

Tutorial: Approximate Message Passing & Replicas for committee machines

Benjamin Aubin[†]

[†] *Institut de Physique Théorique*
CNRS & CEA & Université Paris-Saclay, Saclay, France

Contents

1	Introduction	3
1.1	Model	3
1.1.1	Data	3
1.1.2	Teacher - Planted solution	3
1.1.3	Channel	4
1.1.4	Student	4
1.2	Notations	4
2	Bayesian inference and factor graph	5
3	Replica computation	7
3.1	Heuristic derivation	7
3.1.1	Replica trick	7
3.1.2	Average over the i.i.d input data \mathbf{X}	8
3.1.3	Fourier representation	9
3.2	Replica Symmetric free entropy	11
3.2.1	Replica symmetric ansatz	11
3.2.2	Decomposition of the free entropy functional	13
3.2.3	Consistency conditions $r \rightarrow 0$: $\Theta(1)$ terms	16
3.2.4	Replica trick $r \rightarrow 0$ limit: $\Theta(r)$ terms	17
3.3	Summary	22
3.3.1	Summary - Mismatched case	22
3.3.2	Bayes optimal MMSE estimation	23
3.4	Fixed point equations	23
3.4.1	Mismatched setting	23
3.4.2	Bayes-optimal estimation	24

4	Belief Propagation and Approximate Message Passing	25
4.1	Factor graph and joint probability distribution	25
4.2	Belief Propagation equations	26
4.3	Relaxed Belief Propagation equations	26
4.4	AMP algorithm	29
4.5	State evolution equations of AMP	31
4.5.1	Messages distribution	33
4.5.2	Summary of the SE - mismatched setting	35
4.5.3	Summary of the SE - Bayes-optimal setting	36
5	Consistency between AMP and the replica computation	36

1 Introduction

The purpose of these notes is to provide a short introduction to **Approximate Message Passing (AMP) algorithms** and **replicas** computation, that we illustrate for **committee machines**, that are a simple generalization of **generalized linear models**. For more details, see in particular the references [1, 2, 3, 4, 5, 6, 7].

1.1 Model

We revisit the so-called *teacher-student* scenario from statistical physics [8, 9, 10, 11, 12, 13, 14]. We assume that the synthetic dataset is generated by a *teacher* model.

1.1.1 Data

We assume that the matrix of data inputs $\mathbf{X} \in \mathbb{R}^{n \times d}$, that contains n d -dimensional samples, is drawn i.i.d with distribution $P_{\mathbf{x}}(\mathbf{X}) = \prod_{i,\mu=1}^{d,n} P_{\mathbf{x}}(x_{\mu i})$.

For simplicity, we will consider them to be i.i.d Gaussian with zero mean and unit variance: $\forall \mu \in \llbracket n \rrbracket, \mathbf{x}_{\mu} \sim \mathcal{N}_{\mathbf{x}}(\mathbf{0}, \mathbf{I}_d)$.

1.1.2 Teacher - Planted solution

Moreover, we draw a **planted solution** \mathbf{W}^* from a separable distribution (across the input dimension) $P_{\mathbf{w}^*}(\mathbf{W}^*) = \prod_{i=1}^d P_{\mathbf{w}_i^*}(\mathbf{w}_i^*)$. Therefore the weights could possibly be correlated across the second dimension.

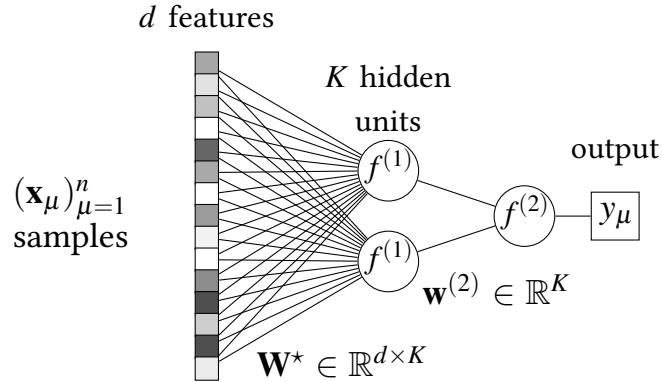


Figure 1 – Illustration of the *committee machine*: it is one of the simplest models belonging to the considered model class (1), and on which we focus to illustrate our results. It is a two-layers neural network with sign activation functions $f^{(1)}, f^{(2)} = \text{sign}$ and weights $\mathbf{w}^{(2)}$ fixed to unity. It is represented for $K = 2$.

1.1.3 Channel

The labels are then generated by multiplying the data matrix \mathbf{X} and the planted weights \mathbf{W}^* and passing their product to a **channel** φ_{out} , acting component-wise, according to:

$$\mathbf{y} = \varphi_{\text{out}} \left(\left\{ \frac{1}{\sqrt{d}} \mathbf{X} \mathbf{w}_k^* \right\}_{k=1}^K \right) \quad \text{or} \quad \mathbf{y} \sim \mathbf{P}_{\text{out}} \left(\cdot \mid \left\{ \frac{1}{\sqrt{d}} \mathbf{X} \mathbf{w}_k^* \right\}_{k=1}^K \right), \quad (1)$$

and illustrated in Fig. 1 in the case of a committee machine with fixed weights $\mathbf{w}^{(2)}$ in the second layer.

The function $\varphi_{\text{out}} : \mathbb{R}^K \mapsto \mathbb{R}$ represents a deterministic, or stochastic function associated to a probability distribution \mathbf{P}_{out} , applied component-wise to each sample.

Remark 1.1. Notice that the factor $\frac{1}{\sqrt{d}}$ is present to insure that the variance of the input data is normalized to the unit.

1.1.4 Student

The committee machine estimation problem consists of trying to fit the n input-output observations $\{\mathbf{X}, \mathbf{y}\} \in \mathbb{R}^{n \times d} \times \mathbb{R}^n$.

In the case the **student** has exactly the same architecture, namely $\mathbf{W}^* \in \mathbb{R}^{d \times K}$ and $\mathbf{W} \in \mathbb{R}^{d \times K}$ have the same dimensionality and prior information $\mathbf{P}_{\mathbf{w}^*} = \mathbf{P}_{\mathbf{w}}$, this setting is called the **Bayes-optimal setting**. Otherwise, it is called the **mis-matched setting**.

In these notes, we focus on the **Bayes-optimal setting**, so that the committee machine estimation problem for the student committee machine (within the same hypothesis class and with weights $\mathbf{W} \in \mathbb{R}^{d \times K}$) reduces to learn the teacher rule generated by the ground truth weights \mathbf{W}^* from the supervised dataset $\{\mathbf{X}, \mathbf{y}\}$.

Remark 1.2. Committee machines are a simple vectorized generalization of **GLM**, whose estimation is performed simultaneously with $K \geq 1$.

1.2 Notations

- $\mathbf{X} \sim \mathbf{P}_{\mathbf{x}}(\mathbf{X}) = \prod_{i=1, \mu=1}^{d, n} \mathbf{P}_{\mathbf{x}}(x_{\mu i})$
- $\mathbf{W} \in \mathbb{R}^{d \times K}$ is the matrix of weights of the first layer with prior distribution:

$$\mathbf{P}_{\mathbf{w}}(\mathbf{W}) = \prod_{i=1}^d \mathbf{P}_{\mathbf{w}}(\mathbf{w}_i)$$
- $\mathbf{y} \in \mathbb{R}^n$ is a set of n scalar observations generated according to (1).

- Indices $\mu \in \llbracket n \rrbracket$ and $i \in \llbracket d \rrbracket$ correspond respectively to data samples and input dimensions.
- $\alpha \equiv \frac{n}{d}$ with $d, n \rightarrow \infty$ and in the limit where $\alpha = \Theta(1)$
- ξ will denote later a Gaussian vectorial variable such that $\xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$

2 Bayesian inference and factor graph

Both **Minimum Mean-Squared Error (MMSE)** and **Maximum A Posteriori (MAP)** estimations boil down to the analysis of the posterior distribution $P(\mathbf{W}|\mathbf{y}, \mathbf{X})$ expressed by the Bayes rule

$$P(\mathbf{W}|\mathbf{y}, \mathbf{X}) = \frac{\mathbb{P}(\mathbf{y}|\mathbf{W}, \mathbf{X}) \mathbb{P}(\mathbf{W})}{P(\mathbf{y}, \mathbf{X})} = \frac{P_{\text{out}}(\mathbf{y}|\mathbf{W}, \mathbf{X}) P_{\text{w}}(\mathbf{W})}{\mathcal{Z}_d(\{\mathbf{y}, \mathbf{X}\})}. \quad (2)$$

The joint distribution is also called the *partition function* $P(\mathbf{y}, \mathbf{X}) \equiv \mathcal{Z}_d(\{\mathbf{y}, \mathbf{X}\})$. To connect with the statistical physics formalism, we introduce the corresponding Hamiltonian, for separable distributions $P_{\text{out}}, P_{\text{w}}$ along one dimension, by

$$\begin{aligned} \mathcal{H}_d(\mathbf{W}, \{\mathbf{y}, \mathbf{X}\}) &= -\log P_{\text{out}}(\mathbf{y}|\mathbf{W}, \mathbf{X}) - \log P_{\text{w}}(\mathbf{W}), \\ &= -\sum_{\mu=1}^n \log P_{\text{out}}(y_{\mu}|\mathbf{W}, \mathbf{x}_{\mu}) - \sum_{i=1}^d P_{\text{w}}(\mathbf{w}_i). \end{aligned}$$

The spin variables represent the weights of the model $\mathbf{W} \in \mathbb{R}^{d \times K}$ and they interact through the random dataset $\{\mathbf{y}, \mathbf{X}\}$ that plays the role of the quenched exchange interactions. However here, the interactions are *fully connected*, meaning that each variable $\mathbf{w}_i \in \mathbb{R}^K$ is connected to every other spin $\{\mathbf{w}_j\}_{j \in \partial j \setminus i}$ as represented in the factor graph in Fig. 2.

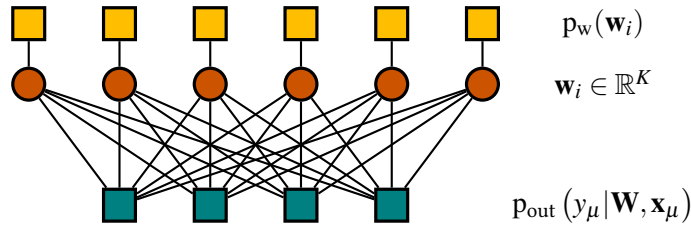


Figure 2 – Factor graph corresponding to the committee machines hypothesis class. The vectorial variables to infer \mathbf{w}_i are fully connected through the quenched disorder $\mathbf{y} \sim P_{\text{out}}^*(.)$ and each variable follow a one-body interaction with a separable prior distribution $p_{\text{w}}(\mathbf{w}_i)$.

The partition function at inverse temperature β is therefore defined by

$$\begin{aligned}\mathcal{Z}_d(\{\mathbf{y}, \mathbf{X}\}; \beta) &\equiv P(\mathbf{y}, \mathbf{X}) = \int_{\mathbb{R}^{d \times K}} d\mathbf{W} e^{-\beta \mathcal{H}_d(\mathbf{W}, \{\mathbf{y}, \mathbf{X}\})} \\ &= \int_{\mathbb{R}^{d \times K}} d\mathbf{w} e^{\beta (\log P_{\text{out}}(\mathbf{y}|\mathbf{W}, \mathbf{X}) + \log P_{\mathbf{w}}(\mathbf{W}))} \\ &= \int_{\mathbb{R}^{d \times K}} d\mathbf{w} p_{\text{out}}(\mathbf{y}|\mathbf{W}, \mathbf{X}) p_{\mathbf{w}}(\mathbf{W}),\end{aligned}\tag{3}$$

and can be exactly mapped to Bayesian estimation for $\beta = 1$. In the context of ERM, MAP estimation can be analyzed by taking the limit $\beta \rightarrow \infty$.

In the considered modern high-dimensional regime with $d \rightarrow \infty$, $n \rightarrow \infty$, $\alpha = n/d = \Theta(1)$ and $K = \Theta(1)$, we are interested in computing the *free entropy* Φ averaged over the input data \mathbf{X} and teacher weights \mathbf{W}^* , or equivalently over the output labels \mathbf{y} generated from it, defined as

$$\Phi(\alpha) \equiv \lim_{d \rightarrow \infty} \frac{1}{d} \mathbb{E}_{\mathbf{y}, \mathbf{X}} [\log \mathcal{Z}_d(\mathbf{y}, \mathbf{X})].\tag{4}$$

Remark 2.1. • The distribution $P(\mathbf{w}|\mathbf{y}; \mathbf{X})$ is often intractable in the limit $d \rightarrow \infty$

- Moreover it is still hard to sample efficiently in the high-dimensional regime

To circumvent this issue:

- We can focus only on marginals $P(\mathbf{w}_i; \mathbf{y}, \mathbf{X})$ which can be estimated with **Belief Propagation (BP)** equations and **Approximate Message Passing (AMP) algorithms**, as illustrated in Sec. 4.
- We may also try to compute directly the free entropy: we present the replica method in Sec. 3 that allows to compute the above average over the random dataset $\{\mathbf{y}, \mathbf{X}\}$, that plays the role of the quenched disorder in usual spin glasses. We show the computation for the more involved committee machine model class and generalization of the GLM class, only for iid data.
- We try to give an idea how these two methods relate and we show they are complementary and consistent.

Remark 2.2. The cumbersome computation for non iid data can be performed as well and lead to more complex expressions and has been performed in particular in [15] in the case of the GLM.

3 Replica computation

3.1 Heuristic derivation

In this section, we present the heuristic derivation of the replica formula in the context of the *committee machine*. This computation is necessary to properly guess the formula that can then be proven using the adaptive interpolation method. The reader interested in the replica approach to neural networks and the committee machine is invited to look as well to some of the classical papers [16, 17, 18, 19, 20, 21].

We present here the replica computation of the averaged free entropy $\Phi(\alpha)$ in eq. (4) for arbitrary *student* prior and channel distributions P_w, P_{w^*} and P_{out}, P_{out^*} , so that the computation remains valid for both the Bayes-optimal and mismatched settings.

3.1.1 Replica trick

The average in eq. (4) is intractable in general, and the computation relies on the so called *replica trick*, that consists in applying the identity

$$\mathbb{E}_{\mathbf{y}, \mathbf{X}} \left[\lim_{d \rightarrow \infty} \frac{1}{d} \log \mathcal{Z}_d(\mathbf{y}, \mathbf{X}) \right] = \lim_{r \rightarrow 0} \left[\lim_{d \rightarrow \infty} \frac{1}{d} \frac{\partial \log \mathbb{E}_{\mathbf{y}, \mathbf{X}} [\mathcal{Z}_d(\mathbf{y}, \mathbf{X})^r]}{\partial r} \right]. \quad (5)$$

The replica trick has been used in a series of previous works to compute the free energy density of GLM for separable distributions [22] and has been rigorous proven in this case by [23].

Eq. (5) is interesting in the sense that it reduces the intractable average to the computation of the moments of the averaged partition function, which are easier quantities to compute. Note that for $r \in \mathbb{N}$, $\mathcal{Z}_d(\mathbf{y}, \mathbf{X})^r = \prod_{a=1}^r \mathcal{Z}_d(\mathbf{y}, \mathbf{X})$ represents the partition function of r identical non-interacting copies of the initial system, called *replicas*. Taking the quenched average over the disorder will then correlate the replicas, before taking the number of replicas $r \rightarrow 0$. Therefore, we assume there exists an analytical continuation so that $r \in \mathbb{R}$ and the limit is well defined. Finally, notice that we exchanged the order of the limits $r \rightarrow 0$ and $d \rightarrow \infty$. These technicalities are crucial points but are not rigorously justified and we will ignore them in the rest of the computation.

First, in order to decouple the contributions of the channel P_{out} and the prior P_w , we introduce the variable $\mathbf{Z} = \frac{1}{\sqrt{d}} \mathbf{XW}$ and a Dirac-delta integral:

$$\mathcal{Z}_d(\mathbf{y}, \mathbf{X}) = \int_{\mathbb{R}^{n \times K}} d\mathbf{z} p_{out}(\mathbf{y}|\mathbf{Z}) \int_{\mathbb{R}^{d \times K}} d\mathbf{w} p_w(\mathbf{W}) \delta \left(\mathbf{Z} - \frac{1}{\sqrt{d}} \mathbf{XW} \right). \quad (6)$$

Thus the replicated partition function for an integer $r \in \mathbb{N}$ in eq. (5) can be written as

$$\begin{aligned} & \mathbb{E}_{\mathbf{y}, \mathbf{X}} [\mathcal{Z}_d(\mathbf{y}, \mathbf{X})^r] \\ &= \mathbb{E}_{\mathbf{y}, \mathbf{X}} \left[\prod_{a=1}^r \int_{\mathbb{R}^n} d\mathbf{Z}^a p_{\text{out}^a}(\mathbf{y}|\mathbf{Z}^a) \right. \end{aligned} \quad (7)$$

$$\begin{aligned} & \quad \times \int_{\mathbb{R}^{d \times K}} d\mathbf{W}^a p_{\mathbf{w}^a}(\mathbf{W}^a) \delta\left(\mathbf{Z}^a - \frac{1}{\sqrt{d}} \mathbf{X} \mathbf{W}^a\right) \Big] \\ &= \mathbb{E}_{\mathbf{X}} \int_{\mathbb{R}^n} d\mathbf{y} \int_{\mathbb{R}^{n \times K}} d\mathbf{Z}^* p_{\text{out}^*}(\mathbf{y}|\mathbf{Z}^*) \end{aligned} \quad (8)$$

$$\begin{aligned} & \quad \times \int_{\mathbb{R}^{d \times K}} d\mathbf{W}^* p_{\mathbf{w}^*}(\mathbf{W}^*) \delta\left(\mathbf{Z}^* - \frac{1}{\sqrt{d}} \mathbf{X} \mathbf{W}^*\right) \\ & \quad \times \left[\prod_{a=1}^r \int_{\mathbb{R}^{n \times K}} d\mathbf{Z}^a p_{\text{out}^a}(\mathbf{y}|\mathbf{Z}^a) \int_{\mathbb{R}^{d \times K}} d\mathbf{W}^a p_{\mathbf{w}^a}(\mathbf{W}^a) \delta\left(\mathbf{Z}^a - \frac{1}{\sqrt{d}} \mathbf{X} \mathbf{W}^a\right) \right] \\ &= \int_{\mathbb{R}^n} d\mathbf{y} \prod_{a=0}^r \int_{\mathbb{R}^{n \times K}} d\mathbf{Z}^a p_{\text{out}^a}(\mathbf{y}|\mathbf{Z}^a) \int_{\mathbb{R}^{d \times K}} d\mathbf{W}^a p_{\mathbf{w}^a}(\mathbf{W}^a) \end{aligned} \quad (9)$$

$$\times \underbrace{\mathbb{E}_{\mathbf{X}} \prod_{a=0}^r \delta\left(\mathbf{Z}^a - \frac{1}{\sqrt{d}} \mathbf{X} \mathbf{W}^a\right)}_{(I)}.$$

Note that the average over \mathbf{y} is equivalent to the one over the ground truth vector \mathbf{W}^* in the case of a *teacher-student*, which can be conveniently grouped with the other terms by just extending the replica indices and considering it as a new replica \mathbf{W}^0 with index $a = 0$, leading to a total of $r + 1$ replicas.

3.1.2 Average over the i.i.d input data \mathbf{X}

Remains to compute the average over \mathbf{X} in the term (I). We suppose that inputs are drawn from an iid distribution, for example a Gaussian $P_{\mathbf{x}}(\mathbf{x}) = \mathcal{N}_{\mathbf{x}}(\mathbf{0}, \mathbf{I}_d)$. More precisely, for $(i, j) \in \llbracket d \rrbracket^2$, $(\mu, \nu) \in \llbracket n \rrbracket^2$, $\mathbb{E}_{\mathbf{X}} [x_{\mu i} x_{\nu j}] = \delta_{\mu \nu} \delta_{ij}$.

