input margin**

# INVERSE PROBLEMS

- We are given $V = ZA + E$

  Given set of measurements

  Unknown signal

  linear operator

  noise

- Standard technique for recovery

$$\min_{Z} \|V - ZA\|_2 \qquad \text{s.t.} \qquad Z \in \Upsilon$$

$Z$ resides in a low dimensional set $\Upsilon$

- Unconstrained form

$$\min_{Z} \|V - ZA\|_2^2 + \lambda f(Z)$$

Regularization parameter

$f$ is a penalty function

# $\ell_1$ MINIMIZATION CASE

- Unconstrained form
$$\min_{Z} \|V - ZA\|_2^2 + \lambda\|Z\|_1$$

- Can be solved by iterative shrinkage and thresholding technique (ISTA)
$$Z^{t+1} = \psi_{\lambda\mu}\left(Z^t + \mu(V - Z^t A)A^T\right)$$

Soft thresholding operation

$-\lambda\mu$    $\lambda\mu$

$\mu$ is the step size

# ISTA CONVERGENCE

- Reconstruction mean squared error (MSE) as a function of the number of iterations

# *LISTA*

- ISTA

$$Z^{t+1} = \psi_{\lambda\mu}\big(Z^t + \mu(V - Z^t A)A^T\big)$$

- Rewriting ISTA:

$$Z^{t+1} = \psi_{\lambda\mu}\big(Z^t\big(I - \mu A A^T\big) + \mu V A^T\big)$$

- Learned ISTA (LISTA):

$$Z^{t+1} = \psi_{\lambda}\big(Z^t X + V S\big)$$

Learned operators

# LISTA CONVERGENCE

- Replacing $I - \mu AA^T$ and $\mu A^T$ in ISTA with the learned $X$ and $S$ improves convergence [Gregor & LeCun, 2010]



- Extensions to other models [Sprechmann, Bronstein & Sapiro, 2015], [Remez, Litani & Bronstein, 2015], [Tompson, Schlachter, Sprechmann & Perlin, 2016].

# LISTA AS A NEURAL NETWORK

Linear operators

$$S$$

$$V \in \mathbb{R}^d$$

$$X$$

$$+$$

$$\psi$$

$$\hat{Z}$$

$$V = ZA + E$$

linear operator

noise

$$Z \in \Upsilon$$

$\psi$ is a non-linear operation

An estimate of $Z$

# $\ell_0$-MINIMIZATION

Iterative hard thresholding algorithm (IHT)

$$I - \mu AA^T$$

$\mu$ is the step size

$$V \in \mathbb{R}^d \longrightarrow \boxed{\mu A^T} \longrightarrow + \longrightarrow \boxed{\psi} \longrightarrow \hat{Z}$$

$$V = ZA + E$$

**Z is a k–sparse vecotr**

$\psi$ is the hard thresholding operation: keeps the largest $k$ entries

A $k$-sparse estimate of $Z$. Aim at solving

$$\min_{\tilde{Z}} \lVert V - \tilde{Z}A \rVert$$

$$s.t \; \lVert \tilde{Z} \rVert_0 \leq k$$

[Blumensath & Davies, 2009]

# $\ell_1$-MINIMIZATION

Projected gradient descent algorithm for $\ell_1$ minimization

$$\boxed{I - \mu A A^T}$$

$\mu$ is the step size

$V \in \mathbb{R}^d$ $\longrightarrow$ $\boxed{\mu A^T}$ $\longrightarrow$ $\bigoplus$ $\longrightarrow$ $\boxed{\psi}$ $\longrightarrow$ $\widehat{Z}$

$V = ZA + E$

$\|Z\|_1 \leq R$

$\psi$ projects onto the $\ell_1$ ball

Estimate of $Z$. Aim at solving

$$\min_{\tilde{Z}} \left\| V - \tilde{Z} A \right\|$$
$$s.t \ \left\| \tilde{Z} \right\|_1 \leq R$$

$R$

$-R$ $\qquad$ $R$

$-R$

# UNCONSTRAINED $\ell_1$-MINIMIZATION

Iterative soft thresholding algorithm (ISTA)

$$I - \mu A A^T$$

$\mu$ is the step size

$$V \in \mathbb{R}^d \longrightarrow \boxed{\mu A^T} \longrightarrow \oplus \longrightarrow \boxed{\psi} \longrightarrow \hat{Z}$$

