# Progress and hurdles
# in the statistical mechanics of deep learning

July 23rd 2020

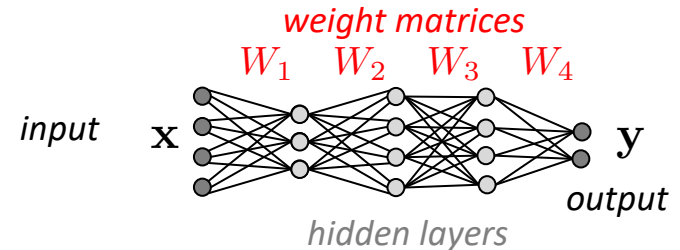Marylou Gabrié (NYU, Flatiron Institute)

**Collaborators:**

Andre Manoel (Owkin), Jean Barbier (ICTP Trieste), Clément Luneau (EPFL),
Nicolas Macris (EPFL), Florent Krzakala (ENS Paris), Lenka Zdeborova (IPHT Saclay)

# Understanding machine learning with deep neural nets

*weight matrices*
$W_1 \quad W_2 \quad W_3 \quad W_4$

▷ **Supervised learning with neural networks**

*input*   $\mathbf{x}$     $\mathbf{y}$   *output*

*hidden layers*

   ▷   training data   $\mathcal{D} = \{\mathbf{y}^{(k)}, \mathbf{x}^{(k)}\}_{k=1}^{P}$

   ▷   fit with class of parametrized functions   $\mathbf{y} = f(\mathbf{W}_L f(\mathbf{W}_{L-1} \ldots f(\mathbf{W}_1 \mathbf{x})))$
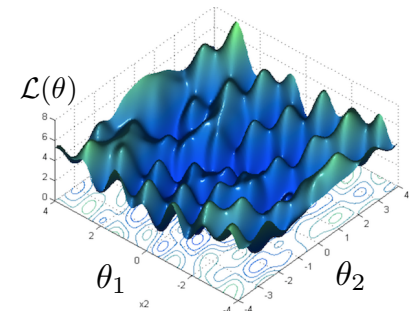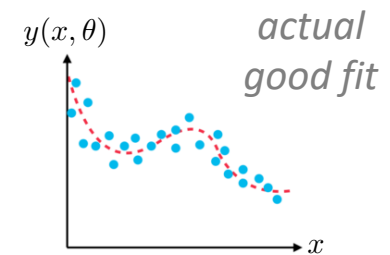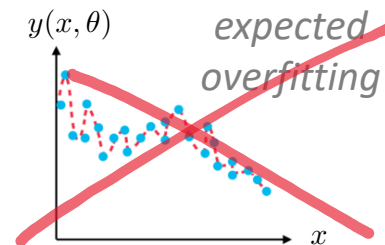
▷ **Impressive performances (automatic vision, natural language processing etc.)**

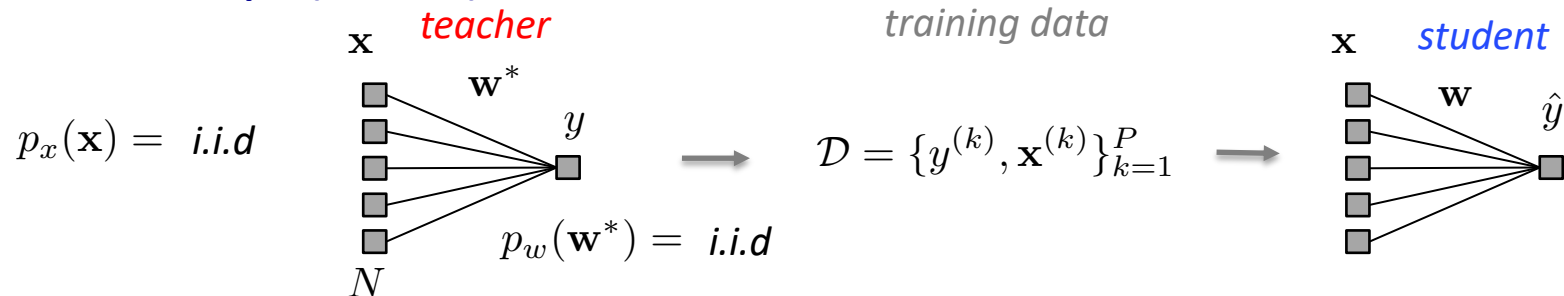▷ **Interesting properties**

   ▷   universal approximators

$y(x,\theta)$    *expected overfitting*

$y(x,\theta)$    *actual good fit*

$x$     $x$

   ▷   not prone to overfitting

   ▷   train with local descent algorithm despite of non-convexity

$\mathcal{L}(\theta)$

$\theta_1 \qquad \theta_2$

# Statistical mechanics of learning, initiated in the 80s

▷ **Focus on simple (solvable) models**



$p_x(\mathbf{x}) =$ *i.i.d*

$p_w(\mathbf{w}^*) =$ *i.i.d*

$\mathcal{D} = \{y^{(k)}, \mathbf{x}^{(k)}\}_{k=1}^{P}$

▷ **Consider the Bayesian posterior statistics** $p(\mathbf{w}|\mathcal{D}) = \dfrac{p(\mathcal{D}|\mathbf{w})\, p_w(\mathbf{w})}{p(\mathcal{D})}$