By definition, the average in (I) defines the probability density $p_{\mathbf{Z}^a}(\mathbf{Z}^a)$ and as $\forall k \in \llbracket K \rrbracket, \forall \mu \in \llbracket n \rrbracket$, $z_{\mu k}^a = \frac{1}{\sqrt{d}} \sum_{i=1}^d x_{\mu i} w_{ik}^a$ is the sum of iid random variables, the CLT insures that in the thermodynamic limit $d \rightarrow \infty$, $z_{\mu k}^a$ follows a Gaussian multivariate distribution, with first moments given by:

$$\begin{aligned} \mathbb{E}_{\mathbf{X}} [z_{\mu k}^a] &= \frac{1}{\sqrt{d}} \sum_{i=1}^d \mathbb{E}_{\mathbf{X}} [x_{\mu i}] w_{ik}^a = 0 \\ \mathbb{E}_{\mathbf{X}} [z_{\mu k}^a z_{\nu k'}^b] &= \frac{1}{d} \sum_{ij} \mathbb{E}_{\mathbf{X}} [x_{\mu i} x_{\nu j}] w_{ik}^a w_{jk'}^b = \frac{1}{d} \sum_{ij} \delta_{ij} w_{ik}^a w_{jk'}^b \delta_{\mu \nu} \equiv \delta_{\mu \nu} Q_{bk'}^{ak}. \end{aligned} \quad (10)$$

Notice that averaging over the quenched disorder introduced correlations between replicas, which were initially independent, described by the symmetric

overlap matrix $\{Q_{bk'}^{ak}\}_{kk'}$ of size $(r+1)K \times (r+1)K$. This matrix order parameter measures the correlations between the replicated matrices $\{\mathbf{W}^a\}_{a=0}^r$ and is formally defined by

$$\mathbf{Q}(\{\mathbf{W}^a\}_{a=0}^r) \equiv \left(\frac{1}{d} \sum_{i=1}^d w_{ik}^a w_{ik'}^b \right)_{k,k'=1..K}^{a,b=0..r},$$

such that $\forall(a,b) \in \llbracket 0:r \rrbracket^2$, $\mathbf{Q}^{ab} \in \mathbb{R}^{K \times K}$. Therefore, again by the CLT, in the limit $d \rightarrow \infty$, the hidden variable $\mathbf{Z}^a \in \mathbb{R}^{n \times K}$ converges in distribution to the multivariate distribution

$$\mathbb{P}_{\mathbf{Z}^a}(\mathbf{Z}^a | \mathbf{Q}) = \exp \left[-\frac{1}{2} \sum_{\mu=1}^n \sum_{a,b=0}^r \sum_{k,k'=1}^K z_{\mu k}^a z_{\mu k'}^b (\mathbf{Q}^{-1})_{kk'}^{ab} \right] / (\det(2\pi\mathbf{Q}))^{\frac{n}{2}}. \quad (11)$$

Inserting this back in the replicated partition function finally writes

$$\mathbb{E}_{\mathbf{y}, \mathbf{X}} [\mathcal{Z}_d(\mathbf{y}, \mathbf{X})^r] = \int_{\mathbb{R}^n} d\mathbf{y} \prod_{a=0}^r \int_{\mathbb{R}^{n \times K}} d\mathbf{Z}^a \mathbb{P}_{\text{out}^a}(\mathbf{y} | \mathbf{Z}^a) \mathbb{P}_{\mathbf{Z}^a}(\mathbf{Z}^a | \mathbf{Q}) \int_{\mathbb{R}^{d \times K}} d\mathbf{W}^a \mathbb{P}_{\mathbf{W}^a}(\mathbf{W}^a) \quad (12)$$

3.1.3 Fourier representation

Next we introduce the change of variable for the new order parameter \mathbf{Q}^{ab} with a Dirac- δ distribution and its Fourier representation. For a variable $x \in \mathbb{R}$, the distribution $\delta(x)$ can be written as an integral over a purely imaginary parameter \hat{x} :

$$\delta(x) = \frac{1}{2i\pi} \int_{i\mathbb{R}} d\hat{x} e^{-\hat{x}x}.$$

Applying the above identity to the change of variable, we obtain

$$\begin{aligned} 1 &= \int_{\mathbb{R}^{(K \times (r+1))^2}} d\mathbf{Q} \prod_{0 \leq a \leq b \leq r, 1 \leq k, k' \leq K} \delta \left(dQ_{kk'}^{ab} - \sum_{i=1}^d w_{ik}^a w_{ik'}^b \right) \\ &\propto \int_{\mathbb{R}^{(K \times (r+1))^2}} d\mathbf{Q} d\hat{\mathbf{Q}} \exp \left(-d \sum_{a=0}^r \sum_{k,k'}^K Q_{kk'}^{aa} \hat{Q}_{kk'}^{aa} - \frac{d}{2} \sum_{a \neq b}^r \sum_{k,k'}^K Q_{kk'}^{ab} \hat{Q}_{kk'}^{ab} \right) \\ &\quad \times \exp \left(\sum_{a=0}^r \sum_{k,k'}^K \hat{Q}_{kk'}^{aa} w_k^a w_{k'}^a + \frac{1}{2} \sum_{a \neq b}^r \sum_{k,k'}^K \hat{Q}_{kk'}^{ab} w_k^a w_{k'}^b \right), \end{aligned} \quad (13)$$

that involves a new ad-hoc purely imaginary matrix parameter $\hat{\mathbf{Q}} \in i\mathbb{R}^{(K \times (r+1))^2}$. Finally, multiplying the replicated partition function by 1, using the Cauchy theorem and rotating the integration, it becomes an integral over the symmetric

matrices $\mathbf{Q} \in \mathbb{R}^{(K \times r+1)^2}$ and $\hat{\mathbf{Q}} \in \mathbb{R}^{(K \times r+1)^2}$

$$\begin{aligned}
& \mathbb{E}_{\mathbf{y}, \mathbf{X}} [\mathcal{Z}_d(\mathbf{y}, \mathbf{X})^r] \\
&= \int \int_{\mathbb{R}^{(K \times r+1)^2}} d\mathbf{Q} d\hat{\mathbf{Q}} \exp \left(-d \sum_{a=0}^r \sum_{k,k'}^K \mathcal{Q}_{kk'}^{aa} \hat{\mathcal{Q}}_{kk'}^{aa} - \frac{d}{2} \sum_{a \neq b}^r \sum_{k,k'}^K \mathcal{Q}_{kk'}^{ab} \hat{\mathcal{Q}}_{kk'}^{ab} \right) \\
&\quad \times \exp \left(\sum_{a=0}^r \sum_{k,k'}^K \hat{\mathcal{Q}}_{kk'}^{aa} w_k^a w_{k'}^a + \frac{1}{2} \sum_{a \neq b}^r \sum_{k,k'}^K \hat{\mathcal{Q}}_{kk'}^{ab} w_k^a w_{k'}^b \right) \\
&\int_{\mathbb{R}^n} d\mathbf{y} \prod_{a=0}^r \int_{\mathbb{R}^{n \times K}} d\mathbf{Z}^a p_{\text{out}^a}(\mathbf{y} | \mathbf{Z}^a) p_{\mathbf{Z}^a}(\mathbf{Z}^a | \mathbf{Q}) \int_{\mathbb{R}^{d \times K}} d\mathbf{W}^a p_{\mathbf{W}^a}(\mathbf{W}^a) \\
&= \int \int_{\mathbb{R}^{(K \times r+1)^2}} d\mathbf{Q} d\hat{\mathbf{Q}} \exp \left(-d \sum_{a=0}^r \sum_{k,k'}^K \mathcal{Q}_{kk'}^{aa} \hat{\mathcal{Q}}_{kk'}^{aa} - \frac{d}{2} \sum_{a \neq b}^r \sum_{k,k'}^K \mathcal{Q}_{kk'}^{ab} \hat{\mathcal{Q}}_{kk'}^{ab} \right) \\
&\quad \times \exp \left(\sum_{a=0}^r \sum_{k,k'}^K \hat{\mathcal{Q}}_{kk'}^{aa} w_k^a w_{k'}^a + \frac{1}{2} \sum_{a \neq b}^r \sum_{k,k'}^K \hat{\mathcal{Q}}_{kk'}^{ab} w_k^a w_{k'}^b \right) \\
&\left[\int_{\mathbb{R}} d\mathbf{y} \prod_{a=0}^r \int_{\mathbb{R}^K} d\mathbf{z}^a p_{\text{out}^a}(\mathbf{y} | \mathbf{z}^a) p_{\mathbf{Z}^a}(\mathbf{z}^a | \mathbf{Q}) \right]^n \left[\prod_{a=0}^r \int_{\mathbb{R}^K} d\mathbf{w}^a p_{\mathbf{W}^a}(\mathbf{w}^a) \right]^d \\
&\simeq \int \int_{\mathbb{R}^{(K \times r+1)^2}} d\mathbf{Q} d\hat{\mathbf{Q}} e^{d\Phi^{(r)}(\mathbf{Q}, \hat{\mathbf{Q}})},
\end{aligned}$$

where in the last step, we used a Laplace method [24] and omitted the sub-leading factors in the thermodynamic limit $d \rightarrow \infty$ to evaluate it as a function of the free entropy potential defined by

$$\begin{aligned}
\Phi^{(r)}(\mathbf{Q}, \hat{\mathbf{Q}}) &\equiv - \sum_{a=0}^r \sum_{k,k'}^K \mathcal{Q}_{kk'}^{aa} \hat{\mathcal{Q}}_{kk'}^{aa} - \frac{1}{2} \sum_{a \neq b}^r \sum_{k,k'}^K \mathcal{Q}_{kk'}^{ab} \hat{\mathcal{Q}}_{kk'}^{ab} \\
&\quad + \log \Psi_{\mathbf{W}}^{(r)}(\hat{\mathbf{Q}}) + \alpha \log \Psi_{\text{out}}^{(r)}(\mathbf{Q}), \\
\Psi_{\mathbf{W}}^{(r)}(\hat{\mathbf{Q}}) &\equiv \prod_{a=0}^r \int_{\mathbb{R}^K} d\mathbf{w}^a p_{\mathbf{W}^a}(\mathbf{w}^a) \\
&\quad \times \exp \left(\sum_{a=0}^r \sum_{k,k'}^K \hat{\mathcal{Q}}_{kk'}^{aa} w_k^a w_{k'}^a + \frac{1}{2} \sum_{a \neq b}^r \sum_{k,k'}^K \hat{\mathcal{Q}}_{kk'}^{ab} w_k^a w_{k'}^b \right), \\
\Psi_{\text{out}}^{(r)}(\mathbf{Q}) &\equiv \prod_{a=0}^r \int_{\mathbb{R}} d\mathbf{y} \int_{\mathbb{R}^K} d\mathbf{z}^a p_{\text{out}^a}(\mathbf{y} | \mathbf{z}^a) p_{\mathbf{Z}^a}(\mathbf{z}^a | \mathbf{Q}),
\end{aligned} \tag{14}$$

and where we decoupled the variable $\mathbf{Z}^a \in \mathbb{R}^{n \times K}$ and $\mathbf{W}^a \in \mathbb{R}^{d \times K}$ along the rows

$$\begin{aligned} p_{\text{out}^a}(\mathbf{y}|\mathbf{Z}^a) &= \prod_{\mu=1}^n p_{\text{out}^a}(y_\mu|\mathbf{z}_\mu^a), \text{ with } \mathbf{z}_\mu^a \in \mathbb{R}^K, \\ p_{\mathbf{Z}^a}(\mathbf{Z}^a|\mathbf{Q}) &= \prod_{\mu=1}^n p(\mathbf{z}_\mu^a|\mathbf{Q}), \\ p_{\mathbf{W}^a}(\mathbf{W}^a) &= \prod_{i=1}^d p_{\mathbf{W}^a}(\mathbf{w}_i^a), \text{ with } \mathbf{w}_i^a \in \mathbb{R}^K, \\ p_{\mathbf{Z}^a}(\mathbf{Z}^a|\mathbf{Q}) &= \exp \left[-\frac{1}{2} \sum_{a,b=0}^r \sum_{k,k'=1}^K z_k^a z_{k'}^b (\mathbf{Q}^{-1})_{kk'}^{ab} \right] / (\det(2\pi\mathbf{Q}))^{\frac{1}{2}}. \end{aligned}$$

Note that the averaged replicated partition function of this fully connected model can be expressed as a saddle point equation only because distributions $P_{\text{out}}, P_{\text{out}^*}$ and $P_{\mathbf{W}}, P_{\mathbf{W}^*}$ are separable so that a pre-factor scaling with the system size d dominates the exponential distribution. Finally, switching the two limits $r \rightarrow 0$ and $d \rightarrow \infty$, the quenched free entropy Φ simplifies as a saddle point equation

$$\Phi(\alpha) = \mathbf{extr}_{\mathbf{Q}, \hat{\mathbf{Q}}} \left\{ \lim_{r \rightarrow 0} \frac{\partial \Phi^{(r)}(\mathbf{Q}, \hat{\mathbf{Q}})}{\partial r} \right\}, \quad (15)$$

over symmetric matrices $\mathbf{Q} \in \mathbb{R}^{(K \times r+1)^2}$ and $\hat{\mathbf{Q}} \in \mathbb{R}^{(K \times r+1)^2}$.

To summarize, we managed to get rid of the original high-dimensional integrals and replace them by an optimization in the space of matrices, which, in this form, is still intractable. We not only have to search in the space of $(r+1) \times (r+1)$ matrices to find the extremiser of $\Phi^{(r)}$, but we also need to compute the limit $r \rightarrow 0^+$. In the following we will assume a simple Ansatz for these matrices in order to first obtain an analytic expression in r before taking the derivative with respect to r .

3.2 Replica Symmetric free entropy

Our goal is to express the functional $\Phi^{(r)}(\mathbf{Q}, \hat{\mathbf{Q}})$ appearing in the free entropy as an analytical function of r , in order to perform the replica trick.

3.2.1 Replica symmetric ansatz

To do so, we will assume that the extremum of $\Phi^{(r)}$ is attained at a point in $\mathbf{Q}, \hat{\mathbf{Q}}$ space such that a *replica symmetry* property is verified. More concretely, we

assume:

$$\begin{aligned}
\exists \mathbf{Q} \in \mathbb{R}^{K \times K} \text{ s.t. } & \forall a \in \llbracket 0 : r \rrbracket \quad \forall (k, k') \in \llbracket K \rrbracket^2 \quad Q_{kk'}^{aa} = Q_{kk'}, \\
& \exists \mathbf{Q}^* \in \mathbb{R}^{K \times K} \text{ s.t. } & \forall (k, k') \in \llbracket K \rrbracket^2 \quad Q_{kk'}^{00} = Q_{kk'}^*, \\
\exists \mathbf{q} \in \mathbb{R}^{K \times K} \text{ s.t. } & \forall (a < b) \in \llbracket 0 : r \rrbracket^2 \quad \forall (k, k') \in \llbracket K \rrbracket^2 \quad Q_{kk'}^{ab} = q_{kk'}, \\
\exists \mathbf{m} \in \mathbb{R}^{K \times K} \text{ s.t. } & \forall a \in \llbracket 0 : r \rrbracket \quad \forall (k, k') \in \llbracket K \rrbracket^2 \quad Q_{kk'}^{0a} = m_{kk'},
\end{aligned} \tag{16}$$

and similarly for the ad-hoc parameter

$$\begin{aligned}
\exists \hat{\mathbf{Q}} \in \mathbb{R}^{K \times K} \text{ s.t. } & \forall a \in \llbracket 0 : r \rrbracket \quad \forall (k, k') \in \llbracket K \rrbracket^2 \quad \hat{Q}_{kk'}^{aa} = -\frac{1}{2} \hat{Q}_{kk'}, \\
& \exists \hat{\mathbf{Q}}^* \in \mathbb{R}^{K \times K} \text{ s.t. } & \forall (k, k') \in \llbracket K \rrbracket^2 \quad \hat{Q}_{kk'}^{00} = \hat{Q}_{kk'}^*, \\
\exists \hat{\mathbf{q}} \in \mathbb{R}^{K \times K} \text{ s.t. } & \forall (a < b) \in \llbracket 0 : r \rrbracket^2 \quad \forall (k, k') \in \llbracket K \rrbracket^2 \quad \hat{Q}_{kk'}^{ab} = \hat{q}_{kk'}, \\
\exists \hat{\mathbf{m}} \in \mathbb{R}^{K \times K} \text{ s.t. } & \forall a \in \llbracket 0 : r \rrbracket \quad \forall (k, k') \in \llbracket K \rrbracket^2 \quad \hat{Q}_{kk'}^{0a} = \hat{m}_{kk'}.
\end{aligned} \tag{17}$$

The factor $-\frac{1}{2}$ is not necessary bu useful to recover commonly used formulations. This Ansatz can be represented by symmetric RS matrices $\mathbf{Q}^{(\text{rs})} \in \mathbb{R}^{(K \times r + 1)^2}$ and $\hat{\mathbf{Q}}^{(\text{rs})} \in \mathbb{R}^{(K \times r + 1)^2}$

$$\mathbf{Q}^{(\text{rs})} = \begin{pmatrix} \mathbf{Q}^* & \mathbf{m} & \cdots & \mathbf{m} \\ \mathbf{m}^\top & \mathbf{Q} & \mathbf{q} & \cdots \\ \vdots & \mathbf{q} & \ddots & \mathbf{q} \\ \mathbf{m}^\top & \cdots & \mathbf{q} & \mathbf{Q} \end{pmatrix} \text{ and } \hat{\mathbf{Q}}^{(\text{rs})} = \begin{pmatrix} \hat{\mathbf{Q}}^* & \hat{\mathbf{m}} & \cdots & \hat{\mathbf{m}} \\ \hat{\mathbf{m}}^\top & -\frac{1}{2} \hat{\mathbf{Q}} & \hat{\mathbf{q}} & \cdots \\ \vdots & \hat{\mathbf{q}} & \ddots & \hat{\mathbf{q}} \\ \hat{\mathbf{m}}^\top & \cdots & \hat{\mathbf{q}} & -\frac{1}{2} \hat{\mathbf{Q}} \end{pmatrix}, \tag{18}$$

where the *overlap* parameters may be reinterpreted as the scalar product between the replicas

$$\forall (a, b) \in \llbracket r \rrbracket^2, \quad \mathbf{q} = \frac{1}{d} \mathbf{W}^{a\top} \mathbf{W}^b,$$

the self-overlap of each replica

$$\forall a \in \llbracket r \rrbracket, \quad \mathbf{Q} = \frac{1}{d} \mathbf{W}^{a\top} \mathbf{W}^a,$$

the scalar product with the ground truth

$$\forall a \in \llbracket r \rrbracket, \quad \mathbf{m} = \frac{1}{d} \mathbf{W}^{*\top} \mathbf{W}^a,$$

and the second moment of the ground truth distribution

$$\mathbf{Q}^* = \frac{1}{d} \mathbf{W}^{*\top} \mathbf{W}^*.$$

The above Ansatz simplifies in the scalar GLM case with $K = 1$ to $q = \frac{1}{d} \mathbf{w}^a \cdot \mathbf{w}^b$ for $a \neq b$, a norm $Q = \frac{1}{d} \|\mathbf{w}^a\|_2^2$, an overlap with the ground truth $m = \frac{1}{d} \mathbf{w}^a \cdot \mathbf{w}^*$ and a second moment $Q^* = \frac{1}{d} \|\mathbf{w}^*\|_2^2$.

3.2.2 Decomposition of the free entropy functional

Let's compute separately the terms involved in the functional $\Phi^{(r)}(\mathbf{Q}, \hat{\mathbf{Q}})$ in (14) by applying this Ansatz: the first is a trace term, the second term $\Psi_w^{(r)}$ depends on the prior distributions P_w, P_{w^*} and finally the third term $\Psi_{\text{out}}^{(r)}$ depends on the channel distributions $P_{\text{out}^*}, P_{\text{out}}$.

Trace term The trace term in (14) can be easily computed at the RS fixed point and takes the following form

$$\begin{aligned} & - \sum_{a=0}^r \sum_{k,k'}^K \mathcal{Q}_{kk'}^{aa} \hat{\mathcal{Q}}_{kk'}^{aa} - \frac{1}{2} \sum_{a \neq b}^r \sum_{k,k'}^K \mathcal{Q}_{kk'}^{ab} \hat{\mathcal{Q}}_{kk'}^{ab} \Big|_{\text{rs}} \\ & = -\text{Tr}(\mathbf{Q}^* \hat{\mathbf{Q}}^*) + \frac{1}{2} r \text{Tr}(\mathbf{Q} \hat{\mathbf{Q}}) - r \text{Tr}(\mathbf{m} \hat{\mathbf{m}}) - \frac{r(r-1)}{2} \text{Tr}(\mathbf{q} \hat{\mathbf{q}}), \end{aligned} \quad (19)$$

and taking the derivative and the limit $r \rightarrow 0$ we obtain

$$\begin{aligned} & \lim_{r \rightarrow 0} \partial_r \left(- \sum_{a=0}^r \sum_{k,k'}^K \mathcal{Q}_{kk'}^{aa} \hat{\mathcal{Q}}_{kk'}^{aa} - \frac{1}{2} \sum_{a \neq b}^r \sum_{k,k'}^K \mathcal{Q}_{kk'}^{ab} \hat{\mathcal{Q}}_{kk'}^{ab} \right) \Big|_{\text{rs}} \\ & = \frac{1}{2} \text{Tr}(\mathbf{Q} \hat{\mathbf{Q}}) - \text{Tr}(\mathbf{m} \hat{\mathbf{m}}) + \frac{1}{2} \text{Tr}(\mathbf{q} \hat{\mathbf{q}}) \end{aligned} \quad (20)$$

Prior integral $\Psi_w^{(r)}$ Evaluated at the RS fixed point the quadratic form reads

$$\begin{aligned} & \sum_{a=0}^r \sum_{k,k'}^K \hat{\mathcal{Q}}_{kk'}^{aa} w_k^a w_{k'}^a + \frac{1}{2} \sum_{a \neq b}^r \sum_{k,k'}^K \hat{\mathcal{Q}}_{kk'}^{ab} w_k^a w_{k'}^b \\ & = \mathbf{w}^{*\top} \hat{\mathbf{Q}}^* \mathbf{w}^* + \sum_{a=1}^r \mathbf{w}^{*\top} \hat{\mathbf{m}} \mathbf{w}^a - \frac{1}{2} \sum_{a=1}^r \mathbf{w}^{a\top} \hat{\mathbf{Q}} \mathbf{w}^a + \frac{1}{2} \sum_{1 \leq a \neq b \leq r} \mathbf{w}^{a\top} \hat{\mathbf{q}} \mathbf{w}^b \\ & = \mathbf{w}^{*\top} \hat{\mathbf{Q}}^* \mathbf{w}^* + \sum_{a=1}^r \mathbf{w}^{*\top} \hat{\mathbf{m}} \mathbf{w}^a - \frac{1}{2} \sum_{a=1}^r \mathbf{w}^{a\top} (\hat{\mathbf{Q}} + \hat{\mathbf{q}}) \mathbf{w}^a + \frac{1}{2} \left(\sum_{a=1}^r \mathbf{w}^a \right)^\top \hat{\mathbf{q}} \left(\sum_{a=1}^r \mathbf{w}^a \right). \end{aligned}$$

Using a Hubbard-Stratonovich transformation:

Proposition 3.1 (Hubbard-Stratonovich transformation). *For $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and a symmetric positive definite matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, for all $\mathbf{x} \in \mathbb{R}^d$*

$$\mathbb{E}_{\boldsymbol{\xi}} \exp \left(\boldsymbol{\xi}^\top \mathbf{A}^{1/2} \mathbf{x} \right) = (2\pi)^{-d/2} \int_{\mathbb{R}^d} d\boldsymbol{\xi} e^{-\frac{1}{2} \mathbf{x}^\top \mathbf{x} + \boldsymbol{\xi}^\top \mathbf{A}^{1/2} \mathbf{x}} = e^{\frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x}}, \quad (21)$$

the prior integral can be further simplified

$$\begin{aligned} & \Psi_w^{(r)}(\hat{\mathbf{Q}}) \Big|_{\text{rs}} = \int_{\mathbb{R}^{(r+1) \times K}} d\mathbf{W} p_w(\mathbf{W}) e^{\sum_{a=0}^r \sum_{k,k'}^K \hat{\mathcal{Q}}_{kk'}^{aa} w_k^a w_{k'}^a + \frac{1}{2} \sum_{a \neq b}^r \sum_{k,k'}^K \hat{\mathcal{Q}}_{kk'}^{ab} w_k^a w_{k'}^b} \\ & = \mathbb{E}_{\boldsymbol{\xi}, \mathbf{w}^* \sim P_{w^*}} \left[e^{\mathbf{w}^{*\top} \hat{\mathbf{Q}}^* \mathbf{w}^*} \mathbb{E}_{\mathbf{w} \sim P_w} \left[e^{\left(\mathbf{w}^\top \hat{\mathbf{m}} \mathbf{w}^* - \frac{1}{2} \mathbf{w}^\top (\hat{\mathbf{Q}} + \hat{\mathbf{q}}) \mathbf{w} + \mathbf{w}^\top \hat{\mathbf{q}}^{1/2} \boldsymbol{\xi} \right)} \right]^r \right]. \end{aligned} \quad (22)$$

Detailed computation of $(\mathbf{Q}^{(\text{rs})})^{-1}$ Let us focus on the matrix $\mathbf{Q}^{(\text{rs})}$ involved in the expression of $\Psi_{\text{out}}^{(r)}$ in (14). Given the replica symmetric ansatz, we must compute the determinant and inverse matrix of the matrix $\mathbf{\Sigma} \equiv \mathbf{Q}^{(\text{rs})}$.

