

Uncertainty quantification with deep ensembles : a toy model

clarte.lucas

July 2024

1 Introduction

1.1 Related work

Bootstrap resampling is a cornerstone of classical statistics [citation needed], and its application to modern machine learning techniques has been more recently explored [?].

2 Setting

Consider training data $\mathcal{D} = (\mathbf{x}_i, y_i)_{i=1}^n$ where the input follows the distribution $\mathbf{x} \sim \mathcal{N}(0, I_d)$ and the labels are generated by a teacher vector \mathbf{w}_* , such that $y = f(\mathbf{w}_*^\top \mathbf{x}_i)$ where f is a function that can be random. For instance, in the case of linear regression, one can use $f(z) = z + \sqrt{\Delta}\varepsilon, \varepsilon \sim \mathcal{N}(0, 1)$. For binary classification, one can use $p(y = 1 | \mathbf{w}_*^\top \mathbf{x})$. The function f and the teacher \mathbf{w}_* will typically be unknown to the statistician, and the goal is to recover \mathbf{w}_* using limited data. Putting the context of overparametrized models, we consider now an ensemble of random features $\mathbf{v}_1, \dots, \mathbf{v}_K$ generated by

$$\mathbf{v}_i^j = \sigma_i(F^j \mathbf{x}_i), \quad F^j \in \mathbb{R}^{p \times d} \quad (1)$$

Where the components of F_i are i.i.d $\mathcal{N}(0, 1/d)$. We will define $\mathcal{D}^j = (\mathbf{v}_i^j, y_i)_{i=1}^n$.

Empirical risk minimizer We will first consider frequentist estimators that minimize the Ridge loss on the random features \mathbf{v} :

$$\hat{\mathbf{w}}^j = \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w} | \mathcal{D}^j) = \sum_i \frac{|y_i - \mathbf{w}^\top \mathbf{v}_i^j|^2}{2} + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (2)$$

Bayes-optimal estimator Given data $\mathcal{D} = (\mathbf{v}_i, y_i)_{i=1}^n$, the Bayes-optimal estimator computes the posterior distribution

$$p(y | \mathbf{v}, \mathcal{D}) = \int d\mathbf{w} p(\mathbf{w} | \mathcal{D}) p(y | \mathbf{w}) \quad (3)$$

By definition, the Bayes-optimal estimator is the one with the best accuracy, given \mathcal{D} .

Pseudo-Bayes Computing the Bayes-optimal estimator eq. (3) is practically impossible as it requires the knowledge of the prior and likelihood function of the data-generating process, which is inaccessible in practice. Then, an alternative is the *pseudo-Bayes estimator* that samples the following distribution associated to the loss.

$$p(\mathbf{w} | \mathcal{D}^j) \propto e^{-\beta \mathcal{L}(\mathbf{w})} \quad (4)$$

where β is the inverse temperature that controls the variance of the distribution. In the context of Bayesian neural network, $\beta > 1$ (respectively $\beta < 1$) correspond to *cold posteriors* (respectively *hot posteriors*).

3 Uncertainty quantification with ensembles

Resampling overparametrized models Resampling techniques such as bootstrap or subsampling are a classical tools in classical statistics to estimate the variance of an estimator and quantify uncertainty. However, in the high-dimensional regime, applying these techniques is not straightforward [citation needed]. One way to model resampling is to reframe it as a weighted empirical risk minimization :

$$\hat{\mathbf{w}}^j = \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w} | \mathcal{D}^j, \mathbf{p}^j) = \sum_i p_i^j \frac{|y_i - \mathbf{w}^\top \mathbf{v}_i^j|^2}{2} + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (5)$$

where the weights $p_i^j \in \mathbb{N}$ follow a multinomial distribution for the Bootstrap, or a Bernoulli distribution for subsampling. Previous work [citation needed] has analyzed the performance of bootstrap and resampling for linear regression, but not in the context of overparametrized models.

3.1 Regression

In the context of regression, when working with an ensemble of estimators (2) $\hat{\mathbf{w}}^1, \dots, \hat{\mathbf{w}}^K$ obtained either by resampling, using different random features or sampling a Bayesian estimator, the standard way to quantify uncertainty is through the prediction variance across the ensemble. For a single sample \mathbf{x} with associated random features $\mathbf{v}^1, \dots, \mathbf{v}^K$ (potentially equal), the prediction variance for \mathbf{x} is

$$\hat{\sigma}^2(\mathbf{x}) := \text{Var}_{\mathbf{w}^j} \left[\hat{\mathbf{w}}^{j\top} \mathbf{v}^j \right] \quad (6)$$

Rescaling in classification For classification tasks, a popular way to improve uncertainty quantification is to rescale the confidence function without changing the predicted class. A popular way is *temperature scaling* that uses the following confidence :

$$f(\mathbf{v}) = \sigma(\mathbf{w}_*^\top \mathbf{v} / T) \quad (7)$$

where T is a scalar factor learned by minimizing the loss on some validation data. Plenty of alternatives to temperature scaling have been developed [citation needed], but all rely on rescaling the confidence function by a few factors.

