

## Supplementary information for

### *Learning Gaussian Mixtures with Generalised Linear Models: Precise Asymptotics in High-dimensions*

**Bruno Loureiro, Gabriele Sicuro, Cédric Gerbelot, Alessandro Pacco, Florent Krzakala, Lenka Zdeborová**

#### A Proof of the main results

This appendix presents the proof of the main technical result, Theorem 1. Throughout the whole proof, we assume that the set of conditions from Sec. 2 is verified.

##### A.1 Required background

In this Section, we give an overview of the main concepts and tools on approximate message passing algorithms which will be required for the proof.

We start with some definitions that commonly appear in the approximate message-passing literature, see e.g. [33, 36, 37]. The main regularity class of functions we will use is that of pseudo-Lipschitz functions, which roughly amounts to functions with polynomially bounded first derivatives. We include the required scaling w.r.t. the dimensions in the definition for convenience.

**Definition 1** (Pseudo-Lipschitz function). *For  $k, K \in \mathbb{N}^*$  and any  $n, m \in \mathbb{N}^*$ , a function  $\phi: \mathbb{R}^{n \times K} \rightarrow \mathbb{R}^{m \times K}$  is called a pseudo-Lipschitz of order  $k$  if there exists a constant  $L(k, K)$  such that for any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{n \times K}$ ,*

$$\frac{\|\phi(\mathbf{x}) - \phi(\mathbf{y})\|_F}{\sqrt{m}} \leq L(k, K) \left( 1 + \left( \frac{\|\mathbf{x}\|_F}{\sqrt{n}} \right)^{k-1} + \left( \frac{\|\mathbf{y}\|_F}{\sqrt{n}} \right)^{k-1} \right) \frac{\|\mathbf{x} - \mathbf{y}\|_F}{\sqrt{n}} \quad (14)$$

where  $\|\bullet\|_F$  denotes the Frobenius norm. Since  $K$  will be kept finite, it can be absorbed in any of the constants.

For example, the function  $f: \mathbb{R}^n \rightarrow \mathbb{R}, \mathbf{x} \mapsto \frac{1}{n} \|\mathbf{x}\|_2^2$  is pseudo-Lipshitz of order 2.

**Moreau envelopes and Bregman proximal operators** — In our proof, we will also frequently use the notions of Moreau envelopes and proximal operators, see e.g. [47, 48]. These elements of convex analysis are often encountered in recent works on high-dimensional asymptotics of convex problems, and more detailed analysis of their properties can be found for example in [12, 31]. For the sake of brevity, we will only sketch the main properties of such mathematical objects, referring to the cited literature for further details. In this proof, we will mainly use proximal operators acting on sets of real matrices endowed with their canonical scalar product. Furthermore, proximals will be defined with matrix valued parameters in the following way: for a given convex function  $f: \mathbb{R}^{d \times K} \rightarrow \mathbb{R}$ , a given matrix  $\mathbf{X} \in \mathbb{R}^{d \times K}$  and a given symmetric positive definite matrix  $\mathbf{V} \in \mathbb{R}^{K \times K}$  with bounded spectral norm, we will consider operators of the type

$$\operatorname{argmin}_{\mathbf{T} \in \mathbb{R}^{d \times K}} \left\{ f(\mathbf{T}) + \frac{1}{2} \operatorname{tr} \left( (\mathbf{T} - \mathbf{X}) \mathbf{V}^{-1} (\mathbf{T} - \mathbf{X})^\top \right) \right\} \quad (15)$$

This operator can either be written as a standard proximal operator by factoring the matrix  $\mathbf{V}^{-1}$  in the arguments of the trace:

$$\operatorname{Prox}_{f(\bullet \mathbf{V}^{1/2})}(\mathbf{X} \mathbf{V}^{-1/2}) \mathbf{V}^{1/2} \in \mathbb{R}^{d \times K} \quad (16)$$

or as a Bregman proximal operator [64] defined with the Bregman distance induced by the strictly convex, coercive function (for positive definite  $\mathbf{V}$ )

$$\mathbf{X} \mapsto \frac{1}{2} \operatorname{tr}(\mathbf{X} \mathbf{V}^{-1} \mathbf{X}^\top) \quad (17)$$

which justifies the use of the Bregman resolvent

$$\operatorname{argmin}_{\mathbf{T} \in \mathbb{R}^{d \times K}} \left\{ f(\mathbf{T}) + \frac{1}{2} \operatorname{tr} \left( (\mathbf{T} - \mathbf{X}) \mathbf{V}^{-1} (\mathbf{T} - \mathbf{X})^\top \right) \right\} = (\operatorname{Id} + \partial f(\bullet) \mathbf{V})^{-1}(\mathbf{X}) \quad (18)$$

Many of the usual or similar properties to that of standard proximal operators (i.e. firm non-expansiveness, link with Moreau/Bregman envelopes,...) hold for Bregman proximal operators defined with the function (17), see e.g. [64, 65]. In particular, we will be using the equivalent notion to firmly nonexpansive operators for Bregman proximity operators, called *D-firm* operators. Consider the Bregman proximal defined with a differentiable, strictly convex, coercive function  $g : \mathcal{X} \rightarrow \mathbb{R}$ , where  $\mathcal{X}$  is a given input Hilbert space. Let  $T$  be the associated Bregman proximal of a given convex function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , i.e., for any  $\mathbf{x} \in \mathcal{X}$

$$T(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y} \in \mathcal{X}} \{f(\mathbf{x}) + D_g(\mathbf{x}, \mathbf{y})\} \quad (19)$$

Then  $T$  is *D-firm*, meaning it verifies

$$\langle T\mathbf{x} - T\mathbf{y}, \nabla g(T\mathbf{x}) - \nabla g(T\mathbf{y}) \rangle \leq \langle T\mathbf{x} - T\mathbf{y}, \nabla g(\mathbf{x}) - \nabla g(\mathbf{y}) \rangle \quad (20)$$

for any  $\mathbf{x}, \mathbf{y}$  in  $\mathcal{X}$ .

**Gaussian concentration** — Gaussian concentration properties are at the root of this proof. Such properties are reviewed in more detail, for example, in [12, 37]. We refer the interested reader to this set of works for a detailed and complete discussion.

**Notations** — For any set of matrices  $\{\mathbf{A}_k \in \mathbb{R}^{n_k \times d_k}\}_{k \in [K]}$  we will use the following notation:

$$\begin{bmatrix} \mathbf{A}_1 & & & \\ & \mathbf{A}_2 & (*) & \\ & (*) & \ddots & \\ & & & \mathbf{A}_K \end{bmatrix} \equiv [\mathbf{A}_k]_{k=1}^K \in \mathbb{R}^{(\sum_{k=1}^K n_k) \times (\sum_{k=1}^K d_k)} \quad (21)$$

where the terms denoted by  $(*)$  will be zero most of the time.

For a given function  $\phi : \mathbb{R}^{d \times K} \rightarrow \mathbb{R}^{d \times K}$ , we write :

$$\phi(\mathbf{X}) = \begin{bmatrix} \phi^1(\mathbf{X}) \\ \vdots \\ \phi^d(\mathbf{X}) \end{bmatrix} \in \mathbb{R}^{d \times K} \quad (22)$$

where each  $\phi^i : \mathbb{R}^{d \times K} \rightarrow \mathbb{R}^K$ . We then write the  $K \times K$  Jacobian

$$\frac{\partial \phi^i}{\partial \mathbf{X}_j}(\mathbf{X}) = \begin{bmatrix} \frac{\partial \phi_1^i(\mathbf{X})}{\partial X_{j1}} & \dots & \frac{\partial \phi_1^i(\mathbf{X})}{\partial X_{jK}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \phi_K^i(\mathbf{X})}{\partial X_{j1}} & \dots & \frac{\partial \phi_K^i(\mathbf{X})}{\partial X_{jK}} \end{bmatrix} \in \mathbb{R}^{K \times K} \quad (23)$$

For a given matrix  $\mathbf{Q} \in \mathbb{R}^{K \times K}$ , we write  $\mathbf{Z} \in \mathbb{R}^{n \times K} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q} \otimes \mathbf{I}_n)$  to denote that the lines of  $\mathbf{Z}$  are sampled i.i.d. from  $\mathcal{N}(\mathbf{0}, \mathbf{Q})$ . Note that this is equivalent to saying that  $\mathbf{Z} = \tilde{\mathbf{Z}}\mathbf{Q}^{1/2}$  where  $\tilde{\mathbf{Z}} \in \mathbb{R}^{n \times K}$  is an i.i.d. standard normal random matrix. The notation  $\stackrel{\text{P}}{\simeq}$  denotes convergence in probability.

**Approximate message-passing** — Approximate message-passing algorithms are a statistical physics inspired family of iterations which can be used to solve high dimensional inference problems [66]. One of the central objects in such algorithms are the so called *state evolution equations*, a low-dimensional recursion equations which allow to exactly compute the high dimensional distribution of the iterates of the sequence. In this proof we will use a specific form of matrix-valued approximate message-passing iteration with non-separable non-linearities. In its full generality, the validity of the state evolution equations in this case is an extension of the works of [36, 37] included in [67]. Consider a sequence Gaussian matrices  $\mathbf{A}(n) \in \mathbb{R}^{n \times d}$  with i.i.d. Gaussian entries,  $A_{ij}(n) \sim \mathcal{N}(0, 1/d)$ . For each  $n, d \in \mathbb{N}$ , consider two sequences of pseudo-Lipschitz functions

$$\{\mathbf{h}_t : \mathbb{R}^{n \times K} \rightarrow \mathbb{R}^{n \times K}\}_{t \in \mathbb{N}} \quad \{\mathbf{e}_t : \mathbb{R}^{d \times K} \rightarrow \mathbb{R}^{d \times K}\}_{t \in \mathbb{N}} \quad (24)$$

initialized on  $\mathbf{u}^0 \in \mathbb{R}^{d \times K}$  in such a way that the limit

$$\lim_{d \rightarrow \infty} \frac{1}{d} \|\mathbf{e}_0(\mathbf{u}^0)^\top \mathbf{e}_0(\mathbf{u}^0)\|_{\text{F}} \quad (25)$$

exists and it is finite, and recursively define:

$$\mathbf{u}^{t+1} = \mathbf{A}^\top \mathbf{h}_t(\mathbf{v}^t) - \mathbf{e}_t(\mathbf{u}^t) \langle \mathbf{h}'_t \rangle^\top \quad (26)$$

$$\mathbf{v}^t = \mathbf{A} \mathbf{e}_t(\mathbf{u}^t) - \mathbf{h}_{t-1}(\mathbf{v}^{t-1}) \langle \mathbf{e}'_t \rangle^\top \quad (27)$$

where the dimension of the iterates are  $\mathbf{u}^t \in \mathbb{R}^{d \times K}$  and  $\mathbf{v}^t \in \mathbb{R}^{n \times K}$ . The terms in brackets are defined as:

$$\langle \mathbf{h}'_t \rangle = \frac{1}{d} \sum_{i=1}^n \frac{\partial \mathbf{h}_t^i}{\partial \mathbf{v}_i}(\mathbf{v}^t) \in \mathbb{R}^{K \times K} \quad \langle \mathbf{e}'_t \rangle = \frac{1}{d} \sum_{i=1}^d \frac{\partial \mathbf{e}_t^i}{\partial \mathbf{u}_i}(\mathbf{u}^t) \in \mathbb{R}^{K \times K} \quad (28)$$

We define now the *state evolution recursion* on two sequences of matrices  $\{\mathbf{Q}_{r,s}\}_{s,r \geq 0}$  and  $\{\hat{\mathbf{Q}}_{r,s}\}_{s,r \geq 1}$  initialized with  $\mathbf{Q}_{0,0} = \lim_{d \rightarrow \infty} \frac{1}{d} \mathbf{e}_0(\mathbf{u}^0)^\top \mathbf{e}_0(\mathbf{u}^0)$ :

$$\mathbf{Q}_{t+1,s} = \mathbf{Q}_{s,t+1} = \lim_{d \rightarrow \infty} \frac{1}{d} \mathbb{E} \left[ \mathbf{e}_s(\hat{\mathbf{Z}}^s)^\top \mathbf{e}_{t+1}(\hat{\mathbf{Z}}^{t+1}) \right] \in \mathbb{R}^{K \times K} \quad (29)$$

$$\hat{\mathbf{Q}}_{t+1,s+1} = \hat{\mathbf{Q}}_{s+1,t+1} = \lim_{d \rightarrow \infty} \frac{1}{d} \mathbb{E} \left[ \mathbf{h}_s(\mathbf{Z}^s)^\top \mathbf{h}_t(\mathbf{Z}^t) \right] \in \mathbb{R}^{K \times K} \quad (30)$$

where  $(\mathbf{Z}^0, \dots, \mathbf{Z}^{t-1}) \sim \mathcal{N}(\mathbf{0}, \{\mathbf{Q}_{r,s}\}_{0 \leq r,s \leq t-1} \otimes \mathbf{I}_n)$ ,  $(\hat{\mathbf{Z}}^1, \dots, \hat{\mathbf{Z}}^t) \sim \mathcal{N}(\mathbf{0}, \{\hat{\mathbf{Q}}_{r,s}\}_{1 \leq r,s \leq t} \otimes \mathbf{I}_d)$ . Then the following holds

**Theorem 4.** *In the setting of the previous paragraph, for any sequence of pseudo-Lipschitz functions  $\phi_n : (\mathbb{R}^{n \times K} \times \mathbb{R}^{d \times K})^t \rightarrow \mathbb{R}$ , for  $n, d \rightarrow \infty$ :*

$$\phi_n(\mathbf{u}^0, \mathbf{v}^0, \mathbf{u}^1, \mathbf{v}^1, \dots, \mathbf{v}^{t-1}, \mathbf{u}^t) \stackrel{\text{P}}{\simeq} \mathbb{E} \left[ \phi_n \left( \mathbf{u}^0, \mathbf{Z}^0, \hat{\mathbf{Z}}^1, \mathbf{Z}^1, \dots, \mathbf{Z}^{t-1}, \hat{\mathbf{Z}}^t \right) \right] \quad (31)$$

where  $(\mathbf{Z}^0, \dots, \mathbf{Z}^{t-1}) \sim \mathcal{N}(\mathbf{0}, \{\mathbf{Q}_{r,s}\}_{0 \leq r,s \leq t-1} \otimes \mathbf{I}_n)$ ,  $(\hat{\mathbf{Z}}^1, \dots, \hat{\mathbf{Z}}^t) \sim \mathcal{N}(\mathbf{0}, \{\hat{\mathbf{Q}}_{r,s}\}_{1 \leq r,s \leq t} \otimes \mathbf{I}_n)$ .