We use tensor products notations \otimes for matrices living in $\mathbb{R}^{r \times r} \otimes \mathbb{R}^{K \times K}$. Let us recall some basic properties of the tensor product of two matrices $\mathbf{A} \in \mathbb{R}^{r \times r}$ and $\mathbf{B} \in \mathbb{R}^{K \times K}$:

$$(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1} \quad (23)$$

$$\det(\mathbf{A} \otimes \mathbf{B}) = \det(\mathbf{A})^K \det(\mathbf{B})^r \quad (24)$$

$$\forall \mathbf{C}, \mathbf{D} \quad (\mathbf{C} \otimes \mathbf{D})(\mathbf{A} \otimes \mathbf{B}) = (\mathbf{CA}) \otimes (\mathbf{DB}) \quad (25)$$

$$\text{Sp}(\mathbf{A} \otimes \mathbf{B}) = \{\lambda_i \mu_j \mid \lambda_i \in \text{Sp}(\mathbf{A}), \mu_j \in \text{Sp}(\mathbf{B})\} \quad (26)$$

Next we denote the *reduced* matrix

$$\mathbf{\Sigma}_r \equiv \begin{pmatrix} \mathbf{Q} & \mathbf{q} & \cdots & \mathbf{q} \\ \mathbf{q} & \mathbf{Q} & \cdots & \mathbf{q} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q} & \mathbf{q} & \cdots & \mathbf{Q} \end{pmatrix} = \mathbf{I}_r \otimes (\mathbf{Q} - \mathbf{q}) + \mathbf{J}_r \otimes \mathbf{q}, \quad (27)$$

where \mathbf{J}_r is the rank-1 matrix with all elements equal to 1. Finally, some basic linear algebra formulas allow us to reduce the problem of computing $\det(\mathbf{\Sigma})$ and $\mathbf{\Sigma}^{-1}$ to those of $\mathbf{\Sigma}_r$.

Calculation of \mathbf{S} We define the $K \times K$ matrix :

$$\mathbf{S} \equiv \mathbf{Q}^* - (\mathbf{1}_r^\top \otimes \mathbf{m}) \mathbf{\Sigma}_r^{-1} (\mathbf{1}_r \otimes \mathbf{m}^\top) \quad (28)$$

Let us compute \mathbf{S} :

$$\begin{aligned} \mathbf{S} &\equiv \mathbf{Q}^* - (\mathbf{1}_r^\top \otimes \mathbf{m}) \mathbf{\Sigma}_r^{-1} (\mathbf{1}_r \otimes \mathbf{m}^\top) \\ &= \mathbf{Q}^* - r \mathbf{m} (\mathbf{Q} - \mathbf{q})^{-1} \mathbf{m}^\top + r^2 \mathbf{m} (\mathbf{Q} + (r-1)\mathbf{q})^{-1} \mathbf{q} (\mathbf{Q} - \mathbf{q})^{-1} \mathbf{m}^\top \\ &= \mathbf{Q}^* - r \mathbf{m} (\mathbf{I}_K - n(\mathbf{Q} + (r-1)\mathbf{q})^{-1} \mathbf{q}) (\mathbf{Q} - \mathbf{q})^{-1} \mathbf{m}^\top \\ &= \mathbf{Q}^* - r \mathbf{m} (\mathbf{Q} + (r-1)\mathbf{q})^{-1} (\mathbf{Q} + (r-1)\mathbf{q} - r\mathbf{q}) (\mathbf{Q} - \mathbf{q})^{-1} \mathbf{m}^\top \\ &= \mathbf{Q}^* - r \mathbf{m} (\mathbf{Q} + (r-1)\mathbf{q})^{-1} \mathbf{m}^\top \end{aligned}$$

Calculation of $\mathbf{\Sigma}_r^{-1}$ $\mathbf{\Sigma}_r$ can be easily analyzed, thanks to a slight generalization of the Sherman-Morrison formula.

$$\mathbf{\Sigma}_r = \mathbf{I}_r \otimes (\mathbf{Q} - \mathbf{q}) + \mathbf{J}_r \otimes \mathbf{q} = \mathbf{I}_r \otimes (\mathbf{Q} - \mathbf{q}) + \mathbf{1}_r \mathbf{1}_r^\top \otimes \mathbf{q} \quad (29)$$

So that the determinant reads

$$\begin{aligned} \det(\mathbf{\Sigma}_r) &= \det(\mathbf{Q} - \mathbf{q})^r \det(\mathbf{I}_K + (\mathbf{1}_r^\top \otimes \mathbf{q})(\mathbf{I}_r \otimes (\mathbf{Q} - \mathbf{q}))^{-1} (\mathbf{1} \otimes \mathbf{I}_K)) \\ &= \det(\mathbf{Q} - \mathbf{q})^r \det(\mathbf{I}_K + r\mathbf{q}(\mathbf{Q} - \mathbf{q})^{-1}) \\ \det(\mathbf{\Sigma}_r) &= \det(\mathbf{Q} - \mathbf{q})^{r-1} \det(\mathbf{Q} + (r-1)\mathbf{q}), \end{aligned} \quad (30)$$

and the inverse

$$\mathbf{\Sigma}_r^{-1} = \mathbf{I}_r \otimes (\mathbf{Q} - \mathbf{q})^{-1} \quad (31)$$

$$\begin{aligned} & - (\mathbf{I}_r \otimes (\mathbf{Q} - \mathbf{q})^{-1}) (\mathbf{1}_r \otimes \mathbf{I}_K) (\mathbf{I}_K + r\mathbf{q}(\mathbf{Q} - \mathbf{q})^{-1})^{-1} (\mathbf{1}_r^\top \otimes \mathbf{q}) (\mathbf{I}_r \otimes (\mathbf{Q} - \mathbf{q})^{-1}) \\ & = \mathbf{I}_r \otimes (\mathbf{Q} - \mathbf{q})^{-1} - \mathbf{J}_r \otimes ((\mathbf{Q} + (r-1)\mathbf{q})^{-1} \mathbf{q}(\mathbf{Q} - \mathbf{q})^{-1}) \end{aligned} \quad (32)$$

Determinant of $\mathbf{Q}^{(\text{rs})}$ Therefore we obtain:

$$\begin{aligned} \det(\mathbf{\Sigma}) &= \det(\mathbf{Q}^{(\text{rs})}) = \det(\mathbf{\Sigma}_r) \det(\mathbf{S}) \\ &= \det(\mathbf{Q} - \mathbf{q})^{r-1} \det(\mathbf{Q} + (r-1)\mathbf{q}) \det(\mathbf{Q}^* - r\mathbf{m}(\mathbf{Q} + (r-1)\mathbf{q})^{-1} \mathbf{m}^\top) \end{aligned} \quad (33)$$

Inverse of $\mathbf{Q}^{(\text{rs})}$ The inverse of the block matrix reads:

$$\mathbf{\Sigma}^{-1} = \begin{pmatrix} \mathbf{S}^{-1} & -\mathbf{S}^{-1}(\mathbf{1}_r^\top \otimes \mathbf{m})\mathbf{\Sigma}_r^{-1} \\ -\mathbf{\Sigma}_r^{-1}(\mathbf{1}_r \otimes \mathbf{m}^\top)\mathbf{S}^{-1} & \mathbf{\Sigma}_r^{-1} + \mathbf{\Sigma}_r^{-1}(\mathbf{1}_r \otimes \mathbf{m}^\top)\mathbf{S}^{-1}(\mathbf{1}_r^\top \otimes \mathbf{m})\mathbf{\Sigma}_r^{-1} \end{pmatrix}$$

where each block can be computed separately. The top diagonal simply reads

$$\mathbf{S}^{-1} = (\mathbf{Q}^* - r\mathbf{m}(\mathbf{Q} + (r-1)\mathbf{q})^{-1} \mathbf{m}^\top)^{-1}$$

and the off-diagonal blocks read:

$$-\mathbf{\Sigma}_r^{-1}(\mathbf{1}_r \otimes \mathbf{m}^\top)\mathbf{S}^{-1} = \quad (34)$$

$$-\mathbf{1} \otimes (\mathbf{Q} + (r-1)\mathbf{q})^{-1} \mathbf{m}^\top (\mathbf{Q}^* - r\mathbf{m}(\mathbf{Q} + (r-1)\mathbf{q})^{-1} \mathbf{m}^\top)^{-1}$$

$$-\mathbf{S}^{-1}(\mathbf{1}_r^\top \otimes \mathbf{m})\mathbf{\Sigma}_r^{-1} = \quad (35)$$

$$-\mathbf{1}_r^\top \otimes ((\mathbf{Q}^* - r\mathbf{m}(\mathbf{Q} + (r-1)\mathbf{q})^{-1} \mathbf{m}^\top)^{-1} \mathbf{m}(\mathbf{Q} + (r-1)\mathbf{q})^{-1})]$$

The bottom diagonal block is a bit more involved :

$$\begin{aligned} & \mathbf{\Sigma}_r^{-1} + \mathbf{\Sigma}_r^{-1}(\mathbf{1}_r \otimes \mathbf{m}^\top)\mathbf{S}^{-1}(\mathbf{1}_r^\top \otimes \mathbf{m})\mathbf{\Sigma}_r^{-1} \\ & = \mathbf{I}_r \otimes (\mathbf{Q} - \mathbf{q})^{-1} - \mathbf{J}_r \otimes \mathbf{A} \end{aligned} \quad (36)$$

with

$$\mathbf{A} \equiv (\mathbf{Q} + (r-1)\mathbf{q})^{-1} \mathbf{q}(\mathbf{Q} - \mathbf{q})^{-1} \quad (37)$$

$$-(\mathbf{Q} + (r-1)\mathbf{q})^{-1} \mathbf{m}^\top (\mathbf{Q}^* - r\mathbf{m}(\mathbf{Q} + (r-1)\mathbf{q})^{-1} \mathbf{m}^\top)^{-1} \mathbf{m}(\mathbf{Q} + (r-1)\mathbf{q})^{-1}$$

Finally we observe that $\mathbf{\Sigma}^{-1}$ is also a block matrix of the form

$$\mathbf{\Sigma}^{-1} = (\mathbf{Q}^{(\text{rs})})^{-1} = \begin{bmatrix} \tilde{\mathbf{Q}}^* & \tilde{\mathbf{m}} & \cdots & \tilde{\mathbf{m}} \\ \tilde{\mathbf{m}}^\top & \tilde{\mathbf{Q}} & \tilde{\mathbf{q}} & \cdots \\ \vdots & \tilde{\mathbf{q}} & \ddots & \tilde{\mathbf{q}} \\ \tilde{\mathbf{m}}^\top & \cdots & \tilde{\mathbf{q}} & \tilde{\mathbf{Q}} \end{bmatrix} \quad (38)$$

with

$$\begin{aligned}
\tilde{\mathbf{Q}}^* &= (\mathbf{Q}^* - r\mathbf{m}(\mathbf{Q} + (r-1)\mathbf{q})^{-1}\mathbf{m}^\top)^{-1} \\
\tilde{\mathbf{m}} &= -(\mathbf{Q}^* - r\mathbf{m}(\mathbf{Q} + (r-1)\mathbf{q})^{-1}\mathbf{m}^\top)^{-1}\mathbf{m}(\mathbf{Q} + (r-1)\mathbf{q})^{-1} \\
\tilde{\mathbf{Q}} &= (\mathbf{Q} - \mathbf{q})^{-1} - (\mathbf{Q} + (r-1)\mathbf{q})^{-1}\mathbf{q}(\mathbf{Q} - \mathbf{q})^{-1} \\
&\quad + (\mathbf{Q} + (r-1)\mathbf{q})^{-1}\mathbf{m}^\top \\
&\quad \times (\mathbf{Q}^* - r\mathbf{m}(\mathbf{Q} + (r-1)\mathbf{q})^{-1}\mathbf{m}^\top)^{-1}\mathbf{m}(\mathbf{Q} + (r-1)\mathbf{q})^{-1} \\
\tilde{\mathbf{q}} &= \tilde{\mathbf{Q}} - (\mathbf{Q} - \mathbf{q})^{-1}
\end{aligned}$$

Channel integral $\Psi_{\text{out}}^{(r)}$ The quadratic form in $\mathbf{p}_{\mathbf{Z}^a}(\mathbf{z}^a|\mathbf{Q}^{(\text{rs})})$ reads

$$\begin{aligned}
& -\frac{1}{2} \sum_{a,b} \sum_{k,k'} z_k^a z_{k'}^b (\mathbf{Q}^{-1})_{kk'}^{ab} \\
& = -\frac{1}{2} \mathbf{z}^{*\top} \tilde{\mathbf{Q}}^* \mathbf{z} - \sum_{a=1}^r \mathbf{z}^{*\top} \tilde{\mathbf{m}} \mathbf{z}^a \\
& \quad - \frac{1}{2} \sum_{a=1}^r \mathbf{z}^{a\top} (\tilde{\mathbf{Q}} - \tilde{\mathbf{q}}) \mathbf{z}^a - \frac{1}{2} \left(\sum_a^r \mathbf{z}^a \right)^\top \tilde{\mathbf{q}} \left(\sum_a^r \mathbf{z}^a \right),
\end{aligned}$$

and using another Gaussian transformation, see Prop. 3.1, we finally obtain

$$\begin{aligned}
\Psi_{\text{out}}^{(r)}(\mathbf{Q}) \Big|_{\text{rs}} &= \int d\mathbf{y} \int_{\mathbb{R}^{(r+1) \times K}} d\mathbf{Z} \mathbf{p}_{\text{out}}(\mathbf{y}|\mathbf{Z}) \mathbf{p}(\mathbf{Z}|\mathbf{Q}) \\
&= \int d\mathbf{y} \int_{\mathbb{R}^{(r+1) \times K}} d\mathbf{Z} \mathbf{p}_{\text{out}}(\mathbf{y}|\mathbf{Z}) e^{-\frac{1}{2} \sum_{a,b=0}^r \sum_{k,k'=1}^K z_k^a z_{k'}^b (\mathbf{Q}^{-1})_{kk'}^{ab}} / \left(\det(2\pi\mathbf{Q})^{(\text{rs})} \right)^{\frac{1}{2}} \\
&= \int d\mathbf{y} \mathbb{E}_{\boldsymbol{\xi}} e^{-\frac{1}{2} \log(\det(2\pi\mathbf{Q}^{(\text{rs})}))} \times \int d\mathbf{z}^* \mathbf{p}_{\text{out}^*}(\mathbf{y}|\mathbf{z}^*) e^{-\frac{1}{2} \mathbf{z}^{*\top} \tilde{\mathbf{Q}}^* \mathbf{z}^*} \\
&\quad \times \left[\int d\mathbf{z} \mathbf{p}_{\text{out}}(\mathbf{y}|\mathbf{z}) \exp \left(-\mathbf{z}^{*\top} \tilde{\mathbf{m}} \mathbf{z} - \frac{1}{2} \mathbf{z}^\top (\tilde{\mathbf{Q}} - \tilde{\mathbf{q}}) \mathbf{z} + \mathbf{z}^\top (-\tilde{\mathbf{q}})^{1/2} \boldsymbol{\xi} \right) \right]^r,
\end{aligned} \tag{39}$$

with $\det(\mathbf{Q}^{(\text{rs})})$ given by (33).

3.2.3 Consistency conditions $r \rightarrow 0$: $\Theta(1)$ terms

It remains to take the limit $r \rightarrow 0^+$ of the expressions for $\Psi_{\text{w}}^{(r)}$ and $\Psi_{\text{out}}^{(r)}$ that are now analytical in r . First, our assumptions must be consistent and thus we need to check the consistency conditions in the limit $r \rightarrow 0$. Indeed, if $\Phi^{(r)}$ is finite we could obtain divergence taking the limit $\lim_{r \rightarrow 0} \frac{1}{r} \Phi^{(r)} = \infty$. Therefore to avoid such divergence, we must at least impose that $\lim_{r \rightarrow 0} \Phi^{(r)} = 0$:

$$\lim_{r \rightarrow 0} \Phi^{(r)}(\mathbf{Q}, \hat{\mathbf{Q}}) = -\text{Tr}(\mathbf{Q}^* \hat{\mathbf{Q}}^*) + \log \Psi_{\text{w}}^0(\hat{\mathbf{Q}}^*) + \alpha \log \Psi_{\text{out}}^0(\mathbf{Q}^*)$$

with

$$\begin{aligned}\Psi_{\mathbf{w}}^0(\hat{\mathbf{Q}}^*) &\equiv \mathbb{E}_{\mathbf{w}^*} \exp(\mathbf{w}^{*\top} \hat{\mathbf{Q}}^* \mathbf{w}^*), \\ \Psi_{\text{out}}^0(\mathbf{Q}^*) &\equiv \int_{\mathbb{R}} dy \int d\mathbf{z}^* p_{\text{out}^*}(y|\mathbf{z}^*) \mathcal{N}_{\mathbf{z}^*}(\mathbf{0}, \mathbf{Q}^*) = 1.\end{aligned}$$

Taking the saddle point equations over \mathbf{Q}^* and $\hat{\mathbf{Q}}^*$, imposing the consistency condition $\lim_{r \rightarrow 0} \Phi^{(r)}(\mathbf{Q}, \hat{\mathbf{Q}}) = 0$, we finally obtain

$$\hat{\mathbf{Q}}^* = \mathbf{0} \quad \text{and} \quad \mathbf{Q}^* = \mathbb{E}_{\mathbf{w}^*} [\mathbf{w}^{*\top} \mathbf{w}^*]. \quad (40)$$

3.2.4 Replica trick $r \rightarrow 0$ limit: $\Theta(r)$ terms

Imposing the conditions (40) avoids divergence in the replica trick, and we can therefore proceed with the $\Theta(r)$ terms.

Prior integral $\Psi_{\mathbf{w}}^{(r)}$ The limit $r \rightarrow 0$ and the derivative of the logarithm of the prior integral (22) can be trivially computed

$$\begin{aligned}& \lim_{r \rightarrow 0} \partial_r \log \Psi_{\mathbf{w}}^{(r)}(\hat{\mathbf{Q}}) \Big|_{\text{rs}} \\ &= \mathbb{E}_{\boldsymbol{\xi}, \mathbf{w}^*} \\ & \times \log \left[\mathbb{E}_{\mathbf{w}} \exp \left(\left[\mathbf{w}^{*\top} \hat{\mathbf{m}} \mathbf{w} - \frac{1}{2} \mathbf{w}^\top (\hat{\mathbf{Q}} + \hat{\mathbf{q}}) \mathbf{w} + \boldsymbol{\xi}^\top \hat{\mathbf{q}}^{1/2} \mathbf{w} \right] \right) \right],\end{aligned}$$

with $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ and $\mathbf{w}^* \sim P_{\mathbf{w}^*}$. To conclude, we can symmetrize and decouple the *teacher* and *student* expectations $\mathbb{E}_{\mathbf{w}^*}, \mathbb{E}_{\mathbf{w}}$. By performing the change of variable $\boldsymbol{\xi} \leftarrow \boldsymbol{\xi} + \hat{\mathbf{q}}^{-1/2} \hat{\mathbf{m}} \mathbf{w}^*$. We finally obtain

$$\begin{aligned}& \lim_{r \rightarrow 0} \partial_r \log \Psi_{\mathbf{w}}^{(r)}(\hat{\mathbf{Q}}) \Big|_{\text{rs}} \\ &= \mathbb{E}_{\boldsymbol{\xi}, \mathbf{w}^*} \exp \left(-\frac{1}{2} \mathbf{w}^{*\top} \hat{\mathbf{m}}^\top \hat{\mathbf{q}}^{-1} \hat{\mathbf{m}} \mathbf{w}^* + \boldsymbol{\xi}^\top \hat{\mathbf{q}}^{-1/2} \hat{\mathbf{m}} \mathbf{w}^* \right) \\ & \times \log \left[\mathbb{E}_{\mathbf{w}} \exp \left(\left[-\frac{1}{2} \mathbf{w}^\top (\hat{\mathbf{Q}} + \hat{\mathbf{q}}) \mathbf{w} + \boldsymbol{\xi}^\top \hat{\mathbf{q}}^{1/2} \mathbf{w} \right] \right) \right] \\ &\equiv \mathbb{E}_{\boldsymbol{\xi}, \mathbf{w}^*} \mathcal{Z}_{\mathbf{w}^*} \left(\hat{\mathbf{m}} \hat{\mathbf{q}}^{-1/2} \boldsymbol{\xi}, \hat{\mathbf{m}}^\top \hat{\mathbf{q}}^{-1} \hat{\mathbf{m}} \right) \log \mathcal{Z}_{\mathbf{w}} \left(\hat{\mathbf{q}}^{1/2} \boldsymbol{\xi}, \hat{\mathbf{Q}} + \hat{\mathbf{q}} \right),\end{aligned} \quad (41)$$

with the corresponding denoising distribution $Q_{\mathbf{w}}$ and functions $\mathcal{Z}_{\mathbf{w}^*}, \mathcal{Z}_{\mathbf{w}}$ defined by

$$\begin{aligned}Q_{\mathbf{w}}(\mathbf{w}; \boldsymbol{\gamma}, \boldsymbol{\Lambda}) &\equiv \frac{P_{\mathbf{w}}(\mathbf{w})}{\mathcal{Z}_{\mathbf{w}}(\boldsymbol{\gamma}, \boldsymbol{\Lambda})} e^{-\frac{1}{2} \mathbf{w}^\top \boldsymbol{\Lambda} \mathbf{w} + \boldsymbol{\gamma}^\top \mathbf{w}}, \\ \mathcal{Z}_{\mathbf{w}}(\boldsymbol{\gamma}, \boldsymbol{\Lambda}) &\equiv \mathbb{E}_{\mathbf{w} \sim P_{\mathbf{w}}} \left[e^{-\frac{1}{2} \mathbf{w}^\top \boldsymbol{\Lambda} \mathbf{w} + \boldsymbol{\gamma}^\top \mathbf{w}} \right] = \int_{\mathbb{R}^K} d\mathbf{w} p_{\mathbf{w}}(\mathbf{w}) e^{-\frac{1}{2} \mathbf{w}^\top \boldsymbol{\Lambda} \mathbf{w} + \boldsymbol{\gamma}^\top \mathbf{w}},\end{aligned} \quad (42)$$

respectively with distribution $P_{\mathbf{w}^*}$ and $P_{\mathbf{w}}$.