Rescaling the variances in regression Similar to temperature scaling for classification tasks [citation needed], one can rescale the ensemble variance of the prediction. If one uses a likelihood $p(y|\mathbf{x})$ as a Gaussian with mean $\hat{\mu}(\mathbf{x}) = \dots$ and variance $\hat{\sigma}(\mathbf{x}) = \text{Var}(\mathbf{x})$ as defined in (6), we can rescale the prediction variance by a factor T

$$\hat{\sigma}^2(\mathbf{x}) \rightarrow \hat{\sigma}^2(\mathbf{x}) / T \quad (8)$$

A natural way to compute the factor T is such that the average prediction variance is equal to the mean square error on the validation set.

3.2 Classification

In the context of classification,

4 Technical results

Our main technical result is the computation of the correlation $\hat{\mathbf{w}}^{1\top} \hat{\mathbf{w}}^2$ for $\hat{\mathbf{w}}^i$ either computed by empirical risk minimization or defined as the average of a Bayesian posterior, trained on the bootstrap resample or subsampled version of \mathcal{D} . This result is an extension of [citation needed].

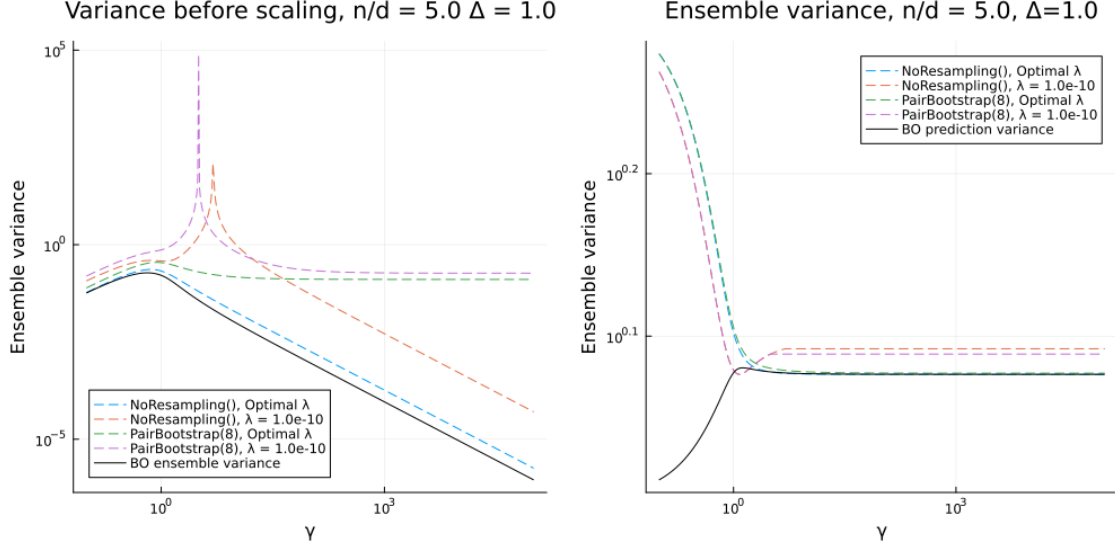


Figure 1: Ensemble variance with and without resampling for random features, at finite overparametrization. Left : before rescaling the variance. Right : after rescaling the variance.

Theorem 4.1. Consider \hat{w}^1, \hat{w}^2 that minimize the loss defined in eq. (5). Then, the variance is given by ?? where the overlaps m, Q that are solutions of the following equations

$$\begin{cases} m &= \dots \\ Q &= \dots \\ V &= \dots \end{cases}, \begin{cases} \hat{m} &= \dots \\ \hat{Q} &= \dots \\ \hat{V} &= \dots \end{cases}$$

theorem 4.1 allows us to compute efficiently the correlation between two estimators in high-dimensions benchmark them in the random features model. For the sake of concreteness, this theorem is formulated for Ridge regression and bootstrap resamples, but we refer to ?? for a more general statement.

Theorem 4.2.

