

AMP-inspired Deep Networks

Phil Schniter



THE OHIO STATE UNIVERSITY

Duke
UNIVERSITY



Collaborators: **Sundeeep Rangan** (NYU), **Alyson Fletcher** (UCLA),
Mark Borgerding (OSU)

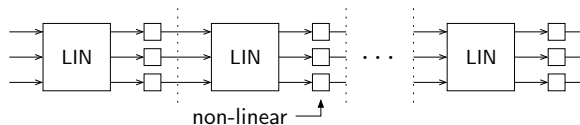
Supported in part by NSF grants IIP-1539960 and CCF-1527162.

Workshop on Information Theory and Applications (ITA) — Feb 2017

Deep Neural Networks

Typical setup:

- Many layers, consisting of (affine) linear stages and scalar nonlinearities.



- Linear stages often constrained (e.g., small convolution kernels).
- Parameters learned by minimizing training error using backprop.

Open questions:

- 1 How should we interpret the learned parameters?
- 2 Can we speed up training?
- 3 Can we design a better network structure?

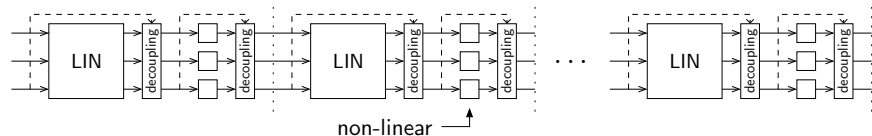
Focus of this talk: Standard Linear Regression

- Consider recovering a vector x from noisy linear observations

$$y = Ax + w,$$

where x is drawn from an iid prior (e.g., sparse¹)

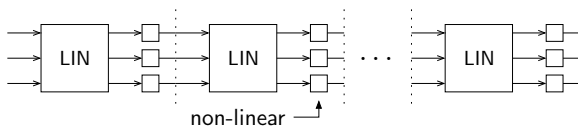
- For this application, we propose a **deep network** that is
 - 1) asymptotically optimal for a large class of A ,
 - 2) interpretable, and
 - 3) easy to train.



¹Gregor/LeCun, Sprechmann/Bronstein/Sapiro, Kamilov/Mansour, Wang/Ling/Huang, Mousavi/Baraniuk, Borgerding/Schniter, etc.

Understanding Deep Networks via Algorithms

- Many **algorithms** have been proposed for high-dimensional inference.
- Often, such algorithms are iterative, where each iteration consists of a linear operation followed by scalar nonlinearities.
- By “**unfolding**” such algorithms, we get deep networks.²



- Can such algorithms help us **design/interpret/train** deep nets?

²Gregor/LeCun, ICML 10.

Algorithmic Approaches to Standard Linear Regression

- Recall goal: recovering/fitting \mathbf{x} from noisy linear observations

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}.$$

- A popular approach is **regularized loss minimization**:

$$\arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \lambda f(\mathbf{x}),$$

where, e.g., $f(\mathbf{x}) = \|\mathbf{x}\|_1$ for the lasso.

- Can also be interpreted as **MAP estimation** of \mathbf{x} under priors

$$\mathbf{x} \sim \exp(-f(\mathbf{x})) \quad \& \quad \mathbf{w} \sim \mathcal{N}(0, \lambda \mathbf{I}).$$

- But often the goal is **minimizing MSE** or inferring marginal posteriors.

High-dimensional MMSE Inference

- High dimensional MMSE inference is difficult in general.
- To simplify things, suppose that
 - 1) x is iid
 - 2) A is large and random.
- The case of iid Gaussian A is well studied, but very restrictive.
- Instead, consider **right-rotationally invariant** A :

$$A = USV^T \text{ with } V \sim \text{Haar and indep of } x.$$

- For this case, the **replica prediction of the MMSE** is³

$$\mathcal{E}(\gamma) = \text{var}\{x|r\}, \quad r = x + \mathcal{N}(0, 1/\gamma), \quad \gamma = R_{A^T A / \sigma^2}(-\mathcal{E}(\gamma))$$

³Tulino/Caire/Verdu/Shamai, IEEE-TIT, 2013.

Achieving MMSE in standard linear regression

- Recently a “**vector approximate message passing**” (VAMP) algorithm has been proposed that iterates linear vector estimation with nonlinear scalar denoising.
- Under large right-rotationally invariant \mathbf{A} and Lipschitz scalar denoisers, VAMP is rigorously characterized by a **scalar state-evolution**.⁴
- Under MMSE scalar denoising, VAMP’s state-evolution fixed-points **agree** with the replica prediction!
- Operating regimes:
 - easy**: VAMP has a **unique** fixed point \Rightarrow attains MMSE
 - hard**: VAMP has **multiple** fixed points \Rightarrow suboptimal
(but so are all other known polynomial-time methods)

⁴Rangan/Schniter/Fletcher, arXiv:1610.03082, 2016.