Prior integral $\Psi_{\text{out}}^{(r)}$ The limit $r \rightarrow 0$ and the derivative of the logarithm of the channel integral (39) is more tricky. First, the limit of the determinant (33) simplifies easily and yields

$$\det(\mathbf{Q}^{(\text{rs})}) \xrightarrow{r \rightarrow 0} \det(\mathbf{Q}^*)$$

and the matrix elements of $(\mathbf{Q}^{(\text{rs})})^{-1}$ in this limit become

$$\begin{aligned}\tilde{\mathbf{Q}}^* &\xrightarrow{r \rightarrow 0} (\mathbf{Q}^*)^{-1}, \\ \tilde{\mathbf{m}} &\xrightarrow{r \rightarrow 0} -(\mathbf{Q}^*)^{-1} \mathbf{m} (\mathbf{Q} - \mathbf{q})^{-1}, \\ \tilde{\mathbf{q}} &\xrightarrow{r \rightarrow 0} -(\mathbf{Q} - \mathbf{q})^{-1} (\mathbf{q} - \mathbf{m}^\top (\mathbf{Q}^*)^{-1} \mathbf{m}) (\mathbf{Q} - \mathbf{q})^{-1}, \\ \tilde{\mathbf{Q}} &\xrightarrow{r \rightarrow 0} \tilde{\mathbf{q}} + (\mathbf{Q} - \mathbf{q})^{-1}.\end{aligned}$$

To take properly the $r \rightarrow 0$ limit, let us first perform the change of variable $\xi \leftarrow \xi - (-\tilde{\mathbf{q}})^{1/2} \tilde{\mathbf{m}}^\top \mathbf{z}^*$ in (39) so that

$$\begin{aligned}\Psi_{\text{out}}^{(r)}(\mathbf{Q}) \Big|_{\text{rs}} &= \int_{\mathbb{R}} dy \mathbb{E}_{\xi} J_0(y, \xi, r) \times J_1(y, \xi, r)^r \\ J_0(y, \xi, r) &= e^{-\frac{1}{2} \log(\det(2\pi \mathbf{Q}^{(\text{rs})}))} \\ &\quad \times \int_{\mathbb{R}^K} d\mathbf{z}^* p_{\text{out}^*}(y|\mathbf{z}^*) e^{-\frac{1}{2} \mathbf{z}^{*\top} (\tilde{\mathbf{Q}}^* - \tilde{\mathbf{m}}^\top \tilde{\mathbf{q}}^{-1} \tilde{\mathbf{m}}) \mathbf{z}^* + \xi^\top (-\tilde{\mathbf{q}})^{1/2} \tilde{\mathbf{m}}^\top \mathbf{z}^*} \\ J_1(y, \xi, r) &= \int_{\mathbb{R}^K} d\mathbf{z} p_{\text{out}}(y|\mathbf{z}) e^{-\frac{1}{2} \mathbf{z}^\top (\tilde{\mathbf{Q}} - \tilde{\mathbf{q}}) \mathbf{z} + \mathbf{z}^\top (-\tilde{\mathbf{q}})^{1/2} \xi}\end{aligned} \quad (43)$$

so that the integral may be written as follows:

$$\begin{aligned}\log(\Psi_{\text{out}}^{(r)}) &= \log \int dy \mathbb{E}_{\xi} J_0(y, \xi, r) J_1(y, \xi, r)^r \\ &= \log \int_{\mathbb{R}} dy \mathbb{E}_{\xi} [J_0(y, \xi, 0) + r \partial_r J_0(y, \xi, 0) + r J_0(y, \xi, 0) \log J_1(y, \xi, 0) + \Theta(r^2)] \\ &= \log \left(1 + r \int_{\mathbb{R}} dy \times \mathbb{E}_{\xi} [\partial_r J_0(y, \xi, 0) + J_0(y, \xi, 0) \log J_1(y, \xi, 0) + \Theta(r^2)] \right) \\ &= r \int_{\mathbb{R}} dy \times \mathbb{E}_{\xi} [\partial_r J_0(y, \xi, 0) + J_0(y, \xi, 0) \log J_1(y, \xi, 0)] + \Theta(r^2)\end{aligned}$$

where we used the consistency condition $\int_{\mathbb{R}} dy \times \mathbb{E}_{\xi} J_0(y, \xi, 0) = 1$. The first term can be evaluated similarly

$$\begin{aligned}\int_{\mathbb{R}} dy \times \mathbb{E}_{\xi} \partial_r J_0(y, \xi, 0) &= \partial_r \left[\int_{\mathbb{R}} dy \times \mathbb{E}_{\xi} J_0(y, \xi, r) \right]_{r=0} \\ &= \partial_r \left[e^{-\frac{1}{2} \log(\det(2\pi \mathbf{Q}^{(\text{rs})}))} \int_{\mathbb{R}^K} d\mathbf{z}^* e^{-\frac{1}{2} \mathbf{z}^{*\top} (\tilde{\mathbf{Q}}^*) \mathbf{z}^*} \right]_{r=0} \\ &= \partial_r \left[e^{-\frac{1}{2} \log(\det(\tilde{\mathbf{Q}}^*) \det(\mathbf{Q}^{(\text{rs})}))} \right]_{r=0} = \dots = 0\end{aligned}$$

Finally we obtain:

$$\lim_{r \rightarrow 0} \partial_r \log \left(\Psi_{\text{out}}^{(r)} \right) = \int_{\mathbb{R}} dy \mathbb{E}_{\xi} J_0(y, \xi) \log J_1(y, \xi) = \lim_{r \rightarrow 0} \int_{\mathbb{R}} dy \mathbb{E}_{\xi} J_0(y, \xi) J_1(y, \xi)^r$$

with

$$\begin{aligned} J_0(y, \xi) &\equiv e^{-\frac{1}{2} \log(\det(2\pi \mathbf{Q}^*))} \int_{\mathbb{R}^K} d\mathbf{z}^* p_{\text{out}^*}(y|\mathbf{z}^*) e^{-\frac{1}{2} \mathbf{z}^{*\top} \mathbf{N} \mathbf{z}^* + \xi^\top \mathbf{M} \mathbf{z}^*} \\ J_1(y, \xi) &\equiv \int_{\mathbb{R}^K} d\mathbf{z} p_{\text{out}}(y|\mathbf{z}) e^{-\frac{1}{2} \mathbf{z}^\top (\mathbf{Q} - \mathbf{q})^{-1} \mathbf{z} + \mathbf{z}^\top \mathbf{P} \xi} \\ &= \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} p_{\text{out}} \left(y | (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{z} \right) e^{\xi^\top \mathbf{P} (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{z}} \end{aligned}$$

where we changed the variable $\mathbf{z} \leftarrow (\mathbf{Q} - \mathbf{q})^{-1/2} \mathbf{z}$ and used the fact that with the second formulation the determinant term at the power r goes away in the limit $r \rightarrow 0$ and with

$$\begin{aligned} \mathbf{P} &\equiv (-\tilde{\mathbf{q}})^{1/2} = (\mathbf{Q} - \mathbf{q})^{-1/2} (\mathbf{q} - \mathbf{m}^\top (\mathbf{Q}^*)^{-1} \mathbf{m})^{1/2} (\mathbf{Q} - \mathbf{q})^{-1/2} \\ \mathbf{M} &\equiv (-\tilde{\mathbf{q}})^{-1/2} \tilde{\mathbf{m}}^\top = -(-\tilde{\mathbf{q}})^{-1/2} (\mathbf{Q} - \mathbf{q})^{-1} \mathbf{m}^\top (\mathbf{Q}^*)^{-1} \\ &= -\mathbf{P}^{-1} (\mathbf{Q} - \mathbf{q})^{-1} \mathbf{m}^\top (\mathbf{Q}^*)^{-1} \\ \mathbf{N} &\equiv \tilde{\mathbf{Q}}^* - \tilde{\mathbf{m}}^\top \tilde{\mathbf{q}}^{-1} \tilde{\mathbf{m}} \\ &= (\mathbf{Q}^*)^{-1} + (\mathbf{Q}^*)^{-1} \mathbf{m} (\mathbf{Q} - \mathbf{q})^{-1} (\mathbf{Q} - \mathbf{q}) (\mathbf{q} - \mathbf{m}^\top (\mathbf{Q}^*)^{-1} \mathbf{m})^{-1} (\mathbf{Q} - \mathbf{q}) (\mathbf{Q} - \mathbf{q})^{-1} \mathbf{m}^\top (\mathbf{Q}^*)^{-1} \\ &= (\mathbf{Q}^*)^{-1} \left[\mathbf{I}_K + \mathbf{m} (\mathbf{q} - \mathbf{m}^\top (\mathbf{Q}^*)^{-1} \mathbf{m})^{-1} \mathbf{m}^\top (\mathbf{Q}^*)^{-1} \right] \\ &= (\mathbf{Q}^*)^{-1} \left[\mathbf{I}_K + [(\mathbf{Q}^*)^{-1} \mathbf{m}^{-\top} \mathbf{q} \mathbf{m}^{-1} - \mathbf{I}_K]^{-1} \right] \\ &= [\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top]^{-1} \end{aligned}$$

Then we rescale $\mathbf{z}^* \leftarrow \mathbf{N}^{1/2} \mathbf{z}^*$ such that

$$\begin{aligned} J_0(y, \xi) &= e^{-\frac{1}{2} \log(\det(2\pi \mathbf{Q}^*))} \int_{\mathbb{R}^K} d\mathbf{z}^* p_{\text{out}^*}(y|\mathbf{z}^*) e^{-\frac{1}{2} \mathbf{z}^{*\top} \mathbf{N} \mathbf{z}^* + \xi^\top \mathbf{M} \mathbf{z}^*} \\ &= e^{-\frac{1}{2} \log(\det(2\pi \mathbf{Q}^*))} e^{\frac{1}{2} \log(\det(\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top))} \\ &\quad \times \mathbb{E}_{\mathbf{z}^* \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} \left[p_{\text{out}^*} \left(y | (\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top)^{1/2} \mathbf{z}^* \right) e^{\xi^\top \mathbf{M} \mathbf{N}^{-1/2} \mathbf{z}^*} \right] \end{aligned}$$

Moreover notice that:

$$\begin{aligned} \mathbf{M} \mathbf{N}^{-1/2} &= \mathbf{P}^{-1} (\mathbf{Q} - \mathbf{q})^{-1} \mathbf{m}^\top (\mathbf{Q}^*)^{-1} [\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top]^{1/2} \\ &= \mathbf{P} \mathbf{P}^{-2} (\mathbf{Q} - \mathbf{q})^{-1} \mathbf{m}^\top (\mathbf{Q}^*)^{-1} [\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top]^{1/2} \\ &= \mathbf{P} (\mathbf{Q} - \mathbf{q}) (\mathbf{q} - \mathbf{m}^\top (\mathbf{Q}^*)^{-1} \mathbf{m})^{-1} (\mathbf{Q} - \mathbf{q}) (\mathbf{Q} - \mathbf{q})^{-1} \mathbf{m}^\top (\mathbf{Q}^*)^{-1} [\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top]^{1/2} \\ &= \mathbf{P} (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{R} \end{aligned}$$

with

$$\mathbf{R} = (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{q}^{-1} \mathbf{m}^\top [\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top]^{-1/2}$$

and notice that $\mathbf{R} = \mathbf{I}_K$ in the Bayes-optimal setting. Therefore:

$$\begin{aligned} J_0(y, \boldsymbol{\xi}) &= e^{-\frac{1}{2} \log(\det(2\pi \mathbf{Q}^*))} e^{\frac{1}{2} \log(\det(\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top))} \\ &\quad \times \mathbb{E}_{\mathbf{z}^* \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} \left[p_{\text{out}^*} \left(y \mid (\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top)^{1/2} \mathbf{z}^* \right) e^{\boldsymbol{\xi}^\top \mathbf{P}(\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{R} \mathbf{z}^*} \right] \end{aligned}$$

As a summary, we have:

$$\lim_{r \rightarrow 0} \partial_r \log \left(\Psi_{\text{out}}^{(r)} \right) = \lim_{r \rightarrow 0} \int_{\mathbb{R}} dy \mathbb{E}_{\boldsymbol{\xi}} J_0(y, \boldsymbol{\xi}) J_1(y, \boldsymbol{\xi})^r$$

with

$$\begin{aligned} J_0(y, \boldsymbol{\xi}) &= e^{-\frac{1}{2} \log(\det(2\pi \mathbf{Q}^*))} e^{\frac{1}{2} \log(\det(\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top))} \\ &\quad \times \mathbb{E}_{\mathbf{z}^* \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} \left[p_{\text{out}^*} \left(y \mid (\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top)^{1/2} \mathbf{z}^* \right) e^{\boldsymbol{\xi}^\top \mathbf{P}(\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{R} \mathbf{z}^*} \right] \\ J_1(y, \boldsymbol{\xi}) &= \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} p_{\text{out}} \left(y \mid (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{z} \right) e^{\boldsymbol{\xi}^\top \mathbf{P}(\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{z}} \end{aligned}$$

We can shift $\mathbf{z} \leftarrow \mathbf{z} - (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{P} \boldsymbol{\xi}$ to obtain

$$\begin{aligned} J_1(y, \boldsymbol{\xi}) &= e^{\frac{1}{2} \boldsymbol{\xi}^\top \mathbf{P}(\mathbf{Q} - \mathbf{q})^{1/2} (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{P} \boldsymbol{\xi}} \\ &\quad \times \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} p_{\text{out}} \left(y \mid (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{z} + (\mathbf{Q} - \mathbf{q}) \mathbf{P} \boldsymbol{\xi} \right) \end{aligned}$$

and $\mathbf{z}^* \leftarrow \mathbf{z}^* - (\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top)^{-1/2} \mathbf{m} \mathbf{q}^{-1} (\mathbf{Q} - \mathbf{q}) \mathbf{P} \boldsymbol{\xi} = \mathbf{z}^* - \mathbf{R}^\top (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{P} \boldsymbol{\xi}$:

$$\begin{aligned} J_0(y, \boldsymbol{\xi}) &= e^{-\frac{1}{2} \log(\det(2\pi \mathbf{Q}^*))} e^{\frac{1}{2} \log(\det(\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top))} \\ &\quad \times \mathbb{E}_{\mathbf{z}^* \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} \left[p_{\text{out}^*} \left(y \mid (\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top)^{1/2} \mathbf{z}^* \right) e^{\boldsymbol{\xi}^\top \mathbf{P}(\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{R} \mathbf{z}^*} \right] \\ &= e^{-\frac{1}{2} \log(\det(2\pi \mathbf{Q}^*))} e^{\frac{1}{2} \log(\det(\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top))} e^{\frac{1}{2} \boldsymbol{\xi}^\top \mathbf{P}(\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{R} \mathbf{R}^\top (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{P} \boldsymbol{\xi}} \\ &\quad \times \mathbb{E}_{\mathbf{z}^* \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} \left[p_{\text{out}^*} \left(y \mid (\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top)^{1/2} \mathbf{z}^* + \mathbf{m} \mathbf{q}^{-1} (\mathbf{Q} - \mathbf{q}) \mathbf{P} \boldsymbol{\xi} \right) \right] \end{aligned}$$

We obtain (taking the limit $r \rightarrow 0$ to the term in front of J_1):

$$\begin{aligned} &\lim_{r \rightarrow 0} \partial_r \log \left(\Psi_{\text{out}}^{(r)} \right) \\ &= \int_{\mathbb{R}} dy \int_{\mathbb{R}^K} d\boldsymbol{\xi} e^{-\frac{1}{2} \log(\det(\mathbf{Q}^*))} e^{\frac{1}{2} \log(\det(\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top))} e^{-\frac{1}{2} \boldsymbol{\xi}^\top [\mathbf{I}_K - \mathbf{P}(\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{R} \mathbf{R}^\top (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{P}] \boldsymbol{\xi}} \\ &\quad \times \mathbb{E}_{\mathbf{z}^* \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} \left[p_{\text{out}^*} \left(y \mid (\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top)^{1/2} \mathbf{z}^* + \mathbf{m} \mathbf{q}^{-1} (\mathbf{Q} - \mathbf{q}) \mathbf{P} \boldsymbol{\xi} \right) \right] \\ &\quad \times \log \left[\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} p_{\text{out}} \left(y \mid (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{z} + (\mathbf{Q} - \mathbf{q}) \mathbf{P} \boldsymbol{\xi} \right) \right] \end{aligned}$$

Computing

$$\begin{aligned}
& \left[\mathbf{I}_K - \mathbf{P}(\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{R} \mathbf{R}^\top (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{P} \right] \\
&= \mathbf{P} \left[\mathbf{P}^{-2} - (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{R} \mathbf{R}^\top (\mathbf{Q} - \mathbf{q})^{1/2} \right] \mathbf{P} \\
&= \mathbf{P}(\mathbf{Q} - \mathbf{q}) \left[(\mathbf{q} - \mathbf{m}^\top (\mathbf{Q}^*)^{-1} \mathbf{m})^{-1} - \mathbf{q}^{-1} \mathbf{m}^\top [\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top]^{-1} \mathbf{m} \mathbf{q}^{-1} \right] (\mathbf{Q} - \mathbf{q}) \mathbf{P} \\
&= \mathbf{P}(\mathbf{Q} - \mathbf{q}) \mathbf{q}^{-1} (\mathbf{Q} - \mathbf{q}) \mathbf{P}
\end{aligned}$$

By finally changing the Gaussian measure $\boldsymbol{\xi} \leftarrow \mathbf{q}^{-1/2} (\mathbf{Q} - \mathbf{q}) \mathbf{P} \boldsymbol{\xi}$ we obtain

$$\begin{aligned}
& \lim_{r \rightarrow 0} \partial_r \log \left(\Psi_{\text{out}}^{(r)} \right) \\
&= \exp(C) \times \int_{\mathbb{R}} dy \mathbb{E}_{\boldsymbol{\xi}} \\
&\times \mathbb{E}_{\mathbf{z}^* \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} \left[p_{\text{out}^*} \left(y \mid (\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top)^{1/2} \mathbf{z}^* + \mathbf{m} \mathbf{q}^{-1/2} \boldsymbol{\xi} \right) \right] \\
&\times \log \left[\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} p_{\text{out}} \left(y \mid (\mathbf{Q} - \mathbf{q})^{1/2} \mathbf{z} + \mathbf{q}^{1/2} \boldsymbol{\xi} \right) \right]
\end{aligned}$$

with

$$\begin{aligned}
C &= -\frac{1}{2} \log(\det(\mathbf{Q}^*)) + \frac{1}{2} \log(\det(\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top)) \\
&\quad - \frac{1}{2} \log \det(\mathbf{P}^2) - \log \det(\mathbf{Q} - \mathbf{q}) + \frac{1}{2} \log \det(\mathbf{q}) \\
&= -\frac{1}{2} \log(\det(\mathbf{Q}^*)) + \frac{1}{2} \log(\det(\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top)) \\
&\quad - \frac{1}{2} \log \det((\mathbf{Q} - \mathbf{q}) (\mathbf{q} - \mathbf{m}^\top (\mathbf{Q}^*)^{-1} \mathbf{m}) (\mathbf{Q} - \mathbf{q})) - \log \det(\mathbf{Q} - \mathbf{q}) + \frac{1}{2} \log \det(\mathbf{q}) \\
&= -\frac{1}{2} \log(\det(\mathbf{Q}^*)) + \frac{1}{2} \log(\det(\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top)) \\
&\quad - \frac{1}{2} \log \det((\mathbf{q} - \mathbf{m}^\top (\mathbf{Q}^*)^{-1} \mathbf{m})) + \frac{1}{2} \log \det(\mathbf{q}) \\
&= -\frac{1}{2} \log(\det(\mathbf{Q}^*)) - \frac{1}{2} \log \det(\mathbf{m} \mathbf{m}^{-1}) - \frac{1}{2} \log \det(\mathbf{q}) + \frac{1}{2} \log(\det(\mathbf{Q}^*)) + \frac{1}{2} \log(\det(\mathbf{q})) \\
&= 0
\end{aligned}$$

We obtain:

$$\begin{aligned}
& \lim_{r \rightarrow 0} \partial_r \log \Psi_{\text{out}}^{(r)}(\mathbf{Q}) \Big|_{\text{rs}} \\
&= \int_{\mathbb{R}} d\mathbf{y} \mathbb{E}_{\boldsymbol{\xi}} \\
&\times \mathbb{E}_{\mathbf{z}^* \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} \left[\text{p}_{\text{out}^*} \left(y \left(\mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top \right)^{1/2} \mathbf{z}^* + \mathbf{m} \mathbf{q}^{-1/2} \boldsymbol{\xi} \right) \right] \\
&\times \log \left[\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} \text{p}_{\text{out}} \left(y \left(\mathbf{Q} - \mathbf{q} \right)^{1/2} \mathbf{z} + \mathbf{q}^{1/2} \boldsymbol{\xi} \right) \right] \\
&= \int_{\mathbb{R}} d\mathbf{y} \mathbb{E}_{\boldsymbol{\xi}} \mathcal{Z}_{\text{out}^*} \left(\mathbf{m} \mathbf{q}^{-1/2} \boldsymbol{\xi}, \mathbf{Q}^* - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top \right) \log \mathcal{Z}_{\text{out}} \left(\mathbf{q}^{1/2} \boldsymbol{\xi}, \mathbf{Q} - \mathbf{q} \right), \tag{44}
\end{aligned}$$

where the denoising distribution \mathbf{Q}_{out} and functions $\mathcal{Z}_{\text{out}^*}$, \mathcal{Z}_{out} are defined by

$$\begin{aligned}
\mathbf{Q}_{\text{out}}(\mathbf{z}; y, \boldsymbol{\omega}, \mathbf{V}) &\equiv \frac{\text{P}_{\text{out}}(y|\mathbf{z})}{\mathcal{Z}_{\text{out}}(y, \boldsymbol{\omega}, \mathbf{V})} \frac{e^{-\frac{1}{2}(\mathbf{z} - \boldsymbol{\omega})^\top \mathbf{V}^{-1}(\mathbf{z} - \boldsymbol{\omega})}}{\sqrt{\det(2\pi\mathbf{V})}}, \\
\mathcal{Z}_{\text{out}}(y, \boldsymbol{\omega}, \mathbf{V}) &\equiv \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)} \left[\text{P}_{\text{out}} \left(y |\mathbf{V}^{1/2} \mathbf{z} + \boldsymbol{\omega}| \right) \right] \\
&= \int_{\mathbb{R}^K} d\mathbf{z} \text{p}_{\text{out}}(y|\mathbf{z}) \frac{e^{-\frac{1}{2}(\mathbf{z} - \boldsymbol{\omega})^\top \mathbf{V}^{-1}(\mathbf{z} - \boldsymbol{\omega})}}{\sqrt{\det(2\pi\mathbf{V})}}. \tag{45}
\end{aligned}$$

for the distributions P_{out^*} , P_{out} .