4.1 Numerics

Ensemble variance In Figure 1, we plot the variance defined in eq. (6). We compare there the ensemble variance of pair bootstrap and ERM without resampling, at optimal regularization and at almost zero regularization ($\lambda = 1e-10$), before (Left) and after rescaling the variance (Right) to match the generalisation error. We see that with rescaling, it makes almost no difference to use bootstrap or not : the determining factor is λ : at optimal λ , the ensemble variance matches the predictive variance of the Bayes-optimal, for the sole reason that at optimal λ , Ridge has almost optimal test error. In other words, good generalisation error induces accurate uncertainty estimation.

Which Bayesian variance are we comparing to ? In the Left panel of fig. 1, we plot the ensemble variance for Bayes posteriors trained on different resamples of F . This ensemble variance is equal to $q_0 - q_1$ using the overlaps from state-evolution. In the Right panel, we plot the variance of the Bayes posterior (which corresponds to the overlap v in the state-evolution equation) for one single random-features model.

Summary It looks like doing resampling does not add much to UQ in this model : using optimal regularization and using post-hoc rescaling of the variance is enough to have good UQ, in the sense that you approximate correctly some notion of variance of the BO estimator. In one case (before the rescaling), the ensemble variance of the ERM estimator is approximately the one of the BO. After rescaling, you approximate the posterior variance.

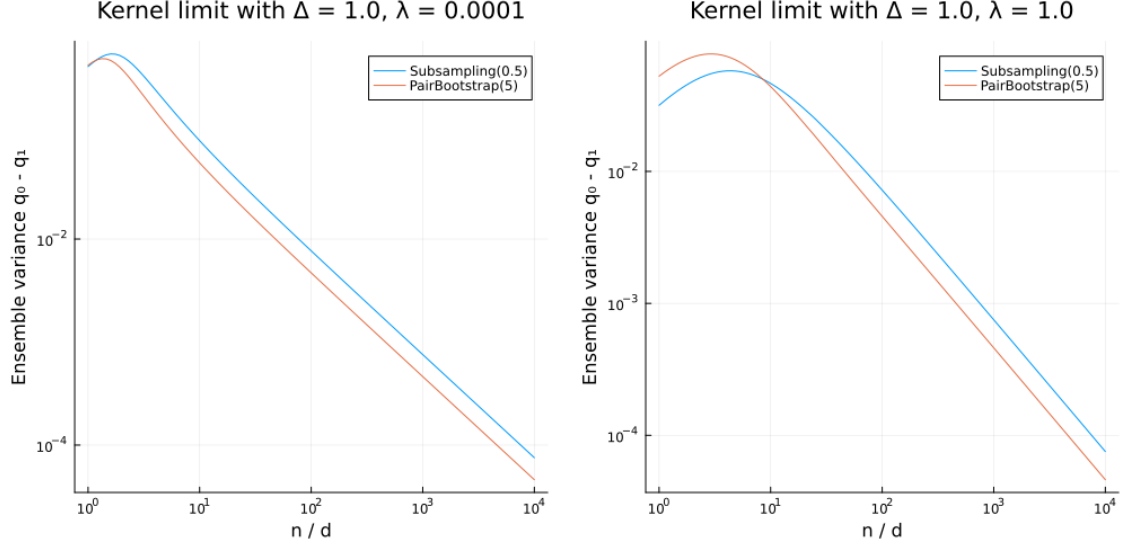


Figure 2: Ensemble variance for subsampling and bootstrap in the kernel limit

Kernel limit At infinite overparametrization, the random features all converge to the same kernel limit independently of the actual realization of the matrix F . For this reason, with resampling the ensemble method converges to 0. On the other hand, in fig. 2, we plot the ensemble variance using bootstrap and subsampling, as a function of the sampling ratio n/d . We observe that in both cases the variance scales as d/n^1

¹This is coherent with what we observed in the bootstrap paper where the variance scales as d/n