**Spatial coupling** As a final premise to our proof, we give the intuition on how to handle a specific form of block random matrix in an AMP sequence. Consider the iteration (26), but this time with a Gaussian matrix defined as:

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 & & & \\ & \mathbf{A}_2 & (0) & \\ & (0) & \ddots & \\ & & & \mathbf{A}_K \end{bmatrix} \in \mathbb{R}^{n \times Kd} \quad (32)$$

where  $\mathbf{A}_k \in \mathbb{R}^{n_k \times d}$  and  $\sum_{k=1}^K n_k = n$ , which leads to the following form for the products between matrices and non-linearities:

$$\mathbf{A}^\top \mathbf{h}_t(\mathbf{v}^t) = \begin{bmatrix} \mathbf{A}_1^\top \mathbf{h}_{1,t}(\mathbf{v}^t) \\ \mathbf{A}_2^\top \mathbf{h}_{2,t}(\mathbf{v}^t) \\ \vdots \\ \mathbf{A}_K^\top \mathbf{h}_{K,t}(\mathbf{v}^t) \end{bmatrix} \in \mathbb{R}^{Kd \times K} \quad \mathbf{A} \mathbf{e}_t(\mathbf{u}^t) = \begin{bmatrix} \mathbf{A}_1 \mathbf{e}_{1,t}(\mathbf{u}^t) \\ \mathbf{A}_2 \mathbf{e}_{2,t}(\mathbf{u}^t) \\ \vdots \\ \mathbf{A}_K \mathbf{e}_{K,t}(\mathbf{u}^t) \end{bmatrix} \in \mathbb{R}^{n \times K} \quad (33)$$

where the blocks  $\mathbf{h}_{k,t}(\mathbf{v}^t) \in \mathbb{R}^{n_k \times K}$ ,  $\mathbf{e}_{k,t}(\mathbf{u}^t) \in \mathbb{R}^{d \times K}$  may depend on their full arguments or only the corresponding blocks depending on their separability. This iteration can be embedded as a subset of the iterates of a larger sequence defined with the full version of the matrix  $\mathbf{A}$  and non-linearities defined as:

$$\begin{aligned} \mathbf{e}_t : \mathbb{R}^{Kd \times K^2} &\rightarrow \mathbb{R}^{Kd \times K^2} \\ \text{generates} &\begin{bmatrix} \mathbf{e}_{1,t}(\bullet) & & & \\ & \mathbf{e}_{2,t}(\bullet) & (0) & \\ & (0) & \ddots & \\ & & & \mathbf{e}_{K,t}(\bullet) \end{bmatrix} \in \mathbb{R}^{Kd \times K^2} \end{aligned} \quad (34)$$

$$\begin{aligned} \mathbf{h}_t : \mathbb{R}^{n \times K^2} &\rightarrow \mathbb{R}^{n \times K^2} \\ \text{generates} &\begin{bmatrix} \mathbf{h}_{1,t}(\bullet) & & & \\ & \mathbf{h}_{2,t}(\bullet) & (0) & \\ & (0) & \ddots & \\ & & & \mathbf{h}_{K,t}(\bullet) \end{bmatrix} \in \mathbb{R}^{n \times K^2} \end{aligned} \quad (35)$$

The original iteration is recovered on the block diagonal of the variables of the iteration. This new setting, however, introduces a richer correlation structure, since each block will be described by a different  $K \times K$  covariance according to the state evolution equations. Formally, the new covariance will be a  $K^2 \times K^2$  block diagonal matrix. Also, the shape of the Onsager term changes from a matrix of size  $K \times K$  to one of size  $K^2 \times K^2$  with a  $K \times (K \times K)$  block diagonal structure.

## A.2 Reformulation of the problem

We start by reformulating problem (2) in a way that can be treated efficiently using an AMP iteration. With respect to the main part of this paper, we will consider the estimator  $\mathbf{W} \in \mathbb{R}^{d \times K}$  instead of  $\mathbb{R}^{K \times d}$ . The normalized (so that the cost does not diverge with the dimension) problem (2) then reads:

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times K}, \mathbf{b} \in \mathbb{R}^K} \frac{1}{d} \left( L \left( \mathbf{Y}, \frac{1}{\sqrt{d}} \mathbf{X} \mathbf{W} + \mathbf{b} \right) + r(\mathbf{W}) \right) \quad (36)$$

where we have introduced the function  $L : \mathbb{R}^{n \times K} \times \mathbb{R}^{n \times K} \rightarrow \mathbb{R}$  acting as

$$\left( \mathbf{Y}, \frac{1}{\sqrt{d}} \mathbf{X} \mathbf{W} + \mathbf{b} \right) \mapsto \sum_{\nu=1}^n \ell \left( \mathbf{y}^\nu, \frac{\mathbf{W} \mathbf{x}^\nu}{\sqrt{d}} + \mathbf{b} \right), \quad (37)$$

the matrix  $\mathbf{Y} \in \mathbb{R}^{n \times K}$  of concatenated one-hot encoded labels, and the matrix of concatenated means  $\mathbf{M} \in \mathbb{R}^{K \times d}$  (in the main we took the transpose  $\mathbf{M} \in \mathbb{R}^{d \times K}$ ). Until further notice, we will drop the scaling  $\frac{1}{d}$  for convenience and study the problem

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times K}, \mathbf{b} \in \mathbb{R}^K} L \left( \mathbf{Y}, \frac{1}{\sqrt{d}} \mathbf{X} \mathbf{W} + \mathbf{b} \right) + r(\mathbf{W}) \quad (38)$$

We will write  $L_k$  the application of  $\ell$  on each row of a sub-block in  $\mathbb{R}^{n_k \times K}$ . Without loss of generality, we can assume that the samples are grouped by clusters in the data matrix, giving the following form for  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , separating the mean part  $\mathbf{Y} \mathbf{M}$  and centered Gaussian part :

$$\mathbf{X} = \mathbf{Y} \mathbf{M} + \tilde{\mathbf{Z}} \boldsymbol{\Sigma} \in \mathbb{R}^{n \times d} \quad (39)$$

where we have introduced the block-diagonal matrix  $\tilde{\mathbf{Z}}$  and the  $Kd \times d$  full-column-rank matrix  $\boldsymbol{\Sigma}$

$$\tilde{\mathbf{Z}} = \begin{bmatrix} \mathbf{Z}_1 & & & & \\ & \mathbf{Z}_2 & (0) & & \\ & & (0) & \ddots & \\ & & & & \mathbf{Z}_K \end{bmatrix} \in \mathbb{R}^{n \times Kd} \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_1^{1/2} \\ \boldsymbol{\Sigma}_2^{1/2} \\ \vdots \\ \boldsymbol{\Sigma}_K^{1/2} \end{bmatrix} \in \mathbb{R}^{Kd \times d}. \quad (40)$$

Here  $(\mathbf{Z}_1, \dots, \mathbf{Z}_K) \in \mathbb{R}^{n_1 \times d} \times \dots \times \mathbb{R}^{n_K \times d}$  are independent, i.i.d. standard normal matrices.

The product between the data matrix and the weights  $\mathbf{W} \in \mathbb{R}^{d \times K}$  then reads:

$$\mathbf{X} \mathbf{W} = \mathbf{Y} \mathbf{M} \mathbf{W} + \tilde{\mathbf{Z}} \boldsymbol{\Sigma} \mathbf{W} = \begin{bmatrix} \mathbf{Y}_1 \mathbf{M} \mathbf{W} + \mathbf{Z}_1 \boldsymbol{\Sigma}_1^{1/2} \mathbf{W} \\ \vdots \\ \mathbf{Y}_K \mathbf{M} \mathbf{W} + \mathbf{Z}_K \boldsymbol{\Sigma}_K^{1/2} \mathbf{W} \end{bmatrix} \in \mathbb{R}^{n \times K} \quad (41)$$

where each  $\mathbf{Y}_k \in \mathbb{R}^{n_k \times d}$  is a  $n_k$  copy of the same label vector. Defining now  $\tilde{\mathbf{W}} = \boldsymbol{\Sigma} \mathbf{W}$ , observe that

$$\tilde{\mathbf{W}} = \boldsymbol{\Sigma} \mathbf{W} \implies \mathbf{W} = \boldsymbol{\Sigma}^+ \tilde{\mathbf{W}}, \quad (42)$$

where

$$\boldsymbol{\Sigma}^+ \equiv \left( \sum_{k=1}^K \boldsymbol{\Sigma}_k \right)^{-1} \boldsymbol{\Sigma}^\top \quad (43)$$

is the pseudo-inverse of the matrix  $\boldsymbol{\Sigma}$ . The optimization problem (2) is thus equivalent to

$$\inf_{\substack{\tilde{\mathbf{W}} \in \mathbb{R}^{Kd \times K} \\ \mathbf{b} \in \mathbb{R}^K}} \sum_{k=1}^K L_k \left( \frac{1}{\sqrt{d}} \mathbf{Y}_k \mathbf{M} \mathbf{W} + \frac{1}{\sqrt{d}} \mathbf{Z}_k \tilde{\mathbf{W}}_k, \mathbf{b} \right) + r(\boldsymbol{\Sigma}^+ \tilde{\mathbf{W}}) \quad (44)$$

Introducing the order parameter  $\mathbf{m} = \frac{1}{\sqrt{d}}\mathbf{M}\mathbf{W} \in \mathbb{R}^{K \times K}$ , we reformulate Eq.(44) as a constrained optimization problem :

$$\begin{aligned} \inf_{\mathbf{m}, \tilde{\mathbf{W}}, \mathbf{b}} \sum_{k=1}^K L_k \left( \frac{1}{\sqrt{d}} \mathbf{Y}_k \mathbf{m} + \frac{1}{\sqrt{d}} \mathbf{Z}_k \tilde{\mathbf{W}}_k \right) + r \left( \Sigma^+ \tilde{\mathbf{W}} \right) \\ \text{s.t. } \frac{1}{\sqrt{d}} \mathbf{M} \Sigma^+ \tilde{\mathbf{W}} = \mathbf{m} \end{aligned} \quad (45)$$

whose Lagrangian form, with dual parameters  $\hat{\mathbf{m}} \in \mathbb{R}^{K \times K}$ , reads

$$\inf_{\mathbf{m}, \tilde{\mathbf{W}}, \mathbf{b}} \sup_{\hat{\mathbf{m}}} \sum_{k=1}^K L_k \left( \mathbf{Y}_k \mathbf{m} + \frac{1}{\sqrt{d}} \mathbf{Z}_k \tilde{\mathbf{W}}_k \right) + r \left( \Sigma^+ \tilde{\mathbf{W}} \right) + \text{tr} \left( \hat{\mathbf{m}}^\top \left( \mathbf{m} - \frac{1}{\sqrt{d}} \mathbf{M} \Sigma^+ \tilde{\mathbf{W}} \right) \right). \quad (46)$$

This is a proper, closed, convex, strictly feasible optimization problem, thus strong duality holds and we can invert the order of the inf-sup to focus on the minimization problem in  $\tilde{\mathbf{W}}$  for fixed  $\mathbf{m}, \hat{\mathbf{m}}, \mathbf{b}$ :

$$\inf_{\tilde{\mathbf{W}} \in \mathbb{R}^{Kd \times K}} \tilde{L} \left( \frac{1}{\sqrt{d}} \tilde{\mathbf{Z}} \tilde{\mathbf{W}} \right) + \tilde{r}(\tilde{\mathbf{W}}) \quad (47)$$

where we defined the loss term

$$\begin{aligned} \tilde{L} : \mathbb{R}^{n \times K} \rightarrow \mathbb{R} \\ \frac{1}{\sqrt{d}} \tilde{\mathbf{Z}} \tilde{\mathbf{W}} \mapsto \sum_{k=1}^K L_k \left( \mathbf{Y}_k \mathbf{m} + \frac{1}{\sqrt{d}} \mathbf{Z}_k \tilde{\mathbf{W}}_k \right) = \sum_{k=1}^K \sum_{i=1}^{n_k} \ell \left( \left[ \mathbf{Y}_k \mathbf{m} + \frac{1}{\sqrt{d}} \mathbf{Z}_k \tilde{\mathbf{W}}_k \right]_i \right) \end{aligned} \quad (48a)$$

and the regularisation term

$$\begin{aligned} \tilde{r} : \mathbb{R}^{Kd \times K} \rightarrow \mathbb{R} \\ \tilde{\mathbf{W}} \mapsto r \left( \Sigma^+ \tilde{\mathbf{W}} \right) + \text{tr} \left( \hat{\mathbf{m}}^\top \left( \mathbf{m} - \frac{1}{\sqrt{d}} \mathbf{M} \Sigma^+ \tilde{\mathbf{W}} \right) \right) \end{aligned} \quad (48b)$$

where  $\Sigma^\top \tilde{\mathbf{W}} = \sum_{k=1}^K \Sigma_k^{1/2} \tilde{\mathbf{W}}_k$  and  $\tilde{\mathbf{Z}} = [\mathbf{Z}_k]_{k=1}^K \in \mathbb{R}^{n \times Kd}$  is an i.i.d. standard normal block diagonal matrix as in Eq. (40).

### A.3 Finding the AMP sequence

We now need to find an AMP iteration relating to  $\tilde{\mathbf{W}}$  that solve the optimization problem in Eq. (47). Although this section is not written as a formal proof, all steps are rigorous. The aim is to give the reader the core intuition on how the AMP iteration is found, otherwise the solution may feel ‘‘parachuted’’. The reader uninterested in the underlying intuition may directly skip to the next section. In order to find the appropriate sequence two key points must be considered :

- the fixed point of the sequence has to match the optimality condition of Eq. (47);
- the update rule of the sequence should have the form Eq. (26) for the state evolution equations to hold.

These two points completely determine the form of the iteration. In the subsequent derivation, we absorb the scaling  $\frac{1}{\sqrt{d}}$  in the matrix  $\tilde{\mathbf{Z}}$ , such that the  $\mathbf{Z}_k \in \mathbb{R}^{n_k \times d}$  have i.i.d.  $\mathcal{N}(0, 1/d)$  elements.