3.3 Summary

3.3.1 Summary - Mismatched case

In the mismatched case, where the teacher and the student have not the same prior distributions, we finally obtain the replica symmetric free entropy Φ_{rs} for the committee machine hypothesis class:

$$\begin{aligned}
\Phi_{\text{rs}}(\alpha) &\equiv \mathbb{E}_{\mathbf{y}, \mathbf{X}} \left[\lim_{d \rightarrow \infty} \frac{1}{d} \log(\mathcal{Z}_d(\mathbf{y}, \mathbf{X})) \right] \\
&= \mathbf{extr}_{\mathbf{Q}, \hat{\mathbf{Q}}, \mathbf{q}, \hat{\mathbf{q}}, \mathbf{m}, \hat{\mathbf{m}}} \left\{ -\text{Tr}(\mathbf{m} \hat{\mathbf{m}}) + \frac{1}{2} \text{Tr}(\mathbf{Q} \hat{\mathbf{Q}}) + \frac{1}{2} \text{Tr}(\mathbf{q} \hat{\mathbf{q}}) \right. \\
&\quad \left. + \Psi_{\text{w}}(\hat{\mathbf{Q}}, \hat{\mathbf{m}}, \hat{\mathbf{q}}) + \alpha \Psi_{\text{out}}(\mathbf{Q}, \mathbf{m}, \mathbf{q}; \boldsymbol{\rho}_{\text{w}^*}) \right\}, \tag{46}
\end{aligned}$$

where $\boldsymbol{\rho}_{\text{w}^*} \equiv \lim_{d \rightarrow \infty} \mathbf{Q}^* = \lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}^*} \frac{1}{d} \mathbf{W}^{*\top} \mathbf{W}^*$ and the channel and prior integrals are defined by

$$\begin{aligned}
\Psi_{\text{w}}(\hat{\mathbf{Q}}, \hat{\mathbf{m}}, \hat{\mathbf{q}}) &\equiv \mathbb{E}_{\boldsymbol{\xi}} \left[\mathcal{Z}_{\text{w}^*} \left(\hat{\mathbf{m}} \hat{\mathbf{q}}^{-1/2} \boldsymbol{\xi}, \hat{\mathbf{m}} \hat{\mathbf{q}}^{-1} \hat{\mathbf{m}} \right) \right. \\
&\quad \left. \times \log \mathcal{Z}_{\text{w}} \left(\hat{\mathbf{q}}^{1/2} \boldsymbol{\xi}, \hat{\mathbf{Q}} + \hat{\mathbf{q}} \right) \right], \tag{47}
\end{aligned}$$

$$\begin{aligned}
\Psi_{\text{out}}(\mathbf{Q}, \mathbf{m}, \mathbf{q}; \boldsymbol{\rho}_{\text{w}^*}) &\equiv \mathbb{E}_{y, \boldsymbol{\xi}} \left[\mathcal{Z}_{\text{out}^*} \left(y, \mathbf{m} \mathbf{q}^{-1/2} \boldsymbol{\xi}, \boldsymbol{\rho}_{\text{w}^*} - \mathbf{m} \mathbf{q}^{-1} \mathbf{m}^\top \right) \right. \\
&\quad \left. \times \log \mathcal{Z}_{\text{out}} \left(y, \mathbf{q}^{1/2} \boldsymbol{\xi}, \mathbf{Q} - \mathbf{q} \right) \right], \tag{48}
\end{aligned}$$

where again $\mathcal{Z}_{\text{out}^*}$, \mathcal{Z}_{w^*} and \mathcal{Z}_{out} , \mathcal{Z}_{w} are defined in (42)- (45) and depend respectively on channel and prior distributions of the *teacher* and *student*.

3.3.2 Bayes optimal MMSE estimation

For MMSE estimation in the Bayes-optimal setting, the student has access to the ground truth distributions of the teacher $P_{\text{out}}(\mathbf{y}|\mathbf{Z}) = P_{\text{out}^*}(\mathbf{y}|\mathbf{Z})$ and $P_{\text{w}}(\mathbf{W}) = P_{\text{w}^*}(\mathbf{W})$, and therefore $\mathcal{Z}_{\text{out}} = \mathcal{Z}_{\text{out}^*}$, $\mathcal{Z}_{\text{w}} = \mathcal{Z}_{\text{w}^*}$. In this idealistic setting, the Nishimori conditions imply that

$$\mathbf{Q} = \mathbf{Q}_{\text{w}^*}, \quad \hat{\mathbf{Q}} = \mathbf{0}, \quad \mathbf{m} = \mathbf{q} \equiv \mathbf{q}_{\text{b}}, \quad \hat{\mathbf{m}} = \hat{\mathbf{q}} \equiv \hat{\mathbf{q}}_{\text{b}}. \quad (49)$$

Therefore the free entropy in eq. (47) simplifies as an optimization problem over the overlaps $\mathbf{q}_{\text{b}}, \hat{\mathbf{q}}_{\text{b}} \in \mathbb{R}^{K \times K}$

$$\Phi_{\text{rs}}^{\text{b}}(\alpha) = \text{extr}_{\mathbf{q}_{\text{b}}, \hat{\mathbf{q}}_{\text{b}}} \left\{ -\frac{1}{2} \text{Tr}(\mathbf{q}_{\text{b}} \hat{\mathbf{q}}_{\text{b}}) + \Psi_{\text{w}}^{\text{b}}(\hat{\mathbf{q}}_{\text{b}}) + \alpha \Psi_{\text{out}}^{\text{b}}(\mathbf{q}_{\text{b}}; \boldsymbol{\rho}_{\text{w}^*}) \right\}, \quad (50)$$

with free entropy terms $\Psi_{\text{w}}^{\text{b}}$ and $\Psi_{\text{out}}^{\text{b}}$ given by

$$\begin{aligned} \Psi_{\text{w}}^{\text{b}}(\hat{\mathbf{q}}) &= \mathbb{E}_{\boldsymbol{\xi}} \left[\mathcal{Z}_{\text{w}^*} \left(\hat{\mathbf{q}}^{1/2} \boldsymbol{\xi}, \hat{\mathbf{q}} \right) \log \mathcal{Z}_{\text{w}^*} \left(\hat{\mathbf{q}}^{1/2} \boldsymbol{\xi}, \hat{\mathbf{q}} \right) \right], \\ \Psi_{\text{out}}^{\text{b}}(\mathbf{q}; \boldsymbol{\rho}_{\text{w}^*}) &= \mathbb{E}_{\mathbf{y}, \boldsymbol{\xi}} \left[\mathcal{Z}_{\text{out}} \left(\mathbf{y}, \mathbf{q}^{1/2} \boldsymbol{\xi}, \boldsymbol{\rho}_{\text{w}^*} - \mathbf{q} \right) \right. \\ &\quad \left. \times \log \mathcal{Z}_{\text{out}} \left(\mathbf{y}, \mathbf{q}^{1/2} \boldsymbol{\xi}, \boldsymbol{\rho}_{\text{w}^*} - \mathbf{q} \right) \right]. \end{aligned} \quad (51)$$

Application to the GLM For the GLM hypothesis class, the same equations are valid if we take $K = 1$ for both the teacher and the student. As a result, we recover the replica symmetric free entropy in the Bayes-optimal setting rigorously proven in [23].

3.4 Fixed point equations

The overlaps parameters, such as \mathbf{m}, \mathbf{q} , play a crucial role since they measure the performances of the statistical estimation. Their behaviours are respectively characterized by the extremization of the free entropy (46) in the mismatched setting and (50) in the Bayes-optimal case. In this section, we give the expressions of the corresponding fixed point equations, whose derivations can be found in [25] Appendix. IV.4-5 for $K = 1$ which can be extended to $K \geq 1$.

3.4.1 Mismatched setting

Extremizing the free entropy eq. (46), we easily obtain the set of six fixed point equations

$$\begin{aligned} \hat{\mathbf{Q}} &= -2\alpha \partial_{\mathbf{Q}} \Psi_{\text{out}}(\mathbf{Q}, \mathbf{m}, \mathbf{q}; \boldsymbol{\rho}_{\text{w}^*}), & \mathbf{Q} &= -2\partial_{\hat{\mathbf{Q}}} \Psi_{\text{w}}(\hat{\mathbf{Q}}, \hat{\mathbf{m}}, \hat{\mathbf{q}}) \\ \hat{\mathbf{q}} &= -2\alpha \partial_{\mathbf{q}} \Psi_{\text{out}}(\mathbf{Q}, \mathbf{m}, \mathbf{q}; \boldsymbol{\rho}_{\text{w}^*}), & \mathbf{q} &= -2\partial_{\hat{\mathbf{q}}} \Psi_{\text{w}}(\hat{\mathbf{Q}}, \hat{\mathbf{m}}, \hat{\mathbf{q}}), \\ \hat{\mathbf{m}} &= \alpha \partial_{\mathbf{m}} \Psi_{\text{out}}(\mathbf{Q}, \mathbf{m}, \mathbf{q}; \boldsymbol{\rho}_{\text{w}^*}), & \mathbf{m} &= \partial_{\hat{\mathbf{m}}} \Psi_{\text{w}}(\hat{\mathbf{Q}}, \hat{\mathbf{m}}, \hat{\mathbf{q}}). \end{aligned} \quad (52)$$

Interestingly, these equations can be reformulated as functions of $\mathcal{Z}_{\text{out}^*}$, \mathcal{Z}_{w^*} and the denoising functions f_{out^*} , f_{w^*} , f_{out} , f_{w} defined by

$$\mathbf{f}_{\text{w}}(\boldsymbol{\gamma}, \boldsymbol{\Lambda}) \equiv \partial_{\boldsymbol{\gamma}} \log(\mathcal{Z}_{\text{w}}(\boldsymbol{\gamma}, \boldsymbol{\Lambda})) = \mathbb{E}_{\mathbf{Q}_{\text{w}}}[\mathbf{w}], \quad (53)$$

$$\partial_{\boldsymbol{\gamma}} \mathbf{f}_{\text{w}}(\boldsymbol{\gamma}, \boldsymbol{\Lambda}) \equiv \mathbb{E}_{\mathbf{Q}_{\text{w}}}[\mathbf{w}\mathbf{w}^{\top}] - \mathbf{f}_{\text{w}}(\boldsymbol{\gamma}, \boldsymbol{\Lambda})^{\otimes 2}, \quad (54)$$

$$\mathbf{f}_{\text{out}}(y, \boldsymbol{\omega}, \mathbf{V}) \equiv \partial_{\boldsymbol{\omega}} \log(\mathcal{Z}_{\text{out}}(y, \boldsymbol{\omega}, \mathbf{V})) = \mathbf{V}^{-1} \mathbb{E}_{\mathbf{Q}_{\text{out}}}[\mathbf{z} - \boldsymbol{\omega}], \quad (55)$$

$$\begin{aligned} \partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out}}(y, \boldsymbol{\omega}, \mathbf{V}) &\equiv \frac{\partial \mathbf{f}_{\text{out}}(y, \boldsymbol{\omega}, \mathbf{V})}{\partial \boldsymbol{\omega}} \\ &= \mathbf{V}^{-1} \mathbb{E}_{\mathbf{Q}_{\text{out}}}[(\mathbf{z} - \boldsymbol{\omega})^{\otimes 2}] \mathbf{V}^{-1} - \mathbf{V}^{-1} - \mathbf{f}_{\text{out}}(y, \boldsymbol{\omega}, \mathbf{V})^{\otimes 2}. \end{aligned} \quad (56)$$

Defining the natural variables $\boldsymbol{\Sigma} = \mathbf{Q} - \mathbf{q}$ and $\hat{\boldsymbol{\Sigma}} = \hat{\mathbf{Q}} + \hat{\mathbf{q}}$ they can reformulated as

$$\begin{aligned} \hat{\mathbf{m}} &= \alpha \mathbb{E}_{y, \boldsymbol{\xi}} \left[\mathcal{Z}_{\text{out}^*} \times \mathbf{f}_{\text{out}^*} \left(y, \mathbf{m} \mathbf{q}^{-1/2} \boldsymbol{\xi}, \boldsymbol{\rho}_{\text{w}^*} - \mathbf{m}^{\top} \mathbf{q}^{-1} \mathbf{m} \right) \right. \\ &\quad \left. \times \mathbf{f}_{\text{out}} \left(y, \mathbf{q}^{1/2} \boldsymbol{\xi}, \boldsymbol{\Sigma} \right)^{\top} \right], \\ \hat{\mathbf{q}} &= \alpha \mathbb{E}_{y, \boldsymbol{\xi}} \left[\mathcal{Z}_{\text{out}^*} \left(y, \mathbf{m} \mathbf{q}^{-1/2} \boldsymbol{\xi}, \boldsymbol{\rho}_{\text{w}^*} - \mathbf{m}^{\top} \mathbf{q}^{-1} \mathbf{m} \right) \mathbf{f}_{\text{out}} \left(y, \mathbf{q}^{1/2} \boldsymbol{\xi}, \boldsymbol{\Sigma} \right)^{\otimes 2} \right], \\ \hat{\boldsymbol{\Sigma}} &= -\alpha \mathbb{E}_{y, \boldsymbol{\xi}} \left[\mathcal{Z}_{\text{out}^*} \left(y, \mathbf{m} \mathbf{q}^{-1/2} \boldsymbol{\xi}, \boldsymbol{\rho}_{\text{w}^*} - \mathbf{m}^{\top} \mathbf{q}^{-1} \mathbf{m} \right) \right. \\ &\quad \left. \times \partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out}} \left(y, \mathbf{q}^{1/2} \boldsymbol{\xi}, \boldsymbol{\Sigma} \right) \right], \\ \mathbf{m} &= \mathbb{E}_{\boldsymbol{\xi}} \left[\mathcal{Z}_{\text{w}^*} \times f_{\text{w}^*} \left(\hat{\mathbf{m}} \hat{\mathbf{q}}^{-1/2} \boldsymbol{\xi}, \hat{\mathbf{m}}^{\top} \hat{\mathbf{q}}^{-1} \hat{\mathbf{m}} \right) \mathbf{f}_{\text{w}} \left(\hat{\mathbf{q}}^{1/2} \boldsymbol{\xi}, \hat{\boldsymbol{\Sigma}} \right) \right], \\ \mathbf{q} &= \mathbb{E}_{\boldsymbol{\xi}} \left[\mathcal{Z}_{\text{w}^*} \left(\hat{\mathbf{m}} \hat{\mathbf{q}}^{-1/2} \boldsymbol{\xi}, \hat{\mathbf{m}}^{\top} \hat{\mathbf{q}}^{-1} \hat{\mathbf{m}} \right) \mathbf{f}_{\text{w}} \left(\hat{\mathbf{q}}^{1/2} \boldsymbol{\xi}, \hat{\boldsymbol{\Sigma}} \right)^2 \right], \\ \boldsymbol{\Sigma} &= \mathbb{E}_{\boldsymbol{\xi}} \left[\mathcal{Z}_{\text{w}^*} \left(\hat{\mathbf{m}} \hat{\mathbf{q}}^{-1/2} \boldsymbol{\xi}, \hat{\mathbf{m}}^{\top} \hat{\mathbf{q}}^{-1} \hat{\mathbf{m}} \right) \partial_{\boldsymbol{\gamma}} \mathbf{f}_{\text{w}} \left(\hat{\mathbf{q}}^{1/2} \boldsymbol{\xi}, \hat{\boldsymbol{\Sigma}} \right) \right], \end{aligned} \quad (57)$$

where we use the abusive notation $\mathbb{E}_y = \int_{\mathbb{R}} dy$.

3.4.2 Bayes-optimal estimation

Extremizing the Bayes-optimal free entropy eq. (50), we easily obtain the set of fixed point equations over the scalar parameters \mathbf{q}_{b} , $\hat{\mathbf{q}}_{\text{b}}$. It can be deduced from eq. (57) using the Nishimori conditions $\mathbf{f}_{\text{w}} = \mathbf{f}_{\text{w}^*}$, $\mathbf{f}_{\text{out}} = \mathbf{f}_{\text{out}^*}$, $\mathbf{m} = \mathbf{q}$, $\boldsymbol{\Sigma} = \boldsymbol{\rho}_{\text{w}^*} - \mathbf{q}$, $\hat{\mathbf{m}} = \hat{\mathbf{q}}$ and $\hat{\boldsymbol{\Sigma}} = \hat{\mathbf{q}}$ that lead to

$$\begin{aligned} \hat{\mathbf{q}}_{\text{b}} &= \alpha \mathbb{E}_{y, \boldsymbol{\xi}} \left[\mathcal{Z}_{\text{out}^*} \left(y, \mathbf{q}_{\text{b}}^{1/2} \boldsymbol{\xi}, \boldsymbol{\rho}_{\text{w}^*} - \mathbf{q}_{\text{b}} \right) \mathbf{f}_{\text{out}^*} \left(y, \mathbf{q}_{\text{b}}^{1/2} \boldsymbol{\xi}, \boldsymbol{\rho}_{\text{w}^*} - \mathbf{q}_{\text{b}} \right)^{\otimes 2} \right], \\ \mathbf{q}_{\text{b}} &= \mathbb{E}_{\boldsymbol{\xi}} \left[\mathcal{Z}_{\text{w}^*} \left(\hat{\mathbf{q}}_{\text{b}}^{1/2} \boldsymbol{\xi}, \hat{\mathbf{q}}_{\text{b}} \right) \mathbf{f}_{\text{w}^*} \left(\hat{\mathbf{q}}_{\text{b}}^{1/2} \boldsymbol{\xi}, \hat{\mathbf{q}}_{\text{b}} \right)^{\otimes 2} \right]. \end{aligned} \quad (58)$$

4 Belief Propagation and Approximate Message Passing

The Approximate Message Passing (AMP) algorithm can be seen as Taylor expansion of the loopy Belief Propagation (BP) approach [1, 26, 27], similar to the so-called TAP equation in spin glass theory [28]. While the behavior of AMP can be rigorously studied [29, 30, 31], it is useful and instructive to see how the derivation can be performed in the framework of BP and the cavity method, as was pioneered in [17] for the single layer problem. The derivation uses the GAMP notations of [32] and follows closely the one of [33]. The computation is presented for the committee machine hypothesis class, which is the vectorized version of the GLM, with $K \geq 1$ vectorial parameters $\mathbf{W} = \{\mathbf{w}_k\}_{k=1}^K \in \mathbb{R}^{d \times K}$.

4.1 Factor graph and joint probability distribution

As a central illustration, we present the instructive derivation of the rBP equations starting with the BP equations in the context of committee machines. We recall the Joint Probability Distribution (JPD)

$$P_d(\mathbf{W}|\mathbf{y}, \mathbf{X}) = \frac{P_{\text{out}}(\mathbf{y}|\mathbf{Z})P_w(\mathbf{W})}{\mathcal{Z}_d(\mathbf{y}, \mathbf{X})} = \frac{\prod_{\mu=1}^n P_{\text{out}}(y_\mu|\mathbf{z}_\mu) \prod_{i=1}^d P_w(\mathbf{w}_i)}{\mathcal{Z}_d(\mathbf{y}, \mathbf{X})}, \quad (59)$$

where we defined $\mathbf{Z} = \frac{1}{\sqrt{d}}\mathbf{X}\mathbf{W} \in \mathbb{R}^{n \times K}$ and we assume that the channel and prior distributions factorize over factors $P_{\text{out}}(y_\mu|\mathbf{z}_\mu)$ and variables $P_w(\mathbf{w}_i)$.