VAMP for linear regression

initialize \mathbf{r}_1, γ_1

for $t = 0, 1, 2, \dots$

$$\hat{\mathbf{x}}_1 \leftarrow (\mathbf{A}^\top \mathbf{A} / \sigma^2 + \gamma_1 \mathbf{I})^{-1} (\mathbf{A}^\top \mathbf{y} / \sigma^2 + \gamma_1 \mathbf{r}_1) \quad \text{LMMSE}$$

$$\alpha_1 \leftarrow \frac{\gamma_1}{N} \text{Tr} \left[(\mathbf{A}^\top \mathbf{A} / \sigma^2 + \gamma_1 \mathbf{I})^{-1} \right] \quad \text{divergence}$$

$$\mathbf{r}_2 \leftarrow \frac{1}{1 - \alpha_1} (\hat{\mathbf{x}}_1 - \alpha_1 \mathbf{r}_1) \quad \text{Onsager correction}$$

$$\gamma_2 \leftarrow \gamma_1 \frac{1 - \alpha_1}{\alpha_1} \quad \text{precision of } \mathbf{r}_2$$

$$\hat{\mathbf{x}}_2 \leftarrow \mathbf{g}(\mathbf{r}_2; \gamma_2) \quad \text{(scalar) denoising}$$

$$\alpha_2 \leftarrow \frac{1}{N} \text{Tr} \left[\frac{\partial \mathbf{g}}{\partial \mathbf{r}}(\mathbf{r}_2; \gamma_2) \right] \quad \text{divergence}$$

$$\mathbf{r}_1 \leftarrow \frac{1}{1 - \alpha_2} (\hat{\mathbf{x}}_2 - \alpha_2 \mathbf{r}_2) \quad \text{Onsager correction}$$

$$\gamma_1 \leftarrow \gamma_2 \frac{1 - \alpha_2}{\alpha_2} \quad \text{precision of } \mathbf{r}_1$$

end

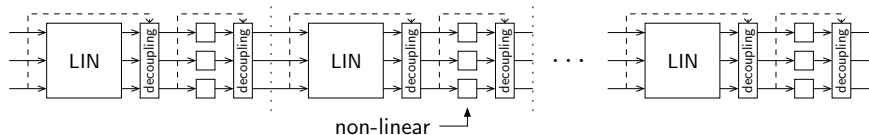
MMSE-VAMP interpreted

```
initialize  $\mathbf{r}_1, \gamma_1$ 
for  $t = 0, 1, 2, \dots$ 
     $\hat{\mathbf{x}}_1 \leftarrow$  MMSE estimate of  $\mathbf{x}$  under
        pseudo-prior  $\mathcal{N}(\mathbf{r}_1, \mathbf{I}/\gamma_1)$  & measurement  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathcal{N}(\mathbf{0}, \sigma^2\mathbf{I})$ 
     $\mathbf{r}_2 \leftarrow$  linear cancellation of  $\mathbf{r}_1$  from  $\hat{\mathbf{x}}_1$ 
     $\hat{\mathbf{x}}_2 \leftarrow$  MMSE estimate of  $\mathbf{x}$  under
        prior  $p(\mathbf{x})$  and pseudo-measurement  $\mathbf{r}_2 = \mathbf{x} + \mathcal{N}(\mathbf{0}, \mathbf{I}/\gamma_2)$ 
     $\mathbf{r}_1 \leftarrow$  linear cancellation of  $\mathbf{r}_2$  from  $\hat{\mathbf{x}}_2$ 
end
```

Linear cancellation “**decouples**” the iterations, so that global MMSE problem can be tackled by solving simpler local MMSE problems.

Unfolding VAMP

Unfolding the VAMP algorithm gives the network⁵



Notice the **two decoupling stages** in each layer.

⁵Borgerding/Schniter'16

Learning the network parameters

After unfolding an algorithm, one can use backpropagation (or similar) to “learn” the optimal network parameters.⁶

- Linear stage: $\hat{x}_1 = Br_1 + Cy$
→ learn (B, C) for each layer.
- Nonlinear stage: $\hat{x}_{2j} = g(r_{2j}) \forall j$
→ learn a scalar function $g(\cdot)$ for each layer.
e.g., spline, piecewise linear, etc.

⁶Gregor/LeCun, ICML 10.

Result of learning

Suppose that the training data $\{\mathbf{y}^{(d)}, \mathbf{x}^{(d)}\}_{d=1}^D$ were constructed using

- 1) iid $x_j^{(d)} \sim p(x)$
- 2) $\mathbf{y}^{(d)} = \mathbf{A}\mathbf{x}^{(d)} + \mathcal{N}(\mathbf{0}, \sigma^2\mathbf{I})$
- 3) right-rotationally invariant \mathbf{A} .

Backpropagation recovers the parameter settings $(\mathbf{B}, \mathbf{C}, f)$ originally prescribed by the VAMP algorithm!