**Resolvent of the loss term** — Going back to problem Eq. (47), its optimality condition will look like :

$$\tilde{\mathbf{Z}}^\top \partial \tilde{L}(\tilde{\mathbf{Z}} \tilde{\mathbf{W}}) + \partial \tilde{r}(\tilde{\mathbf{W}}) = 0 \iff \begin{bmatrix} \mathbf{Z}_1^\top & & & \\ & \mathbf{Z}_2^\top & (0) & \\ & (0) & \ddots & \\ & & & \mathbf{Z}_K^\top \end{bmatrix} \begin{bmatrix} \partial \tilde{L}_1(\mathbf{Z}_1 \tilde{\mathbf{W}}_1) \\ \partial \tilde{L}_2(\mathbf{Z}_2 \tilde{\mathbf{W}}_2) \\ \vdots \\ \partial \tilde{L}_K(\mathbf{Z}_K \tilde{\mathbf{W}}_K) \end{bmatrix} + \partial \tilde{r}(\tilde{\mathbf{W}}) = 0 \quad (49)$$

where each  $\mathbf{Z}_k \in \mathbb{R}^{n_k \times d}$ , and the subdifferential of  $\tilde{L}$  is separable across blocks of size  $n_k \times d$ , and  $\partial\tilde{r}(\tilde{\mathbf{W}}) \in \mathbb{R}^{Kd \times K}$ . Following the intuition of spatial coupling, we introduce the *full* matrix  $\mathbf{Z} \in \mathbb{R}^{n \times Kd}$ , with i.i.d.  $\mathcal{N}(0, 1/d)$  entries. The optimality condition can then be written on the diagonal of a  $Kd \times K^2$  matrix:

$$\mathbf{Z}^\top \begin{bmatrix} \partial\tilde{L}_1(\mathbf{Z}_1 \tilde{\mathbf{W}}_1) & & & & \\ & \partial\tilde{L}_2(\mathbf{Z}_2 \tilde{\mathbf{W}}_2) & (0) & & \\ & (0) & & \ddots & \\ & & & & \partial\tilde{L}_K(\mathbf{Z}_K \tilde{\mathbf{W}}_K) \end{bmatrix} + \begin{bmatrix} \partial\tilde{r}(\tilde{\mathbf{W}})_1 & & & & \\ & \partial\tilde{r}(\tilde{\mathbf{W}})_2 & (0) & & \\ & (0) & & \ddots & \\ & & & & \partial\tilde{r}(\tilde{\mathbf{W}})_K \end{bmatrix} = \mathbf{0} \quad (50)$$

where  $\partial\tilde{r}(\tilde{\mathbf{W}})_k$  represents the  $k$ -th block of the subdifferential of  $\tilde{r}$  which is non-separable across the blocks of  $\tilde{\mathbf{W}}$ . To make the resolvents/proximals appear, we add the argument of the subdifferentials on both sides weighted by a (symmetric) positive definite matrix  $\mathbf{S}_k \in \mathbb{R}^{K \times K}$  which will be used to allow for Onsager correction while respecting the fixed point condition. Using the notation defined in section A.1

$$\begin{aligned} & \left[ \mathbf{Z}_k^\top \partial\tilde{L}_k(\mathbf{Z}_k \tilde{\mathbf{W}}_k) \right]_{k=1}^K + \left[ \partial\tilde{r}(\tilde{\mathbf{W}}) \right]_{k=1}^K = 0 \\ \iff & \left[ \mathbf{Z}_k^\top \partial\tilde{L}_k(\mathbf{Z}_k \tilde{\mathbf{W}}_k) + \mathbf{Z}_k^\top \mathbf{Z}_k \tilde{\mathbf{W}}_k \mathbf{S}_k^{-1} \right]_{k=1}^K + \left[ \partial\tilde{r}(\tilde{\mathbf{W}}) \right]_{k=1}^K = \left[ \mathbf{Z}_k^\top \mathbf{Z}_k \tilde{\mathbf{W}}_k \mathbf{S}_k^{-1} \right]_{k=1}^K \end{aligned} \quad (51)$$

for a given set of positive definite matrices  $\{\mathbf{S}_k\}_{k \in [K]}$ . Again, the reason for introducing different  $\mathbf{S}_k$  on each block is to match the expected structure of the Onsager term. We can introduce the resolvent, formally Bregman resolvent/proximal operator:

$$\mathbf{U}_k \equiv \partial\tilde{L}_k(\mathbf{Z}_k \tilde{\mathbf{W}}_k) \mathbf{S}_k + \mathbf{Z}_k \tilde{\mathbf{W}}_k \iff \mathbf{Z}_k \tilde{\mathbf{W}}_k = \mathbf{R}_{\tilde{L}_k, \mathbf{S}_k}(\mathbf{U}_k) \quad (52)$$

where

$$\begin{aligned} \mathbf{R}_{\tilde{L}_k, \mathbf{S}_k}(\mathbf{U}_k) &= (\text{Id} + \partial\tilde{L}_k(\bullet) \mathbf{S}_k)^{-1}(\mathbf{U}_k) \\ &= \underset{\mathbf{T} \in \mathbb{R}^{n_k \times K}}{\text{argmin}} \left\{ \tilde{L}_k(\mathbf{T}) + \frac{1}{2} \text{tr}((\mathbf{T} - \mathbf{U}_k) \mathbf{S}_k^{-1} (\mathbf{T} - \mathbf{U}_k)^\top) \right\} \\ &= \underset{\mathbf{T} \in \mathbb{R}^{n_k \times K}}{\text{argmin}} \left\{ L_k(\mathbf{T}) + \frac{1}{2} \text{tr}((\mathbf{T} - (\mathbf{Y}_k \mathbf{m} + \mathbf{U}_k)) \mathbf{S}_k^{-1} (\mathbf{T} - (\mathbf{Y}_k \mathbf{m} + \mathbf{U}_k))^\top) \right\} - \mathbf{Y}_k \mathbf{m}. \end{aligned} \quad (53)$$

In the previous expressions  $\partial\tilde{L}_k \in \mathbb{R}^{n_k \times K}$  and  $\mathbf{V}_k \in \mathbb{R}^{K \times K}$ . The following formulation of the optimality condition is reached:

$$\begin{aligned} & \left[ \mathbf{Z}_k^\top \mathbf{U}_k \mathbf{S}_k^{-1} \right]_{k=1}^K + \left[ \partial\tilde{r}(\tilde{\mathbf{W}})_k \right]_{k=1}^K = \left[ \mathbf{Z}_k^\top \mathbf{R}_{\tilde{L}_k, \mathbf{S}_k}(\mathbf{U}_k) \mathbf{S}_k^{-1} \right]_{k=1}^K \\ \iff & \left[ \mathbf{Z}_k^\top (\mathbf{U}_k - \mathbf{R}_{\tilde{L}_k, \mathbf{S}_k}(\mathbf{U}_k)) \mathbf{S}_k^{-1} \right]_{k=1}^K + \left[ \partial\tilde{r}(\tilde{\mathbf{W}})_k \right]_{k=1}^K = 0 \end{aligned} \quad (54)$$

**Resolvent of the regularization term** Determining the block decomposition of the subdifferential of the regularization term is less simple. We would like a block expression in the flavour of:

$$\left[ \partial\tilde{r}(\tilde{\mathbf{W}})_k \right]_{k=1}^K + \left[ \tilde{\mathbf{W}}_k \hat{\mathbf{S}}_k^{-1} \right]_{k=1}^K = \left[ \tilde{\mathbf{W}}_k \hat{\mathbf{S}}_k^{-1} \right]_{k=1}^K \quad (55)$$

At this point it becomes clear that we cannot consider the resolvent as acting on  $\tilde{\mathbf{W}} \in \mathbb{R}^{Kd \times K}$  otherwise there could be only one  $\hat{\mathbf{S}} \in \mathbb{R}^{K \times K}$  and there would be a mismatch with the expected form of the Onsager terms. As specified by the definitions Eq.(48), the subdifferential of  $\tilde{r}$  is acting on the whole block diagonal matrix  $[\tilde{\mathbf{W}}_k]_{k=1}^K$ , by way of summation due to the action of the pseudo-inverse

$\Sigma^+$ . We can thus consider its proximal acting on  $\mathbb{R}^{d \times K^2}$  as  $[\tilde{\mathbf{W}}_1 \tilde{\mathbf{W}}_2 \dots \tilde{\mathbf{W}}_K]$  (note that we could have also worked directly with a block diagonal matrix in  $\mathbb{R}^{Kd \times K^2}$ ). Proceeding in this way, we can directly write our expression as an application parametrized by another set of positive definite matrices  $\{\hat{\mathbf{S}}_k\}_{k \in [K]}$ .

$$\hat{\mathbf{U}} = (\text{Id} + \partial \tilde{r}(\bullet) \hat{\mathbf{S}}) (\tilde{\mathbf{W}}) \quad \tilde{\mathbf{W}} = \mathbf{R}_{\tilde{r}, \hat{\mathbf{S}}}(\hat{\mathbf{U}}) \quad (56)$$

where

$$\begin{aligned} \mathbf{R}_{\tilde{r}, \hat{\mathbf{S}}}(\hat{\mathbf{U}}) &= (\text{Id} + \partial \tilde{r}(\bullet) \hat{\mathbf{S}})^{-1} (\hat{\mathbf{U}}) \\ &= \underset{\mathbf{T} \in \mathbb{R}^{d \times K^2}}{\text{argmin}} \left\{ \tilde{r}(\mathbf{T}) + \frac{1}{2} \text{tr} \left( (\mathbf{T} - \hat{\mathbf{U}}) \hat{\mathbf{S}}^{-1} (\mathbf{T} - \hat{\mathbf{U}})^\top \right) \right\} \end{aligned} \quad (57)$$

where  $\hat{\mathbf{S}} \in \mathbb{R}^{K^2 \times K^2}$  block diagonal, and  $\hat{\mathbf{U}} \in \mathbb{R}^{d \times K^2}$ . This would lead to the equivalent optimality condition for the regularization part:

$$\hat{\mathbf{U}} \hat{\mathbf{S}}^{-1} = \mathbf{R}_{\tilde{r}, \hat{\mathbf{S}}}(\hat{\mathbf{U}}) \hat{\mathbf{S}}^{-1} \iff \left[ \hat{\mathbf{U}}_k \hat{\mathbf{S}}_k^{-1} \right]_{k=1}^K = \left[ \mathbf{R}_{\tilde{r}, \hat{\mathbf{S}}, k}(\hat{\mathbf{U}}) \hat{\mathbf{S}}_k^{-1} \right]_{k=1}^K \quad (58)$$

We now need to figure out the block structure of this resolvent since we want to spread it across a block diagonal matrix. Let  $\mathbf{C} = \sum_{k=1}^K \Sigma_k$ , so that  $\Sigma^+ = \mathbf{C}^{-1} \Sigma^\top$ , and the blocks  $\mathbf{T}_k \in \mathbb{R}^{d \times K}$  are the solution to the minimization problem

$$\begin{aligned} \min_{\{\mathbf{T}_k\}_{k \in [K]} \in (\mathbb{R}^{d \times K})^K} & r(\mathbf{C}^{-1} \sum_{k=1}^K \Sigma_k^{1/2} \mathbf{T}_k) + \frac{1}{2} \text{tr} \left( (\mathbf{T} - \hat{\mathbf{U}}) \hat{\mathbf{S}}^{-1} (\mathbf{T} - \hat{\mathbf{U}})^\top \right) \\ & + \text{tr} \left( \hat{\mathbf{m}}^\top \left( \mathbf{m} - \frac{1}{\sqrt{d}} \mathbf{M} \Sigma^+ \mathbf{T} \right) \right) \end{aligned} \quad (59)$$

Let  $\tilde{\mathbf{T}} = \mathbf{C}^{-1} \sum_{k=1}^K \Sigma_k^{1/2} \mathbf{T}_k \in \mathbb{R}^{d \times K}$ , and the equivalent reformulation as a constraint optimization problem:

$$\begin{aligned} \min_{\substack{\mathbf{T}_k \in [K] \in \mathbb{R}^{d \times K} \\ \tilde{\mathbf{T}} \in \mathbb{R}^{d \times K}}} & r(\tilde{\mathbf{T}}) + \frac{1}{2} \text{tr} \left( (\mathbf{T} - \hat{\mathbf{U}}) \hat{\mathbf{S}}^{-1} (\mathbf{T} - \hat{\mathbf{U}})^\top \right) + \text{tr} \left( \hat{\mathbf{m}}^\top \left( \mathbf{m} - \frac{1}{\sqrt{d}} \mathbf{M} \tilde{\mathbf{T}} \right) \right) \\ \text{s.t.} & \tilde{\mathbf{T}} = \mathbf{C}^{-1} \sum_{k=1}^K \Sigma_k^{1/2} \mathbf{T}_k \end{aligned} \quad (60)$$

This is a feasible convex problem under convex constraint with a strongly convex term, it thus has a unique solution and strong duality holds. Introducing the Lagrange multiplier  $\lambda \in \mathbb{R}^{d \times K}$ , we get the equivalent representation:

$$\begin{aligned} \min_{\substack{\mathbf{T}_k \in [K] \in \mathbb{R}^{d \times K} \\ \tilde{\mathbf{T}} \in \mathbb{R}^{d \times K}}} \max_{\lambda \in \mathbb{R}^{d \times K}} & r(\tilde{\mathbf{T}}) + \sum_{k=1}^K \text{tr} \left( (\mathbf{T}_k - \hat{\mathbf{U}}_k) \hat{\mathbf{S}}_k^{-1} (\mathbf{T}_k - \hat{\mathbf{U}}_k)^\top \right) \\ & + \text{tr} \left( \lambda^\top \left( \tilde{\mathbf{T}} - \mathbf{C}^{-1} \sum_{k=1}^K \Sigma_k^{1/2} \mathbf{T}_k \right) \right) + \text{tr} \left( \hat{\mathbf{m}}^\top \left( \mathbf{m} - \frac{1}{\sqrt{d}} \mathbf{M} \tilde{\mathbf{T}} \right) \right). \end{aligned} \quad (61)$$