The posterior distribution may be represented by the following bipartite factor graph in Fig. 3. In the following, we attach a set of *messages* $\{m_{i \rightarrow \mu}, \tilde{m}_{\mu \rightarrow i}\}_{i=1..n}^{\mu=1..m}$ to

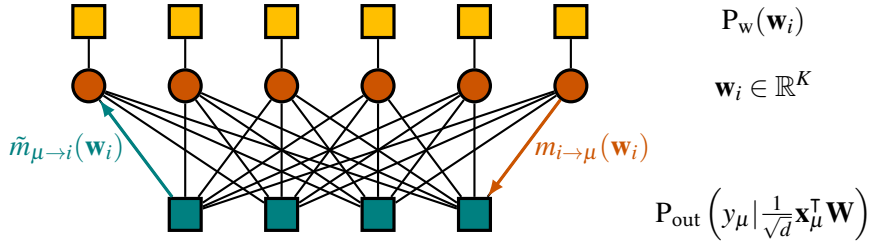


Figure 3 – Factor graph representation of the joint distribution for committee machines.

the edges of this bipartite factor graph. These messages correspond to the marginal probabilities of $\mathbf{w}_i \in \mathbb{R}^K$ if we remove the edges $(i \rightarrow \mu)$ or $(\mu \rightarrow i)$. We define the auxiliary variable $\mathbf{z}_\mu = \frac{1}{\sqrt{d}}\mathbf{x}_\mu^T \mathbf{W} \in \mathbb{R}^K$ which is $\Theta(1)$ thanks to the pre-factor rescaling $1/\sqrt{d}$. This scaling is crucial as it allows the BP equations to hold true even though the factor graph is not tree-like and is instead fully connected with short loops.

4.2 Belief Propagation equations

The BP equations (also called the sum-product equations) for $\mathbf{w}_i = (w_{ik})_{k=1..K} \in \mathbb{R}^K$ on the factor graph Fig. 3 can be formulated as:

$$\begin{aligned} m_{i \rightarrow \mu}^{t+1}(\mathbf{w}_i) &= \frac{1}{\mathcal{Z}_{i \rightarrow \mu}} P_{\mathbf{w}}(\mathbf{w}_i) \prod_{v \neq \mu}^n \tilde{m}_{v \rightarrow i}^t(\mathbf{w}_i) \\ \tilde{m}_{\mu \rightarrow i}^t(\mathbf{w}_i) &= \frac{1}{\mathcal{Z}_{\mu \rightarrow i}} \int_{\mathbb{R}^K} \prod_{j \neq i}^d d\mathbf{w}_j P_{\text{out}} \left(y_{\mu} \middle| \frac{1}{\sqrt{d}} \sum_{j=1}^d x_{\mu j} \mathbf{w}_j \right) m_{j \rightarrow \mu}^t(\mathbf{w}_j), \end{aligned} \quad (60)$$

The BP equations assume that incoming messages are independent. Hence these equations are exact on a tree (no loop), but they remain exact if "correlations decrease fast enough / long loops". We assume in the following that the hypothesis is true in our model.

4.3 Relaxed Belief Propagation equations

The idea of the relaxed BP equations is to simply expand in the limit $d \rightarrow \infty$ the set of $\Theta(d^2)$ messages \tilde{m} of the BP equations in (60) before plugging them in m . Truncating the expansion and keeping only terms of order $\Theta(1/d)$, messages become *Gaussian*. Hence messages are therefore parametrized only by the mean $\hat{\mathbf{w}}_{i \rightarrow \mu}^t$ and the covariance matrix $\hat{\mathbf{C}}_{i \rightarrow \mu}^t$ of the marginal distribution at time t :

$$\begin{aligned} \hat{\mathbf{w}}_{i \rightarrow \mu}^t &\equiv \int_{\mathbb{R}^K} d\mathbf{w}_i m_{i \rightarrow \mu}^t(\mathbf{w}_i) \mathbf{w}_i \\ \hat{\mathbf{C}}_{i \rightarrow \mu}^t &\equiv \int_{\mathbb{R}^K} d\mathbf{w}_i m_{i \rightarrow \mu}^t(\mathbf{w}_i) \mathbf{w}_i \mathbf{w}_i^{\top} - \hat{\mathbf{w}}_{i \rightarrow \mu}^t (\hat{\mathbf{w}}_{i \rightarrow \mu}^t)^{\top} \end{aligned} \quad (61)$$

To decouple the argument of P_{out} , we first by introducing its Fourier transform \hat{P}_{out} according to

$$\begin{aligned} P_{\text{out}} \left(y_{\mu} \middle| \frac{1}{\sqrt{d}} \sum_{j=1}^d x_{\mu j} \mathbf{w}_j \right) &= \frac{1}{(2\pi)^{K/2}} \\ &\times \int_{\mathbb{R}^K} d\boldsymbol{\xi} \exp \left(i\boldsymbol{\xi}^{\top} \left(\frac{1}{\sqrt{d}} \sum_{j=1}^d x_{\mu j} \mathbf{w}_j \right) \right) \hat{P}_{\text{out}}(y_{\mu}, \boldsymbol{\xi}). \end{aligned}$$

Injecting this representation in the BP equations, (60) becomes:

$$\begin{aligned} \tilde{m}_{\mu \rightarrow i}^t(\mathbf{w}_i) &= \frac{1}{(2\pi)^{K/2} \mathcal{Z}_{\mu \rightarrow i}} \int_{\mathbb{R}^K} d\boldsymbol{\xi} \hat{P}_{\text{out}}(y_{\mu}, \boldsymbol{\xi}) \exp \left(i\boldsymbol{\xi}^{\top} \frac{1}{\sqrt{d}} x_{\mu i} \mathbf{w}_i \right) \\ &\times \underbrace{\prod_{j \neq i}^d \int_{\mathbb{R}^K} d\mathbf{w}_j m_{j \rightarrow \mu}^t(\mathbf{w}_j) \exp \left(i\boldsymbol{\xi}^{\top} \frac{1}{\sqrt{d}} x_{\mu j} \mathbf{w}_j \right)}_{\equiv I_j} \end{aligned} \quad (62)$$

In the limit $d \rightarrow \infty$ the term I_j can be easily expanded and expressed using $\hat{\mathbf{w}}$ and $\hat{\mathbf{C}}$ in (61):

$$\begin{aligned} I_j &= \int_{\mathbb{R}^K} d\mathbf{w}_j m'_{j \rightarrow \mu}(\mathbf{w}_j) \exp\left(i \frac{x_{\mu j}}{\sqrt{d}} \boldsymbol{\xi}^\top \mathbf{w}_j\right) \\ &\simeq \exp\left(i \frac{x_{\mu j}}{\sqrt{d}} \boldsymbol{\xi}^\top \hat{\mathbf{w}}'_{j \rightarrow \mu} - \frac{1}{2} \frac{x_{\mu j}^2}{d} \boldsymbol{\xi}^\top \hat{\mathbf{C}}'_{j \rightarrow \mu} \boldsymbol{\xi}\right). \end{aligned}$$

Finally using the inverse Fourier transform:

$$\begin{aligned} \tilde{m}'_{\mu \rightarrow i}(\mathbf{w}_i) &= \frac{1}{(2\pi)^{K/2} \mathcal{Z}_{\mu \rightarrow i}} \int_{\mathbb{R}^K} d\mathbf{z} P_{\text{out}}(y_\mu | \mathbf{z}) \int_{\mathbb{R}^K} d\boldsymbol{\xi} e^{-i \boldsymbol{\xi}^\top \mathbf{z}} e^{i x_{\mu i} \boldsymbol{\xi}^\top \mathbf{w}_i} \\ &\quad \times \prod_{j \neq i}^d \exp\left(i \frac{x_{\mu j}}{\sqrt{d}} \boldsymbol{\xi}^\top \hat{\mathbf{w}}'_{j \rightarrow \mu} - \frac{1}{2} \frac{x_{\mu j}^2}{d} \boldsymbol{\xi}^\top \hat{\mathbf{C}}'_{j \rightarrow \mu} \boldsymbol{\xi}\right) \\ &= \frac{1}{(2\pi)^K \mathcal{Z}_{\mu \rightarrow i}} \int_{\mathbb{R}^K} d\mathbf{z} P_{\text{out}}(y_\mu | \mathbf{z}) \\ &\quad \int_{\mathbb{R}^K} d\boldsymbol{\xi} e^{-i \boldsymbol{\xi}^\top \mathbf{z}} e^{i x_{\mu i} \boldsymbol{\xi}^\top \mathbf{w}_i} e^{i \sum_{j \neq i}^d \frac{x_{\mu j}}{\sqrt{d}} \boldsymbol{\xi}^\top \hat{\mathbf{w}}'_{j \rightarrow \mu}} e^{-\frac{1}{2} \sum_{j \neq i}^d \frac{x_{\mu j}^2}{d} \boldsymbol{\xi}^\top \hat{\mathbf{C}}'_{j \rightarrow \mu} \boldsymbol{\xi}} \\ &= \frac{1}{(2\pi)^K \mathcal{Z}_{\mu \rightarrow i}} \int_{\mathbb{R}^K} d\mathbf{z} P_{\text{out}}(y_\mu | \mathbf{z}) \\ &\quad \times \sqrt{\frac{(2\pi)^K}{\det(\mathbf{V}_{\mu \rightarrow i}^t)}} e^{-\frac{1}{2} \left(\mathbf{z} - \frac{x_{\mu i}}{\sqrt{d}} \mathbf{w}_i - \boldsymbol{\omega}_{\mu \rightarrow i}^t\right)^\top (\mathbf{V}_{\mu \rightarrow i}^t)^{-1} \left(\mathbf{z} - \frac{x_{\mu i}}{\sqrt{d}} \mathbf{w}_i - \boldsymbol{\omega}_{\mu \rightarrow i}^t\right)}, \end{aligned}$$

where we defined the mean and variance, depending on the node i :

$$\boldsymbol{\omega}_{\mu \rightarrow i}^t \equiv \frac{1}{\sqrt{d}} \sum_{j \neq i}^d x_{\mu j} \hat{\mathbf{w}}'_{j \rightarrow \mu}, \quad \mathbf{V}_{\mu \rightarrow i}^t \equiv \frac{1}{d} \sum_{j \neq i}^d x_{\mu j}^2 \hat{\mathbf{C}}'_{j \rightarrow \mu}.$$

Again, in the limit $d \rightarrow \infty$, the term $H_{\mu \rightarrow i}$ can be expanded as

$$\begin{aligned} H_{\mu \rightarrow i} &\simeq e^{-\frac{1}{2} (\mathbf{z} - \boldsymbol{\omega}_{\mu \rightarrow i}^t)^\top (\mathbf{V}_{\mu \rightarrow i}^t)^{-1} (\mathbf{z} - \boldsymbol{\omega}_{\mu \rightarrow i}^t)} \\ &\quad \times \left(1 + \frac{x_{\mu i}}{\sqrt{d}} \mathbf{w}_i^\top (\mathbf{V}_{\mu \rightarrow i}^t)^{-1} (\mathbf{z} - \boldsymbol{\omega}_{\mu \rightarrow i}^t) - \frac{1}{2} \frac{x_{\mu i}^2}{d} \mathbf{w}_i^\top (\mathbf{V}_{\mu \rightarrow i}^t)^{-1} \mathbf{w}_i \right. \\ &\quad \left. + \frac{1}{2} \frac{x_{\mu i}^2}{d} \mathbf{w}_i^\top (\mathbf{V}_{\mu \rightarrow i}^t)^{-1} (\mathbf{z} - \boldsymbol{\omega}_{\mu \rightarrow i}^t) (\mathbf{z} - \boldsymbol{\omega}_{\mu \rightarrow i}^t)^\top (\mathbf{V}_{\mu \rightarrow i}^t)^{-1} \mathbf{w}_i \right). \end{aligned}$$

Putting all pieces together, the message $\tilde{m}_{\mu \rightarrow i}$ can be expressed using definitions of \mathbf{f}_{out} and $\partial_{\omega} \mathbf{f}_{\text{out}}$ in Appendix. ?? . We finally obtain

$$\begin{aligned}
\tilde{m}_{\mu \rightarrow i}^t(\mathbf{w}_i) &\sim \frac{1}{\mathcal{Z}_{\mu \rightarrow i}} \left\{ 1 + \frac{x_{\mu i}}{\sqrt{d}} \mathbf{w}_i^{\top} \mathbf{f}_{\text{out}}(y_{\mu}, \boldsymbol{\omega}_{\mu \rightarrow i}^t, \mathbf{V}_{\mu \rightarrow i}^t) \right. \\
&\quad + \frac{1}{2} \frac{x_{\mu i}^2}{d} \mathbf{w}_i^{\top} \mathbf{f}_{\text{out}} \mathbf{f}_{\text{out}}^{\top}(y_{\mu}, \boldsymbol{\omega}_{\mu \rightarrow i}^t, \mathbf{V}_{\mu \rightarrow i}^t) \mathbf{w}_i \\
&\quad \left. + \frac{1}{2} \frac{x_{\mu i}^2}{d} \mathbf{w}_i^{\top} \partial_{\omega} \mathbf{f}_{\text{out}}(y_{\mu}, \boldsymbol{\omega}_{\mu \rightarrow i}^t, \mathbf{V}_{\mu \rightarrow i}^t) \mathbf{w}_i \right\} \\
&= \frac{1}{\mathcal{Z}_{\mu \rightarrow i}} \left\{ 1 + \mathbf{w}_i^{\top} \mathbf{b}_{\mu \rightarrow i}^t + \frac{1}{2} \mathbf{w}_i^{\top} \mathbf{b}_{\mu \rightarrow i}^t (\mathbf{b}_{\mu \rightarrow i}^t)^{\top} (\mathbf{w}_i) - \frac{1}{2} \mathbf{w}_i^{\top} \mathbf{A}_{\mu \rightarrow i}^t \mathbf{w}_i \right\} \\
&= \sqrt{\frac{\det(\mathbf{A}_{\mu \rightarrow i}^t)}{(2\pi)^K}} e^{-\frac{1}{2} (\mathbf{w}_i^{\top} - (\mathbf{A}_{\mu \rightarrow i}^t)^{-1} \mathbf{b}_{\mu \rightarrow i}^t)^{\top} \mathbf{A}_{\mu \rightarrow i}^t (\mathbf{w}_i^{\top} - (\mathbf{A}_{\mu \rightarrow i}^t)^{-1} \mathbf{b}_{\mu \rightarrow i}^t)}
\end{aligned}$$

with the following definitions of $\mathbf{A}_{\mu \rightarrow i}$ and $\mathbf{b}_{\mu \rightarrow i}$

$$\begin{aligned}
\mathbf{b}_{\mu \rightarrow i}^t &\equiv \frac{x_{\mu i}}{\sqrt{d}} \mathbf{f}_{\text{out}}(y_{\mu}, \boldsymbol{\omega}_{\mu \rightarrow i}^t, \mathbf{V}_{\mu \rightarrow i}^t), \\
\mathbf{A}_{\mu \rightarrow i}^t &\equiv -\frac{x_{\mu i}^2}{d} \partial_{\omega} \mathbf{f}_{\text{out}}(y_{\mu}, \boldsymbol{\omega}_{\mu \rightarrow i}^t, \mathbf{V}_{\mu \rightarrow i}^t).
\end{aligned}$$

The set of BP equations can finally be closed over the Gaussian messages $\{m_{i \rightarrow \mu}\}_{i=1..n}^{\mu=1..n}$ according to

$$\begin{aligned}
m_{i \rightarrow \mu}^{t+1}(\mathbf{w}_i) &= \frac{1}{\mathcal{Z}_{i \rightarrow \mu}} \mathbf{P}_{\mathbf{w}}(\mathbf{w}_i) \prod_{v \neq \mu}^n \sqrt{\frac{\det(\mathbf{A}_{v \rightarrow i}^t)}{(2\pi)^K}} \\
&\quad \times e^{-\frac{1}{2} (\mathbf{w}_i^{\top} - (\mathbf{A}_{v \rightarrow i}^t)^{-1} \mathbf{b}_{v \rightarrow i}^t)^{\top} \mathbf{A}_{v \rightarrow i}^t (\mathbf{w}_i^{\top} - (\mathbf{A}_{v \rightarrow i}^t)^{-1} \mathbf{b}_{v \rightarrow i}^t)}.
\end{aligned}$$

In the end, computing the mean and variance of the product of Gaussians, the messages are updated using $\mathbf{f}_{\mathbf{w}}$ and $\partial_{\gamma} \mathbf{f}_{\mathbf{w}}$, defined in Appendix. ?? , according to

$$\hat{\mathbf{w}}_{i \rightarrow \mu}^{t+1} = \mathbf{f}_{\mathbf{w}}(\boldsymbol{\gamma}_{\mu \rightarrow i}^t, \boldsymbol{\Lambda}_{\mu \rightarrow i}^t), \quad \hat{\mathbf{C}}_{i \rightarrow \mu}^{t+1} = \partial_{\gamma} \mathbf{f}_{\mathbf{w}}(\boldsymbol{\gamma}_{\mu \rightarrow i}^t, \boldsymbol{\Lambda}_{\mu \rightarrow i}^t),$$

with

$$\boldsymbol{\gamma}_{\mu \rightarrow i}^t = \sum_{v \neq \mu}^n \mathbf{b}_{v \rightarrow i}^t, \quad \boldsymbol{\Lambda}_{\mu \rightarrow i}^t = \sum_{v \neq \mu}^n \mathbf{A}_{v \rightarrow i}^t.$$

Summary of the rBP equations In the end, the rBP equations are simply the following set of equations:

$$\begin{aligned}
\hat{\mathbf{w}}_{i \rightarrow \mu}^{t+1} &= \mathbf{f}_w(\boldsymbol{\gamma}_{\mu \rightarrow i}^t, \boldsymbol{\Lambda}_{\mu \rightarrow i}^t), & \hat{\mathbf{C}}_{i \rightarrow \mu}^{t+1} &= \partial_{\boldsymbol{\gamma}} \mathbf{f}_w(\boldsymbol{\gamma}_{\mu \rightarrow i}^t, \boldsymbol{\Lambda}_{\mu \rightarrow i}^t) \\
\boldsymbol{\gamma}_{\mu \rightarrow i}^t &= \sum_{v \neq \mu}^n \mathbf{b}_{v \rightarrow i}^t, & \boldsymbol{\Lambda}_{\mu \rightarrow i}^t &= \sum_{v \neq \mu}^n \mathbf{A}_{v \rightarrow i}^t \\
\mathbf{b}_{\mu \rightarrow i}^t &= \frac{x_{\mu i}}{\sqrt{d}} \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_{\mu \rightarrow i}^t, \mathbf{V}_{\mu \rightarrow i}^t), \\
\mathbf{A}_{\mu \rightarrow i}^t &= -\frac{x_{\mu i}^2}{d} \partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_{\mu \rightarrow i}^t, \mathbf{V}_{\mu \rightarrow i}^t) \\
\boldsymbol{\omega}_{\mu \rightarrow i}^t &= \sum_{j \neq i}^d \frac{x_{\mu j}}{\sqrt{d}} \hat{\mathbf{w}}_{j \rightarrow \mu}^t, & \mathbf{V}_{\mu \rightarrow i}^t &= \sum_{j \neq i}^d \frac{x_{\mu j}^2}{d} \hat{\mathbf{C}}_{j \rightarrow \mu}^t.
\end{aligned} \tag{63}$$

4.4 AMP algorithm

The rBP equations eq. (63) contains $\Theta(d^2)$ messages. However all the messages depend weakly on the target node. The missing message is negligible in the limit $d \rightarrow \infty$, that allows us to expand the rBP around the *full* messages:

$$\begin{aligned}
\boldsymbol{\omega}_\mu^t &\equiv \sum_{j=1}^d \frac{x_{\mu j}}{\sqrt{d}} \hat{\mathbf{w}}_{j \rightarrow \mu}^t, & \mathbf{V}_\mu^t &\equiv \sum_{j=1}^d \frac{x_{\mu j}^2}{d} \hat{\mathbf{C}}_{j \rightarrow \mu}^t \\
\boldsymbol{\gamma}_i^t &\equiv \sum_{\mu=1}^n \mathbf{b}_{\mu \rightarrow i}^t, & \boldsymbol{\Lambda}_i^t &\equiv \sum_{\mu=1}^n \mathbf{A}_{\mu \rightarrow i}^t.
\end{aligned} \tag{64}$$

By completing the sum, we naturally remove the target node dependence and reduce the set of messages to $\Theta(d)$. Let us now perform the expansion of the rBP messages.

