

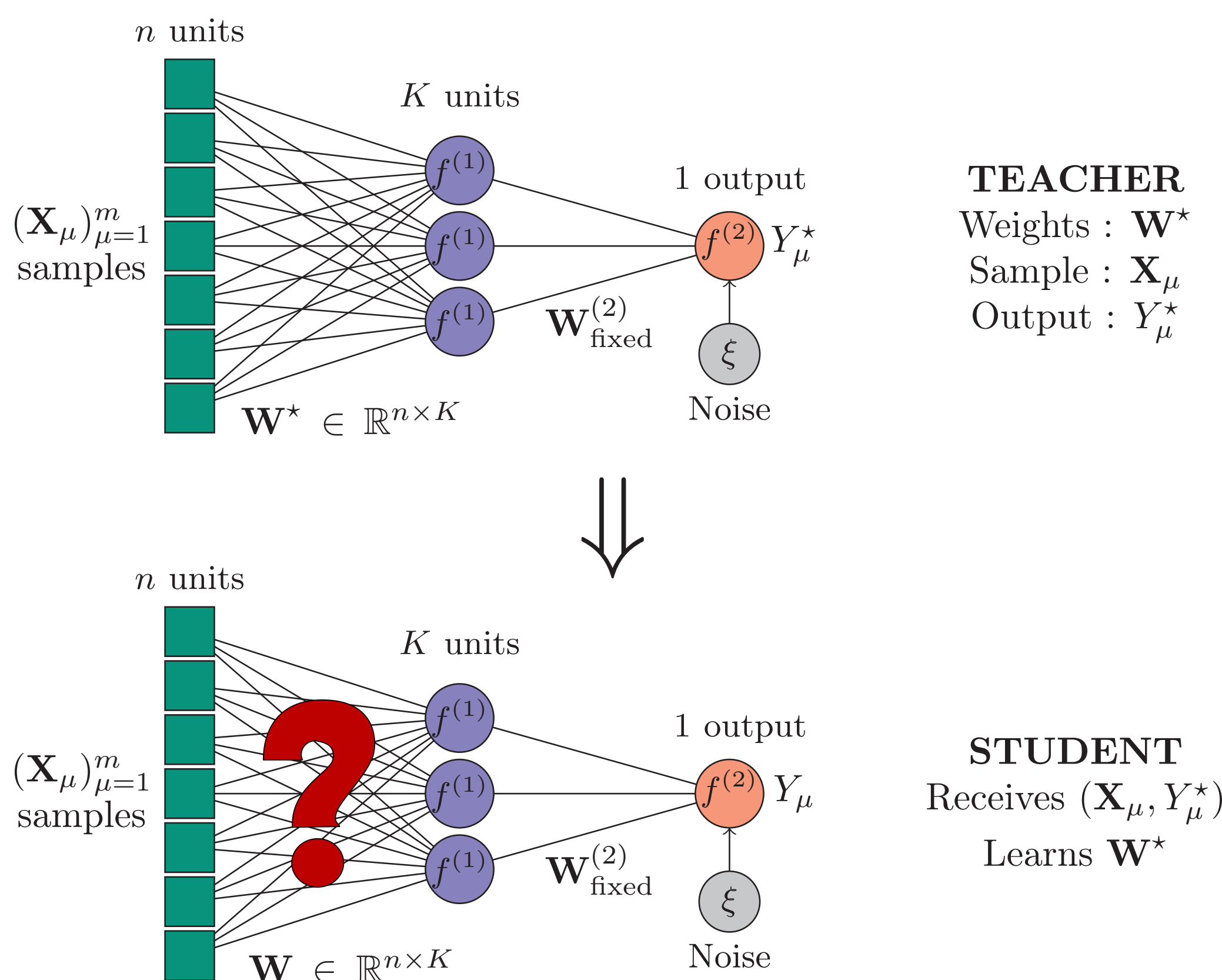
# The committee machine: Computational to statistical gaps in learning a two-layers neural network