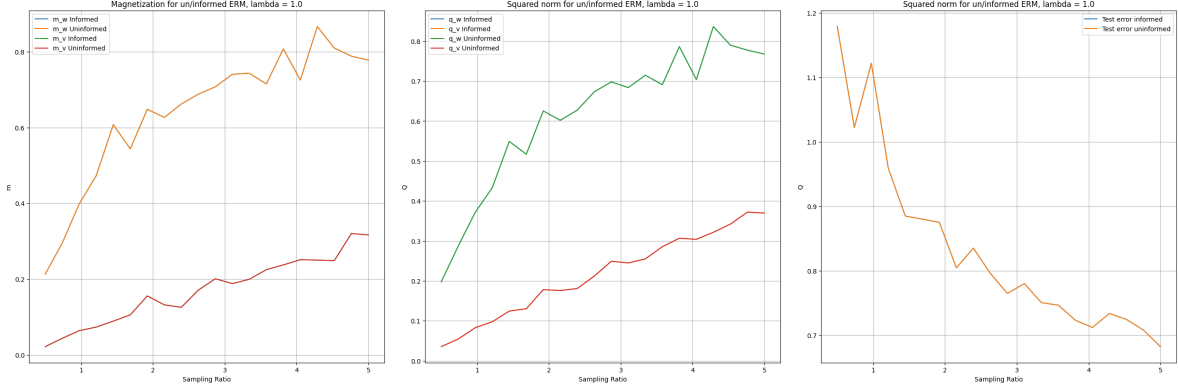


Figure 3: From left to right : plot of m , q and the test error of estimator $\hat{\mathbf{w}}, \hat{\mathbf{v}}$ obtained by minizing (10) for $\lambda = 1$, which corresponds to the MAP estimator.

5 Model for mean-variance estimation

5.1 Setting

The data will be generated by the model

$$y = \mathbf{w}_\star^\top \mathbf{x} + \varepsilon, \varepsilon \sim \mathcal{N}(0, \sigma(\mathbf{v}_\star^\top \mathbf{x})) \quad (9)$$

Where once again the input \mathbf{x} comes from the standard normal distribution. The teacher vectors $\mathbf{w}_\star, \mathbf{v}_\star$ are also sampled from the distribution $\mathcal{N}(\mathbf{0}, I_d)$. We will use the same non-linearity σ to parametrize the noise variance as the original deep ensemble paper, namely the softplus function²

$$\sigma(z) = \log(1 + e^z)$$

We will infer the vectors \mathbf{w}, \mathbf{v} with the empirical risk

$$\mathcal{L}(\mathbf{w}, \mathbf{v}) = \sum_{i=1}^n \frac{(\mathbf{w}^\top \mathbf{x}_i - y_i)^2}{2\sigma(\mathbf{v}^\top \mathbf{x}_i)} + \frac{\lambda}{2} (\|\mathbf{w}\|^2 + \|\mathbf{v}\|^2) \quad (10)$$

The loss \mathcal{L} corresponds to the log-posterior when $\lambda = 1$. The goal is to characterize the performance of the empirical risk minimizer of \mathcal{L} in high-dimensions. First thing to note is that the loss (10) is non-convex, meaning that gradient descent may be stuck in local minima.

5.2 Numerics

Uninformed and informed initialization To test whether the non-convexity of the training loss has an impact on the estimation, we run gradient descent on \mathcal{L} with an *informed* initialization, starting from the teacher $\mathbf{w}_\star, \mathbf{v}_\star$, and an *uninformed* initialization starting from random vectors. In Figure 3, we plot the overlaps and the test error for $\lambda = 1$ for both initializations. For this choice of λ , both estimators coincide. Same for Figure 4 where we take $\lambda = 0.1$.

5.3 State-evolution

We have the loss

$$\ell(y, \mathbf{z}) = \frac{(y - z_1)^2}{2\sigma(z_2)} \quad (11)$$

²Note : one could use other non-linearities, such as the square function.

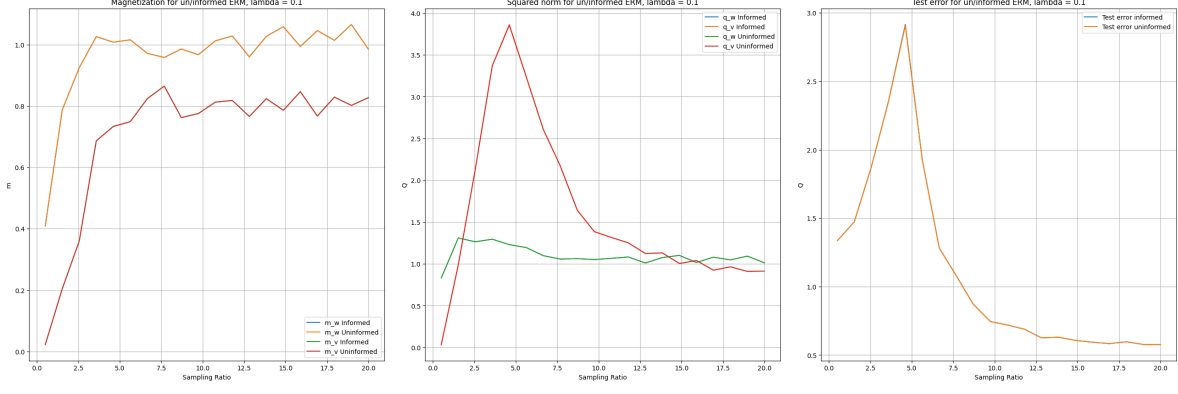


Figure 4: From left to right : plot of m , q and the test error of estimator \hat{w}, \hat{v} obtained by minizing (10) for $\lambda = 0.1$.