Soft thresholding operation

Step size $\mu$ obeys
$$\frac{1}{\mu} \geq \|A\|$$

$- \lambda\mu \qquad \lambda\mu$

Minimizer of
$$\min_{\widetilde{Z}} \|V - \widetilde{Z}A\| + \lambda\|\widetilde{Z}\|_1$$

$$V = ZA + E$$

[Daubechies, Defrise & Mol, 2004], [Beck & Teboulle, 2009]

# ISTA CONVERGENCE

- Reconstruction mean squared error (MSE) as a function of the number of iterations

# LEARNED ISTA (LISTA)

Learned linear operators

$V \in \mathbb{R}^d$

$V = ZA + E$

$Z \in \Upsilon$

$S$

$X$

$+$

$\psi$

$\hat{Z}$

Soft thresholding operation

$-\lambda$      $\lambda$

An estimate of $Z$

[Gregor & LeCun, 2010]

34

# LISTA CONVERGENCE

- Replacing $I - \mu AA^T$ and $\mu A^T$ in ISTA with the learned $X$ and $S$ improves convergence [Gregor & LeCun, 2010]



- Extensions to other models [Sprechmann, Bronstein & Sapiro, 2015], [Remez, Litani & Bronstein, 2015], [Tompson, Schlachter, Sprechmann & Perlin, 2016].

# PROJECTED GRADIENT DESCENT (PGD)

$$I - \mu AA^T$$

$\mu$ is the step size

$$V \in \mathbb{R}^d \rightarrow \mu A^T \rightarrow + \rightarrow \psi \rightarrow \widehat{Z}$$

$$V = ZA + E$$

$$f(Z) \leq R$$

$\psi$ projects onto the set $\Upsilon$

$$f(\widetilde{Z}) \leq R$$

Estimate of $Z$. Aim at solving

$$\min_{\widetilde{Z}} \|V - \widetilde{Z}A\|$$

$$s.t. \ f(\widetilde{Z}) \leq R$$

- Theorem 8: Let $Z \in \mathbb{R}^d$, $f: \mathbb{R}^d \to \mathbb{R}$ a proper function, $f(Z) \leq R$, $C_f(Z)$ the tangent cone of $f$ at point $x$, $A \in \mathbb{R}^{d \times m}$ a random Gaussian matrix and $V = ZA + E$. Then the estimate of PGD at iteration $t$, $\hat{Z}^t$, obeys

$$\left\| \hat{Z}^t - Z \right\| \leq \left( \kappa_f \rho \right)^t \|Z\|,$$

where $\rho = \sup_{U, W \in C_f(Z) \cap \mathcal{B}^d} U\left( I - \mu A A^T \right) W^T$

and $\kappa_f = 1$ if $f$ is convex and $\kappa_f = 2$ otherwise. [Oymak, Recht & Soltanolkotabi, 2016].

# PGD CONVERGENCE RATE

- $\rho = \sup\limits_{U,W \in C_f(Z) \cap \mathcal{B}^d} U\left(I - \mu AA^T\right)W^T$ is the convergence rate of PGD.

- Let $\omega$ be the Gaussian mean width of $C_f(Z) \cap \mathcal{B}^d$.

- If $\mu = \frac{1}{\left(\sqrt{m}+\sqrt{d}\right)^2} \simeq \frac{1}{d}$ then $\rho = 1 - O\left(\frac{\sqrt{m}-\omega}{m+d}\right)$.

- If $\mu = \frac{1}{m}$ then $\rho = O\left(\frac{\omega}{\sqrt{m}}\right)$.

- For the $k$-sparse model $\omega^2 = O\left(k\log(\mathrm{d})\right)$

- For GMM with $k$ Gaussians $\omega^2 = O(k)$.

- How may we cause $\omega$ to become smaller for having a better convergence rate?

# INACCURATE PROJECTION

- PGD iterations projects onto $\Upsilon = \{\tilde{Z}: f(\tilde{Z}) \leq R\}$.

- Smaller $\Upsilon \Rightarrow$ Smaller $\omega$.