*e.g. Bayes optimal estimator (minimum mean square error)*

$$\min_{\hat{\mathbf{w}}} \int d\mathbf{w}\,(\mathbf{w} - \hat{\mathbf{w}})^2\, p_S(\mathbf{w}\mathcal{D}) \longrightarrow \hat{\mathbf{w}}_{\mathrm{MMSE}} = \int d\mathbf{w}\,\mathbf{w}\, p_S(\mathbf{w}|\mathcal{D})$$

▷ **The thermodynamic limit = infinitely large model**
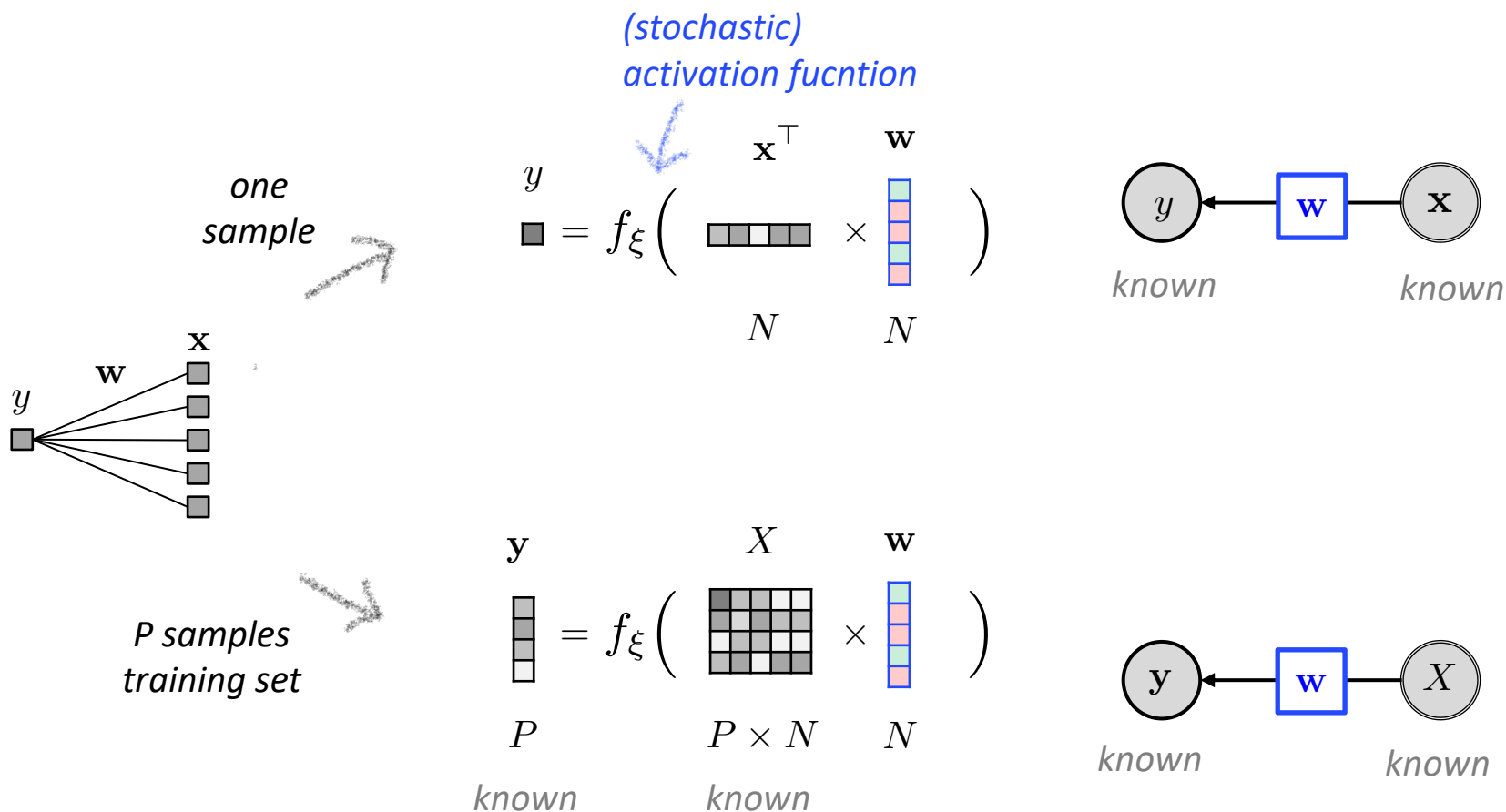**→ typical cases concentrate at the average**

$$N \to \infty \qquad \alpha = P/N$$

**Mean-field tools from the stat. phys. of disordered systems:**

*e.g. Bayes optimal square error* $\quad \mathrm{MMSE}(\alpha) = \displaystyle\lim_{N \to \infty} \frac{1}{N} \int d\mathbf{w}\,(\mathbf{w} - \hat{\mathbf{w}}_{\mathrm{MMSE}})^2\, p_S(\mathbf{w}|\mathcal{D})$