for which we can write the proximal operator (and the corresponding channel function) :

$$\text{prox}(y, \omega, V) = \arg \min_z \ell(y, z) + V^{-1}(z - \omega)/2 \quad (12)$$

$$g(y, \omega, V) = V^{-1}(\text{prox}(y, \omega, V) - \omega) \quad (13)$$

Note that because the contributions of z_1 and z_2 in the loss are not separable, the covariance matrix V is *a priori* not diagonal.

The denoising step is similar as for Bootstrap [reference needed], since the prior is factorized.

6 To put in the appendix : Gaussian Covariate Model

Recall that by Gaussian universality, this model is equivalent in the high-dimensional limit to

$$\mathbf{v}_i = \kappa_i^1 F_i \mathbf{x} + \kappa_i^* z_i \quad (14)$$

where z_i follows the normal distribution. In what follows, we consider the case $K = 2$. Then, the vector $(\mathbf{x}, \mathbf{v}_1, \mathbf{v}_2)$ follows a zero mean Gaussian with covariance :

$$\begin{pmatrix} I_d & \Phi_1 & \Phi_2 \\ \Phi_1 & \Omega_1 & \Phi_{1,2} \\ \Phi_2 & \Phi_{1,2} & \Omega_2 \end{pmatrix}, \quad \begin{cases} \Phi_i &= \kappa_i^1 F_i \\ \Omega_i &= \kappa_i^1 F_i F_i^T + \kappa_i^{*2} \\ \Phi_{1,2} &= \kappa_1^1 \times \kappa_2^1 F_1 F_2^T \end{cases} \quad (15)$$

The random features define a *projected teacher* for each random feature :

$$\mathbf{w}_{*i} = \Omega_i^{-1} \Phi_i^T \boldsymbol{\theta}_* \sim \mathcal{N}(0, \Sigma_i^*), \quad \Sigma_i^* = \Omega_i^{-1} \Phi_i^T \Phi_i \Omega_i^{-1} \quad (16)$$

The label as a function of $\mathbf{w}_{*i}^\top \mathbf{v}_i$ is

$$y = \mathbf{w}_{*i}^\top \mathbf{v}_i / \sqrt{d} + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \Delta + \tau_{add}) \quad (17)$$

with the additional the noise variance $\tau_{add} = 1 - \mathbb{E}_{x \sim MP} \left[\frac{(\kappa^1)^2 x}{(\kappa^1)^2 x + (\kappa^*)^2} \right]$. Here, the MP distribution is the limiting distribution of the spectrum of FF^T because the entries are i.i.d Gaussian.

Bayes-optimal estimator in the GCM To compute the Bayes-optimal estimator, we sample from the distribution

$$p(\mathbf{w}_i | \mathcal{D}) \propto \prod_{\mu=1}^n p(y_\mu | \mathbf{w}^\top \mathbf{v}_{i\mu}) \prod \mathcal{N}(\mathbf{w} | \mathbf{0}, \Sigma_*) \quad (18)$$

In the case of a Gaussian likelihood and Gaussian prior, the Bayes-optimal estimator $\mathbf{w}_{BO,i}$ is just the mean of this distribution. The performance of a single Bayes-optimal estimator has been characterized in our AISTATS paper. The goal here is to compute the correlation between two estimators trained on different random features, coming from two matrices F_1 and F_2 . We thus define the overlap

$$q_1 := \frac{1}{d} \mathbb{E} [\mathbf{w}_1^\top \mathbf{v}_1 \times \mathbf{w}_2^\top \mathbf{v}_2] = \frac{1}{d} \mathbf{w}_1^\top \text{Cov}(\mathbf{v}_1, \mathbf{v}_2) \mathbf{w}_2 = \frac{1}{d} \mathbf{w}_1^\top \Phi_{1,2} \mathbf{w}_2 \quad (19)$$

The computation of this correlation has been done in Loureiro et al. for Ridge estimators. They derive the expression (for Ridge estimator $l = (y - z)^2 / 2\hat{\Delta}$) :

$$\begin{cases} \hat{q}_1 &= \frac{\rho + q_1 - 2m + \Delta + \tau_{add}}{\hat{\Delta} + v} \\ q_1 &= \left(1 + \frac{\hat{q}_1}{\hat{m}^2}\right) \times m^2 \end{cases} \quad (20)$$

where the overlaps $m, q, v, \hat{m}, \hat{q}, \hat{v}$ are the usual solutions of the fixed point equations. For the Bayes-optimal estimators, q_1 also solves (20). In 7, we also derive the expression of q_1 from the state-evolution equations of AMP.

