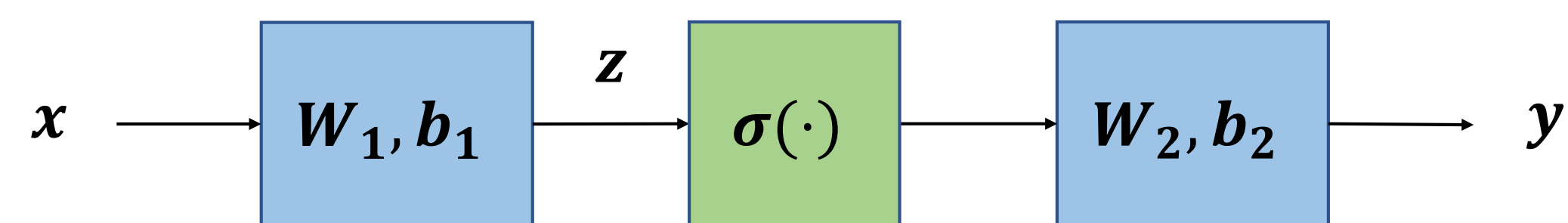


Parthe Pandit, Mojtaba Sahraee-Ardakan, Alyson K. Fletcher  
University of California, Los Angeles

Philip Schniter  
Ohio State University

Sundeeep Rangan  
New York University

## Learning Two-Layer Neural Networks



- Two-layer fully connected neural network
- $z = W_1^T x + b_1, \quad y = W_2^T \sigma(z) + b_2$
- $W_2, b_2$  known
- Goal:** Learn  $W_1, b_1$  from samples  $(x_i, y_i), i = 1, \dots, N$  generated by a ground truth network

## Empirical Risk Minimization

### Regularized empirical risk minimization

- Learn parameters by minimizing regularized empirical risk

$$\min_{W_1, b_1} \phi(W_1) + \sum_{i=1}^N \mathcal{L}(\hat{y}_i, y_i)$$

- $\phi(\cdot)$  is a regularizer
- $\mathcal{L}(\cdot, \cdot)$  is a loss function
- $\hat{y}_i = W_2^T \sigma(z_i) + b_2$  where  $z_i = W_1^T x_i + b_1$

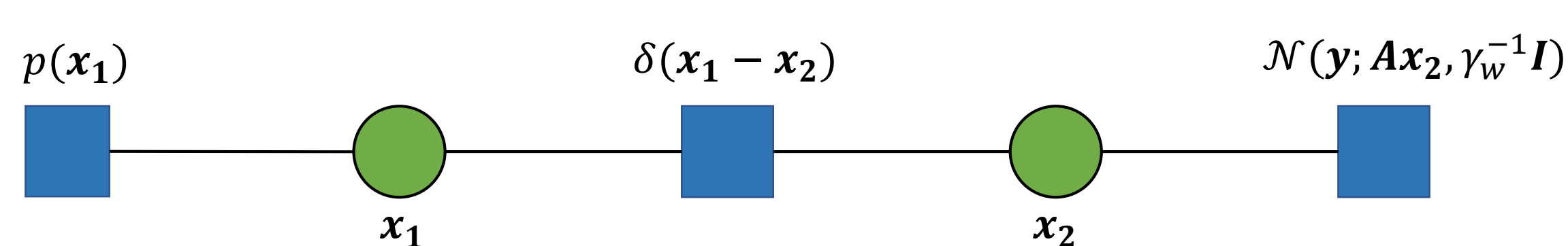
## Main Questions

- Optimization typically performed using some variant of SGD
  - Works well in practice
  - But hard to analyze
  - Error bounds not known to be optimal
  - When does learning succeed?
  - What is the generalization error?

## This Work: Analysis via Message Passing

- Approximate message passing (AMP) algorithms:
  - Efficiently solve inverse problems
  - Originally developed for linear inverse problems
  - Extended to multi-layer networks (ML-VAMP)
- Key property: State evolution provides performance guarantees in high dimensional regime
- Main contributions:
  - Learning parameters of a two-layer network using AMP framework
  - State evolution provides estimation error at each iteration
  - Predicts when learning will or will not work

## Background on Vector AMP



### Linear inverse problem

- Model:  $x \sim p(x), \quad y = Ax + w, \quad w \sim \mathcal{N}(0, \gamma_w^{-1} I)$
- Posterior density factorizes:

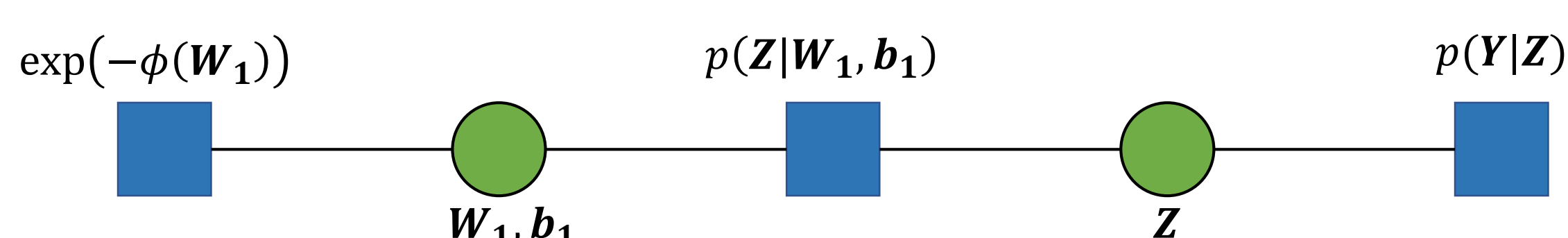
$$p(x|y) \propto \exp\left(-\frac{\gamma_w}{2} \|y - Ax\|^2\right) p(x)$$

- MAP estimation:  $\hat{x} = \operatorname{argmin}_x \frac{\gamma_w}{2} \|y - Ax\|^2 - \log p(x)$
- Apply expectation propagation-like inference over factor graph
  - Gaussian approximation of messages
- In each iteration:
  - Linear Gaussian estimation
  - Estimation with prior; separable for separable priors

### Properties of VAMP

- State evolution
  - Behavior in each iteration exactly explained for large random  $A$
  - Provides MSE in each iteration
- Provably Bayes optimal in certain cases
  - Including non-convex cases
- Relates to ADMM with carefully chosen adaptive step size

### Factor graph of the empirical risk

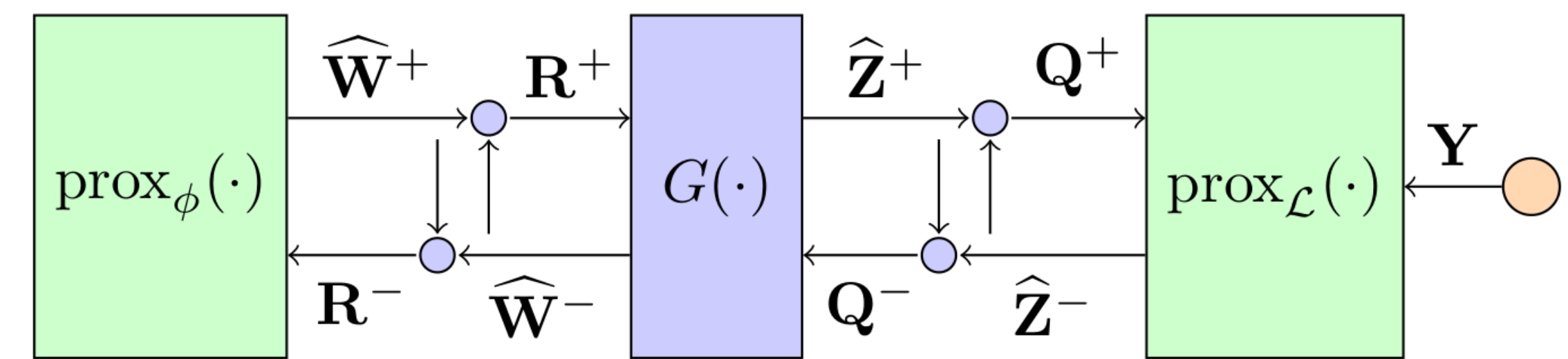


- Apply expectation propagation over this graph to get matrix AMP

### Acknowledgements:

The work of P. Pandit, M. Sahraee-Ardakan, A. K. Fletcher was supported in part by the NSF Grants 1254204 and 1738286, and the Office of Naval Research Grant N00014-15-1-2677. S. Rangan was supported in part by the NSF Grants 1116589, 1302336, and 1547332, NIST, the industrial affiliates of NYU WIRELESS, and the SRC. The work of P. Schniter was supported in part by the National Science Foundation Grant 1716388.