Amit, D. J., Gutfreund, H., & Sompolinsky, H. (1985) Amit, D. J., Hanoch, G., & Sompolinsky, H. (1986) Gardner, E. (1987) E. Gardner, B. Derrida, P. Mottishaw (1987), B. Derrida, J.P. Nadal (1987), Peterson, C., & Anderson, J. R. (1987), W. Krauth, M. Mézard (1987), Gardner, E. (1988), Mézard, M. (1989). Györgyi, G. (1990). Opper, M., & Haussler, D. (1991) H.S. Seung, H. Sompolinsky, N. Tishby (1992), T.L.H. Watkin, A. Rau, M. Biehl. (1993), Monasson, R., Zecchina, R. (1995) etc..

# Starting point:
# Perceptron a.k.a Generalized Linear Model (GLM)



*(stochastic) activation fucntion*

*one sample*

$$y = f_\xi \left( \quad \mathbf{x}^\top \times \mathbf{w} \quad \right)$$

$$N \qquad N$$

*known*      *known*

*P samples training set*

$$\mathbf{y} = f_\xi \left( \quad X \times \mathbf{w} \quad \right)$$

$$P \qquad P \times N \qquad N$$

*known*      *known*

*known*      *known*

# Statistical mechanics of the Perceptron

$p_x(\mathbf{x})$

*i.i.d. Rademacher*

$p_w(\mathbf{w}^*)$

*i.i.d. binary*

$$\mathbf{y} = f_\xi\left(\; X \times \mathbf{w}^* \;\right)$$

$P \qquad P \times N \qquad N$



first order phase transition
perfect generalization

$\mathbf{w} = \mathbf{w}^*$

$\epsilon_g(\alpha)$

SE
Optimal
AMP, N=$10^4$
Logistic, N=$10^4$

1. Teacher-student / planted problem
2. Bayesian optimal

3. Mean-field analysis / typical case

generalization error:

*teacher*    *student*

$$\epsilon_g(N, P) = \mathbb{E}_{\mathbf{w}|\mathcal{D}}\left[(\hat{y} - \hat{y}(\mathbf{w}))^2\right]$$

$$\epsilon_g(\alpha) = \lim_{N \to \infty} \frac{1}{N}\epsilon_g(N, P) \qquad \alpha = P/N$$

✓ Information theoretic analysis with mean-field replica method [1]

✓ Rigorously proven [2]

✓ (Generalized) Approximate Message Passing (AMP) algorithm [3]

✓ State evolution statistical analysis of algorithm performance [3]

[1] Györgyi (1990). *First-order transition to perfect generalization in a neural network with binary synapses*
[2] Barbier, et al. (2018). *Phase Transitions, Optimal Errors and Optimality of Message-Passing in Generalized Linear Models*
[3] Rangan (2011) *Generalized Approximate Message Passing for Estimation with Random Linear Mixing.*

# Mean-Field methods for statistical inference analysis
# The tools

**Information theoretic analysis**

**Non-rigorous computations of asymptotic posterior statistics**

*replica method*

**Algorithms**

**Message passing algorithms
for inference on finite size models**

*belief propagation (BP),
approximate message passing (AMP),
expectation propagation (EP)*

*high temperature expansions
(naïve MF, TAP)*

**Mathematical rigorous proofs
of the conjecture**

*Guerra interpolation,
Adaptive interpolation*

**Statistical analysis of
asymptotic performance
of message passing algorithms**

*state evolution (SE)*

# "Mean-field approximations" in deep learning literature

**- more general than tools above**

**- neglect correlations thanks to randomness in the thermodynamic (large-size) limit**

**\* Analysis of statistical of inference** ➔ Focus of this talk

*non-exhaustive!*

Reviews: - Zdeborová & Krzakala (2016) *Statistical physics of inference: Thresholds and algorithms.*
- Gabrié. (2020) *Mean field inference methods for neural networks.*

**\* Signal propagation in deep neural networks**

- Trainability of very deep network at init.   e.g. Schoenholz et al.(2017*). Deep Information Propagation.*

- Separation of structured data
    e.g. Cohen, et al (2020). *Separability and geometry of object manifolds in deep neural networks.*