7 State evolution for q_1 from AMP

An option is to write an AMP algorithm that estimates the two estimators simultaneously, with the same channel function

$$\dots \quad (21)$$

and the denoisers related to their own prior $\mathcal{N}(0, \Sigma_{\star i})$ where (cf. the AISTATS paper)

$$\Sigma_{\star, i} = \Omega_i^{-1/2} \Phi_i^T \Phi_i \Omega_i^{1/2} \quad (22)$$

which is equivalent of defining a single prior on $(\mathbf{w}_1, \mathbf{w}_2) \in \mathbb{R}^{2 \times p}$ with a factorized form

$$\mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \Sigma_{\star, 1} & 0 \\ 0 & \Sigma_{\star, 2} \end{pmatrix}\right)$$

the denoiser for a prior $\mathcal{N}(0, \Sigma_{\star})$ is

$$f_a(\mathbf{b}, \mathbf{A}) = (\mathbf{A} + \Sigma_{\star}^{-1})^{-1} \mathbf{b} \in \mathbb{R}^p, \quad f_v(\mathbf{b}, \mathbf{A}) = (\mathbf{A} + \Sigma_{\star}^{-1})^{-1} \in \mathbb{R}^p \quad (23)$$

with a diagonal matrix $\mathbf{A} = \text{Diag}(\mathbf{A})$.

Algorithm 1 GAMP for Bayes-optimal on random features \mathbf{v}_{μ}^j for $1 \leq j \leq B$ different random feature matrices F_j .

Input: $(\mathbf{v}_{\mu}^j)_{\mu=1}^n \in \mathbb{R}^{n \times p}$ ($1 \leq j \leq B$), $\mathbf{y} \in \mathbb{R}^n$

Initialize: $\mathbf{g}_{\text{out}\mu}^{(0)} = \mathbf{0}$ for $1 \leq \mu \leq n$, $\mathbf{A}_i^{(0)} = \mathbf{I}_B$ for $1 \leq i \leq d$

Initialize: $\hat{\boldsymbol{\theta}}_i^{(1)} \in \mathbb{R}^B$ and $\hat{\mathbf{C}}_i^{(1)} \in \mathbb{R}^{B \times B}$ for $1 \leq i \leq d$

Repeat for $t = 1, 2, \dots$:

// Update of the means $\boldsymbol{\omega}_{\mu}^j \in \mathbb{R}$ and covariances $\mathbf{V}_{\mu} \in \mathbb{R}^{B \times B}$ for $1 \leq \mu \leq n$:

$$\boldsymbol{\omega}_{\mu}^{j(t)} = \sum_{i=1}^d v_{\mu,i}^j \mathbf{w}_i^{j(t)} - \mathbf{v}_{\mu,i}^{j2} \hat{\mathbf{C}}_i^{j(t)} \mathbf{g}_{\text{out}\mu}^{j(t-1)} \mid V_{\mu}^{j,k(t)} = \sum_{i=1}^d \mathbf{v}_{\mu,i}^j \mathbf{v}_{\mu,i}^k \hat{\mathbf{C}}_i^{j,k(t)} \text{ for } 1 \leq j, k \leq B$$

// Update of $\mathbf{g}_{\text{out}\mu} \in \mathbb{R}^B$ and $\partial_{\omega} \mathbf{g}_{\text{out}\mu} \in \mathbb{R}^{B \times B}$ for $1 \leq \mu \leq n$:

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

// Update of means $\mathbf{b}_i \in \mathbb{R}^B$ and covariances $\mathbf{A}_i \in \mathbb{R}^{B \times B}$ for $1 \leq i \leq d$:

$$\mathbf{A}_i^{j,k(t)} = -\sum_{\mu=1}^n \mathbf{v}_{\mu,i}^j \mathbf{v}_{\mu,i}^k \partial_{\omega} \mathbf{g}_{\text{out}\mu}^{j,k(t)} \mid \mathbf{b}_i^{j(t)} = \mathbf{A}_i^{jj(t)} \hat{\boldsymbol{\theta}}_i^{j(t)} + \sum_{\mu=1}^n \mathbf{v}_{\mu,i}^j \mathbf{g}_{\text{out}\mu}^{jj(t)}$$