The optimality condition for this problem reads:

$$\partial_{\tilde{\mathbf{T}}} : \quad \partial r(\tilde{\mathbf{T}}) + \lambda - \frac{1}{\sqrt{d}} \mathbf{M}^\top \hat{\mathbf{m}} = 0 \quad (62)$$

$$\partial_{\mathbf{T}} : \quad (\mathbf{T}_k - \hat{\mathbf{U}}_k) \hat{\mathbf{S}}_k^{-1} = \Sigma_k^{1/2} \mathbf{C}^{-1} \lambda \quad \forall k \in [K] \quad (63)$$

$$\partial_{\lambda} : \quad \tilde{\mathbf{T}} = \mathbf{C}^{-1} \sum_{k=1}^K \Sigma_k^{1/2} \mathbf{T}_k \quad (64)$$

Using the gradient condition on  $\mathbf{T}$ , we get

$$\sum_{k=1}^K \Sigma_k^{1/2} (\mathbf{T}_k - \hat{\mathbf{U}}_k) \hat{\mathbf{S}}_k^{-1} = \lambda \quad (65)$$

The constraint  $\tilde{\mathbf{T}} = \mathbf{C}^{-1} \sum_{k=1}^K \Sigma_k^{1/2} \mathbf{T}_k$  is solved by  $\mathbf{T}_k = \Sigma_k^{1/2} \tilde{\mathbf{T}}$  which gives the solution for  $\lambda$

$$\lambda = \sum_{k=1}^K \Sigma_k^{1/2} (\Sigma_k^{1/2} \tilde{\mathbf{T}} - \hat{\mathbf{U}}_k) \hat{\mathbf{S}}_k^{-1} = \sum_{k=1}^K \Sigma_k \tilde{\mathbf{T}} \hat{\mathbf{S}}_k^{-1} - \sum_{k=1}^K \Sigma_k^{1/2} \hat{\mathbf{U}}_k \hat{\mathbf{S}}_k^{-1} \quad (66)$$

and prescribes the following form for  $\tilde{\mathbf{T}}$ , as solution to the problem

$$\begin{aligned} \partial r(\tilde{\mathbf{T}}) + \sum_{k=1}^K \Sigma_k \tilde{\mathbf{T}} \hat{\mathbf{S}}_k^{-1} - \sum_{k=1}^K \Sigma_k^{1/2} \hat{\mathbf{U}}_k \hat{\mathbf{S}}_k^{-1} - \frac{1}{\sqrt{d}} \mathbf{M}^\top \hat{\mathbf{m}} &= 0 \\ \iff \operatorname{argmin}_{\tilde{\mathbf{T}}} r(\tilde{\mathbf{T}}) + \frac{1}{2} \sum_{k=1}^K \Sigma_k \tilde{\mathbf{T}} \hat{\mathbf{S}}_k^{-1} \tilde{\mathbf{T}} - \left( \sum_{k=1}^K \Sigma_k^{1/2} \hat{\mathbf{U}}_k \hat{\mathbf{S}}_k^{-1} + \frac{1}{\sqrt{d}} \mathbf{M}^\top \hat{\mathbf{m}} \right) \tilde{\mathbf{T}} & \quad (67) \end{aligned}$$

We then recover  $\mathbf{T}$  from  $\mathbf{T} = \Sigma \tilde{\mathbf{T}}$ . Thus, defining the function

$$\begin{aligned} \eta : \mathbb{R}^{d \times K^2} &\rightarrow \mathbb{R}^{d \times K} \\ \hat{\mathbf{U}} &\mapsto \operatorname{argmin}_{\tilde{\mathbf{T}}} r(\tilde{\mathbf{T}}) + \frac{1}{2} \sum_{k=1}^K \Sigma_k \tilde{\mathbf{T}} \hat{\mathbf{S}}_k^{-1} \tilde{\mathbf{T}} - \left( \sum_{k=1}^K \Sigma_k^{1/2} \hat{\mathbf{U}}_k \hat{\mathbf{S}}_k^{-1} + \frac{1}{\sqrt{d}} \mathbf{M}^\top \hat{\mathbf{m}} \right) \tilde{\mathbf{T}} \quad (68) \end{aligned}$$

the block decomposition of the resolvent for the regularizer reads:

$$\mathbf{R}_{\tilde{\mathbf{r}}, \hat{\mathbf{S}}, k}(\hat{\mathbf{U}}) = \Sigma_k^{1/2} \eta(\hat{\mathbf{U}}) \quad (69)$$

**Matching the optimality condition with the AMP fixed point** The global optimality condition then reads:

$$\left[ \mathbf{Z}_k^\top \left( \mathbf{R}_{\tilde{\mathbf{L}}, k, \mathbf{S}_k}(\mathbf{U}_k) - \mathbf{U}_k \right) \mathbf{S}_k^{-1} \right]_{k=1}^K = \left[ (\hat{\mathbf{U}}_k - \mathbf{R}_{\tilde{\mathbf{r}}, \hat{\mathbf{S}}, k}(\hat{\mathbf{U}})) \hat{\mathbf{S}}_k^{-1} \right]_{k=1}^K \quad (70)$$

$$\left[ \mathbf{Z}_k \mathbf{R}_{\tilde{\mathbf{r}}, \hat{\mathbf{S}}, k}(\hat{\mathbf{U}}) \right]_{k=1}^K = \left[ \mathbf{R}_{\tilde{\mathbf{L}}, k, \mathbf{S}_k}(\mathbf{U}_k) \right]_{k=1}^K \quad (71)$$

where both equations should be satisfied. We can now define update functions based on the previously obtained block decomposition. The fixed point of the matrix-valued AMP Eq.(26) reads:

$$\operatorname{Id} + \mathbf{e}(\mathbf{u}) \langle \mathbf{h}' \rangle^\top = \mathbf{Z}^\top \mathbf{h}(\mathbf{v}) \quad (72)$$

$$\operatorname{Id} + \mathbf{h}(\mathbf{v}) \langle \mathbf{e}' \rangle^\top = \mathbf{Z} \mathbf{e}(\mathbf{u}) \quad (73)$$

Matching this fixed point with the optimality condition Eq.(70) suggests the following mapping:

$$\begin{aligned} \mathbf{h}_k(\mathbf{U}_k) &= \left( \mathbf{R}_{\tilde{\mathbf{L}}, k, \mathbf{S}_k}(\mathbf{U}_k) - \mathbf{U}_k \right) \mathbf{S}_k^{-1}, & \mathbf{S}_k &= \langle \mathbf{e}'_k \rangle, \\ \mathbf{e}_k(\hat{\mathbf{U}}) &= \mathbf{R}_{\tilde{\mathbf{r}}, \hat{\mathbf{S}}, k}(\hat{\mathbf{U}} \hat{\mathbf{S}}), & \hat{\mathbf{S}}_k &= -\langle \mathbf{h}'_k \rangle^{-1}, \end{aligned} \quad (74)$$

where we redefined  $\hat{\mathbf{U}} \equiv \hat{\mathbf{U}} \hat{\mathbf{S}}$  in (56), and the subscripts on the non-linearities are block indexes.



**Lemma 5.** Consider the sequence defined by Eq.(81), for any fixed  $\mathbf{m}, \hat{\mathbf{m}}, \mathbf{b}$ . For any sequences of pseudo-Lipschitz functions  $\phi_{1,n} : \mathbb{R}^{d \times K^2} \rightarrow \mathbb{R}, \phi_{2,n} : \mathbb{R}^{n \times K^2} \rightarrow \mathbb{R}$ , for any  $t \in \mathbb{N}^*$ :

$$\phi_{1,n}(\mathbf{u}_1^t, \dots, \mathbf{u}_K^t) \stackrel{\text{P}}{\simeq} \mathbb{E} \left[ \phi_{1,n}(\mathbf{H}_1(\hat{\mathbf{Q}}_{1,t})^{1/2}, \dots, \mathbf{H}_K(\hat{\mathbf{Q}}_{K,t})^{1/2}) \right] \quad (83)$$

$$\phi_{2,n}(\mathbf{v}_1^t, \dots, \mathbf{v}_K^t) \stackrel{\text{P}}{\simeq} \mathbb{E} \left[ \phi_{2,n}(\mathbf{G}_1(\mathbf{Q}_{1,t})^{1/2}, \dots, \mathbf{G}_K(\mathbf{Q}_{K,t})^{1/2}) \right] \quad (84)$$

where the matrices  $\mathbf{H}_k \in \mathbb{R}^{d \times K}, \mathbf{G}_k \in \mathbb{R}^{n_k \times K}$  are independent matrices with i.i.d. standard normal elements, and at each time step  $t \geq 1$

$$\begin{aligned} \mathbf{Q}_{k,t} &= \lim_{d \rightarrow +\infty} \frac{1}{d} \mathbb{E} \left[ \mathbf{e}_{k,t} (\{\mathbf{H}_k(\hat{\mathbf{Q}}_{k,t})^{1/2}(\hat{\mathbf{V}}_{k,t})^{-1}\}_{k \in [K]})^\top \mathbf{e}_{k,t} (\{\mathbf{H}_k(\hat{\mathbf{Q}}_{k,t})^{1/2}(\hat{\mathbf{V}}_{k,t})^{-1}\}_{k \in [K]}) \right] \\ &\in \mathbb{R}^{K \times K} \end{aligned} \quad (85)$$

$$\hat{\mathbf{Q}}_{k,t} = \lim_{d \rightarrow +\infty} \frac{1}{d} \mathbb{E} \left[ \mathbf{h}_{k,t-1} (\mathbf{G}_k(\mathbf{Q}_{k,t-1})^{1/2})^\top \mathbf{h}_{k,t-1} (\mathbf{G}_k(\mathbf{Q}_{k,t-1})^{1/2}) \right] \in \mathbb{R}^{K \times K} \quad (86)$$

$$\mathbf{V}_{k,t} = \lim_{d \rightarrow +\infty} \frac{1}{d} \sum_{i=1}^d \frac{\partial e_{k,t-1} (\{\mathbf{H}_k(\hat{\mathbf{Q}}_{k,t-1})^{1/2}\}_{k \in [K]})}{\partial (\mathbf{H}_k(\hat{\mathbf{Q}}_{k,t-1})^{1/2})_i} \in \mathbb{R}^{K \times K} \quad (87)$$

$$\hat{\mathbf{V}}_{k,t} = - \lim_{d \rightarrow +\infty} \frac{1}{d} \sum_{i=1}^{n_k} \frac{\partial \mathbf{h}_{k,t} (\mathbf{G}_k(\mathbf{Q}_{k,t})^{1/2})}{\partial (\mathbf{G}_k(\mathbf{Q}_{k,t})^{1/2})_i} \in \mathbb{R}^{K \times K} \quad (88)$$

where the sequence is initialized with  $\hat{\mathbf{V}}_0, \mathbf{e}_0, \mathbf{Q}_{0,0} = \lim_{d \rightarrow \infty} \frac{1}{d} \|\mathbf{e}_0(\mathbf{u}^0)^\top \mathbf{e}_0(\mathbf{u}^0)\|_{\text{F}}$ .

*Proof.* Lemma 5 is a consequence of Theorem 4 whose assumptions have been verified in the paragraph.  $\square$

Note that in Lemma 5, we have directly written the block decomposition of the state evolution corresponding to the iteration Eq. (81), which involves the block diagonal matrices  $\mathbf{Q}_t, \hat{\mathbf{Q}}_t, \mathbf{V}_t, \hat{\mathbf{V}}_t$  which are all in  $\mathbb{R}^{K^2 \times K^2}$ . Using the notations introduced in section A.1

$$\mathbf{V} = [\mathbf{V}_k]_{k=1}^K \quad \hat{\mathbf{V}} = [\hat{\mathbf{V}}_k]_{k=1}^K \quad \mathbf{Q} = [\mathbf{Q}_k]_{k=1}^K \quad \hat{\mathbf{Q}} = [\hat{\mathbf{Q}}_k]_{k=1}^K \quad (89)$$

Also note that we do not use the full state evolution giving the correlations across all time steps, but only use those at equal times  $t$ .

**Trajectories and fixed point of the AMP sequence** Now that we have a sequence with state evolution equations, the following two lemmas link the fixed points of this iteration to any optimal solution of problem Eq.(47).

**Lemma 6.** Consider any fixed point  $\mathbf{V}, \hat{\mathbf{V}}, \mathbf{Q}, \hat{\mathbf{Q}}$  of the state evolution equations from Lemma 5. For any fixed point  $\mathbf{u}^*, \mathbf{v}^*$  of iteration Eq.(81), the quantity

$$\mathbf{R}_{\tilde{\mathbf{r}}, \hat{\mathbf{V}}^{-1}}(\mathbf{u}^* \hat{\mathbf{V}}^{-1}) = \left( \text{Id} + \partial \tilde{\mathbf{r}}(\bullet) \hat{\mathbf{V}}^{-1} \right) (\mathbf{u}^* \hat{\mathbf{V}}^{-1}) \quad (90)$$

is an optimal solution  $\tilde{\mathbf{W}}^*$  of problem Eq.(47). Furthermore

$$\mathbf{R}_{\tilde{\mathbf{L}}, \mathbf{V}}(\mathbf{v}^*) = (\text{Id} + \partial \tilde{\mathbf{L}}(\bullet) \mathbf{V})(\mathbf{v}^*) = \mathbf{Z} \tilde{\mathbf{W}}^* \quad (91)$$

where the block decompositions of each resolvents have been explicitly calculated in section A.3.

*Proof.* Lemma 6 is a direct consequence of the analysis carried out in section A.3.  $\square$

At this point we know the fixed points of the AMP iteration correspond to the optimal solutions of problem Eq.(47). Note that the resolvents/proximals linking the fixed point of the AMP iteration with the solutions of Eq.(47) are Lipschitz continuous, making them acceptable transforms for state evolution observables. However this does not guarantee that the optimal solution is characterized by the fixed point of the state evolution equations. Indeed, we need to show that a converging trajectory can be systematically found for any instance of the problem Eq.(47). This is the purpose of the following lemma.