**\* Role of over-parametrization in trainability with Gradient Descent methods**

- Convergence of  SGD for 2-layers neural networks
    Chizat & Bach (2018), Mei, Montanari & Nguyen (2018), Rotskoff & Vanden-Eijnden (2018)

- Neural Tangent Kernels, Equivalence to Gaussian processes, "Lazy training"
    Jacot et al (2018), Lee et al (2019),  review: Bahri et al (2020) *Statistical Mechanics of Deep Learning*

- Online learning   e.g. Goldt, et al (2019). Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student setup

**\* Gradient Descent algorithms and landscape interactions**

Dauphin et al (2014). Identifying and attacking the saddle point problem in high-dimensional non-convex optimization
Sarao Mannelli & Zdeborova (2020). *Thresholds of descending algorithms in inference problems.*

# From perceptron/GLM with random i.i.d. matrices to deep neural networks ?

*perceptron*



*i.i.d. prior*

*known*          *known*

*i.i.d. entries*

*deep neural network*



*known*          *known*

# From perceptron/GLM with random i.i.d. matrices to deep neural networks ?



*perceptron*

*deep neural network*

*i.i.d. prior*

*known*     *known*

*i.i.d. entries*

*known*   *known*    *known*   *known*

1. **Inference of layers variables in deep networks (with learned weight matrices)**

2. **The challenge of weight inference and structured weights**

# From perceptron/GLM with random i.i.d. matrices to deep neural networks ?



*perceptron*

**w**    **x**

$y$

*i.i.d. prior*

**y** ← **w** ← $X$

*known*    *known*

*i.i.d. entries*

*deep neural network*

$\mathbf{z}$

**y**    $W_2$    $W_1$    **x**

*known*

$Y$ ← $W_2$ ← $Z$ ← $W_1$ ← $X$

*known*    *known*

1. **Inference of layers variables in deep networks (with learned weight matrices)**

2. **The challenge of weight inference and structured weights**

# From perceptron/GLM with random i.i.d. matrices to deep neural networks ?



*perceptron*

*deep neural network*
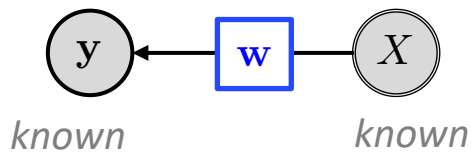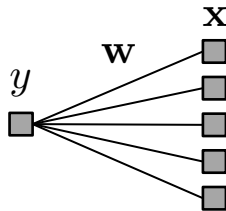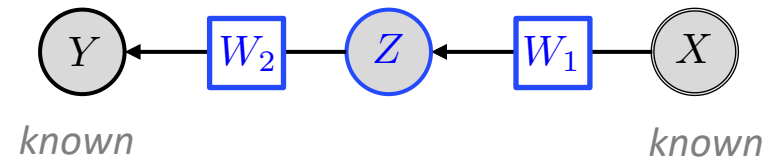
*i.i.d. prior*

*i.i.d. entries*

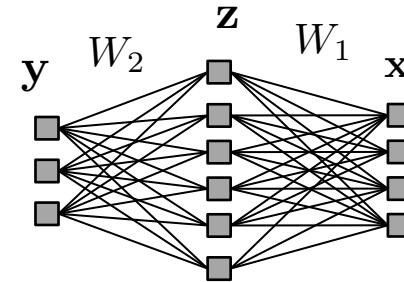1. **Inference of layers variables in deep networks (with learned weight matrices)**

2. **The challenge of weight inference and structured weights**

# Inferring neural networks layer states from output

**Single layer = perceptron / GLM**



$$\mathbf{y} = f_\xi \left( X \times \mathbf{w} \right)$$

$$P \qquad P \times N \qquad N$$

$$\mathbf{y} = f_\xi \left( W \times \mathbf{x} \right)$$

$$P \qquad P \times N \qquad N$$

**Multi-layer ?**



$$\mathbf{y} = f_\xi \left( W_2 \times f_\xi \left( W_1 \times \mathbf{x} \right) \right)$$

*known*     *known*     *known*

# Layers inference in deep neural network with i.i.d weights

*i.i.d. entries*

*i.i.d. prior*

$$\mathbf{y} = f_\xi \Big( W_2 \times f_\xi \big( W_1 \times \mathbf{x} \big) \Big)$$

$M$      $M \times K$      $K \times N$    $N$

*known*      *known*      *known*

$N$  *input dim*
$K$  *hidden dim*
$M$  *output dim*

**decompose the problem in sub-problems**   $N \to \infty$    $\alpha_1 = K/N$    $\alpha_2 = M/N$

$\mathbf{z}$  *hidden layer state*

$$\mathbf{y} = f_\xi \big( W_2 \times \mathbf{z} \big) \approx GLM$$

$K$

$$\mathbf{z} = f_\xi \big( W_1 \times \mathbf{x} \big) \approx GLM$$

✓ Multi-Layer AMP (arbitrary depth) [1]