## Matrix AMP



### Algorithm 1 Matrix Approximate Message Passing

**Require:**  $K_{\text{it}}$ ,  $\operatorname{prox}_\phi(\cdot)$ ,  $\operatorname{prox}_\mathcal{L}(\cdot)$ , and linear solver:  $\_, \_ = G(\cdot)$

- Select initial  $R_0^- \in \operatorname{Range}(\mathbf{X})$  and  $\Gamma_0^- \succ 0$ .
- for**  $k = 0, 1, \dots, K_{\text{it}}$  **do**
- // Forward updates
- $\widehat{W}^+ \leftarrow \operatorname{prox}_\phi(R^-, \Gamma^-) = \operatorname{argmin}_W \phi(W) + \frac{1}{2} \operatorname{tr}((W - R^-) \Gamma^- (W - R^-)^\top)$
- $\Phi^+ \leftarrow \langle \nabla_R \operatorname{prox}_\phi(R^-, \Gamma^-) \rangle^{-1} \Gamma^-, \Gamma^+ \leftarrow \Phi^+ - \Gamma^-$
- $R^+ \leftarrow (\widehat{W}^+ \Phi^+ - R^- \Gamma^-) (\Gamma^+)^{-1}$
- $\_, \widehat{Z}^+ \leftarrow G(R^+, Q^-, \Gamma^+, \Upsilon^-)$
- $\Lambda^+ \leftarrow \langle \nabla_Q G(R^+, Q^-, \Gamma^+, \Upsilon^-) \rangle^{-1} \Upsilon^-, \Upsilon^+ \leftarrow \Lambda^+ - \Upsilon^-$
- $Q^+ \leftarrow (\widehat{Z}^+ \Lambda^+ - Q^- \Upsilon^-) (\Upsilon^+)^{-1}$
- // Backward updates
- $\widehat{Z}^- \leftarrow \operatorname{prox}_\mathcal{L}(Q^+, \Upsilon^+) = \operatorname{argmin}_Z \mathcal{L}(Z|Y) + \frac{1}{2} \operatorname{tr}((Z - Q^+) \Upsilon^+ (Z - Q^+)^\top)$
- $\Lambda^- \leftarrow \langle \nabla_Q \operatorname{prox}_\mathcal{L}(Q^+, \Upsilon^+) \rangle^{-1} \Upsilon^+, \Upsilon^- \leftarrow \Lambda^- - \Upsilon^+$
- $Q^- \leftarrow (\widehat{Z}^- \Lambda^- - Q^+ \Upsilon^+) (\Upsilon^-)^{-1}$
- $\widehat{W}^-, \_ \leftarrow G(R^+, Q^-, \Gamma^+, \Upsilon^-)$
- $\Phi^- \leftarrow \langle \nabla_R G(R^+, Q^-, \Gamma^+, \Upsilon^-) \rangle^{-1} \Gamma^+, \Gamma^- \leftarrow \Phi^- - \Gamma^+$
- $R^- \leftarrow (\widehat{W}^- \Phi^- - R^+ \Gamma^+) (\Gamma^-)^{-1}$
- end for**

$$\widehat{W}, \widehat{Z} \leftarrow G(R, Q, \Gamma, \Upsilon) \text{ solves } \mathbf{X}^\top \mathbf{X} \widehat{W} \Upsilon + \widehat{W} \Gamma = \mathbf{X}^\top \mathbf{Q} \Upsilon + \mathbf{R} \Gamma, \quad \text{and} \quad \widehat{Z} = \mathbf{X} \widehat{W}$$

## High-Dimensional Analysis

If a sequence of problems indexed by  $N$  is solved by Algorithm 1, with  $\lim_{N \rightarrow \infty} \frac{N_{\text{in}}}{N} = \beta \in (0, \infty)$

- Data:**  $X(N) = USV^T$  with  $U, V$  are square Haar distributed and  $S_n(N)$  are bounded
- True weights:**  $W^*(N) \in \mathbb{R}^{N_{\text{in}} \times N_{\text{hid}}}$  such that rows are i.i.d. samples of  $\mathcal{W} \in \mathbb{R}^{1 \times N_{\text{hid}}}$
- Prox:** Act row-wise, have symmetric Jacobians which are uniformly Lipschitz