Benjamin Aubin<sup>2</sup>, Antoine Maillard<sup>1</sup>, Jean Barbier<sup>1,3,4</sup>, Florent Krzakala<sup>1</sup>, Nicolas Macris<sup>3</sup> & Lenka Zdeborová<sup>2</sup>  
<sup>1</sup>Ecole Normale Supérieure (Paris) <sup>2</sup>Institut de Physique Théorique (Saclay) <sup>3</sup>Ecole Polytechnique Fédérale de Lausanne <sup>4</sup>International Center for Theoretical Physics (Trieste)

## Motivation and outline

- Understand the *generalization effectiveness* of neural networks.
- Understand the *typical case scenario*, and locate *phase transitions* in the generalization abilities: when can the network learn efficiently ?
- Traditional bounds are *worst case* bounds, based on the value of the VC dimension and Rademacher complexity.
- ⇒ Complementary approach: we are able to compute the *optimal* generalization error in the typical case, and to compute the corresponding solution in *polynomial* time.
- Drawbacks: we need an i.i.d.-sampled dataset, and the optimal algorithm is slower than SGD.

## A general teacher-student model



- Limit  $n \rightarrow \infty$  with  $\lim_{n \rightarrow \infty} \frac{m}{n} = \alpha \in (0, \infty)$  and  $K = \Theta_n(1)$ .
- Gaussian i.i.d. samples  $\mathbf{X}_\mu \sim \mathcal{N}(0, Id)$ .
- $\forall i \in [1, n]$ , the weight vector  $\mathbf{W}_i = \{W_{il}\}_{l=1}^K$  has a Bayesian prior  $P_0$ , with zero mean and well-defined covariance matrix  $\rho$ .
- We look at the *typical case*, by averaging over the weights of the teacher and the noise.

At a non-rigorous level, these models were previously investigated by the theoretical physics literature e.g. [1], [2].

### Notations and definitions

- $\mathcal{S}_K^+$  : real positive symmetric  $K \times K$  matrices.
- $\mathcal{S}_K^+(\rho) = \{q \in \mathcal{S}_K^+ \text{ s.t. } q \in \mathcal{S}_K^+ \wedge (\rho - q) \in \mathcal{S}_K^+\}$ .
- For  $(\mathbf{V}, \mathbf{U}) \in \mathbb{R}^K$ ,  $\mathbf{z}_{q,\mathbf{V}}(\mathbf{U}) \equiv q^{1/2}\mathbf{V} + (\rho - q)^{1/2}\mathbf{U}$ .
- For  $\mathbf{x} \in \mathbb{R}^d$ ,  $\mathcal{D}\mathbf{x} = \frac{1}{(2\pi)^{d/2}} e^{-\frac{1}{2}\|\mathbf{x}\|^2} d\mathbf{x}$ .
- We write  $Y_\mu \sim P_{\text{out}}\left(\cdot \mid \frac{1}{\sqrt{n}}\mathbf{X}_\mu\mathbf{W}\right)$
- Two auxiliary functions for  $r \in \mathcal{S}_K^+$  and  $q \in \mathcal{S}_K^+(\rho)$ :

$$\mathcal{I}_{P_0}(r) \equiv - \int \mathcal{D}\mathbf{Z} dP_0(\mathbf{W}_0) \ln \int dP_0(\mathbf{W}) e^{-\frac{1}{2}\|r^{1/2}(\mathbf{W}_0 - \mathbf{W}) + \mathbf{Z}\|^2}$$

$$\mathcal{I}_{\text{out}}^{(\rho)}(q) \equiv \int d\hat{\mathbf{Y}} \mathcal{D}\mathbf{V} P_{\text{out}}(\hat{\mathbf{Y}}|\rho^{1/2}\mathbf{V}) \ln P_{\text{out}}(\hat{\mathbf{Y}}|\rho^{1/2}\mathbf{V})$$

$$- \int d\hat{\mathbf{Y}} \mathcal{D}\mathbf{V} \mathcal{D}\mathbf{u}^* P_{\text{out}}(\hat{\mathbf{Y}}|\mathbf{z}_{q,\mathbf{V}}(\mathbf{u}^*)) \ln \int d\mathbf{u} P_{\text{out}}(\hat{\mathbf{Y}}|\mathbf{z}_{q,\mathbf{V}}(\mathbf{u}))$$