- The learned linear stages are MMSE under pseudo-prior $\mathbf{x} \sim \mathcal{N}(\mathbf{r}_1, \mathbf{I}/\gamma_1)$ & measurement $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathcal{N}(\mathbf{0}, \sigma^2\mathbf{I})$
- The learned scalar nonlinearities are MMSE under prior $p(x_j)$ and pseudo-measurements $r_{2j} = x_j + \mathcal{N}(0, 1/\gamma_2)$

\rightsquigarrow *This deep network is interpretable!*

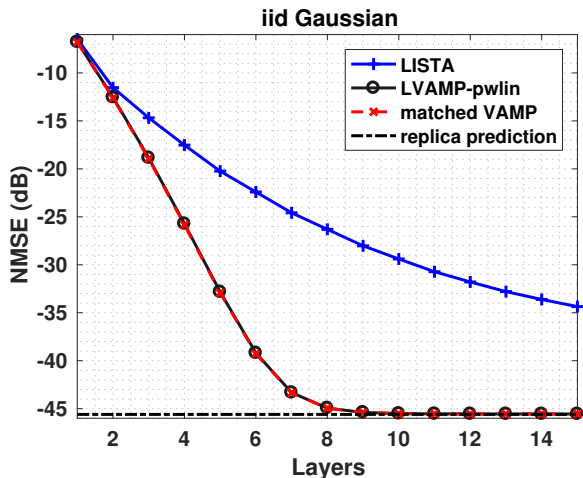
Implications for training

Due to the decoupling stages...

- each linear stage should be **locally** MSE optimal, so
↷ the linear params (\mathbf{B}, \mathbf{C}) can be learned locally in each layer!
- each non-linear stage should be **locally** MSE optimal, so
↷ the nonlinear function $g(\cdot)$ can be learned locally in each layer!

This deep network is easy to train!

Example with iid Gaussian A



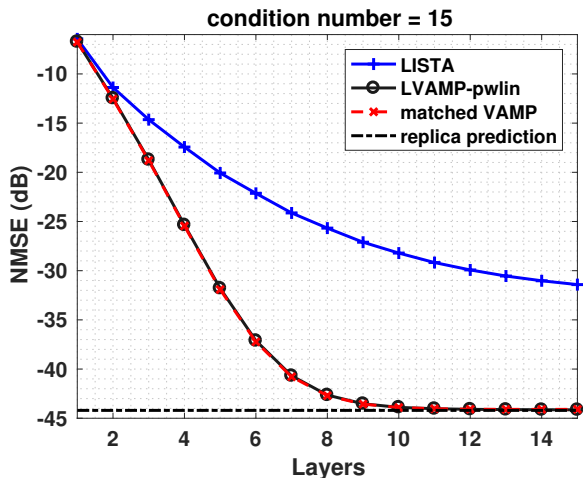
$$n = 1024$$
$$m/n = 0.5$$

$$A \sim \text{iid } \mathcal{N}(0, 1)$$

$$x \sim \text{Bernoulli-Gaussian}$$
$$\Pr\{x \neq 0\} = 0.1$$

$$\text{SNR} = 40 \text{ dB}$$

Example with non-iid Gaussian A



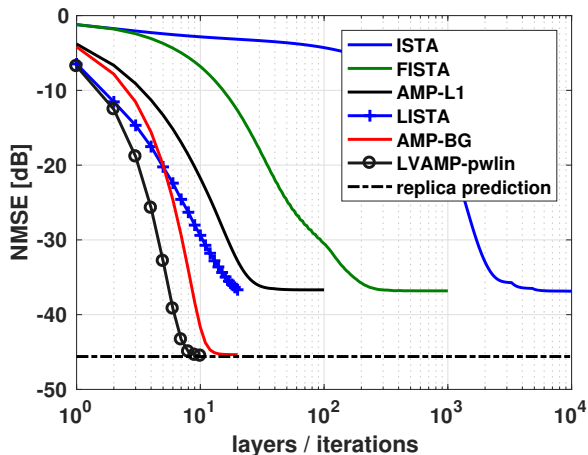
$$n = 1024$$
$$m/n = 0.5$$

$$A = USV^T$$
$$U, V \sim \text{Haar}$$
$$s_n/s_{n-1} = \phi \quad \forall n$$

$$x \sim \text{Bernoulli-Gaussian}$$
$$\Pr\{x \neq 0\} = 0.1$$

$$\text{SNR} = 40 \text{ dB}$$

Deep nets vs algorithms



$n = 1024$
 $m/n = 0.5$

$\mathbf{A} \sim \text{iid } \mathcal{N}(0, 1)$

$x \sim \text{Bernoulli-Gaussian}$
 $\Pr\{x \neq 0\} = 0.1$

SNR = 40 dB

Conclusions

- Our goal is to understand the **design and interpretation of deep nets**.
- For this talk, we restricted our focus to the problem of (e.g., sparse) **linear regression**.
- We proposed a deep net that is
 - 1) **asymptotically MSE-optimal** (for iid \mathbf{x} and RRI \mathbf{A})
 - 2) **interpretable**: ... LMMSE/decoupling/NL-MMSE/decoupling ...
 - 3) **locally trainable**.
- The proposed network is obtained by “**unfolding**” the VAMP algorithm and **learning** its parameters.
- In ongoing work, we are extending these ideas **beyond** linear regression.