$\Rightarrow$ Faster convergence as
$\rho = 1 - O\left(\frac{\sqrt{m}-\omega}{m+d}\right)$ or $O\left(\frac{\omega}{\sqrt{m}}\right)$

$\Upsilon$

$f(\widetilde{\mathbf{Z}}) \leq R$

- Let us assume that our signal belongs to a smaller set $\widehat{\Upsilon} = \{\tilde{Z}: \hat{f}(\tilde{Z}) \leq R\}$ with $\widehat{\omega} \ll \omega$.

- Ideally, we would like to project onto $\widehat{\Upsilon}$ instead of $\Upsilon$.

$\widehat{\Upsilon}$

$\hat{f}(\widetilde{\mathbf{Z}}) \leq R$

- This will lead to faster convergence.

- What if such a projection is not feasible?

# INACCURATE PROJECTION

- We will estimate the projection onto $\widehat{\Upsilon}$ by
    - A linear projection $P$
    - Followed by a projection onto $\Upsilon$

- Assumptions:
    - $\|\mathscr{P}_{\Upsilon}(ZP) - Z\| \leq \epsilon$

$$\widehat{f}(\tilde{Z}) \leq R$$

Projection of the target vector $Z$
onto P and then onto $\Upsilon$

# INACCURATE PGD (IPGD)

$$(I - \mu AA^T)P$$

$\mu$ is the step size

$$V \in \mathbb{R}^d \longrightarrow \mu A^T P \longrightarrow + \longrightarrow \psi \longrightarrow \widehat{Z}$$

$$V = ZA + E$$

$\widehat{\Upsilon}$

$$\hat{f}(Z) \leq R$$

$\psi$ projects onto the set $\Upsilon$

$$f(Z) \leq R$$

Estimate of $Z$. Aim at solving

$$\min_{\widetilde{Z}} \|V - \widetilde{Z}A\|$$

$$s.t. \; \hat{f}(\widetilde{Z}) \leq R$$

# THEORY FOR IPGD

- Theorem 9: Let $Z \in \mathbb{R}^d$, $f: \mathbb{R}^d \to \mathbb{R}$ a proper convex* function, $f(Z) \leq R$, $\hat{C}_f(Z)$ the tangent cone of $f$ at point $Z$, $A \in \mathbb{R}^{d \times m}$ a random Gaussian matrix and $V = ZA + E$. Then the estimate of IPGD at iteration $t$, $\hat{Z}^t$, obeys

$$\left\| \hat{Z}^t - Z \right\| \leq \left( (\rho_P)^t + \frac{1 - (\rho_P)^t}{1 - \rho_P} \tilde{\epsilon} \right) \|Z\|,$$

where $\rho_p = \sup_{U,W \in C_f(Z) \cap \mathcal{B}^d} UP\left( I - \mu AA^T \right) PW^T$

and $\tilde{\epsilon} = (2 + \rho_p)\epsilon$.
[Giryes, Eldar, Bronstein & Sapiro, 2016]

*We have a version of this theorem also when $f$ is non-proper or non-convex function

# CONVERGENCE RATE COMPARISON

- PGD convergence:

$$(\rho)^t$$

- IPGD convergence:

$$(\rho_P)^t + \frac{1 - (\rho_P)^t}{1 - \rho_P}(2 + \rho_p)\epsilon$$

$$\underset{(a)}{\cong} (\rho_P)^t + \epsilon \underset{(b)}{\cong} (\rho_P)^t \underset{(c)}{\ll} (\rho)^t$$

(a) $\epsilon$ is negligible compared to $\rho_P$

(b) For small values of $t$ (early iterations).

(c) Faster convergence as $\rho_P \ll \rho$ (because $\omega_p \ll \omega$).

# MODEL BASED COMPRESSED SENSING

- $\widehat{\Upsilon}$ is the set of sparse vectors with sparsity patterns that obey a tree structure.

- Projecting onto $\widehat{\Upsilon}$ improves convergence rate compared to projecting onto the set of sparse vectors $\Upsilon$ [Baraniuk et al., 2010].

- The projection onto $\widehat{\Upsilon}$ is more demanding than onto $\Upsilon$.

- Note that the probability of selecting atoms from lower tree levels is smaller than upper ones.

- $P$ will be a projection onto certain tree levels – zeroing the values at lower levels.

# MODEL BASED COMPRESSED SENSING



Non-zeros picked entries has zero mean random Gaussian distribution with variance:
- 1 at first two levels
- $0.5^2$ at the third level
- $0.2^2$ at the rest of the levels

# SPECTRAL COMPRESSED SENSING

- $\widehat{Y}$ is the set of vectors with sparse representation in a 2-times redundant DCT dictionary such that:



- We set $P$ to be a pooling-like operation that keeps in each window of size 3 only the largest value.