## Main theoretical result

### Theorem: Replica-symmetric formula

- (i) The normalized *mutual information*  $i_n \equiv \frac{1}{n} I(\mathbf{W}; \mathbf{Y}|\mathbf{X})$  converges to :

$$i_\infty = \inf_{r \in \mathcal{S}_K^+} \sup_{q \in \mathcal{S}_K^+(\rho)} \left\{ \mathcal{I}_{P_0}(r) + \alpha \mathcal{I}_{\text{out}}^{(\rho)}(q) - \frac{1}{2} \text{Tr}[r(\rho - q)] \right\}$$

Let us call  $(r^*, q^*)$  the extremizers in this formula.

- (ii) The *Bayes-optimal* generalization error

$$\epsilon_g^{(n)} \equiv \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{W}^*} \left[ \left( \mathbb{E}_{\mathbf{W}|\mathbf{X}} [Y(\mathbf{X}\mathbf{W})] - Y^*(\mathbf{X}\mathbf{W}^*) \right)^2 \right]$$

converges as  $n \rightarrow \infty$  to a limit  $\epsilon_g(q^*)$  that only depends on  $q^*$ .

$q^*$  can be interpreted as the *overlap matrix* between the weights of the teacher and the weights of the student : “ $q^* = \frac{1}{n} \mathbf{W}^\top \mathbf{W}^*$ ”.

## Sketch of proof

We use an adaptive version of Guerra’s interpolation [3], developed in [4]. Extension of these techniques applied in general linear estimation [5].

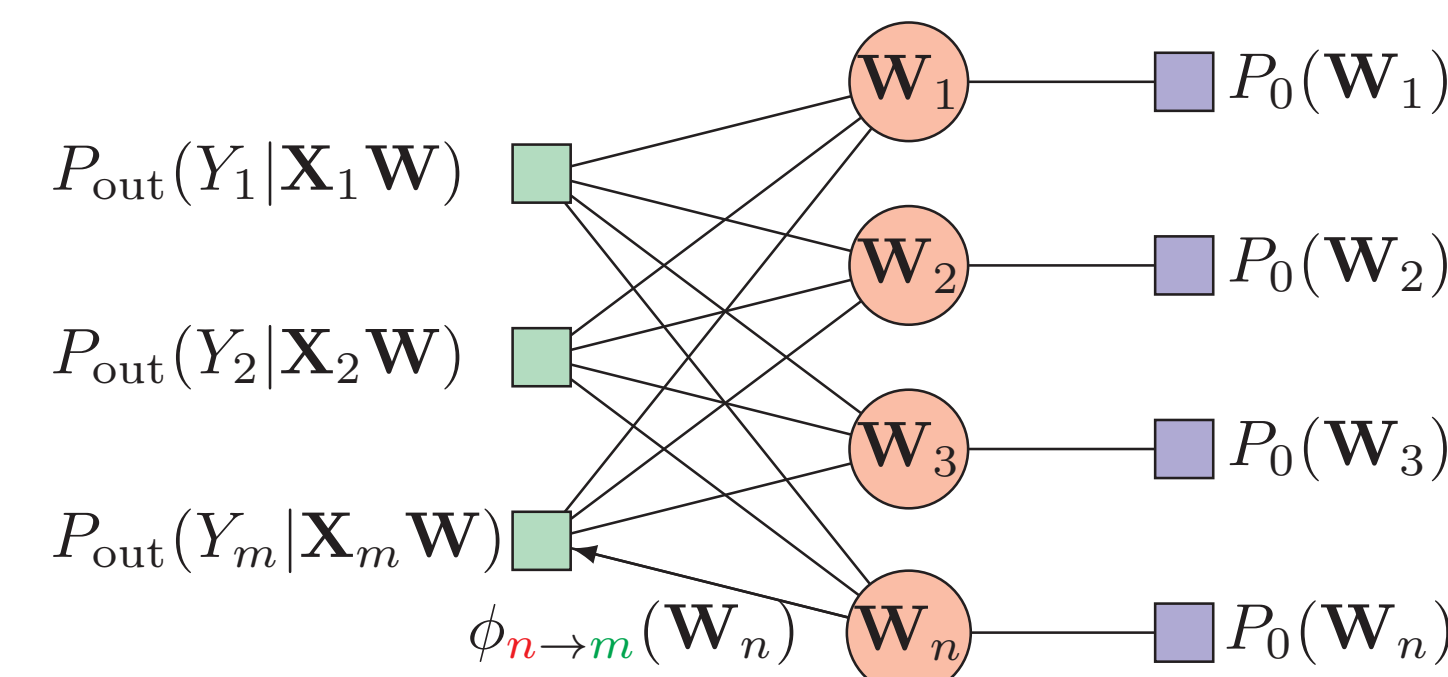
- (1) Choose functions  $q(t), r(t)$ . Two estimations problems:

- (a)  $\mathbf{Y}_0 = r(t)^{1/2}\mathbf{W}_0 + \mathbf{Z}$   $i_a(t=1) = \mathcal{I}_{P_0}(r(1))$   
(b)  $\hat{\mathbf{Y}} \sim P_{\text{out}}(\cdot|\mathbf{z}_{q(t),\mathbf{V}}(\mathbf{u}^*))$   $i_b(t=1) = \mathcal{I}_{\text{out}}^{(\rho)}(q(1))$

- (2) Interpolate between the original problem and (a) + (b). The interpolated m.i. verifies  $i_n(t=0) = i_n$ ;  $i_n(t=1) = i_a + i_b$ . So  $i_n = i_a + i_b - \int_0^1 i'_n(t) dt$ . Choose a *smart path*  $\{q(t), r(t)\}$  to conclude.

## Approximate message-passing algorithm (1)

*Factor graph* representation of the interactions:



- The “message”  $\phi_{i \rightarrow \mu}(\mathbf{W}_i)$  is the marginal probability of  $\mathbf{W}_i$  in the absence of the node  $\mu$ . One can write *belief propagation* [6] iterative equations on the  $\phi_{i \rightarrow \mu}(\mathbf{W}_i)$ .
- Gaussian approximation for messages  $\Rightarrow$  AMP [7] [8].
- Different from SGD : **we do not optimize a cost function !**
- AMP is conjectured (and often shown) to perform the optimal learning for inference problems of this class.**
- Rigorous track of AMP via *state evolution* iterations:**

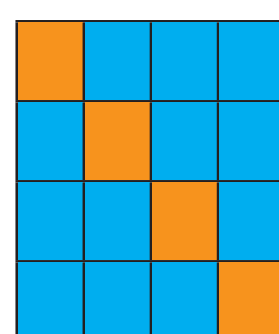
$$q_{\text{AMP}}^{t+1} = \rho - 2\partial_r \mathcal{I}_{P_0}(r_{\text{AMP}}^t), \quad r_{\text{AMP}}^{t+1} = -2\alpha \partial_q \mathcal{I}_{\text{out}}^{(\rho)}(q_{\text{AMP}}^t).$$

It is the variational condition of the replica symmetric-formula !

## A special case: the committee machine

$$Y_\mu^* = \text{sign} \left[ \frac{1}{\sqrt{K}} \sum_{l=1}^K \text{sign} \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n X_{\mu i} W_{il}^* \right) \right].$$