**Lemma 7.** Consider iteration Eq.(81), where the parameters  $\mathbf{Q}, \hat{\mathbf{Q}}, \mathbf{V}, \hat{\mathbf{V}}$  are initialized at any fixed point of the state evolution equations of Lemma 5. For any sequence initialized with  $\hat{\mathbf{V}}_0 = \hat{\mathbf{V}}$  and  $\mathbf{u}^0$  such that

$$\lim_{d \rightarrow \infty} \frac{1}{d} \mathbf{e}_0(\mathbf{u}^0)^\top \mathbf{e}_0(\mathbf{u}^0) = \mathbf{Q} \quad (92)$$

the following holds

$$\lim_{t \rightarrow \infty} \lim_{d \rightarrow \infty} \frac{1}{\sqrt{d}} \|\mathbf{u}^t - \mathbf{u}^*\|_{\text{F}} = 0 \quad \lim_{t \rightarrow \infty} \lim_{d \rightarrow \infty} \frac{1}{\sqrt{d}} \|\mathbf{v}^t - \mathbf{v}^*\|_{\text{F}} = 0 \quad (93)$$

*Proof.* The proof of Lemma 7 is deferred to subsection A.7.  $\square$

Note that the  $\mathbf{G}$  defined here is not the same as the  $\mathbf{G}$  in the replica computation. Combining the lemmas 5, 6 and 7 with the pseudo-Lipschitz property, we have reached the following lemma

**Lemma 8.** For any fixed  $\mathbf{m}, \hat{\mathbf{m}}, \mathbf{b}$ , consider the fixed point  $(\mathbf{Q}, \hat{\mathbf{Q}}, \mathbf{V}, \hat{\mathbf{V}})$  of the state evolution equations from Lemma. 5. Then, for any sequences of pseudo-Lipschitz functions  $\phi_{1,n} : \mathbb{R}^{d \times K^2} \rightarrow \mathbb{R}, \phi_{2,n} : \mathbb{R}^{n \times K} \rightarrow \mathbb{R}$ , for  $n, d \rightarrow \infty$

$$\phi_{1,n}(\tilde{\mathbf{W}}^*) \stackrel{\text{P}}{\simeq} \mathbb{E} \left[ \phi_{1,n} \left( R_{\hat{\mathbf{r}}, \hat{\mathbf{V}}}(\mathbf{H} \hat{\mathbf{Q}}^{1/2} \hat{\mathbf{V}}^{-1}) \right) \right] \quad (94)$$

$$\phi_{2,n}(\mathbf{Z} \tilde{\mathbf{W}}^*) \stackrel{\text{P}}{\simeq} \mathbb{E} \left[ \phi_{2,n} \left( R_{\hat{\mathbf{L}}, \mathbf{V}}(\mathbf{G} \mathbf{Q}^{1/2}) \right) \right] \quad (95)$$

where we remind that  $\mathbf{G} = [\mathbf{G}_k]_{k=1}^K, \mathbf{H} = [\mathbf{H}_k]_{k=1}^K$  are block diagonal i.i.d. standard normal matrices as in Lemma 5, and  $\mathbf{Q} = [\mathbf{Q}_k]_{k=1}^K, \hat{\mathbf{Q}} = [\hat{\mathbf{Q}}_k]_{k=1}^K$  are the  $K^2 \times K^2$  block diagonal covariances.

*Proof.* Lemma 8 is a consequence of Lemmas 5,6,7 and applying the pseudo-Lipschitz property along with the fact that the iterates of the AMP have bounded norm using the state evolution and that the estimator also has bounded norm (feasibility assumption). Note that, for a generically non-strictly convex problem, being close to the zero gradient condition does not guarantee being close to the estimator. This is further discussed in Appendix A.5.  $\square$

Note that the resolvents are implicitly acting on the block diagonals of their arguments. At this point we are quite close to Theorem 1(details for the exact matching will be given later), but we are missing the equations on  $\mathbf{m}, \hat{\mathbf{m}}, \mathbf{b}$ .

**Fixed point equations for  $\mathbf{m}, \hat{\mathbf{m}}, \mathbf{b}$**  We drop the dependence on the bias term  $\mathbf{b}$  as its solution is very similar to the one for  $\mathbf{m}, \hat{\mathbf{m}}$ . To obtain the equations for  $\mathbf{m}, \hat{\mathbf{m}}$ , we go back to the complete optimization problem

$$\begin{aligned} \inf_{\mathbf{m}, \tilde{\mathbf{W}}, \mathbf{b}} \sup_{\hat{\mathbf{m}}} L(\mathbf{Y}_k \mathbf{m} + \mathbf{Z}_k \tilde{\mathbf{W}}_k) + r(\Sigma^+ \tilde{\mathbf{W}}) \\ + \text{tr} \left( \hat{\mathbf{m}}^\top \left( \mathbf{m} - \frac{1}{\sqrt{d}} \mathbf{M} \Sigma^+ \tilde{\mathbf{W}} \right) \right) \end{aligned} \quad (96)$$

where we can use strong duality to write the equivalent form

$$\begin{aligned} \inf_{\mathbf{m}, \mathbf{b}} \sup_{\hat{\mathbf{m}}} L(\mathbf{Y}_k \mathbf{m} + \mathbf{Z}_k \tilde{\mathbf{W}}_k^*) + r(\Sigma^+ \tilde{\mathbf{W}}) \\ + \text{tr} \left( \hat{\mathbf{m}}^\top \left( \mathbf{m} - \frac{1}{\sqrt{d}} \mathbf{M} \Sigma^+ \tilde{\mathbf{W}}^* \right) \right) \end{aligned} \quad (97)$$

The gradients w.r.t.  $\mathbf{m}, \hat{\mathbf{m}}$  then read:

$$\partial \hat{\mathbf{m}} = \mathbf{m} - \frac{1}{\sqrt{d}} \mathbf{M} \Sigma^+ \tilde{\mathbf{W}}^* \quad (98)$$

$$\partial \mathbf{m} = \hat{\mathbf{m}} + \partial_{\mathbf{m}} L(\mathbf{Y} \mathbf{m} + \mathbf{Z} \tilde{\mathbf{W}}^*) \quad (99)$$

Uniform convergence of derivatives and conditions for the dominated convergence theorem are verified using similar arguments as in [12, Lemma 12]. We can thus invert limits and derivatives, and expectations and derivatives. To facilitate taking the derivative  $\partial_{\mathbf{m}}$ , we use Lemma 8 (assuming the normalized loss function is pseudo-Lipschitz, which is a very loose assumption verified by most machine learning losses) to obtain, reintroducing the scaling  $1/d$

$$\frac{1}{d}L(\mathbf{Y}\mathbf{m} + \mathbf{Z}\tilde{\mathbf{W}}^*) \xrightarrow{d \rightarrow \infty} \frac{1}{d}\mathbb{E} \left[ L(\mathbf{Y}\mathbf{m} + \mathbf{R}_{\bar{L},\mathbf{V}}(\mathbf{G}\mathbf{Q}^{1/2})) \right] \quad (100)$$

Using the block decomposition from Eq.(53), the blocks  $(\mathbf{R}_{\bar{L},\mathbf{V}}(\mathbf{G}\mathbf{Q}^{1/2}))_k \in \mathbb{R}^{n_k \times K}$  are given by:

$$\operatorname{argmin}_{\mathbf{T} \in \mathbb{R}^{n_k \times K}} \left\{ L_k(\mathbf{T}) + \frac{1}{2} \operatorname{tr} \left( (\mathbf{T} - (\mathbf{Y}_k \mathbf{m} + \mathbf{G}_k \mathbf{Q}_k^{1/2})) \mathbf{V}_k^{-1} (\mathbf{T} - (\mathbf{Y}_k \mathbf{m} + \mathbf{G}_k \mathbf{Q}_k^{1/2}))^\top \right) \right\} - \mathbf{Y}_k \mathbf{m} \quad (101)$$

Using a block diagonal representation, we can write:

$$\begin{aligned} \frac{1}{d}L(\mathbf{Y}\mathbf{m} + \mathbf{R}_{\bar{L},\mathbf{V}}(\mathbf{G}\mathbf{Q}^{1/2})) &= \frac{1}{d}L(\mathbf{R}_{L,\mathbf{V}}(\mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2})) \\ &= \frac{1}{d}\mathcal{M}_{L,\mathbf{V}}(\mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2}) - \\ &\frac{1}{2d} \operatorname{tr} \left( (\mathbf{R}_{L,\mathbf{V}}(\mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2}) - (\mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2})) \mathbf{V}^{-1} (\mathbf{R}_{L,\mathbf{V}}(\mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2}) - (\mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2}))^\top \right) \end{aligned} \quad (102)$$

where we have introduced the Bregman-envelope [65] with respect to the distance Eq. (17)

$$\begin{aligned} \mathcal{M}_{L,\mathbf{V}}(\mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2}) &= \\ \min_{\mathbf{T}} \left\{ L(\mathbf{T}) + \frac{1}{2} \operatorname{tr} \left( (\mathbf{T} - (\mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2})) \mathbf{V}^{-1} (\mathbf{T} - (\mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2}))^\top \right) \right\} \end{aligned} \quad (103)$$

Then, using the state evolution equations from Lemma 5 and Stein's lemma, we can write:

$$\frac{1}{d}L(\mathbf{Y}\mathbf{m} + \mathbf{R}_{\bar{L},\mathbf{V}}(\mathbf{G}\mathbf{Q}^{1/2})) = \frac{1}{d}\mathcal{M}_{L,\mathbf{V}}(\mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2}) - \frac{1}{2} \operatorname{tr}(\mathbf{V}^\top \mathbf{Q}) \quad (104)$$

Taking the gradient w.r.t.  $\mathbf{m}$  using the expression for the derivative of a Bregman envelope [65], we get:

$$\partial_{\mathbf{m}} L(\mathbf{Y}\mathbf{m} + \mathbf{R}_{\bar{L},\mathbf{V}}(\mathbf{G}\mathbf{Q}^{1/2})) = \frac{1}{d} \mathbf{Y}^\top \left( \mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2} - \mathbf{R}_{L,\mathbf{V}}(\mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2}) \right) \mathbf{V}^{-1} \quad (105)$$

which prescribes, using Lemma 8

$$\hat{\mathbf{m}} \stackrel{\text{P}}{\simeq} \frac{1}{d} \mathbf{Y}^\top \left( \mathbf{R}_{L,\mathbf{V}}(\mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2}) - \mathbf{Y}\mathbf{m} + \mathbf{G}\mathbf{Q}^{1/2} \right) \mathbf{V}^{-1} \quad (106)$$

For  $\mathbf{m}$ , we use the block decomposition from Eq.(67), which simplifies the pseudo-inverse  $\Sigma^+$  in Eq. (98) to give, using Lemma 8 again

$$\mathbf{m} \stackrel{\text{P}}{\simeq} \frac{1}{\sqrt{d}} \mathbf{M} \boldsymbol{\eta} (\mathbf{H} \hat{\mathbf{Q}}^{1/2} \hat{\mathbf{V}}^{-1}) \quad (107)$$

where the function  $\boldsymbol{\eta}$  acts on the block diagonal and is defined by Eq.(68). Using those results and the definition of  $\tilde{\mathbf{W}}$ , the solution  $\mathbf{W}^*$  and the quantity  $\mathbf{X}\mathbf{W}^*$  are characterized, in the pseudo-Lipschitz sense of Theorem 1, by the fixed point of the system of equations (the first four equations are meant

for all  $1 \leq k \leq K$ ):

$$\mathbf{Q}_k = \lim_{d \rightarrow +\infty} \frac{1}{d} \mathbb{E} \left[ \mathbf{e}_k(\{\mathbf{H}_k(\hat{\mathbf{Q}}_k)^{1/2} \hat{\mathbf{V}}_k^{-1}\}_{k \in [K]})^\top \mathbf{e}_k(\{\mathbf{H}_k(\hat{\mathbf{Q}}_k)^{1/2} \hat{\mathbf{V}}_k^{-1}\}_{k \in [K]}) \right] \in \mathbb{R}^{K \times K} \quad (108)$$

$$\hat{\mathbf{Q}}_k = \lim_{d \rightarrow +\infty} \frac{1}{d} \mathbb{E} \left[ \mathbf{h}_k(\mathbf{G}_k \mathbf{Q}_k^{1/2})^\top \mathbf{h}_k(\mathbf{G}_k \mathbf{Q}_k^{1/2}) \right] \in \mathbb{R}^{K \times K} \quad (109)$$

$$\mathbf{V}_k = \lim_{d \rightarrow +\infty} \frac{1}{d} \sum_{i=1}^d \mathbb{E} \left[ \frac{\partial \mathbf{e}_k(\{\mathbf{H}_k(\hat{\mathbf{Q}}_k)^{1/2}\}_{k \in [K]})}{\partial (\mathbf{H}_k(\hat{\mathbf{Q}}_k)^{1/2})_i} \right] \in \mathbb{R}^{K \times K} \quad (110)$$

$$\hat{\mathbf{V}}_k = - \lim_{d \rightarrow +\infty} \frac{1}{d} \sum_{i=1}^{n_k} \mathbb{E} \left[ \frac{\partial \mathbf{h}_{k,t}(\mathbf{G}_k(\mathbf{Q}_{k,t})^{1/2})}{\partial (\mathbf{G}_k(\mathbf{Q}_k)^{1/2})_i} \right] \in \mathbb{R}^{K \times K} \quad (111)$$

$$\mathbf{m} = \frac{1}{\sqrt{d}} \mathbb{E} \left[ \mathbf{M} \boldsymbol{\eta} (\mathbf{H} \hat{\mathbf{Q}}^{1/2} \hat{\mathbf{V}}^{-1}) \right] \in \mathbb{R}^{K \times K} \quad (112)$$

$$\hat{\mathbf{m}} = \frac{1}{d} \mathbf{Y}^\top \left( \mathbf{R}_{L,\mathbf{V}} (\mathbf{Y} \mathbf{m} + \mathbf{G} \mathbf{Q}^{1/2}) - \mathbf{Y} \mathbf{m} + \mathbf{G} \mathbf{Q}^{1/2} \right) \mathbf{V}^{-1} \in \mathbb{R}^{K \times K} \quad (113)$$

Using the explicit form of the different functions given in section A.3 and Stein's lemma for the derivatives, these equations match those of Theorem 1. This completes the proof.