✓ Corresponding state evolution (SE) [1]

✓ Replica free energy (mmse, entropy) [1]

✓ Rigorously proven for 2 layers [2, 3]

[1] Manoel et al (2017) *Multi-layer generalized linear estimation.*
[2] Gabrié et al (2018) *Entropy and mutual information in models of deep neural networks.*
[3] Reeves (2018) *Additivity of Information in Multilayer Networks via Additive Gaussian Noise Transforms.*

# Layers inference in deep neural network with weight matrices with correlations



*i.i.d. prior*

*rotationally invariant*

arbitrary diagonal

i.i.d. from Haar

✓ Multi-Layer Vector-AMP [1]

✓ Corresponding state evolution [1]

✓ Replica free energy (mmse, entropy) [2, 3]
   (extension of single layer formula by [4])

✗ Proof ?

[1] Fletcher et al (2018) *Inference in deep networks in high dimensions.*
[2] Gabrié et al (2018) *Entropy and mutual information in models of deep neural networks.*
[3] Reeves (2018) *Additivity of Information in Multilayer Networks via Additive Gaussian Noise Transforms.*
[4] Shinzato & Kabashima (2009) Learning from correlated patterns by simple perceptrons

# Explicit weight learning, empirical verification

**Learning the weight matrices while remaining rotationally inv.?**

**e.g. with gradient descent**

- Initialize Gaussian i.i.d W matrices
- Singular value decomposition
- Only learn spectrum (N degrees of freedom instead of $N^2$)

orthogonal    diagonal    orthogonal

$$W_\ell = U_\ell \times S_\ell \times V_\ell^\mathsf{T}$$

fixed    updated    fixed

**Numerical verification?**

- Linear networks trained
- Gaussian inputs

$Y \leftarrow W_3 \leftarrow Z_2 \leftarrow W_2 \leftarrow Z_1 \leftarrow W_1 \leftarrow X$



Replica correct
with learned matrices

[1] Gabrié et al (2018) *Entropy and mutual information in models of deep neural networks.*

# From perceptron/GLM with random i.i.d. matrices to deep neural networks ?



*perceptron*

**w**  **x**

$y$

*i.i.d. prior*

**y** ← **w** ← $X$

*known*  *known*
*i.i.d. entries*

*deep neural network*

**z**

**y**  $W_2$  $W_1$  **x**

*i.i.d. entries*

$Y$ ← $W_2$ ← $Z$ ← $W_1$ ← $X$

*known*  *known*  *known*  *known*
*learned diag.*  *learned diag.*

1. **Inference of layers variables in deep networks (with learned weight matrices)**

2. **The challenge of weight inference and structured weights**

# From GLM with random i.i.d. matrices to deep neural networks ?



*perceptron*

*deep neural network*

*known*      *known*      *known*      *known*

1. **Inference of layers variables in deep networks (with learned weight matrices)**

2. **The challenge of weight inference and structured weights**

# Weight inference in deep neural networks decomposed

$$Y = f_\xi \left( W_2 \times f_\xi \left( W_1 \times X \right) \right)$$

$M \times P$ *known*    $M \times K$    $K \times N$    $N \times P$ *known*

$N$   *input dim*
$K$   *hidden dim*
$M$   *output dim*
$P$   *sample size*

**First idea: decompose the inference in sub-problems**

**(alike Multi layer - AMP)**

*hidden layer states*
*over the $P$ samples*

$Y$   $W_2$   $Z$

▷   $Y = f_\xi \left( W_2 \times Z \right)$

*known*    $K \times P$

≈ *matrix factorization*
*with rank $K$*

$Z$   $W_1$   $X$

▷   $Z = f_\xi \left( W_1 \times X \right)$

*known*

≈ *P x GLMs*

# Scaling of the size of the hidden layer?

$$Y = f_\xi \left( W_2 \times Z \right) \approx \text{matrix factorization with rank } K$$

$M \times P$ *known*   $M \times K$   $K \times P$

$N \to \infty$

▷ $K = O(1)$

- "low-rank matrix factorization": good mean field understanding [1, 2]

- finite number of hidden units, committee machines: great body of work! [3, 4, 5, 6, ..]