Symmetry  $\Rightarrow q^* =$

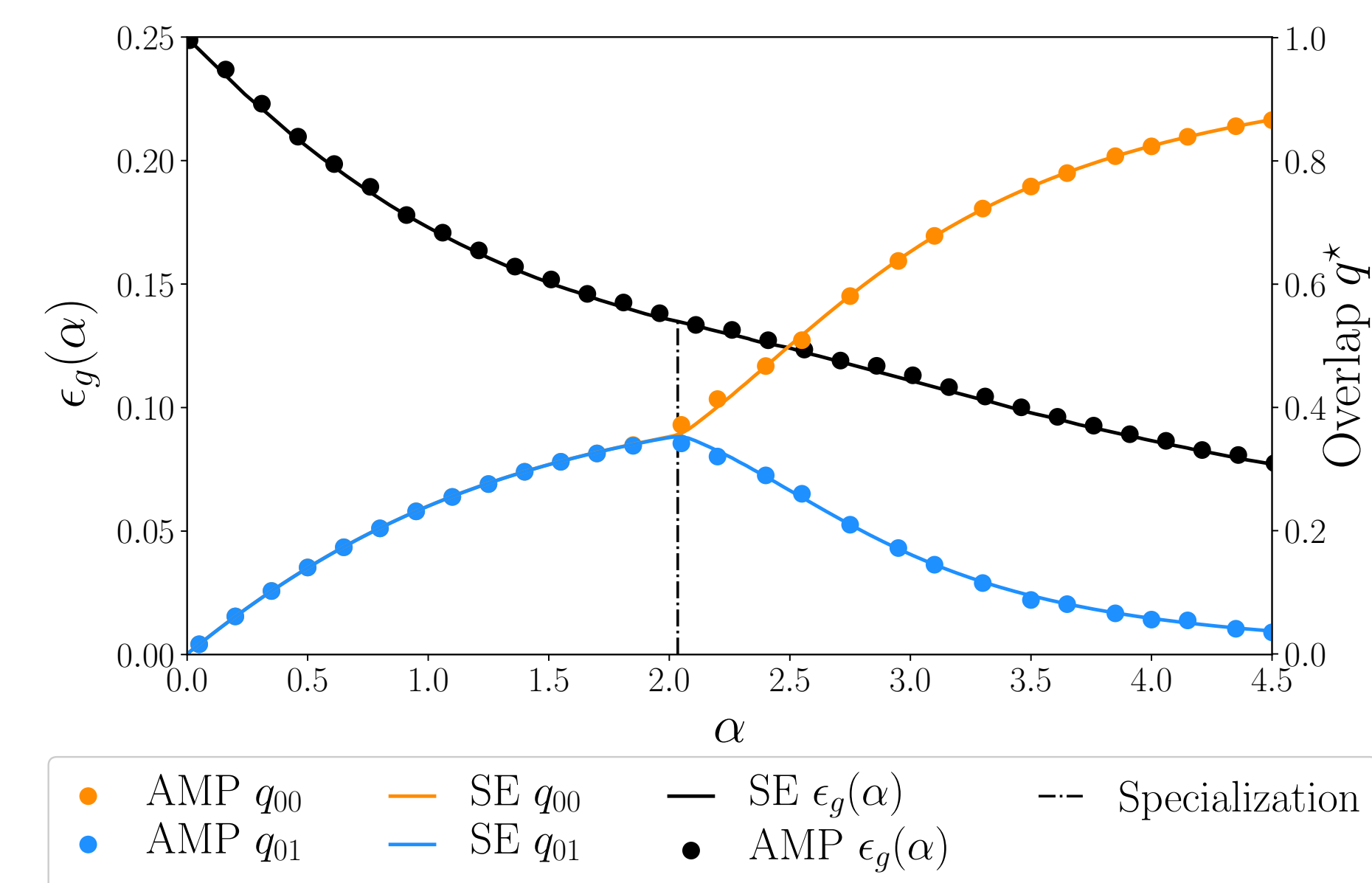


## Approximate message-passing algorithm (2)

### AMP iterative algorithm

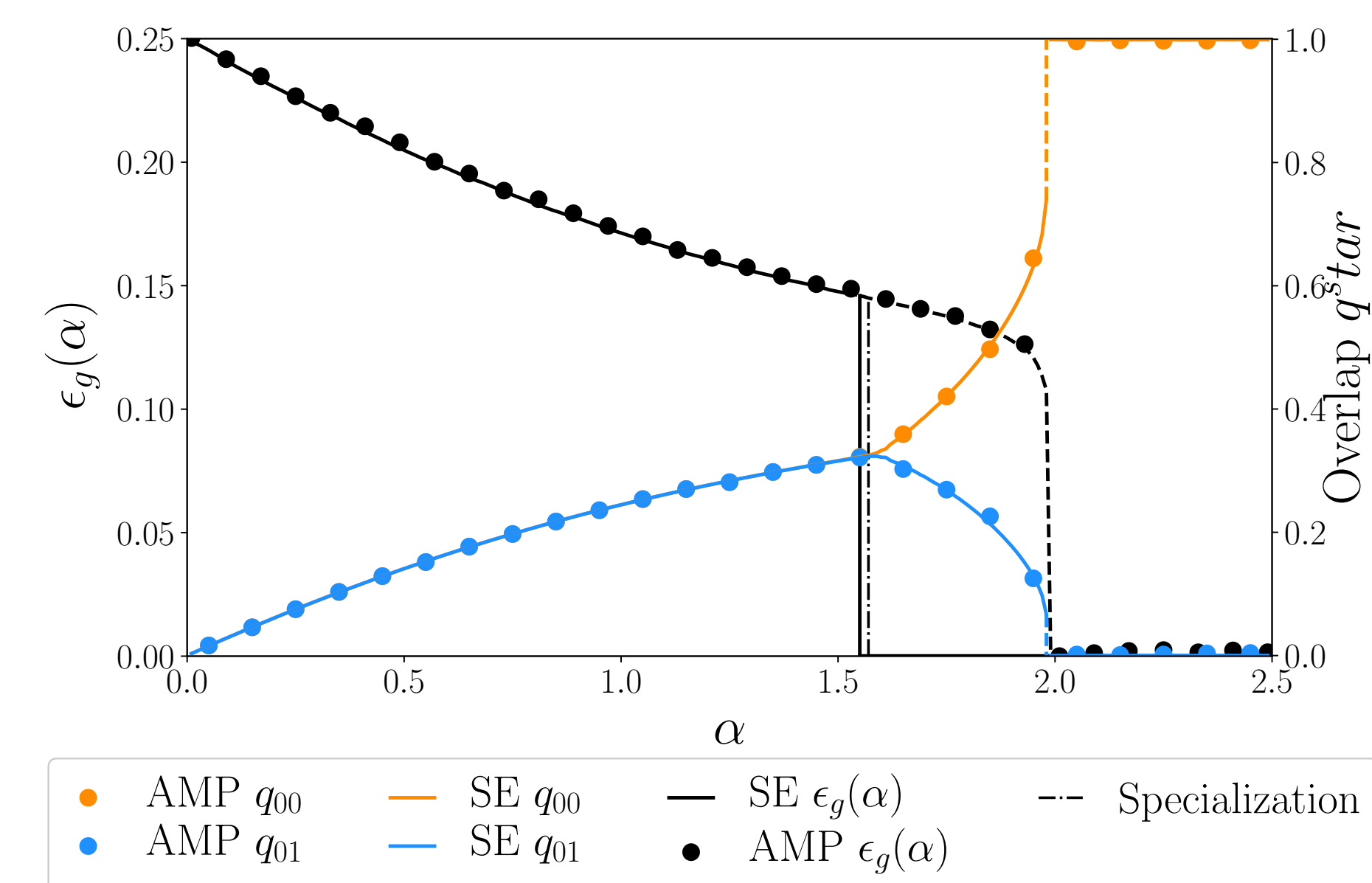
- Input:  $\mathbf{Y} \in \mathbb{R}^m$ ,  $\mathbf{X} \in \mathbb{R}^{m \times n}$ .
- Iteration  $t \rightarrow t+1$ :  
 $\hat{\mathbf{W}}_\mu^t = \sum_{i=1}^n \left[ X_{\mu i} \hat{\mathbf{W}}_i^t - X_{\mu i}^2 (\Sigma_i^{t-1})^{-1} \hat{C}_i^t \Sigma_i^{t-1} \mathbf{g}_\mu^{t-1} \right]$   
 $V_\mu^t = \sum_{i=1}^n X_{\mu i}^2 \hat{C}_i^t$   
 $\mathbf{g}_\mu^t = G(Y_\mu, \omega_\mu^t, V_\mu^t)$  ;  $h_\mu^t = H(Y_\mu, \omega_\mu^t, V_\mu^t)$   
 $\mathbf{T}_i^t = \Sigma_i^t \left( \sum_{\mu=1}^m X_{\mu i} \mathbf{g}_\mu^t - X_{\mu i}^2 h_\mu^t \hat{\mathbf{W}}_i^t \right)$   
 $\Sigma_i^t = - \left( \sum_{\mu=1}^m X_{\mu i}^2 h_\mu^t \right)^{-1}$   
 $\hat{\mathbf{W}}_i^{t+1} = f_W(\mathbf{T}_i^t, \Sigma_i^t)$  ;  $\hat{C}_i^{t+1} = f_C(\mathbf{T}_i^t, \Sigma_i^t)$   
 $[G, H, f_W \text{ and } f_C \text{ are simple functions of } P_0 \text{ and } P_{\text{out}}.]$
- Output:  $\{\hat{\mathbf{W}}_i, \hat{C}_i\}$  (mean and variance of the weights).