Partial covariance \mathbf{f}_w : $\boldsymbol{\Lambda}_{\mu \rightarrow i}^t$

$$\begin{aligned}
\boldsymbol{\Lambda}_{\mu \rightarrow i}^t &= \sum_{v \neq \mu}^n \mathbf{A}_{v \rightarrow i}^t = \sum_{v=1}^n \mathbf{A}_{v \rightarrow i}^t - \mathbf{A}_{\mu \rightarrow i}^t \\
&= \boldsymbol{\Lambda}_i^t - \mathbf{A}_{\mu \rightarrow i}^t = \boldsymbol{\Lambda}_i^t + \Theta\left(\frac{1}{d}\right).
\end{aligned}$$

Partial mean \mathbf{f}_w : $\boldsymbol{\gamma}_{\mu \rightarrow i}^t$

$$\boldsymbol{\gamma}_{\mu \rightarrow i}^t = \sum_{v \neq \mu}^n \mathbf{b}_{v \rightarrow i}^t = \sum_{v=1}^n \mathbf{b}_{v \rightarrow i}^t - \mathbf{b}_{\mu \rightarrow i}^t = \boldsymbol{\gamma}_i^t - \mathbf{b}_{\mu \rightarrow i}^t + \Theta\left(\frac{1}{d}\right).$$

Mean $\hat{\mathbf{w}}_{i \rightarrow \mu}^{t+1}$ update

$$\begin{aligned}
\hat{\mathbf{w}}_{i \rightarrow \mu}^{t+1} &= \mathbf{f}_w(\boldsymbol{\gamma}_{\mu \rightarrow i}^t, \boldsymbol{\Lambda}_{\mu \rightarrow i}^t) = \mathbf{f}_w(\boldsymbol{\gamma}_i^t - \mathbf{b}_{\mu \rightarrow i}^t, \boldsymbol{\Lambda}_i^t) + \Theta\left(\frac{1}{d}\right) \\
&= \mathbf{f}_w(\boldsymbol{\gamma}_i^t, \boldsymbol{\Lambda}_i^t) - \partial_{\boldsymbol{\gamma}} \mathbf{f}_w(\boldsymbol{\gamma}_i^t, \boldsymbol{\Lambda}_i^t) \mathbf{b}_{\mu \rightarrow i}^t + \Theta\left(\frac{1}{d}\right) \\
&= \hat{\mathbf{w}}_i^{t+1} - \hat{\mathbf{C}}_i^{t+1} \mathbf{b}_{\mu \rightarrow i}^t + \Theta\left(\frac{1}{d}\right) \\
&= \hat{\mathbf{w}}_i^{t+1} - \frac{x_{\mu i}}{\sqrt{d}} \hat{\mathbf{C}}_i^{t+1} \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^t, \mathbf{V}_\mu^t) + \Theta\left(\frac{1}{d}\right).
\end{aligned}$$

where we defined the prior updates

$$\hat{\mathbf{w}}_i^{t+1} \equiv \mathbf{f}_w(\boldsymbol{\gamma}_i^t, \boldsymbol{\Lambda}_i^t), \quad \hat{\mathbf{C}}_i^{t+1} \equiv \partial_{\boldsymbol{\gamma}} \mathbf{f}_w(\boldsymbol{\gamma}_i^t, \boldsymbol{\Lambda}_i^t),$$

and used the fact that $\mathbf{b}_{\mu \rightarrow i}^t \simeq \frac{x_{\mu i}}{\sqrt{d}} \hat{\mathbf{C}}_i^{t+1} \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^t, \mathbf{V}_\mu^t)$ by expanding the equation over $\mathbf{b}_{\mu \rightarrow i}^t$ in (63).

Covariance $\hat{\mathbf{C}}_{i \rightarrow \mu}^{t+1}$ update

$$\begin{aligned}
\hat{\mathbf{C}}_{i \rightarrow \mu}^{t+1} &= \partial_{\boldsymbol{\gamma}} \mathbf{f}_w(\boldsymbol{\gamma}_{\mu \rightarrow i}^t, \boldsymbol{\Lambda}_{\mu \rightarrow i}^t) \\
&\simeq \partial_{\boldsymbol{\gamma}} \mathbf{f}_w(\boldsymbol{\gamma}_i^t, \boldsymbol{\Lambda}_i^t) + \Theta\left(\frac{1}{\sqrt{d}}\right) = \hat{\mathbf{C}}_i^{t+1} + \Theta\left(\frac{1}{\sqrt{d}}\right).
\end{aligned}$$

Channel update function $\mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_{\mu \rightarrow i}^t, \mathbf{V}_{\mu \rightarrow i}^t)$

$$\begin{aligned}
\mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_{\mu \rightarrow i}^t, \mathbf{V}_{\mu \rightarrow i}^t) &= \mathbf{f}_{\text{out}}\left(y_\mu, \boldsymbol{\omega}_\mu^t - \frac{x_{\mu i}}{\sqrt{d}} \hat{\mathbf{w}}_{i \rightarrow \mu}^t, \mathbf{V}_\mu^t - \frac{x_{\mu i}^2}{d} \hat{\mathbf{C}}_{i \rightarrow l}^t\right) \\
&= \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^t, \mathbf{V}_\mu^t) - \frac{x_{\mu i}}{\sqrt{d}} \partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^t, \mathbf{V}_\mu^t) \underbrace{\hat{\mathbf{w}}_{i \rightarrow \mu}^t}_{= \hat{\mathbf{w}}_i^t + \Theta\left(\frac{1}{\sqrt{d}}\right)} + \Theta\left(\frac{1}{d}\right) \\
&= \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^t, \mathbf{V}_\mu^t) - \frac{x_{\mu i}}{\sqrt{d}} \partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^t, \mathbf{V}_\mu^t) \hat{\mathbf{w}}_i^t + \Theta\left(\frac{1}{d}\right).
\end{aligned}$$

Covariance $\mathbf{f}_{\text{out}}: \mathbf{V}_\mu^t$

$$\mathbf{V}_\mu^t \equiv \sum_{j=1}^d \frac{x_{\mu j}^2}{d} \hat{\mathbf{C}}_{j \rightarrow \mu}^t = \sum_{j=1}^d \frac{x_{\mu j}^2}{d} \hat{\mathbf{C}}_{j \rightarrow \mu}^t + \Theta\left(\frac{1}{d^{3/2}}\right).$$

Mean \mathbf{f}_{out} : $\boldsymbol{\omega}_\mu^t$

$$\begin{aligned}\boldsymbol{\omega}_\mu^t &= \sum_{i=1}^d \frac{x_{\mu i}}{\sqrt{d}} \hat{\mathbf{w}}_{i \rightarrow \mu}^t \\ &= \sum_{i=1}^d \frac{x_{\mu i}}{\sqrt{d}} \left(\hat{\mathbf{w}}_i^t - x_{\mu i} \hat{\mathbf{C}}_i^t \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^{t-1}, \mathbf{V}_\mu^{t-1}) + \Theta\left(\frac{1}{d}\right) \right) \\ &= \sum_{i=1}^d \frac{x_{\mu i}}{\sqrt{d}} \hat{\mathbf{w}}_i^t - \sum_{i=1}^d \frac{x_{\mu i}^2}{d} \hat{\mathbf{C}}_i^t \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^{t-1}, \mathbf{V}_\mu^{t-1}) + \Theta\left(\frac{1}{d^{3/2}}\right).\end{aligned}$$

Covariance \mathbf{f}_w : $\boldsymbol{\Lambda}_i^t$

$$\begin{aligned}\boldsymbol{\Lambda}_i^t &\equiv \sum_{\mu=1}^n \mathbf{A}_{\mu \rightarrow i}^t = \sum_{\mu=1}^n -\frac{x_{\mu i}^2}{d} \partial \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_{\mu \rightarrow i}^t, \mathbf{V}_{\mu \rightarrow i}^t) \\ &= \sum_{\mu=1}^n -\frac{x_{\mu i}^2}{d} \partial \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^t, \mathbf{V}_\mu^t) + \Theta\left(\frac{1}{d^{3/2}}\right).\end{aligned}$$

Mean \mathbf{f}_w : $\boldsymbol{\gamma}_i^t$

$$\begin{aligned}\boldsymbol{\gamma}_i^t &= \sum_{\mu=1}^n \mathbf{b}_{\mu \rightarrow i}^t = \sum_{\mu=1}^n \frac{x_{\mu i}}{\sqrt{d}} \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_{\mu \rightarrow i}^t, \mathbf{V}_{\mu \rightarrow i}^t) \\ &= \sum_{\mu=1}^n \frac{x_{\mu i}}{\sqrt{d}} \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^t, \mathbf{V}_\mu^t) \\ &\quad - \frac{x_{\mu i}^2}{d} \partial \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^t, \mathbf{V}_\mu^t) \hat{\mathbf{w}}_i^t + \Theta\left(\frac{1}{d^{3/2}}\right).\end{aligned}$$

Summary - AMP algorithm

We finally obtain the AMP algorithm as a reduced set of $\Theta(d)$ messages in Algo. 1.

4.5 State evolution equations of AMP

In this section we derive the behavior of the AMP algorithm in Algo. 1 in the thermodynamic limit $d \rightarrow \infty$. This average asymptotic behavior can be tracked with some overlap parameters at time t , \mathbf{m}^t , \mathbf{q}^t , $\boldsymbol{\Sigma}^t$, that respectively measure the correlation of the AMP estimator with the ground truth, the norms of student and teacher weights, the estimator variance and the second moment of the teacher network $\boldsymbol{\rho}_{w^*}$, defined by

$$\begin{aligned}\mathbf{m}^t &\equiv \mathbb{E} \lim_{d \rightarrow \infty} \frac{1}{d} \hat{\mathbf{W}}^{t\top} \hat{\mathbf{W}}^*, & \mathbf{q}^t &\equiv \mathbb{E} \lim_{d \rightarrow \infty} \frac{1}{d} \hat{\mathbf{W}}^{t\top} \hat{\mathbf{W}}^t, \\ \boldsymbol{\Sigma}^t &\equiv \mathbb{E} \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=1}^d \hat{\mathbf{C}}_i^t, & \boldsymbol{\rho}_{w^*} &\equiv \mathbb{E} \lim_{d \rightarrow \infty} \frac{1}{d} \mathbf{W}^{*\top} \mathbf{W}^*,\end{aligned}\tag{65}$$

Input: vector $\mathbf{y} \in \mathbb{R}^n$ and matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$.

Initialize: $\hat{\mathbf{w}}_i, \mathbf{f}_{\text{out},\mu} \in \mathbb{R}^K$ and $\hat{\mathbf{V}}_i, \partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out},\mu} \in \mathbb{R}^{K \times K}$ for $1 \leq i \leq d$ and $1 \leq \mu \leq n$ at $t = 0$.

repeat

Channel: Update the mean $\boldsymbol{\omega}_\mu \in \mathbb{R}^K$ and variance $V_\mu \in \mathbb{R}^{K \times K}$:

$$\mathbf{V}_\mu^t = \sum_{i=1}^d \frac{x_{\mu i}^2}{d} \hat{\mathbf{C}}_i^t$$

$$\boldsymbol{\omega}_\mu^t = \sum_{i=1}^d \frac{x_{\mu i}}{\sqrt{d}} \hat{\mathbf{w}}_i^t - \mathbf{V}_\mu^t \mathbf{f}_{\text{out},\mu}^{t-1},$$

Update $\mathbf{f}_{\text{out},\mu}$ and $\partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out},\mu}$:

$$\mathbf{f}_{\text{out},\mu}^t = \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^t, \mathbf{V}_\mu^t), \partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out},\mu}^t = \partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out}}(y_\mu, \boldsymbol{\omega}_\mu^t, \mathbf{V}_\mu^t)$$

Prior: Update the mean $\boldsymbol{\gamma}_i \in \mathbb{R}^K$ and variance $\boldsymbol{\Lambda}_i \in \mathbb{R}^{K \times K}$:

$$\boldsymbol{\Lambda}_i^t = \sum_{\mu=1}^n -\frac{x_{\mu i}^2}{d} \partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out},\mu}$$

$$\boldsymbol{\gamma}_i^t = \sum_{\mu=1}^n \frac{x_{\mu i}}{\sqrt{d}} \mathbf{f}_{\text{out},\mu}^t + \boldsymbol{\Lambda}_i^t \hat{\mathbf{w}}_i^t,$$

Update the estimated marginals $\hat{\mathbf{w}}_i \in \mathbb{R}$ and $\hat{\mathbf{C}}_i \in \mathbb{R}^+$:

$$\hat{\mathbf{w}}_i^{t+1} = \mathbf{f}_w(\boldsymbol{\gamma}_i^t, \boldsymbol{\Lambda}_i^t), \hat{\mathbf{C}}_i^{t+1} = \partial_{\boldsymbol{\gamma}} \mathbf{f}_w(\boldsymbol{\gamma}_i^t, \boldsymbol{\Lambda}_i^t)$$

$t \leftarrow t + 1$

until Convergence on $\hat{\mathbf{w}}_i, \hat{\mathbf{C}}_i$.

Output: $\{\hat{\mathbf{w}}_i\}_{i=1}^d$ and $\{\hat{\mathbf{C}}_i\}_{i=1}^d$.

Algorithm 1: Approximate Message Passing algorithm for committee machines.

where the expectation is over ground truth signals \mathbf{W}^* and input data \mathbf{X} . The aim is to derive the asymptotic behavior of these overlap parameters, called SE. The idea is simply to compute the overlap distributions starting with the set of rBP equations in (63).

4.5.1 Messages distribution

In order to get the asymptotic behavior of the overlap parameters, we first need to compute the distribution of \mathbf{W}^{t+1} and, as a result, of the mean $\mathbf{y}_{\mu \rightarrow i}^t$ and covariance $\mathbf{\Lambda}_{\mu \rightarrow i}^t$. Recalling that under the BP assumption incoming messages are independent, the messages $\mathbf{\omega}_{\mu \rightarrow i}^t$ and \mathbf{z}_μ are the sum of independent variables and follow Gaussian distributions. However, these two variables are correlated and we need to compute correctly the covariance matrix.

To compute it, we will make use of different ingredients. First, we recall that in the T-S scenario, the output has been generated by a teacher such that $\forall \mu \in \llbracket n \rrbracket$, $y_\mu = \varphi_{\text{out}^*} \left(\frac{1}{\sqrt{d}} \mathbf{x}_\mu^\top \mathbf{W}^* \right)$. By convenience, we define $\mathbf{z}_\mu \equiv \frac{1}{\sqrt{d}} \mathbf{x}_\mu^\top \mathbf{W}^* = \frac{1}{\sqrt{d}} \sum_{i=1}^d x_{\mu i} \mathbf{w}_i^*$ and $z_{\mu \rightarrow i} \equiv \frac{1}{\sqrt{d}} \sum_{j \neq i}^d x_{\mu j} \mathbf{w}_j^*$. Second, in the case the input data are i.i.d Gaussian, we have $\mathbb{E}_{\mathbf{X}}[x_{\mu i}] = 0$ and $\mathbb{E}_{\mathbf{X}}[x_{\mu i}^2] = 1$.

Partial mean $\mathbf{f}_{\text{out}}: \mathbf{\omega}_{\mu \rightarrow i}^t$ Let's compute the first two moments, using expansions of the rBP equations (63):

$$\begin{aligned} \mathbb{E}[\mathbf{\omega}_{\mu \rightarrow i}^t] &= \frac{1}{\sqrt{d}} \sum_{j \neq i}^d \mathbb{E}_{\mathbf{X}}[x_{\mu j}] \mathbb{E}[\hat{\mathbf{w}}_{j \rightarrow \mu}^t] = \mathbf{0}, \\ \mathbb{E}[\mathbf{\omega}_{\mu \rightarrow i}^t (\mathbf{\omega}_{\mu \rightarrow i}^t)^\top] &= \frac{1}{d} \sum_{j \neq i}^d \mathbb{E}_{\mathbf{X}}[x_{\mu j}^2] \mathbb{E}[\hat{\mathbf{w}}_{j \rightarrow \mu}^t (\hat{\mathbf{w}}_{j \rightarrow \mu}^t)^\top] \\ &= \frac{1}{d} \sum_{i=1}^d \mathbb{E}_{\mathbf{X}}[x_{\mu i}^2] \mathbb{E}[\hat{\mathbf{w}}_i^t (\hat{\mathbf{w}}_i^t)^\top] + \Theta(d^{-3/2}) \xrightarrow{d \rightarrow \infty} \mathbf{q}^t. \end{aligned}$$

Hidden variable \mathbf{z}_μ Let us compute the first moments of the hidden variable \mathbf{z}_μ :

$$\begin{aligned} \mathbb{E}[\mathbf{z}_\mu] &= \frac{1}{\sqrt{d}} \sum_{i=1}^d \mathbb{E}_{\mathbf{X}}[x_{\mu i}] \mathbb{E}_{\mathbf{W}^*}[\mathbf{w}_i^*] = \mathbf{0}, \\ \mathbb{E}[\mathbf{z}_\mu \mathbf{z}_\mu^\top] &= \frac{1}{d} \sum_{i=1}^d \mathbb{E}_{\mathbf{X}}[x_{\mu i}^2] \mathbb{E}_{\mathbf{W}^*}[\mathbf{w}_i^* (\mathbf{w}_i^*)^\top] \xrightarrow{d \rightarrow \infty} \boldsymbol{\rho}_{\mathbf{w}^*}. \end{aligned}$$

Correlation between \mathbf{z}_μ and $\boldsymbol{\omega}_{\mu \rightarrow i}^t$ The cross correlation is given by

$$\begin{aligned}\mathbb{E}[\boldsymbol{\omega}_{\mu \rightarrow i}^t \mathbf{z}_\mu^\top] &= \frac{1}{d} \sum_{j \neq i, k=1}^d \mathbb{E}_{\mathbf{X}}[x_{\mu j} x_{\mu k}] \mathbb{E}_{\mathbf{W}^*}[\hat{\mathbf{w}}_{j \rightarrow \mu}^t (\mathbf{w}_k^*)^\top] \\ &= \frac{1}{d} \sum_{j \neq i}^d \mathbb{E}_{\mathbf{W}^*}[\hat{\mathbf{w}}_{j \rightarrow \mu}^t (\mathbf{w}_j^*)^\top] = \frac{1}{d} \sum_i^d \mathbb{E}_{\mathbf{W}^*}[\hat{\mathbf{w}}_i^t (\mathbf{w}_i^*)^\top] + \Theta(d^{-3/2}) \\ &\xrightarrow{d \rightarrow \infty} \mathbf{m}^t.\end{aligned}$$

Hence asymptotically the random vector $(\mathbf{z}_\mu, \boldsymbol{\omega}_{\mu \rightarrow i}^t)$ follow a multivariate Gaussian distribution with covariance matrix $\mathbf{Q}^t = \begin{bmatrix} \boldsymbol{\rho}_{\mathbf{w}^*} & \mathbf{m}^t \\ \mathbf{m}^t & \mathbf{q}^t \end{bmatrix} \in \mathbb{R}^{(2K) \times (2K)}$.

Partial variance \mathbf{f}_{out} : $\mathbf{V}_{\mu \rightarrow i}$ concentrates around its mean:

$$\mathbb{E}[\mathbf{V}_{\mu \rightarrow i}^t] = \frac{1}{d} \sum_{j \neq i}^d \mathbb{E}_{\mathbf{X}}[x_{\mu j}^2] \hat{\mathbf{C}}_{j \rightarrow \mu}^t = \frac{1}{d} \sum_i^d \hat{\mathbf{C}}_i^t + \Theta(d^{-3/2}) \xrightarrow{d \rightarrow \infty} \boldsymbol{\Sigma}^t.$$

Ad-hoc overlaps Let us define some other ad-hoc order parameters, that will appear in the following:

$$\begin{aligned}\hat{\mathbf{q}}^t &\equiv \alpha \mathbb{E}_{\boldsymbol{\omega}, \mathbf{z}}[\mathbf{f}_{\text{out}}(\boldsymbol{\varphi}_{\text{out}^*}(\mathbf{z}), \boldsymbol{\omega}, \boldsymbol{\Sigma}^t)^{\otimes 2}], \\ \hat{\mathbf{m}}^t &\equiv \alpha \mathbb{E}_{\boldsymbol{\omega}, \mathbf{z}}[\partial_{\mathbf{z}} \mathbf{f}_{\text{out}}(\boldsymbol{\varphi}_{\text{out}^*}(\mathbf{z}), \boldsymbol{\omega}, \boldsymbol{\Sigma}^t)], \\ \hat{\boldsymbol{\chi}}^t &\equiv \alpha \mathbb{E}_{\boldsymbol{\omega}, \mathbf{z}}[-\partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out}}(\boldsymbol{\varphi}_{\text{out}^*}(\mathbf{z}), \boldsymbol{\omega}, \boldsymbol{\Sigma}^t)].\end{aligned}\tag{66}$$

Partial mean \mathbf{f}_{w} : $\boldsymbol{\gamma}_{\mu \rightarrow i}^t$ Using the expression $y_v = \boldsymbol{\varphi}_{\text{out}^*}(\mathbf{z}_{v \rightarrow i} + \frac{1}{\sqrt{d}} x_{vi} \mathbf{w}_i^*)$ and expanding $\boldsymbol{\gamma}_{\mu \rightarrow i}^t$, we obtain

$$\begin{aligned}\boldsymbol{\gamma}_{\mu \rightarrow i}^t &= \sum_{v \neq \mu}^n \mathbf{b}_{v \rightarrow i}^t = \sum_{v \neq \mu}^n \frac{x_{vi}}{\sqrt{d}} \mathbf{f}_{\text{out}}(y_v, \boldsymbol{\omega}_{v \rightarrow i}^t, \mathbf{V}_{v \rightarrow i}^t) \\ &= \frac{1}{\sqrt{d}} \sum_{v \neq \mu}^n x_{vi} \mathbf{f}_{\text{out}}(\boldsymbol{\varphi}_{\text{out}^*}(\mathbf{z}_{v \rightarrow i}), \boldsymbol{\omega}_{v \rightarrow i}^t, \mathbf{V}_{v \rightarrow i}^t) \\ &\quad + \frac{1}{d} \sum_{v \neq \mu}^n x_{vi}^2 \partial_{\mathbf{z}} \mathbf{f}_{\text{out}}(\boldsymbol{\varphi}_{\text{out}^*}(\mathbf{z}_{v \rightarrow i}), \boldsymbol{\omega}_{v \rightarrow i}^t, \mathbf{V}_{v \rightarrow i}^t) \mathbf{w}_i^*.\end{aligned}$$

Thus, taking the average

$$\begin{aligned}\mathbb{E}[\boldsymbol{\gamma}_{\mu \rightarrow i}^t] &= \mathbf{0} + \frac{1}{d} \sum_{v \neq \mu}^n \mathbb{E}_{\mathbf{z}, \boldsymbol{\omega}} [\partial_{\mathbf{z}} \mathbf{f}_{\text{out}}(\varphi_{\text{out}^*}(\mathbf{z}_{v \rightarrow i}), \boldsymbol{\omega}_{v \rightarrow i}^t, \mathbf{V}_{v \rightarrow i}^t)] \mathbf{w}_i^* \\ &\xrightarrow{d \rightarrow \infty} \hat{\mathbf{m}}^t \mathbf{w}_i^*, \\ \mathbb{E}[(\boldsymbol{\gamma}_{\mu \rightarrow i}^t)^{\otimes 2}] &= \frac{1}{d} \sum_{v \neq \mu}^n \mathbb{E}_{\mathbf{z}, \boldsymbol{\omega}} [\mathbf{f}_{\text{out}}(\varphi_{\text{out}^*}(\mathbf{z}_{v \rightarrow i}), \boldsymbol{\omega}_{v \rightarrow i}^t, \mathbf{V}_{v \rightarrow i}^t)^{\otimes 2}] \\ &\xrightarrow{d \rightarrow \infty} \hat{\mathbf{q}}^t.\end{aligned}$$

Hence $\boldsymbol{\gamma}_{\mu \rightarrow i}^t \sim \hat{\mathbf{m}}^t \mathbf{w}_i^* + (\hat{\mathbf{q}}^t)^{1/2} \boldsymbol{\xi}$ with $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$.