// Update of the estimated marginals $\hat{\boldsymbol{\theta}}_i \in \mathbb{R}^B$ and $\hat{\mathbf{C}}_i \in \mathbb{R}^{B \times B}$ for $1 \leq i \leq d$:

$$\hat{\boldsymbol{\theta}}_i^{j(t+1)} = \mathbf{f}_a(\mathbf{b}_i^{j(t)}, \mathbf{A}_i^{jj(t)}) \mid \hat{\mathbf{C}}_i^{j(t+1)} = \partial_{\mathbf{b}} \mathbf{f}_a(\mathbf{b}_i^{j(t)}, \mathbf{A}_i^{jj(t)})$$

Until convergence

Output: $\hat{\boldsymbol{\theta}}_1, \dots, \hat{\boldsymbol{\theta}}_d$ and $\hat{\mathbf{C}}_1, \dots, \hat{\mathbf{C}}_d$

Note on 1 : we know that the matrix $\mathbf{V} \in \mathbb{R}^{2 \times 2}$ is diagonal because the distribution that we are sampling from is factorized between \mathbf{w}^1 and \mathbf{w}^2 (so conditionned on the data, the two estimators are independent).

7.1 State-evolution equations

To compute the state-evolution, we consider the relaxed-BP which is equivalent to AMP in high-dimensions. We define :

$$\begin{cases} \mathbf{b}_{\mu \rightarrow i}^{j(t)} &= \sum_{\nu \neq \mu} \mathbf{v}_{\nu,i}^j \mathbf{g}_{\text{out}}(\varphi_0(z_{\nu}), \boldsymbol{\omega}_{\nu \rightarrow i}^{j(t)}, \mathbf{V}_{\nu \rightarrow i}^{j(t)}) \\ \mathbf{A}_{\mu \rightarrow i}^{j(t)} &= -\sum_{\nu \neq \mu} \mathbf{v}_{\nu,i}^{j2} \partial_{\omega} \mathbf{g}_{\text{out}}(\varphi_0(z_{\nu}), \boldsymbol{\omega}_{\nu \rightarrow i}^{j(t)}, \mathbf{V}_{\nu \rightarrow i}^{j(t)}) \end{cases}$$

Concentration of \mathbf{V}_{μ}^j $\mathbb{E}[\mathbf{V}_{\mu}^j] = \sum_{k \neq i} \mathbb{E}[\mathbf{v}_k^{j2}] \hat{\mathbf{c}}_k^j \rightarrow \frac{1}{p} \sum_{p=1} \hat{\mathbf{c}}^j$ as the expectation of \mathbf{v}_k^{j2} is simply $\frac{1}{d}$.

Limitation of distribution of \mathbf{b} For $j = 1, 2$, we have asymptotically

$$\mathbb{E} \left[b_{\mu \rightarrow i}^j \right] \rightarrow \mathbf{w}_{\star i}^j \hat{m}_j \quad (24)$$

Where

$$\hat{m}_j = n/p \mathbb{E}_{z, \omega} \left[\partial_z g_{out}^j \right] \quad (25)$$

Note that if the random matrices are sampled from the same distribution, then the channel functions are the same i.e $\hat{m}^1 = \hat{m}^2$.

The covariance is

$$\mathbb{E} \left[b_{\mu \rightarrow i}^j b_{\mu \rightarrow i}^k \right] \rightarrow \Phi_{j, k, i, i} \times n/p \mathbb{E} \left[g_{out}^j \times g_{out}^k \right] \quad (26)$$

Note that the correlation is indexed by i because of the non-i.i.d covariance matrix between the two student samples, but the matrix $\hat{\mathbf{Q}} = \frac{n}{p} \mathbb{E} \left[g_{out}^1 \times g_{out}^2 \right]$ is independent from the index.

In the end, for $1 \leq i \leq d$, the vector $\mathbf{b}_i \in \mathbb{R}^2$ has the distribution

$$\mathbf{b}_i = \mathbf{w}_{\star i} \otimes \hat{\mathbf{m}} + \boldsymbol{\xi}_i, \boldsymbol{\xi}_i \sim \mathcal{N} \left(\mathbf{0}, \begin{pmatrix} \hat{\mathbf{Q}}_{1,1} & \Phi_{i,i}^{1,2} \times \hat{\mathbf{Q}}_{1,2} \\ \Phi_{i,i}^{1,2} \times \hat{\mathbf{Q}}_{1,2} & \hat{\mathbf{Q}}_{2,2} \end{pmatrix} \right) \quad (27)$$