## Small number of hidden neurons (1)



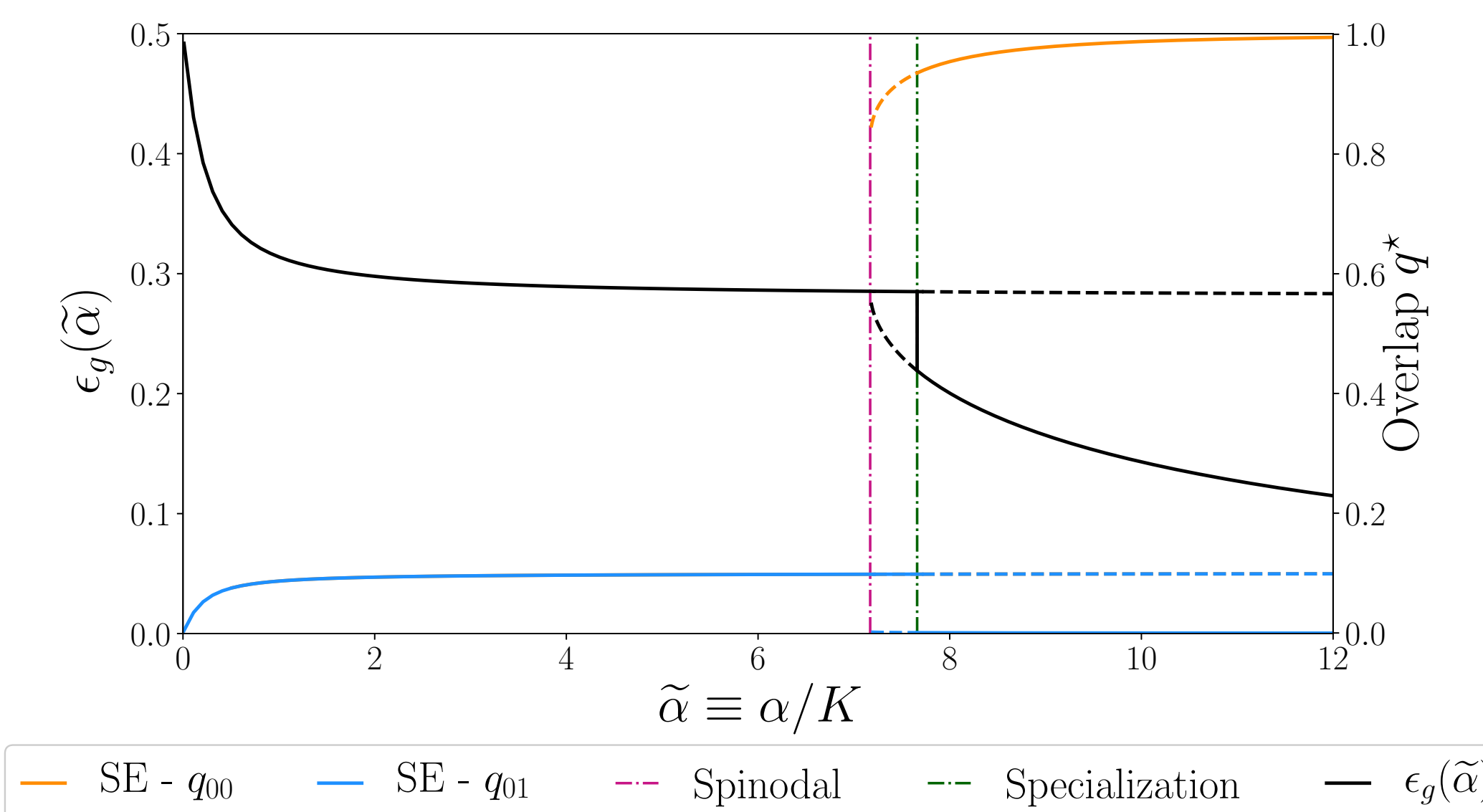
- $K=2$  with *Gaussian prior* on the weights.
- Specialization* transition at  $\alpha_c \simeq 2$  : for  $\alpha \leq \alpha_c$ , the student believes the data is linearly separable, i.e.e that the teacher has only one hidden unit. At  $\alpha \geq \alpha_c$ , the student “learns” to differentiate the neurons of the teacher.
- Specialization is here a *second order* phase transition.
- Perfect generalization reached at  $\alpha \gg 1$ .