### A.5 On the strict convexity assumption

If the optimization problem defining  $\mathbf{W}^*$  is strictly convex, there is only one minimizer and the provided proof is enough. Additionally it is shown in [68] that for any loss function that is strictly convex in its argument and penalized with the  $\ell_1$  norm, provided the data is sampled from a continuous distribution, the solution is unique with probability one regardless of the rank of the design matrix. Thus finding a point verifying the optimality condition of (47) is also enough to complete the proof. For generic convex (non-strictly) problems a more careful analysis could be performed in the same spirit as the one of [51]. Empirically the result still holds.

### A.6 On the uniqueness of the solution to the fixed point equations (108)

It is possible to reconstruct Bregman envelopes on problem (47) for the loss and regularization as we have done for the loss in the previous section. We can then show that the fixed point equations (108) are the optimality condition of a convex-concave problem involving both Bregman envelopes and linear combinations of the order parameters. In the same spirit as [12, 49], this problem should be asymptotically strictly convex. This is supported by the simulations presented in the experiments sections but left as an assumption in the main paper.

### A.7 Proof of Lemma 7

This proof follows a similar argument to the one used to control the trajectory of the AMP studied in [50]. Note that, because of the way the AMP is initialized using the fixed point of the state evolution equations, for any  $t \geq 1$  the following holds:

$$\lim_{d \rightarrow +\infty} \frac{1}{d} \mathbb{E} \left[ \mathbf{e}(\mathbf{u}^t)^\top \mathbf{e}(\mathbf{u}^t) \right] = \mathbf{Q} \in \mathbb{R}^{K^2 \times K^2} \quad (114)$$

$$\lim_{d \rightarrow +\infty} \frac{1}{d} \mathbb{E} \left[ \mathbf{h}(\mathbf{v}^t)^\top \mathbf{h}(\mathbf{v}^t) \right] = \hat{\mathbf{Q}} \in \mathbb{R}^{K^2 \times K^2} \quad (115)$$

where

$$\mathbf{e}(\mathbf{u}^t) = (\text{Id} + \partial \tilde{r}(\bullet) \hat{\mathbf{V}}^{-1})^{-1} (\mathbf{u}^t \hat{\mathbf{V}}^{-1}) \quad \mathbf{h}(\mathbf{v}^t) = \left( (\text{Id} + \partial \tilde{L}(\bullet) \mathbf{V})^{-1} (\mathbf{v}^t) - \mathbf{v}^t \right) \mathbf{V}^{-1} \quad (116)$$

then the limit we are looking for reads:

$$\begin{aligned} \lim_{d \rightarrow \infty} \frac{1}{d} \|\mathbf{u}^t - \mathbf{u}^{t-1}\|_F^2 &= \lim_{d \rightarrow \infty} 2(\hat{\mathbf{Q}} - \frac{1}{d} \text{tr}((\mathbf{u}^t)^\top \mathbf{u}^{t-1})) \\ \lim_{d \rightarrow \infty} \frac{1}{d} \|\mathbf{v}^t - \mathbf{v}^{t-1}\|_F^2 &= 2(\mathbf{Q} - \frac{1}{d} \text{tr}((\mathbf{v}^t)^\top \mathbf{v}^{t-1})) \end{aligned} \quad (117)$$

We thus need to study the correlation between successive iterates. At each time step, denote  $(\hat{C}_t, C_t)$  in  $\mathbb{R}^{K^2 \times K^2}$  the correlation matrices between iterates at times  $t, t-1$  describing the Gaussian fields respectively associated to  $\mathbf{u}^t, \mathbf{v}^t$  i.e.,

$$\lim_{d \rightarrow \infty} \frac{1}{d} \text{tr}((\mathbf{u}^t)^\top \mathbf{u}^{t-1}) = \hat{C}_t \quad \lim_{d \rightarrow \infty} \frac{1}{d} \text{tr}((\mathbf{v}^t)^\top \mathbf{v}^{t-1}) = C_t \quad (118)$$

we can then write the block diagonal Gaussian fields  $\hat{\mathbf{Z}}^t, \hat{\mathbf{Z}}^{t-1}, \mathbf{Z}^t, \mathbf{Z}^{t-1}$  in  $\mathbb{R}^{Kd \times K^2}$  and in the following way

$$\hat{\mathbf{Z}}^t \sim \mathbf{H}(\hat{C}_t)^{1/2} + \mathbf{H}'(\hat{Q} - \hat{C}_t)^{1/2} \quad (119)$$

$$\hat{\mathbf{Z}}^{t-1} \sim \mathbf{H}(\hat{C}_t)^{1/2} + \mathbf{H}''(\hat{Q} - \hat{C}_t)^{1/2} \quad (120)$$

$$\mathbf{Z}^t \sim \mathbf{G}(C_t)^{1/2} + \mathbf{G}'(Q - C_t)^{1/2} \quad (121)$$

$$\mathbf{Z}^{t-1} \sim \mathbf{G}(C_t)^{1/2} + \mathbf{G}''(Q - C_t)^{1/2} \quad (122)$$

where the matrices  $\mathbf{H}, \mathbf{H}', \mathbf{H}''$  are in  $\mathbb{R}^{Kd \times K^2}$ ,  $\mathbf{G}, \mathbf{G}', \mathbf{G}''$  are in  $\mathbb{R}^{n \times K^2}$  and all have i.i.d. standard normal elements. The recursion describing the evolution of these correlations then reads :

$$C_{t+1} = \frac{1}{d} \mathbb{E} \left[ e(\mathbf{H}\hat{C}_t^{1/2} + \mathbf{H}'(\hat{Q} - \hat{C}_t)^{1/2})^\top e(\mathbf{H}\hat{C}_t^{1/2} + \mathbf{H}''(\hat{Q} - \hat{C}_t)^{1/2}) \right] \quad (123)$$

$$\hat{C}_t = \frac{1}{d} \mathbb{E} \left[ \mathbf{h}(\mathbf{G}C_t^{1/2} + \mathbf{G}'(Q - C_t)^{1/2})^\top \mathbf{h}(\mathbf{G}C_t^{1/2} + \mathbf{G}''(Q - C_t)^{1/2}) \right] \quad (124)$$

Integrating out the independent  $\mathbf{H}', \mathbf{H}''$  first, we get

$$C_{t+1} = \int_{\mathbb{R}^{Kd \times K^2}} d\mu(\mathbf{H}) \mathbf{I}(\mathbf{H})^\top \mathbf{I}(\mathbf{H}) \quad (125)$$

where  $\mathbf{I}(\mathbf{H}) = \int_{\mathbb{R}^{Kd \times K^2}} d\mu(\mathbf{H}') e(\mathbf{H}\hat{C}_t^{1/2} + \mathbf{H}'(\hat{Q} - \hat{C}_t)^{1/2})$ . So  $C^t$  is symmetric positive definite, assuming the resolvents aren't trivial. The same argument applied to  $\hat{C}^t$  shows it is also symmetric positive definite. From [64], the operators

$$(Id + \partial \tilde{r}(\bullet) \hat{V}^{-1})^{-1}(\bullet) \quad \left( Id + \partial \tilde{L}(\bullet) \mathbf{V} \right)^{-1}(\bullet) \quad (126)$$

are *D-firm* w.r.t. the Bregman distances induced by the differentiable, strictly convex functions  $\frac{1}{2} \text{tr}(X \hat{V} X^\top)$  and  $\frac{1}{2} \text{tr}(X \mathbf{V}^{-1} X^\top)$  respectively. Recall

$$e(\mathbf{u}^t) = (Id + \partial \tilde{r}(\bullet) \hat{V}^{-1})^{-1}(\mathbf{u}^t \hat{V}^{-1}) \quad \mathbf{h}(\mathbf{v}^t) = \left( \left( Id + \partial \tilde{L}(\bullet) \mathbf{V} \right)^{-1}(\mathbf{v}^t) - \mathbf{v}^t \right) \mathbf{V}^{-1} \quad (127)$$

Then, using the definition of *D-firm*

$$\langle e(\hat{\mathbf{Z}}^t) - e(\hat{\mathbf{Z}}^{t-1}), (e(\hat{\mathbf{Z}}^t) - e(\hat{\mathbf{Z}}^{t-1})) \hat{V} \rangle \leq \langle e(\hat{\mathbf{Z}}^t) - e(\hat{\mathbf{Z}}^{t-1}), (\hat{\mathbf{Z}}^t - \hat{\mathbf{Z}}^{t-1}) \hat{V}^{-1} \hat{V} \rangle \quad (128)$$

then, adding the normalization by  $\frac{1}{d}$ , using the representation Eq.(119-122), taking expectations and applying the matrix form of Stein's lemma, see for example [67] Lemma 12, we get:

$$\text{tr}((Q - C_{t+1}) \hat{V}) \leq \text{tr}((\hat{Q} - \hat{C}_t) \mathbf{V}) \quad (129)$$

Using a similar argument on  $\mathbf{h}$ , we get

$$\text{tr}((\hat{Q} - \hat{C}_t) \mathbf{V}) \leq \text{tr}((Q - C_t) \hat{V}) \quad (130)$$

and

$$\text{tr}(C_{t+1} \hat{V}) \geq \text{tr}(C_t \hat{V}) \quad (131)$$

thus the sequence  $\text{tr}(C_{t+1} \hat{V})$  is a bounded (above) monotone (increasing) sequence, and therefore converges. Since  $\hat{V}$  is positive definite and given the iteration defining  $C_{t+1}$  from  $C_t$ , any fixed point of this iteration is a fixed point of  $\text{tr}(C_t \hat{V})$ . Assuming there is only one fixed point to the set of self-consistent equations Eq.(8) (see previous section), the proof is complete. (A similar argument can be carried out on  $\hat{C}_t$ ).

## B Replica computation

### B.1 Setting of the problem

In this Section we give a full derivation of the results in Theorem 1 and Theorem 2 by means of the replica approach, a standard method developed in the realm of statistical physics of disordered systems [69]. In the general computation, we will consider the classification problem of  $K$  clusters, assuming a dataset  $\{(\mathbf{x}^\nu, \mathbf{y}^\nu)\}_{\nu \in [n]}$  of  $n$  independent datapoints where, as in the main text, the labels  $\mathbf{y}$  takes value in a set of  $K$  elements,  $\mathbf{y}^\nu \in \{\mathbf{e}_k\}_k$ , with  $\mathbf{e}_k \in \mathbb{R}^L$ . The elements of the dataset are independently generated by a mixture density in the form

$$P(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^K \mathbb{I}(\mathbf{y} = \mathbf{e}_k) \rho_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad \sum_{k=1}^K \rho_k = 1. \quad (132)$$

We will perform our classification task searching for a set of parameters  $(\mathbf{W}^*, \mathbf{b}^*)$  that will allow us to construct an estimator. The parameters will be chosen by minimising an empirical risk function in the form

$$\mathcal{R}(\mathbf{W}, \mathbf{b}) \equiv \sum_{\nu=1}^n \ell \left( \mathbf{y}^\nu, \frac{\mathbf{W} \mathbf{x}^\nu}{\sqrt{d}} + \mathbf{b} \right) + \lambda r(\mathbf{W}), \quad (133)$$

i.e., they are given by

$$(\mathbf{W}^*, \mathbf{b}^*) \equiv \underset{\mathbf{W} \in \mathbb{R}^{L \times d}, \mathbf{b} \in \mathbb{R}^L}{\operatorname{argmin}} \mathcal{R}(\mathbf{W}, \mathbf{b}). \quad (134)$$

We will say that  $\mathbf{W} \in \mathbb{R}^{L \times d}$  and  $\mathbf{b} \in \mathbb{R}^L$  are the weights and bias to be learned respectively,  $\ell$  is a convex loss function with respect to its second argument, and  $r$  is a regularisation function whose strength is tuned by the parameter  $\lambda \geq 0$ . Finally, we will assume that a classifier  $\varphi: \mathbb{R}^L \rightarrow \{\mathbf{e}_k\}_k$  is given, such that, once  $(\mathbf{W}^*, \mathbf{b}^*)$  are obtained, a new point  $\mathbf{x}$  is assigned to the label

$$\mathbf{x} \mapsto \varphi \left( \frac{\mathbf{W}^* \mathbf{x}}{\sqrt{d}} + \mathbf{b}^* \right) \in \{\mathbf{e}_k\}_k. \quad (135)$$

The described setting is slightly more general than the one given in Theorem 1. As a consequence of the fact that we choose  $L$ -dimensional labels, the order parameters that appear in the computation are  $L$  dimensional vectors or  $L \times L$  matrices. A typical ‘‘high-dimensional encoding’’ is the one-hot encoding convention adopted in Theorem 1, where  $L = K$  and  $\{\mathbf{e}_k\}_k \subset \mathbb{R}^K$  is the canonical basis of  $\mathbb{R}^K$ . In this case, the adopted classifier is

$$\varphi(\mathbf{x}) \equiv \hat{\mathbf{y}}(\mathbf{x}), \quad \hat{y}_k(\mathbf{x}) = \mathbb{I}(\max_{\kappa} x_{\kappa} = x_k). \quad (136)$$

Assuming *scalar* labels  $\{e_k\}_k \in \mathbb{R}$ , we deal with scalar order parameters. For example, in the case of binary classification ( $K = 2$ ) it is common to adopt  $L = 1$  and  $\{e_1, e_2\} = \{+1, -1\}$ . In this case  $\varphi(x) = \operatorname{sign}(x)$ , see also Section C.2.