[1] Lesieur et al (2016), *MMSE of probabilistic low-rank matrix estimation: Universality with respect to the output channel*
[2] Lesieur et al (2017), *Constrained Low-rank Matrix Estimation: Phase Transitions, Approximate Message Passing and Applications*
[3] Aubin et al (2018). *The committee machine: Computational to statistical gaps in learning a two-layers neural network*
[4] Monasson et ql (2004). *Learning and Generalization Theories of Large Committee-Machines*
[5] Schwarze & Hertz (1993). *Generalization in Fully Connected Committee Machines.*
[6] Schwarze (1993). *Learning a Rule in a Multilayer Neural-Network.*

# Phase transitions for committee machines



$$y = \text{sign}\left( w_2 \times \text{sign}\left( W_1 \times X \right) \right)$$

$N$ *input dim*
$K$ *hidden dim*
$P$ *sample size*

$P$ *known*    $K$ *known*    $K \times N$   $N \times P$ *known*

K = 2
binary weights

| | | |
|---|---|---|
| AMP $q_{00}$ | — SE $q_{00}$ | — SE $\epsilon_g(\alpha)$ |
| AMP $q_{01}$ | — SE $q_{01}$ | • AMP $\epsilon_g(\alpha)$ |

✓ Committee-AMP [1]

✓ Corresponding state evolution [1]

✓ Replica free energy (mmse) [2, 3, 4]

✓ Proof [1]

*teacher*     *student*

$q_{00} = \text{overlap}\{(W_1^*)_{0,.} ; (W_1)_{0,.}\}$

$q_{01} = \text{overlap}\{(W_1^*)_{0,.} ; (W_1)_{1,.}\}$



*specialization*

*perfect generalization*

[1] Aubin et al (2018). *The committee machine: Computational to statistical gaps in learning a two-layers neural network*
[2] Monasson et ql (2004). *Learning and Generalization Theories of Large Committee-Machines*
[3] Schwarze & Hertz (1993). *Generalization in Fully Connected Committee Machines.*
[4] Schwarze (1993). *Learning a Rule in a Multilayer Neural-Network.*

# Scaling of the size of the hidden layer?

$$Y = f_\xi \left( W_2 \times Z \right) \qquad \approx \textit{matrix factorization with rank } K$$

$M \times P$
*known*

$M \times K \qquad K \times P$

$N \to \infty$

▷ $K = O(1)$

- "low-rank matrix factorization": good mean field understanding [1, 2]

- finite number of hidden units, committee machines: great body of work! [3, 4, 5, 6, ..]

▷ $K = O(N)$

- "high-rank matrix factorization": mean-field analysis?

- number of hidden units scaling like the inputs

[1] Lesieur et al (2016), *MMSE of probabilistic low-rank matrix estimation: Universality with respect to the output channel*
[2] Lesieur et al (2017), *Constrained Low-rank Matrix Estimation: Phase Transitions, Approximate Message Passing and Applications*
[3]Aubin et al (2018). *The committee machine: Computational to statistical gaps in learning a two-layers neural network*
[4] Monasson et ql (2004). *Learning and Generalization Theories of Large Committee-Machines*
[5] Schwarze & Hertz  (1993). *Generalization in Fully Connected Committee Machines.*
[6] Schwarze (1993). *Learning a Rule in a Multilayer Neural-Network.*

# Structured weights inference $K = O(N)$

$$Y = f_\xi \left( W_2 \times Z \right)$$

$M \times P$
*known*

$M \times K$

$K \times P$

*# parameters* $O(N^2)$

**Second idea: learn structured simpler weights**

$$\triangleright \quad Y = f_\xi \left( S_2 \times \tilde{W}_2 \times Z \right)$$

$M \times P$
*known*

$M \times M$

$M \times K$
*known*

$K \times P$

*# parameters* $O(N)$

▷ Also used in deep learning literature:

- Speed / memory concerns: e.g. ACDC layers [1], Ensemble learning [2]
- Theoretical papers:  e.g. Porcupine networks [3], Replica entropy [4]

▷ Signal processing literature: a.k.a. Blind Calibration

[1] Moczulski et al (2015), *ACDC: A Structured Efficient Linear Layer*
[2] Wen et al (2020), *BatchEnsemble: An Alternative Approach to Efficient Ensemble and Lifelong Learning*
[3] Feizi et al (2016) *Porcupine Neural Networks: (Almost) All Local Optima are Global*
[4] Gabrié et al  (2018), *Entropy and mutual information in models of deep neural networks*

# Blind calibration mean field analysis

**Simultaneous recovery of input signal and "calibration variables"**

*i.i.d. prior*    *i.i.d. entries*    *i.i.d. prior*

$$Y = f_\xi \left( S \times \tilde{W} \times X \right)$$

| | | | |
|---|---|---|---|
| $M \times P$ | $M \times M$ | $M \times N$ | $N \times P$ |
| *known* | *to be calibrated* | *known* | |

$N$  *input dim*

$M$  *output dim*

$P$  *sample size*

✓  Calibration - AMP algorithm [1, 2]

✓  Corresponding state evolution [3]

✓  Replica free energy [3]

✗  Rigorous proof

[1] Schulke C. et al (2013), *Blind Calibration in Compressed Sensing using Message Passing Algorithms*
[2] Schulke C. et al (2016), *Blind sensor calibration using approximate message passing*
[3] Gabrié M. et al (2020), *Blind calibration for compressed sensing: State evolution and an online algorithm*