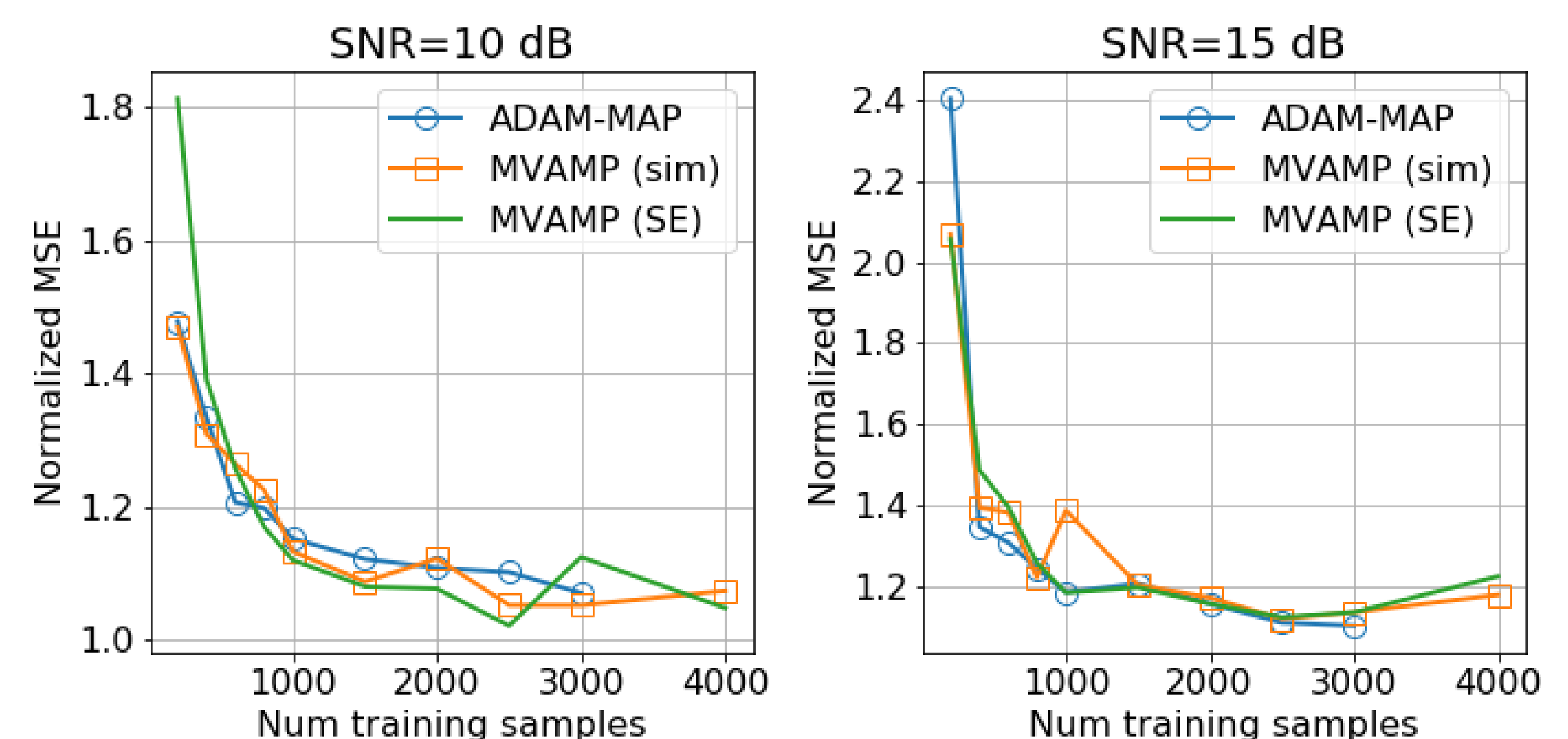
Then there exists deterministic  $\bar{\Gamma}_t$  and  $\bar{\Sigma}_t$  obtained recursively,  $\mathcal{W} \perp Z \sim \mathcal{N}(0, I_{N_{\text{hid}}})$  such that for all iterations  $t$

$$(\widehat{W}_{t,n}^+, W_{n,*}^*) \xrightarrow{d} (\operatorname{prox}_\phi(\mathcal{W} + Z \bar{\Sigma}_t, \bar{\Gamma}_t), \mathcal{W})$$

- Joint distribution of estimate and true vectors converge
- Recursive formula for  $(\widehat{W}_{t,n}^+, W_{n,*}^*)$
- Provides exact prediction of parameter and generation error

## Predicting Performance for a Synthetic Network

- Two-layer ground truth network
  - $N_{\text{in}} = 100, d = 4, N_{\text{out}} = 1$
  - Input: iid Gaussian with variance  $1/N_{\text{in}}$
  - Noise added to output to get different SNR levels
- Key takeaways:
  - ADAM optimizer achieves similar results to Matrix AMP
  - State Evolution accurately predicts performance of Matrix AMP



### References:

- Pandit, P., Sahraee, M., Rangan, S., & Fletcher, A. K. (2019). Asymptotics of map inference in deep networks. IEEE International Symposium on Information Theory (ISIT), 842-846
- Rangan, S., Schniter, P., & Fletcher, A. K. (2019). Vector approximate message passing. IEEE Transactions on Information Theory.
- Bayati, M., & Montanari, A. (2011). The dynamics of message passing on dense graphs, with applications to compressed sensing. IEEE Transactions on Information Theory, 57(2), 764-785.