Partial covariance \mathbf{f}_w : $\boldsymbol{\Lambda}_{\mu \rightarrow i}^t$

$$\begin{aligned}\boldsymbol{\Lambda}_{\mu \rightarrow i}^t &= \sum_{v \neq \mu}^n \mathbf{A}_{v \rightarrow i}^t = -\frac{1}{d} \sum_{v \neq \mu}^n x_{\mu i}^2 \partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out}}(y_v, \boldsymbol{\omega}_{v \rightarrow i}^t, \mathbf{V}_{v \rightarrow i}^t) \\ &= -\frac{1}{d} \sum_{v \neq \mu}^n x_{\mu i}^2 \partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out}}(\varphi_{\text{out}^*}(\mathbf{z}_{v \rightarrow i}), \boldsymbol{\omega}_{v \rightarrow i}^t, \mathbf{V}_{v \rightarrow i}^t) + \Theta(d^{-3/2})\end{aligned}$$

and taking the average

$$\begin{aligned}\mathbb{E}[\boldsymbol{\Lambda}_{\mu \rightarrow i}^t] &= -\frac{1}{d} \sum_{v \neq \mu}^n \mathbb{E}_{\mathbf{z}, \boldsymbol{\omega}} [\partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out}}(\varphi_{\text{out}^*}(\mathbf{z}_{v \rightarrow i}), \boldsymbol{\omega}_{v \rightarrow i}^t, \mathbf{V}_{v \rightarrow i}^t)] \\ &\xrightarrow{d \rightarrow \infty} \hat{\boldsymbol{\chi}}^t,\end{aligned}$$

so that in the thermodynamic limit $\boldsymbol{\Lambda}_{\mu \rightarrow i}^t \sim \hat{\boldsymbol{\chi}}^t$.

4.5.2 Summary of the SE - mismatched setting

Using the definition of the overlaps in (65) at time $t+1$ and the message distributions, we finally obtain the set of SE equations of the AMP algorithm in Algo. 1 in the mismatched setting:

$$\begin{aligned}\mathbf{m}^{t+1} &\equiv \mathbb{E} \lim_{d \rightarrow \infty} \frac{1}{d} \hat{\mathbf{W}}^{t+1 \top} \hat{\mathbf{W}}^* = \mathbb{E}_{\mathbf{w}^*, \boldsymbol{\xi}} \left[\mathbf{f}_w \left(\hat{\mathbf{m}}^t \mathbf{w}^* + (\hat{\mathbf{q}}^t)^{1/2} \boldsymbol{\xi}, \hat{\boldsymbol{\chi}}^t \right) \mathbf{w}^{* \top} \right], \\ \mathbf{q}^{t+1} &\equiv \mathbb{E} \lim_{d \rightarrow \infty} \frac{1}{d} \hat{\mathbf{W}}^{t+1 \top} \hat{\mathbf{W}}^{t+1} = \mathbb{E}_{\mathbf{w}^*, \boldsymbol{\xi}} \left[\mathbf{f}_w \left(\hat{\mathbf{m}}^t \mathbf{w}^* + (\hat{\mathbf{q}}^t)^{1/2} \boldsymbol{\xi}, \hat{\boldsymbol{\chi}}^t \right)^{\otimes 2} \right], \\ \boldsymbol{\Sigma}^{t+1} &\equiv \mathbb{E} \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=1}^d \hat{\mathbf{C}}_i^{t+1} = \mathbb{E}_{\mathbf{w}^*, \boldsymbol{\xi}} \left[\partial_{\boldsymbol{\gamma}} \mathbf{f}_w \left(\hat{\mathbf{m}}^t \mathbf{w}^* + (\hat{\mathbf{q}}^t)^{1/2} \boldsymbol{\xi}, \hat{\boldsymbol{\chi}}^t \right) \right],\end{aligned}\tag{67}$$

and

$$\begin{aligned}
\hat{\mathbf{q}}^t &= \alpha \int_{\mathbb{R}^K} \int_{\mathbb{R}^K} d\boldsymbol{\omega} d\mathbf{z} \mathcal{N}_{(\mathbf{z}, \boldsymbol{\omega})}(\mathbf{0}_{2K}, \mathbf{Q}^t) \mathbf{f}_{\text{out}}(\varphi_{\text{out}^*}(\mathbf{z}), \boldsymbol{\omega}, \boldsymbol{\Sigma}^t)^{\otimes 2} \\
\hat{\mathbf{m}}^t &= \alpha \int_{\mathbb{R}^K} \int_{\mathbb{R}^K} d\boldsymbol{\omega} d\mathbf{z} \mathcal{N}_{(\mathbf{z}, \boldsymbol{\omega})}(\mathbf{0}_{2K}, \mathbf{Q}^t) \partial_{\mathbf{z}} \mathbf{f}_{\text{out}}(\varphi_{\text{out}^*}(\mathbf{z}), \boldsymbol{\omega}, \boldsymbol{\Sigma}^t), \\
\hat{\boldsymbol{\chi}}^t &= -\alpha \int_{\mathbb{R}^K} \int_{\mathbb{R}^K} d\boldsymbol{\omega} d\mathbf{z} \mathcal{N}_{(\mathbf{z}, \boldsymbol{\omega})}(\mathbf{0}_{2K}, \mathbf{Q}^t) \partial_{\boldsymbol{\omega}} \mathbf{f}_{\text{out}}(\varphi_{\text{out}^*}(\mathbf{z}), \boldsymbol{\omega}, \boldsymbol{\Sigma}^t).
\end{aligned} \tag{68}$$

with $\mathbf{Q}^t = \begin{bmatrix} \boldsymbol{\rho}_{\mathbf{w}^*} & \mathbf{m}^t \\ \mathbf{m}^t & \mathbf{q}^t \end{bmatrix} \in \mathbb{R}^{(2K) \times (2K)}$.

4.5.3 Summary of the SE - Bayes-optimal setting

In the Bayes-optimal setting, the student $\mathbf{P}_{\mathbf{w}} = \mathbf{P}_{\mathbf{w}^*}$ and $\mathbf{P}_{\text{out}} = \mathbf{P}_{\text{out}^*}$, so that we have $\mathbf{f}_{\mathbf{w}} = \mathbf{f}_{\mathbf{w}^*}$ and $\mathbf{f}_{\text{out}} = \mathbf{f}_{\text{out}^*}$. Moreover, the Nishimori conditions, recalled in Appendix. ??, imply that

$$\mathbf{m}^t = \mathbf{q}^t \equiv \mathbf{q}_{\mathbf{b}}^t, \quad \hat{\mathbf{q}}^t = \hat{\mathbf{m}}^t = \hat{\boldsymbol{\chi}}^t \equiv \hat{\mathbf{q}}_{\mathbf{b}}^t, \quad \boldsymbol{\Sigma}^t = \boldsymbol{\rho}_{\mathbf{w}^*} - \mathbf{q}^t.$$

Therefore the set of SE equations simplify and reduce to

$$\begin{aligned}
\mathbf{q}_{\mathbf{b}}^{t+1} &= \mathbb{E}_{\mathbf{w}^*, \boldsymbol{\xi}} \left[\mathbf{f}_{\mathbf{w}^*} \left(\hat{\mathbf{q}}_{\mathbf{b}}^t \mathbf{w}^* + (\hat{\mathbf{q}}_{\mathbf{b}}^t)^{1/2} \boldsymbol{\xi}, \hat{\mathbf{q}}_{\mathbf{b}}^t \right)^{\otimes 2} \right] \\
\hat{\mathbf{q}}_{\mathbf{b}}^t &= \alpha \int_{\mathbb{R}^K} \int_{\mathbb{R}^K} d\boldsymbol{\omega} d\mathbf{z} \mathcal{N}_{(\mathbf{z}, \boldsymbol{\omega})}(\mathbf{0}_{2K}, \mathbf{Q}_{\mathbf{b}}^t) \mathbf{f}_{\text{out}^*}(\varphi_{\text{out}^*}(\mathbf{z}), \boldsymbol{\omega}, \boldsymbol{\rho}_{\mathbf{w}^*} - \mathbf{q}_{\mathbf{b}}^t)^{\otimes 2}
\end{aligned} \tag{69}$$

with the simplified covariance matrix $\mathbf{Q}_{\mathbf{b}}^t = \begin{bmatrix} \boldsymbol{\rho}_{\mathbf{w}^*} & \mathbf{q}_{\mathbf{b}}^t \\ \mathbf{q}_{\mathbf{b}}^t & \mathbf{q}_{\mathbf{b}}^t \end{bmatrix}$.

5 Consistency between AMP and the replica computation

Very surprisingly, the SE of the AMP algorithm can be obtained in a convoluted and more rapid way. It turns out that in the Bayes-optimal setting, AMP performs a gradient ascent on the RS free entropy in (50). Meaning that at convergence, and under good initialization, the AMP overlaps are given by the saddle point equations of the RS free entropy $\Phi^{(\text{rs})}$. To see this, we shall start performing the change of variable $\boldsymbol{\xi} \leftarrow \boldsymbol{\xi} + (\hat{\mathbf{q}}_{\mathbf{b}}^t)^{1/2} \mathbf{w}^*$ in (69) so that we directly obtain the first equation of (58) with the corresponding time indices

$$\mathbf{q}_{\mathbf{b}}^{t+1} = \mathbb{E}_{\mathbf{w}^*, \boldsymbol{\xi}} \left[\mathcal{Z}_{\mathbf{w}^*} \left((\hat{\mathbf{q}}_{\mathbf{b}}^t)^{1/2} \boldsymbol{\xi}, \hat{\mathbf{q}}_{\mathbf{b}}^t \right) \mathbf{f}_{\mathbf{w}^*} \left((\hat{\mathbf{q}}_{\mathbf{b}}^t)^{1/2} \boldsymbol{\xi}, \hat{\mathbf{q}}_{\mathbf{b}}^t \right)^{\otimes 2} \right]. \tag{70}$$

Moreover in this setting, we notice that variables $\boldsymbol{\omega}_{\mu \rightarrow i}^t$ and $\mathbf{z}_\mu - \boldsymbol{\omega}_{\mu \rightarrow i}^t$ become independent since

$$\begin{aligned}\mathbb{E} \left[\boldsymbol{\omega}_{\mu \rightarrow i}^t (\mathbf{z}_\mu - \boldsymbol{\omega}_{\mu \rightarrow i}^t)^\top \right] &\xrightarrow{d \rightarrow \infty} \mathbf{m}^t - \mathbf{q}^t = \mathbf{q}_b^t - \mathbf{q}_b^t = \mathbf{0}, \\ \mathbb{E} \left[\boldsymbol{\omega}_{\mu \rightarrow i}^t (\boldsymbol{\omega}_{\mu \rightarrow i}^t)^\top \right] &\xrightarrow{d \rightarrow \infty} \mathbf{q}_b^t, \\ \mathbb{E} \left[(\mathbf{z}_\mu - \boldsymbol{\omega}_{\mu \rightarrow i}^t) (\mathbf{z}_\mu - \boldsymbol{\omega}_{\mu \rightarrow i}^t)^\top \right] &\xrightarrow{d \rightarrow \infty} \boldsymbol{\rho}_{w^*} - \mathbf{q}_b^t,\end{aligned}$$

so that the multivariate Gaussian distribution factorize to

$$\mathcal{N}_{(\mathbf{z}, \boldsymbol{\omega})}(\mathbf{0}, \mathbf{Q}_b^t) = \mathcal{N}_{\boldsymbol{\omega}}(\mathbf{0}_K, \mathbf{q}_b^t) \mathcal{N}_{\mathbf{z}}(\boldsymbol{\omega}, \boldsymbol{\rho}_{w^*} - \mathbf{q}_b^t).$$

Using $\mathbf{P}_{\text{out}^*}(y|\mathbf{z}) = \delta(y - \varphi_{\text{out}^*}(\mathbf{z}))$ the second equation of (69) becomes

$$\begin{aligned}\hat{\mathbf{q}}^t &= \alpha \int_{\mathbb{R}^K} \int_{\mathbb{R}^K} d\boldsymbol{\omega} d\mathbf{z} \mathcal{N}_{(\mathbf{z}, \boldsymbol{\omega})}(\mathbf{0}_{2K}, \mathbf{Q}_b^t) \mathbf{f}_{\text{out}^*}(\varphi_{\text{out}^*}(\mathbf{z}), \boldsymbol{\omega}, \boldsymbol{\rho}_{w^*} - \mathbf{q}_b^t)^{\otimes 2} \\ &= \alpha \int_{\mathbb{R}} dy \int_{\mathbb{R}^K} d\boldsymbol{\omega} \mathcal{N}_{\boldsymbol{\omega}}(\mathbf{0}_K, \mathbf{q}_b^t) \\ &\quad \times \int_{\mathbb{R}^K} d\mathbf{z} \mathbf{P}_{\text{out}^*}(y|\mathbf{z}) \mathcal{N}_{\mathbf{z}}(\boldsymbol{\omega}; \boldsymbol{\rho}_{w^*} - \mathbf{q}_b^t) \mathbf{f}_{\text{out}^*}(y, \boldsymbol{\omega}, \boldsymbol{\rho}_{w^*} - \mathbf{q}_b^t)^{\otimes 2} \\ &= \alpha \int_{\mathbb{R}} dy \int_{\mathbb{R}^K} d\boldsymbol{\xi} \mathcal{N}_{\boldsymbol{\xi}}(\mathbf{0}; \mathbf{I}_K) \int_{\mathbb{R}^K} d\mathbf{z} \mathbf{P}_{\text{out}^*}(y|\mathbf{z}) \\ &\quad \times \mathcal{N}_{\mathbf{z}}\left((\mathbf{q}_b^t)^{1/2} \boldsymbol{\xi}; \boldsymbol{\rho}_{w^*} - \mathbf{q}_b^t\right) \mathbf{f}_{\text{out}^*}(y, (\mathbf{q}_b^t)^{1/2} \boldsymbol{\xi}, \boldsymbol{\rho}_{w^*} - \mathbf{q}_b^t)^{\otimes 2} \\ &\quad \text{(Change of variable } \boldsymbol{\xi} \leftarrow (\mathbf{q}_b^t)^{-1/2} \boldsymbol{\omega}^t) \\ &= \alpha \int_{\mathbb{R}} dy \mathbb{E}_{\boldsymbol{\xi}} \mathcal{Z}_{\text{out}^*}\left(y, (\mathbf{q}_b^t)^{1/2} \boldsymbol{\xi}, \boldsymbol{\rho}_{w^*} - \mathbf{q}_b^t\right) \\ &\quad \times \mathbf{f}_{\text{out}^*}\left(y, (\mathbf{q}_b^t)^{1/2} \boldsymbol{\xi}, \boldsymbol{\rho}_{w^*} - \mathbf{q}_b^t\right),\end{aligned}$$

which is exactly the second fixed point equation of the RS free entropy (58).

References

- [1] Marc Mézard, Giorgio Parisi, and Miguel Virasoro. *Spin glass theory and beyond: An Introduction to the Replica Method and Its Applications*, volume 9. World Scientific Publishing Company, 1987.
- [2] Tommaso Castellani and Andrea Cavagna. Spin-glass theory for pedestrians. *Journal of Statistical Mechanics: Theory and Experiment*, (5):215–266, 2005.
- [3] Jean Barbier, Florent Krzakala, Nicolas Macris, Léo Miolane, and Lenka Zdeborová. Phase Transitions, Optimal Errors and Optimality of Message-Passing in Generalized Linear Models. pages 1–59, 2017.
- [4] JS Yedidia and WT Freeman. Understanding belief propagation and its generalizations. *Exploring artificial intelligence in the new millennium*, pages 239 – 269, 2001.
- [5] Florent Krzakala, Marc Mézard, Francois Sausset, Yifan Sun, and Lenka Zdeborová. Probabilistic reconstruction in compressed sensing: Algorithms, phase diagrams, and threshold achieving matrices. *Journal of Statistical Mechanics: Theory and Experiment*, 2012(8):1–42, 2012.
- [6] Lenka Zdeborová and Florent Krzakala. Statistical physics of inference: thresholds and algorithms. *Advances in Physics*, 65(5):453–552, 2016.
- [7] Thibault Lesieur, Florent Krzakala, and Lenka Zdeborová. Constrained low-rank matrix estimation: Phase transitions, approximate message passing and applications. *Journal of Statistical Mechanics: Theory and Experiment*, 2017(7), 2017.
- [8] S Patarnello and P Carnevali. Learning networks of neurons with boolean logic. *Europhysics Letters (EPL)*, 4(4):503–508, aug 1987.
- [9] Elizabeth Gardner and Bernard Derrida. Three unfinished works on the optimal storage capacity of networks. *Journal of Physics A: Mathematical and General*, 22(12):1983, 1989.
- [10] Naftali Tishby, Esther Levin, and Sara A. Solla. Consistent inference of probabilities in layered networks: Predictions and generalization. In Anon, editor, *IJCNN Int Jt Conf Neural Network*, pages 403–409. Publ by IEEE, December 1989. IJCNN International Joint Conference on Neural Networks ; Conference date: 18-06-1989 Through 22-06-1989.
- [11] H. Sompolinsky, N. Tishby, and H. S. Seung. Learning from examples in large neural networks. *Physical Review Letters*, 65(13):1683–1686, 1990.
- [12] Sebastian Seung, Haim Sompolinsky, and Naftali Tishby. Statistical mechanics of learning from examples. *Physical Review A*, 45(8):6056, 1992.

- [13] Timothy LH Watkin, Albrecht Rau, and Michael Biehl. The statistical mechanics of learning a rule. *Reviews of Modern Physics*, 65(2):499, 1993.
- [14] G. Györgyi. Techniques of replica symmetry breaking and the storage problem of the McCulloch-Pitts neuron. *Physics Report*, 342(4-5):263–392, 2001.
- [15] Yoshiyuki Kabashima. Inference from correlated patterns: a unified theory for perceptron learning and linear vector channels. In *Journal of Physics: Conference Series*, volume 95, page 012001. IOP Publishing, 2008.
- [16] Elizabeth Gardner and Bernard Derrida. Optimal storage properties of neural network models. *Journal of Physics A: Mathematical and general*, 21(1):271, 1988.
- [17] Marc Mézard. The space of interactions in neural networks: Gardner’s computation with the cavity method. *Journal of Physics A: Mathematical and General*, 22(12):2181, 1989.
- [18] Henry Schwarze and John Hertz. Generalization in a large committee machine. *EPL (Europhysics Letters)*, 20(4):375, 1992.
- [19] Henry Schwarze and John Hertz. Generalization in fully connected committee machines. *EPL (Europhysics Letters)*, 21(7):785, 1993.
- [20] Henry Schwarze. Learning a rule in a multilayer neural network. *Journal of Physics A: Mathematical and General*, 26(21):5781, 1993.
- [21] Rémi Monasson and Riccardo Zecchina. Learning and generalization theories of large committee-machines. *Modern Physics Letters B*, 9(30):1887–1897, 1995.
- [22] Florent Krzakala, Marc Mézard, Francois Sausset, Yifan Sun, and Lenka Zdeborová. Probabilistic reconstruction in compressed sensing: algorithms, phase diagrams, and threshold achieving matrices. *Journal of Statistical Mechanics: Theory and Experiment*, 2012(08):P08009, 2012.
- [23] Jean Barbier, Florent Krzakala, Nicolas Macris, Léo Miolane, and Lenka Zdeborová. Optimal errors and phase transitions in high-dimensional generalized linear models. *Proceedings of the National Academy of Sciences*, 116(12):5451–5460, 2019.
- [24] R. Wong. *Asymptotic Approximations of Integrals Computer Science and Scientific Computing*. 1989.
- [25] Benjamin Aubin, Florent Krzakala, Yue M Lu, and Lenka Zdeborová. Generalization error in high-dimensional perceptrons: Approaching bayes error with convex optimization. In *Advances in Neural Information Processing Systems* 33. 2020.

- [26] Marc Mézard and Andrea Montanari. *Information, physics, and computation*. Oxford University Press, 2009.
- [27] Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- [28] David J Thouless, Philip W Anderson, and Robert G Palmer. Solution of ‘solvable model of a spin glass’. *Philosophical Magazine*, 35(3):593–601, 1977.
- [29] Mohsen Bayati and Andrea Montanari. The dynamics of message passing on dense graphs, with applications to compressed sensing. *IEEE Transactions on Information Theory*, 57(2):764–785, 2011.
- [30] Adel Javanmard and Andrea Montanari. State evolution for general approximate message passing algorithms, with applications to spatial coupling. *Information and Inference: A Journal of the IMA*, 2(2):115–144, 2013.
- [31] Mohsen Bayati, Marc Lelarge, Andrea Montanari, et al. Universality in polytope phase transitions and message passing algorithms. *The Annals of Applied Probability*, 25(2):753–822, 2015.
- [32] Sundeep Rangan. Generalized approximate message passing for estimation with random linear mixing. In *Information Theory Proceedings (ISIT), 2011 IEEE International Symposium on*, pages 2168–2172. IEEE, 2011.
- [33] Lenka Zdeborová and Florent Krzakala. Statistical physics of inference: Thresholds and algorithms. *Advances in Physics*, 65(5):453–552, 2016.