### B.2 Gibbs minimisation

The problem stated in Section 1 is formulated as an optimisation problem. We can tackle such optimisation problem introducing a Gibbs measure over the weights  $(\mathbf{W}, \mathbf{b})$ , namely

$$\mu_{\beta}(\mathbf{W}, \mathbf{b}) \propto e^{-\beta \mathcal{R}(\mathbf{W}, \mathbf{b})} = \underbrace{e^{-\beta r(\mathbf{W})}}_{P_w(\mathbf{W})} \prod_{\nu=1}^n \underbrace{\exp \left[ -\beta \ell \left( \mathbf{y}^\nu, \frac{\mathbf{W} \mathbf{x}^\nu}{\sqrt{d}} + \mathbf{b} \right) \right]}_{P_y(\mathbf{y} | \mathbf{W}, \mathbf{b})}. \quad (137)$$

The parameter  $\beta > 0$  is introduced for convenience: in the  $\beta \rightarrow +\infty$  limit, the Gibbs measure concentrates on the values  $(\mathbf{W}^*, \mathbf{b}^*)$  which minimize the empirical risk  $\mathcal{R}(\mathbf{W}, \mathbf{b})$  and are therefore the goal of the learning process. The functions  $P_y$  and  $P_w$  can be interpreted as a (unnormalised) likelihood and prior distribution respectively. Our analysis will go through the computation of the average free energy density associated to such Gibbs measure, i.e.,

$$f_{\beta} = - \lim_{\substack{n, d \rightarrow +\infty \\ n/d = \alpha}} \mathbb{E}_{\{(\mathbf{x}, \mathbf{y})\}} \left[ \frac{\ln \mathcal{Z}_{\beta}}{d} \right], \quad (138)$$

where  $\mathbb{E}_{\{(\mathbf{x}, \mathbf{y})\}}[\bullet]$  is the average over the training dataset, and we have introduced the partition function

$$\mathcal{Z}_\beta \equiv \int e^{-\beta \mathcal{R}(\mathbf{W}, \mathbf{b})} d\mathbf{W} \quad (139)$$

To perform the computation of such quantity, we use the so-called replica method, i.e., we compute

$$-\lim_{\substack{n, d \rightarrow +\infty \\ n/d = \alpha}} \mathbb{E}_{\{(\mathbf{x}, \mathbf{y})\}} \left[ \frac{\ln \mathcal{Z}_\beta}{d\beta} \right] = \lim_{\substack{n, d \rightarrow +\infty \\ n/d = \alpha}} \lim_{s \rightarrow 0} \frac{1 - \mathbb{E}_{\{(\mathbf{x}, \mathbf{y})\}}[\mathcal{Z}_\beta^s]}{sd\beta}, \quad (140)$$

### B.3 Replica approach

We proceed in our calculation considering the bias vector assuming no prior on  $\mathbf{b}$ , which will play a role of an extra parameter. The equations for the bias  $\mathbf{b}$  will be derived extremising with respect to it the final result for the free energy. We need to evaluate

$$\mathbb{E}_{\{(\mathbf{x}, \mathbf{y})\}}[\mathcal{Z}_\beta^s] = \prod_{a=1}^s \int d\mathbf{W}^a P_w(\mathbf{W}^a) \left( \sum_k \rho_k \mathbb{E}_{\mathbf{x}|\mathbf{y}=e_k} \left[ \prod_{a=1}^s P_y \left( e_k \left| \frac{\mathbf{W}^a \mathbf{x}}{\sqrt{d}} + \mathbf{b} \right. \right) \right] \right)^n. \quad (141)$$

Let us take the inner average introducing a new variable  $\boldsymbol{\eta}$ ,

$$\begin{aligned} \mathbb{E}_{\mathbf{x}|\mathbf{y}=e_k} \left[ \prod_{a=1}^s P_y \left( e_k \left| \frac{\mathbf{W}^a \mathbf{x}}{\sqrt{d}} + \mathbf{b} \right. \right) \right] &= \prod_{a=1}^s \int d\boldsymbol{\eta}^a P_y(e_k | \boldsymbol{\eta}^a) \mathbb{E}_{\mathbf{x}} \left[ \prod_{a=1}^s \delta \left( \boldsymbol{\eta}^a - \frac{\mathbf{W}^a \mathbf{x}}{\sqrt{d}} + \mathbf{b} \right) \right] \\ &= \prod_{a=1}^s \int d\boldsymbol{\eta}^a P_y(e_k | \boldsymbol{\eta}^a) \mathcal{N} \left( \boldsymbol{\eta} \left| \frac{\mathbf{W}^a \boldsymbol{\mu}_k}{\sqrt{d}} - \mathbf{b}; \frac{\mathbf{W}^a \boldsymbol{\Sigma}_k \mathbf{W}^{b\top}}{d} \right. \right). \end{aligned} \quad (142)$$

We can write then

$$\begin{aligned} \mathbb{E}_{\{(\mathbf{x}, \mathbf{y})\}}[\mathcal{Z}_\beta^s] &= \\ &= \prod_{a=1}^n \int d\mathbf{W}^a P_w(\mathbf{W}^a) \left( \sum_k \rho_k \prod_{a=1}^s \int d\boldsymbol{\eta}^a P_y(e_k | \boldsymbol{\eta}^a) \mathcal{N} \left( \boldsymbol{\eta}; \frac{\mathbf{W}^a \boldsymbol{\mu}_k}{d} + \mathbf{b}; \frac{\mathbf{W}^a \boldsymbol{\Sigma}_k \mathbf{W}^{b\top}}{d} \right) \right)^n \\ &= \left( \prod_{k=1}^K \prod_{a \leq b} \iint \frac{d\mathbf{Q}_k^{ab} d\hat{\mathbf{Q}}_k^{ab}}{(2\pi)^{L^2/2}} \right) \left( \prod_k \prod_a \int \frac{d\mathbf{m}_k^a d\hat{\mathbf{m}}_k^a}{(2\pi)^{L/2}} \right) e^{-d\beta\Phi^{(s)}}. \end{aligned} \quad (143)$$

where we introduced the *order parameters*

$$\mathbf{Q}_k^{ab} = \frac{\mathbf{W}^a \boldsymbol{\Sigma}_k \mathbf{W}^{b\top}}{d} \in \mathbb{R}^{L \times L}, \quad a, b = 1, \dots, s, \quad (144)$$

$$\mathbf{m}_k^a = \frac{\mathbf{W}^a \boldsymbol{\mu}_k}{\sqrt{d}} \in \mathbb{R}^L, \quad a = 1, \dots, s, \quad (145)$$

and the replicated free-energy

$$\begin{aligned} \beta\Phi^{(s)}(\mathbf{Q}, \mathbf{m}, \hat{\mathbf{Q}}, \hat{\mathbf{m}}, \mathbf{b}) &= \sum_{k=1}^K \sum_a \hat{\mathbf{m}}_k^{a\top} \mathbf{m}_k^a + \sum_{k=1}^K \sum_{a \leq b} \text{tr} \left[ \hat{\mathbf{Q}}_k^{ab\top} \mathbf{Q}_k^{ab} \right] \\ &\quad - \frac{1}{d} \ln \prod_{a=1}^s \int P_w(\mathbf{W}^a) d\mathbf{W}^a \prod_k \left( \prod_{a \leq b} e^{\text{tr}[\hat{\mathbf{Q}}_k^{ab\top} \mathbf{W}^a \boldsymbol{\Sigma}_k \mathbf{W}^{b\top}]} \prod_a e^{\sqrt{d} \hat{\mathbf{m}}_k^{a\top} \mathbf{W}^a \boldsymbol{\mu}_k} \right) \\ &\quad - \alpha \ln \sum_k \rho_k \prod_{a=1}^s \int d\boldsymbol{\eta}^a P_y(e_k | \boldsymbol{\eta}^a) \mathcal{N}(\boldsymbol{\eta} | \mathbf{m}_k^a + \mathbf{b}, \mathbf{Q}_k^{ab}). \end{aligned} \quad (146)$$

At this point, the free energy  $f_\beta$  should be computed extremising with respect to all the order parameters by virtue of the Laplace approximation (in addition to  $\mathbf{b}$ ),

$$f_\beta = \lim_{s \rightarrow 0} \text{Extr}_{\{\mathbf{m}, \mathbf{Q}, \hat{\mathbf{m}}, \hat{\mathbf{Q}}\}, \mathbf{b}} \frac{\Phi^{(s)}(\mathbf{Q}, \mathbf{m}, \hat{\mathbf{Q}}, \hat{\mathbf{m}}, \mathbf{b})}{s}. \quad (147)$$

However, the convexity of the problem allows us to make an important simplification.

**Replica symmetric ansatz** — Before taking the  $s \rightarrow 0$  limit we make the assumptions

$$\begin{aligned} \mathbf{Q}_k^{aa} &= \begin{cases} \mathbf{R}_k, & a = b \\ \mathbf{Q}_k & a \neq b \end{cases} & \hat{\mathbf{Q}}_k^{aa} &= \begin{cases} -\frac{1}{2}\mathbf{R}_k, & a = b \\ \hat{\mathbf{Q}}_k & a \neq b \end{cases} \\ \mathbf{m}_k^a &= \mathbf{m}_k & \hat{\mathbf{m}}_k^a &= \hat{\mathbf{m}}_k \quad \forall a \end{aligned} \quad (148)$$

This ansatz is justified by the fact that we are assuming  $\ell$  and  $r$  to be convex, and  $\lambda > 0$ . In this case, the problem admit one solution only that, therefore, coincide with the replica symmetric solution, in which overlaps between two replicas do not depend on the chosen replicas. By means of the replica symmetric hypothesis, we can write

$$\mathbf{Q}_k^{ab} \mapsto \mathbf{Q}_k \equiv \mathbf{I}_{s,s} \otimes (\mathbf{R}_k - \mathbf{Q}_k) + \mathbf{1}_s \otimes \mathbf{Q}_k. \quad (149)$$

The inverse matrix is therefore

$$\mathbf{Q}_k^{-1} = \mathbf{1}_s \otimes (\mathbf{R}_k - \mathbf{Q}_k)^{-1} - \mathbf{I}_{s,s} \otimes [(\mathbf{R}_k + (s-1)\mathbf{Q}_k)^{-1} \mathbf{Q}_k (\mathbf{R}_k - \mathbf{Q}_k)^{-1}], \quad (150)$$

whereas

$$\begin{aligned} \det \mathbf{Q}_k &= \det(\mathbf{R}_k - \mathbf{Q}_k)^{s-1} \det(\mathbf{R}_k + (s-1)\mathbf{Q}_k) \\ &= 1 + s \ln \det(\mathbf{R}_k - \mathbf{Q}_k) + s \operatorname{tr} [(\mathbf{R}_k - \mathbf{Q}_k)^{-1} \mathbf{Q}_k] + o(s). \end{aligned} \quad (151)$$

If we denote  $\mathbf{V}_k \equiv \mathbf{R}_k - \mathbf{Q}_k$

$$\begin{aligned} &\ln \sum_k \rho_k \prod_{a=1}^s \int d\boldsymbol{\eta}^a P_y(e_k | \boldsymbol{\eta}^a) \mathcal{N}(\boldsymbol{\eta} | \mathbf{m}_k^a + \mathbf{b}, \mathbf{Q}_k^{ab}) \\ &= s \sum_k \rho_k \mathbb{E}_{\boldsymbol{\xi}} \ln \left( \int \frac{d\boldsymbol{\eta} P_y(e_k | \boldsymbol{\eta})}{\sqrt{\det(2\pi \mathbf{V}_k)}} e^{-\frac{1}{2}(\boldsymbol{\eta} - \mathbf{m}_k - \mathbf{b} - \mathbf{Q}_k^{1/2} \boldsymbol{\xi})^\top \mathbf{V}_k^{-1} (\boldsymbol{\eta} - \mathbf{b} - \mathbf{m}_k - \mathbf{Q}_k^{1/2} \boldsymbol{\xi})} \right) + o(s) \\ &= s \sum_k \rho_k \mathbb{E}_{\boldsymbol{\xi}} \left[ \ln Z(e_k, \mathbf{m}_k + \mathbf{b} + \mathbf{Q}_k^{1/2} \boldsymbol{\xi}, \mathbf{V}_k) \right] + o(s), \end{aligned} \quad (152)$$

with  $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_L)$  is a normally distributed vector and we have introduced the function

$$Z(e_k, \mathbf{m}, \mathbf{V}) \equiv \int \frac{d\boldsymbol{\eta} P_y(e_k | \boldsymbol{\eta})}{\sqrt{\det(2\pi \mathbf{V})}} e^{-\frac{1}{2}(\boldsymbol{\eta} - \mathbf{m})^\top \mathbf{V}^{-1} (\boldsymbol{\eta} - \mathbf{m})} \quad (153)$$

On the other hand, denoting by  $\hat{\mathbf{V}}_k = \hat{\mathbf{R}}_k + \hat{\mathbf{Q}}_k$ ,

$$\begin{aligned} &\frac{1}{d} \ln \prod_{a=1}^s \left( \int P_w(\mathbf{W}^a) d\mathbf{W}^a \prod_k e^{-\frac{1}{2} \operatorname{tr} [\hat{\mathbf{V}}_k^\top \mathbf{W}^a \boldsymbol{\Sigma}_k (\mathbf{W}^a)^\top] + \sqrt{d} \hat{\mathbf{m}}_k^\top \mathbf{W}^a \boldsymbol{\mu}_k} \prod_{b,k} e^{\frac{1}{2} \operatorname{tr} [\hat{\mathbf{Q}}_k \mathbf{W}^a \boldsymbol{\Sigma}_k (\mathbf{W}^b)^\top]} \right) = \\ &= \frac{s}{d} \mathbb{E}_{\boldsymbol{\Xi}} \ln \left[ \int P_w(\mathbf{W}) d\mathbf{W} \prod_k \exp \left( -\frac{\operatorname{tr} [\hat{\mathbf{V}}_k^\top \mathbf{W} \boldsymbol{\Sigma}_k \mathbf{W}^\top]}{2} + \sqrt{d} \hat{\mathbf{m}}_k^\top \mathbf{W} \boldsymbol{\mu}_k + \boldsymbol{\Xi}_k \odot \sqrt{\hat{\mathbf{Q}}_k \otimes \boldsymbol{\Sigma}_k \odot \mathbf{W}} \right) \right] \\ &\quad + o(s). \end{aligned} \quad (154)$$