# Numerical results for sparse priors

**Example sparse signal recovery:**

*output dim / input dim* $\quad \alpha = M/N$

*input sparsity* $\quad \rho$

*naive count* $\quad \alpha_{\min} = \rho \dfrac{P}{P-1}$

$$Y = f_\xi \left( S \times \tilde{W} \times X \right)$$

$M \times P$    $S$    $\tilde{W}$    $X$

*i.i.d. entries*

$N \times P$

*i.i.d. prior*      *i.i.d. $\rho$ -sparse*

Cal- AMP reconstruction errors ($P = 2$)

$S$ error      $X$ error

Cal- AMP State evolution



*Cal-AMP reconstructs efficiently with a finite number of samples*

*Good agreement SE and Cal-AMP*

[1] Gabrié M. et al (2020), *Blind calibration for compressed sensing: State evolution and an online algorithm*

# Statistical mechanics of online learning

$$Y = f_\xi \left( S \times \tilde{W} \times X \right)$$

$M \times P$      $N \times P$

*observation*    $\hat{S} \; \hat{\mathbf{x}}^{(1)}$      $\hat{S} \; \hat{\mathbf{x}}^{(2)}$

*approx. posterior*       *approx. posterior*

$$p_S(S) \; p_x(\mathbf{x}) \xrightarrow{\mathbf{y}^{(1)}} \hat{p}(S, \mathbf{x} \,|\, \mathbf{y}^{(1)}) \xrightarrow{\mathbf{y}^{(2)}} \hat{p}(S, \mathbf{x} \,|\, \mathbf{y}^{(2)}) \xrightarrow{\mathbf{y}^{(3)}} etc.$$