## Small number of hidden neurons (2)



- $K=2$  with *Rademacher prior* on the weights.
- Existence of a *specialization* transition at  $\alpha_c \simeq 1.5$ .
- Perfect generalization* transition: the perfect generalization solution exists for  $\alpha \simeq 1.5$ , but it is found by AMP for  $\alpha \geq 2$ .
- ⇒ Existence of a *computational gap* / *hard phase*.

## Large number of hidden neurons



- $K \gg 1$  with *Gaussian prior* on the weights.
- First order* specialization transition.
- At all  $\tilde{\alpha}$ , AMP is stuck in the local minimum of the non-specialized solution with high generalization error.
- It is a stable local minima at least up to  $\alpha \sim K^2$  !
- Very large computational gap* / *hard phase* !

## Conclusion and perspectives

- Proof of the heuristic replica formula for a general model of a two-layers neural network with one hidden layer.
- Exact/Rigorous computation of the Bayes optimal error and optimal learning to unveil computational gaps.
- Evidence for a specialization phase transition at both small and large  $K$  and different priors on the weights.
- Algorithmic evidence for a large ‘hard’ phase with  $K \gg 1$ , and for small  $K$  with binary weights.
- Algorithmic evidence of a perfect generalization phase transition in large binary networks.
- Side result : these transitions do not appear in linear networks.
- (Some) remaining questions :
  - What happens for  $K$  diverging with  $n$ ? Even at the (non-rigorous) replica level, this is a challenging question.
  - What if the student has a different architecture than the teacher ?
  - What if we learn the second layer ? And a deep network ?

## References

- T. L. H. Watkin, A. Rau, and M. Biehl. Apr. 1993.
- H. S. Seung, H. Sompolinsky, and N. Tishby. Apr. 1992.
- F. Guerra. eprint: [arXiv:cond-mat/0205123](https://arxiv.org/abs/cond-mat/0205123).
- J. Barbier and N. Macris. eprint: [arXiv:1705.02780](https://arxiv.org/abs/1705.02780).
- J. Barbier et al. eprint: [arXiv:1708.03395](https://arxiv.org/abs/1708.03395).
- M. Mezard and A. Montanari. *Information, physics, and computation*. Oxford University Press, 2009.
- D. L. Donoho, A. Maleki, and A. Montanari. eprint: [arXiv:0907.3574](https://arxiv.org/abs/0907.3574).
- S. Rangan. eprint: [arXiv:1010.5141](https://arxiv.org/abs/1010.5141).