# SPECTRAL COMPRESSED SENSING

# SPECTRAL COMPRESSED SENSING

- $\widehat{Y}$ is the set of vectors with sparse representation in a 4-times redundant DCT dictionary such that:



- We set $P$ to be a pooling-like operation that keeps in each window of size 5 only the largest value.

# SPECTRAL COMPRESSED SENSING

# LEARNING THE PROJECTION

- If we have no explicit information about $\widehat{\Upsilon}$ it might be desirable to learn the projection.

- Instead of learning $P$, it is possible to replace $(I - \mu A A^T)P$ and $\mu A^T P$ with two learned matrices $S$ and $X$ respectively.

- This leads to a very similar scheme to the one of LISTA and provides a theoretical foundation for the success of LISTA.

# LEARNED IPGD

Learned linear operators

$V \in \mathbb{R}^d$

$X$

$+$

$\psi$

$S$

$\widehat{Z}$

$V = ZA + E$

$\widehat{Y}$

$\hat{f}(Z) \leq R$

$\psi$ projects onto the set $Y$

$f(Z) \leq R$

Estimate of $Z$. Aim at solving

$$\min_{\widetilde{Z}} \|V - \widetilde{Z}A\|$$

$$s.t. \ \hat{f}(\widetilde{Z}) \leq R$$

# LISTA

Learned linear operators

$S$

$V \in \mathbb{R}^d$ → $X$ → $+$ → $\psi$ → $\widehat{Z}$

$V = ZA + E$

$\widehat{Y}$

$\hat{f}(Z) \leq R$

$\psi$ is a proximal mapping.

$$\psi(U) = \underset{\widetilde{Z} \in \mathbb{R}^d}{\text{argmin}} \|U - \widetilde{Z}\| + \lambda f(\widetilde{Z})$$

Estimate of $Z$. Aim at solving

$$\underset{\widetilde{Z}}{min} \|V - \widetilde{Z}A\| + \lambda \hat{f}(\widetilde{Z})$$

# LISTA MIXTURE MODEL

- Approximation of the projection onto $\widehat{\Upsilon}$ with one linear projection may not be accurate enough.

- This requires more LISTA layers/iterations.

- Instead, one may use several LISTA networks, where each approximates a different part of $\widehat{\Upsilon}$

- Training multiple LISTA networks accelerate the convergence further.

# LISTA MIXTURE MODEL

# RELATED WORKS

- In [Bruna et al. 2016] it is shown that a learning may give a gain due to better preconditioning of *A.*

- In [Xin et al. 2016] a relation to the restricted isometry property (RIP) is drawn

- In [Borgerding & Schniter, 2016] a connection is drawn to approximate message passing (AMP).

- All these works consider only the sparsity case

# Take Home Message

DNN may solve optimization problems

Generalization error depends on the DNN input margin

# ACKNOWLEDGEMENTS



Guillermo Sapiro
Duke University

Yonina C. Eldar
Technion

Alex M. Bronstein
Technion

Miguel Rodrigues
UCL

Jure Sokolic
UCL

# QUESTIONS?

WEB.ENG.TAU.AC.IL/~RAJA

# GAUSSIAN MEAN WIDTH

- Gaussian mean width:

$$\omega(\Upsilon) = E \sup_{V,W \in \Upsilon} \langle V - W, g \rangle, \qquad g \sim N(0, I).$$

The width of the set $\Upsilon$ in the direction of $g$:

# MEASURE FOR LOW DIMENSIONALITY

- Gaussian mean width:
$$\omega(\Upsilon) = E \sup_{V,W \in \Upsilon} \langle V - W, g \rangle, \qquad g \sim N(0, I).$$

- $\omega^2(\Upsilon)$ is a measure for the dimensionality of the data.

- Examples:

If $\Upsilon \subset \mathbb{B}^d$ is a Gaussian Mixture Model with $k$ Gaussians then
$$\omega^2(\Upsilon) = O(k)$$

If $\Upsilon \subset \mathbb{B}^d$ is a data with $k$-sparse representations then
$$\omega^2(\Upsilon) = O(k \log d)$$