*prior*      *updated prior*      *updated prior*

▷ **Streaming AMP for GLM [1], for blind calibration [2]**

**Numerical results:**

**Example of sparse signal recovery**

$X$   *i.i.d. $\rho$ -sparse*

$\alpha = M/N$

    *output dim / input dim*



[1] Manoel et al. (2018). Streaming Bayesian inference: Theoretical limits and mini-batch approximate message-passing
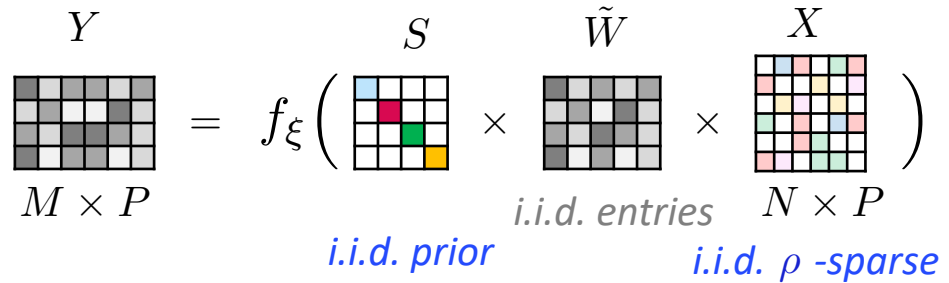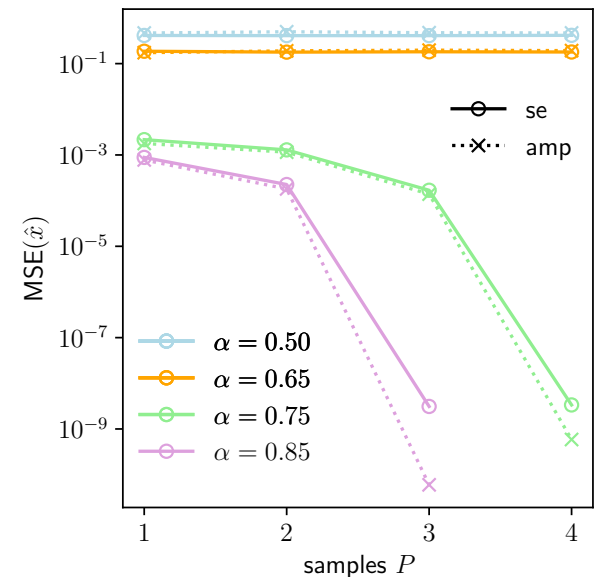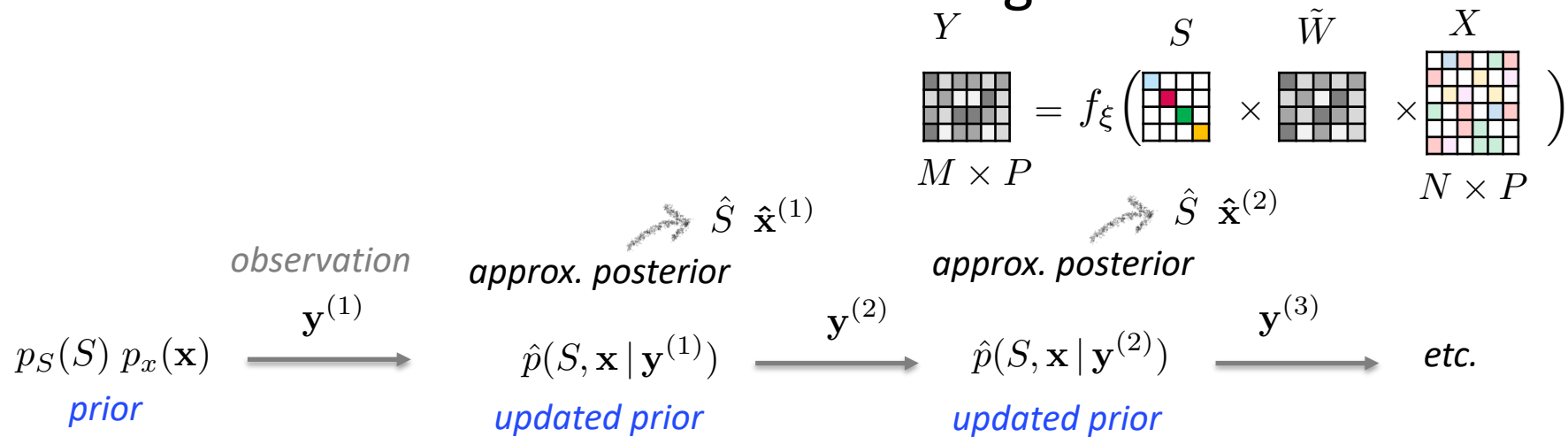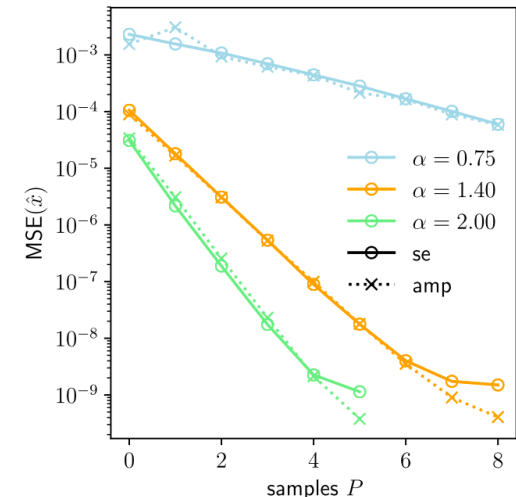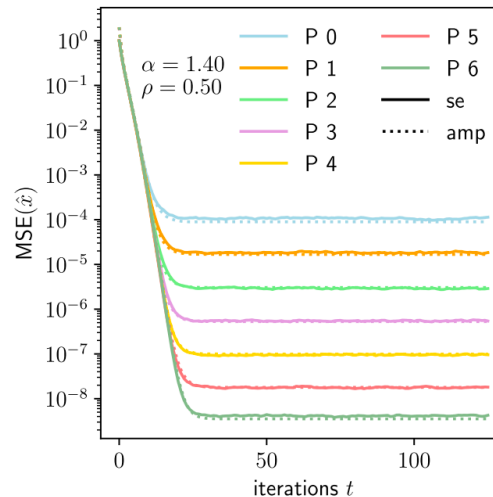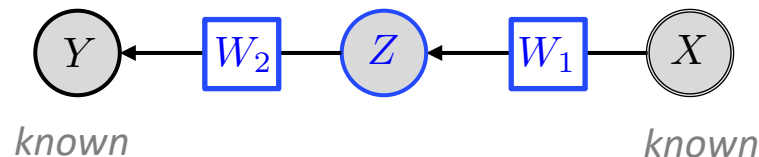[2] Gabrié M. et al (2020), *Blind calibration for compressed sensing: State evolution and an online algorithm*

# Perspectives for weight inference in deep NNs

$$Y \quad S \quad \tilde{W} \quad X$$



$$Y = f_\xi \left( S \times \tilde{W} \times X \right)$$

**Weight inference in hidden layers for the stat mech of deep learning (offline/batch and online/mini-batch)**

▷ Perspective: Combine Cal-AMP in layers to infer structured weights in NNs (extensive number of hidden units!)

▷ Challenge: Back to the teacher-student scenario?



$$Y \longleftarrow W_2 \longleftarrow Z \longleftarrow W_1 \longleftarrow X$$

*known*         *known*

# Perspectives for mean-field methods for inference and information/computational thresholds

▷ **More and more complex matrix ensembles (weights, data)**

▷ **Combining solutions to more complex models**

▷ **Great open source package for algorithms**        sphinxteam / **tramp**

Tutorial review:

 Gabrié (2020), *Mean field inference methods for neural networks* – arXiv/1911.00890

Software:

 Baker et al (2020), *Compositional Inference with Tree Approximate Message Passing*

# Thank you!