Computation of the overlaps Consider the magnetization $\mathbf{m} \in \mathbb{R}^2$. Asymptotically, we have that³

$$\mathbf{m} \in \mathbb{R}^2 = \frac{1}{p} \mathbb{E}_{\mathbf{w}_{\star}, \boldsymbol{\xi}} \left[\begin{pmatrix} \mathbf{f}_a \left(\hat{m}_1 \mathbf{w}_{\star}^1 + (\boldsymbol{\xi})_1, \mathbf{A}^1 \right)^\top \mathbf{w}_{\star}^1 \\ \mathbf{f}_a \left(\hat{m}_2 \mathbf{w}_{\star}^2 + (\boldsymbol{\xi})_2, \mathbf{A}^2 \right)^\top \mathbf{w}_{\star}^2 \end{pmatrix} \right] \quad (28)$$

The expression of e.g. the first coordinate of the expectation is

$$\mathbb{E} \left[\mathbf{f}_a \left(\hat{m}_1 \mathbf{w}_{\star}^1 + (\boldsymbol{\xi})_1, \hat{\mathbf{V}}_1 \right)^\top \mathbf{w}_{\star}^1 \right] = \mathbb{E} \left[\mathbf{w}_{\star}^{1\top} (\mathbf{A} + \Sigma_{\star,1}^{-1})^{-1} (\mathbf{w}_{\star i}^1 \hat{m}_1 + \boldsymbol{\xi}) \right] \quad (29)$$

which equates

$$m^1 = \frac{1}{p} \hat{m}^1 \times \mathbb{E} \left[\text{Tr}(\mathbf{w}_{\star}^1 \mathbf{w}_{\star}^{1\top} (\mathbf{A} + \Sigma_{\star,1}^{-1})^{-1}) \right] = \hat{m}^1 \times \text{Tr}(\Sigma_{\star,1} (\mathbf{A} + \Sigma_{\star,1}^{-1})^{-1}) \quad (30)$$

and we recover the result of previous work.

Consider now the overlap $\mathbf{Q} \in \mathbb{R}^{2 \times 2}$. We have for $i \neq j$ (the case $i = j$ has been covered before in our AISTATS paper)

$$Q_{1,2} = \mathbb{E}_{\mathbf{w}_{\star}, \boldsymbol{\xi}} \left[\mathbf{f}_a \left(\mathbf{w}_{\star}^1 \hat{m}^1 + \boldsymbol{\xi}^1, \hat{\mathbf{V}}^1 \right)^\top \mathbf{f}_a \left(\mathbf{w}_{\star}^2 \hat{m}^2 + \boldsymbol{\xi}^2, \hat{\mathbf{V}}^2 \right) \right] \quad (31)$$

$$= \hat{m}^1 \hat{m}^2 \times \mathbb{E} \left[\mathbf{w}_{\star}^{1\top} \left(\hat{\mathbf{V}}^1 + \Sigma_{\star,1}^{-1} \right)^{-1} \left(\hat{\mathbf{V}}^2 + \Sigma_{\star,2}^{-1} \right)^{-1} \mathbf{w}_{\star}^2 \right] + \mathbb{E} \left[\boldsymbol{\xi}^{1\top} \left(\hat{\mathbf{V}}^1 + \Sigma_{\star,1}^{-1} \right)^{-1} \left(\hat{\mathbf{V}}^2 + \Sigma_{\star,2}^{-1} \right)^{-1} \boldsymbol{\xi}^2 \right] \quad (32)$$

$$= \hat{m}^1 \hat{m}^2 \times \text{Tr} \left(\text{Cov}(\mathbf{w}_{\star}^1, \mathbf{w}_{\star}^2) \left(\hat{\mathbf{V}}^1 + \Sigma_{\star,1}^{-1} \right)^{-1} \left(\hat{\mathbf{V}}^2 + \Sigma_{\star,2}^{-1} \right)^{-1} \right) + \hat{Q}_{1,2} \times \text{Tr} \left(\text{Cov}(\boldsymbol{\xi}^1, \boldsymbol{\xi}^2) \left(\hat{\mathbf{V}}^1 + \Sigma_{\star,1}^{-1} \right)^{-1} \left(\hat{\mathbf{V}}^2 + \Sigma_{\star,2}^{-1} \right)^{-1} \right) \quad (33)$$

The last line seems coherent with the expression found by Loureiro et al. when the covariance is $\frac{1}{\lambda} I_d$

Overlap V The overlap V is given by

$$V^j = \frac{1}{p} \mathbb{E} \left[(\mathbf{A} + \Sigma_{\star,j}^{-1})^{-1} \right] \quad (34)$$

³The $1/p$ factor comes from the scalar product with the input vector.