In the expression above we have used the tensorial product  $\hat{\mathbf{Q}} \otimes \boldsymbol{\Sigma} = (\hat{Q}_{kk'} \Sigma_{ij})_{ki, k'j'}$ . Given a matrix  $\mathbf{B} \in \mathbb{R}^{L \times d}$  and the tensors  $\mathbf{A}, \mathbf{A}' \in \mathbb{R}^{L \times d} \otimes \mathbb{R}^{L \times d}$ , we denote  $(\mathbf{B} \odot \mathbf{A})_{ki} \equiv \sum_{k'i'} B_{k'i'} A_{k'i' ki} \in \mathbb{R}^{L \times d}$ ,  $(\mathbf{A} \odot \mathbf{B})_{ki} \equiv \sum_{k'i'} A_{ki k'i'} B_{k'i'}$  and  $(\mathbf{A} \odot \mathbf{A}')_{ki k'i'} = \sum_{\kappa j} A_{ki \kappa j} A_{\kappa j k'i'}$ . In this way, we define  $\sqrt{\mathbf{A}}$  as the tensor such that  $\mathbf{A} = \sqrt{\mathbf{A}} \odot \sqrt{\mathbf{A}}$ . Finally, we have also introduced a set of  $k$  matrices  $\boldsymbol{\Xi}_k \in \mathbb{R}^{L \times d}$  with i.i.d. random Gaussian entries with zero mean and variance 1, and the average over them  $\mathbb{E}_{\boldsymbol{\Xi}}[\bullet]$ . Therefore, the (replicated) replica symmetric free-energy is given by

$$\begin{aligned} \lim_{s \rightarrow 0} \frac{\beta}{s} \Phi_{\text{RS}}^{(s)} &= \sum_{k=1}^K \hat{\mathbf{m}}_k^\top \mathbf{m}_k + \frac{1}{2} \sum_{k=1}^K \operatorname{tr} [\hat{\mathbf{V}}_k^\top \mathbf{Q}_k] - \frac{1}{2} \sum_{k=1}^K \operatorname{tr} [\hat{\mathbf{Q}}_k^\top \mathbf{V}_k] - \frac{1}{2} \sum_{k=1}^K \operatorname{tr} [\hat{\mathbf{V}}_k^\top \mathbf{V}_k] \\ &\quad - \alpha \beta \Psi_{\text{out}}(\mathbf{m}, \mathbf{Q}, \mathbf{V}) - \beta \Psi_w(\hat{\mathbf{m}}, \hat{\mathbf{Q}}, \hat{\mathbf{V}}) \end{aligned} \quad (155)$$

where we have defined two contributions

$$\Psi_{\text{out}}(\mathbf{m}, \mathbf{Q}, \mathbf{V}) \equiv \beta^{-1} \sum_k \rho_k \mathbb{E}_{\xi_k} \ln Z(e_k, \boldsymbol{\omega}_k, \mathbf{V}_k) \quad (156)$$

$$\Psi_w(\hat{\mathbf{m}}, \hat{\mathbf{Q}}, \hat{\mathbf{V}}) \equiv \frac{1}{\beta d} \mathbb{E}_{\xi} \ln \left( \int P_w(\mathbf{W}) d\mathbf{W} \prod_k e^{-\frac{\text{tr}[\hat{\mathbf{V}}_k^{\top} \mathbf{W} \boldsymbol{\Sigma}_k \mathbf{W}^{\top}]}{2} + \sqrt{d} \hat{\mathbf{m}}_k^{\top} \mathbf{W} \boldsymbol{\mu}_k + \boldsymbol{\Xi}_k \odot \sqrt{\hat{\mathbf{Q}}_k \otimes \boldsymbol{\Sigma}_k} \odot \mathbf{W}} \right) \quad (157)$$

and introduced, for future convenience,

$$\boldsymbol{\omega}_k \equiv \mathbf{m}_k + \mathbf{b} + \mathbf{Q}_k^{1/2} \boldsymbol{\xi}_k. \quad (158)$$

Note that we have separated the contribution coming from the chosen loss (the so-called *channel* part  $\Psi_{\text{out}}$ ) from the contribution depending on the regularisation (the *prior* part  $\Psi_w$ ). To write down the saddle-point equations in the  $\beta \rightarrow +\infty$  limit, let us first rescale our order parameters as  $\hat{\mathbf{m}}_k \mapsto \beta \hat{\mathbf{m}}_k$ ,  $\hat{\mathbf{Q}}_k \mapsto \beta^2 \hat{\mathbf{Q}}_k$ ,  $\hat{\mathbf{V}}_k \mapsto \beta \hat{\mathbf{V}}_k$  and  $\mathbf{V}_k \mapsto \beta^{-1} \mathbf{V}_k$ . For  $\beta \rightarrow +\infty$  the channel part is

$$\Psi_{\text{out}}(\mathbf{m}, \mathbf{Q}, \mathbf{V}) = - \sum_k \rho_k \mathbb{E}_{\xi} \left[ \mathcal{M}_{\ell(e_k, \mathbf{V}_k^{1/2} \bullet)} \left( \mathbf{V}_k^{-1/2} \boldsymbol{\omega}_k \right) \right]. \quad (159)$$

Here and in the following the quantity

$$\mathcal{M}_{f(\bullet)}(\mathbf{u}) \equiv \min_{\mathbf{v} \in \text{domain}(\mathbf{v})} \left[ \frac{1}{2} \|\mathbf{v} - \mathbf{u}\|_{\text{F}}^2 + f(\mathbf{v}) \right] \quad (160)$$

is the Moreau envelope of  $f: \text{domain}(\mathbf{v}) \rightarrow \mathbb{R}$ , whereas  $\|\bullet\|_{\text{F}}$  is the Frobenius norm. We can write the contribution  $\Psi_{\text{out}}$  in terms of a proximal

$$\mathbf{h}_k = \mathbf{V}_k^{1/2} \text{Prox}_{\ell(e_k, \mathbf{V}_k^{1/2} \bullet)}(\mathbf{V}_k^{-1/2} \boldsymbol{\omega}_k) \equiv \mathbf{V}_k^{1/2} \arg \min_{\mathbf{u} \in \mathbb{R}^L} \left[ \frac{1}{2} \|\mathbf{u} - \mathbf{V}_k^{-1/2} \boldsymbol{\omega}_k\|_{\text{F}}^2 + \ell(e_k, \mathbf{V}_k^{1/2} \mathbf{u}) \right]. \quad (161)$$

as

$$\Psi_{\text{out}}(\mathbf{m}, \mathbf{Q}, \mathbf{V}) = - \sum_k \rho_k \mathbb{E}_{\xi} \left[ \frac{1}{2} \|\mathbf{V}_k^{-1/2} \mathbf{h}_k - \mathbf{V}_k^{-1/2} \boldsymbol{\omega}_k\|_{\text{F}}^2 + \ell(e_k, \mathbf{h}_k) \right] \quad (162)$$

A similar expression can be obtained for  $\Psi_w$ . Defining

$$\mathbf{A} = \left( \sum_k \hat{\mathbf{V}}_k \otimes \boldsymbol{\Sigma}_k \right)^{-1}, \quad \mathbf{B} = \sqrt{d} \sum_k \boldsymbol{\mu}_k \hat{\mathbf{m}}_k^{\top} + \sum_k \boldsymbol{\Xi}_k \odot \sqrt{\hat{\mathbf{Q}}_k \otimes \boldsymbol{\Sigma}_k}. \quad (163)$$

$\Psi_w$  can be written as

$$\begin{aligned} \Psi_w(\hat{\mathbf{m}}, \hat{\mathbf{Q}}, \hat{\mathbf{V}}) &= \frac{1}{2d} \mathbb{E}_{\xi} [ \mathbf{B} \odot \mathbf{A} \odot \mathbf{B} ] \\ &+ \frac{1}{\beta d} \mathbb{E}_{\xi} \ln \left[ \int d\mathbf{W} \exp \left( -\frac{\beta}{2} \|\mathbf{A}^{-1/2} \odot \mathbf{W} - \mathbf{A}^{1/2} \odot \mathbf{B}\|_{\text{F}}^2 - \beta r(\mathbf{W}) \right) \right]. \end{aligned} \quad (164)$$

It follows that, for  $\beta \rightarrow +\infty$ ,

$$\Psi_w(\hat{\mathbf{m}}, \hat{\mathbf{Q}}, \hat{\mathbf{V}}) = \frac{1}{2d} \mathbb{E}_{\xi} [ \mathbf{B} \odot \mathbf{A} \odot \mathbf{B} ] - \frac{1}{d} \mathbb{E}_{\xi} \left[ \mathcal{M}_{r(\mathbf{A}^{1/2} \odot \bullet)}(\mathbf{A}^{1/2} \odot \mathbf{B}) \right]. \quad (165)$$

As before, let us introduce the proximal

$$\mathbf{G} = \mathbf{A}^{1/2} \odot \text{Prox}_{r(\mathbf{A}^{1/2} \odot \bullet)}(\mathbf{A}^{1/2} \odot \mathbf{B}) \in \mathbb{R}^{L \times d} \quad (166)$$

We can rewrite the prior contribution  $\Psi_w$  as

$$\Psi_w(\hat{\mathbf{m}}, \hat{\mathbf{Q}}, \hat{\mathbf{V}}) = \frac{1}{2d} \mathbb{E}_{\xi} [ \mathbf{B} \odot \mathbf{A} \odot \mathbf{B} ] - \frac{1}{d} \mathbb{E}_{\xi} \left[ \frac{\|\mathbf{A}^{-1/2} \odot \mathbf{G} - \mathbf{A}^{1/2} \odot \mathbf{B}\|_{\text{F}}^2}{2} + r(\mathbf{G}) \right]. \quad (167)$$

The parallelism between the two contributions is evident, aside from the different dimensionality of the involved objects. The replica symmetric free energy in the  $\beta \rightarrow +\infty$  limit is computed extremising with respect to the introduced order parameters,

$$f_{\text{RS}} = \text{Extr}_{\substack{\mathbf{m}, \mathbf{Q}, \mathbf{V}, \mathbf{b} \\ \hat{\mathbf{m}}, \hat{\mathbf{Q}}, \hat{\mathbf{V}}}} \left[ \sum_{k=1}^K \hat{\mathbf{m}}_k^\top \mathbf{m}_k + \frac{1}{2} \sum_{k=1}^K \text{tr} [\hat{\mathbf{V}}_k^\top \mathbf{Q}_k] - \frac{1}{2} \sum_{k=1}^K \text{tr} [\hat{\mathbf{Q}}_k^\top \mathbf{V}_k] \right. \\ \left. - \frac{1}{2} \sum_{k=1}^K \text{tr} [\hat{\mathbf{V}}_k^\top \mathbf{V}_k] - \alpha \Psi_{\text{out}}(\mathbf{m}, \mathbf{Q}, \mathbf{V}) - \Psi_w(\hat{\mathbf{m}}, \hat{\mathbf{Q}}, \hat{\mathbf{V}}) \right]. \quad (168)$$

To do so, we have to write down a set of saddle-point equations and solve them.

**Saddle-point equations** — The saddle-point equations are derived straightforwardly from the obtained free energy extremising with respect to all parameters. A first set of equations is obtained from  $\Psi_{\text{out}}$  as<sup>1</sup>

$$\hat{\mathbf{Q}}_k = \alpha \rho_k \mathbb{E}_\xi [\mathbf{f}_k \mathbf{f}_k^\top], \quad (169a)$$

$$\hat{\mathbf{V}}_k = -\alpha \rho_k \mathbf{Q}_k^{-1/2} \mathbb{E}_\xi [\mathbf{f}_k \xi^\top], \quad (169b)$$

$$\hat{\mathbf{m}}_k = \alpha \rho_k \mathbb{E}_\xi [\mathbf{f}_k], \quad (169c)$$

$$\mathbf{b} = \sum_k \rho_k \mathbb{E}_\xi [\mathbf{h}_k - \mathbf{m}_k] \iff \sum_k \rho_k \mathbb{E}_\xi [\mathbf{V}_k \mathbf{f}_k] = \mathbf{0}. \quad (169d)$$

where for brevity we have denoted

$$\mathbf{f}_k \equiv \mathbf{V}_k^{-1}(\mathbf{h}_k - \boldsymbol{\omega}_k). \quad (170)$$

Similarly, the saddle-point equations from  $\Psi_{\text{out}}$  are

$$\mathbf{V}_k = \frac{1}{d} \mathbb{E}_\Xi \left[ \left( \mathbf{G} \odot \left( \hat{\mathbf{Q}}_k \otimes \boldsymbol{\Sigma}_k \right)^{-1/2} \odot \left( \mathbf{I}_k \otimes \boldsymbol{\Sigma}_k \right) \right) \boldsymbol{\Xi}_k^\top \right] \quad (171a)$$

$$\mathbf{Q}_k = \frac{1}{d} \mathbb{E}_\xi [\mathbf{G} \boldsymbol{\Sigma}_k \mathbf{G}^\top] \quad (171b)$$

$$\mathbf{m}_k = \frac{1}{\sqrt{d}} \mathbb{E}_\xi [\mathbf{G} \boldsymbol{\mu}_k]. \quad (171c)$$

To obtain the replica symmetric free energy, therefore, the given set of equation has to be solved, and the result then plugged in Eq. (168). No further simplification can be obtained in the most general setting. We will explore however some simple (but important) applications in Appendix C. Before going on, however, it is important to express the relevant quantities for learning, i.e., the training and generalization errors, in terms of the obtained order parameters.

#### B.4 Training and test errors

The order parameters introduced to solve the problem allow us to reach our ultimate goal of computing the average errors of the learning process. We will start from the estimation of the training loss. The complication in computing this quantity is that the order parameters found in the learning process are, of course, correlated with the dataset used for the learning itself. We need to compute

$$\epsilon_\ell \equiv \frac{1}{n} \sum_{\nu=1}^n \ell \left( \mathbf{y}^\nu, \frac{\mathbf{W}^* \mathbf{x}^\nu}{\sqrt{d}} + \mathbf{b}^* \right) \quad (172)$$

in the  $n \rightarrow +\infty$  limit. Denoting for brevity  $\ell_k(\mathbf{x}) \equiv \ell(e_k, \mathbf{x})$ , the best way to proceed is to observe that  $\mathbb{E}_{\{(\mathbf{y}^\nu, \mathbf{x}^\nu)\}_\nu} [\mathcal{R}(\mathbf{W}^*, \mathbf{b}^*)] = -\lim_{\beta \rightarrow +\infty} \mathbb{E}_{\{(\mathbf{y}^\nu, \mathbf{x}^\nu)\}_\nu} [\partial_\beta \ln \mathcal{Z}_\beta] = \lambda \mathbb{E}_{\{(\mathbf{y}^\nu, \mathbf{x}^\nu)\}_\nu} [r(\mathbf{W}^*)] + \epsilon_\ell$ , where

$$\epsilon_\ell = -\lim_{\beta \rightarrow +\infty} \partial_\beta (\beta \Psi_{\text{out}}) = \lim_{\beta \rightarrow +\infty} \sum_k \rho_k \int \ell_k(\boldsymbol{\eta}) \frac{e^{-\frac{\beta}{2}(\boldsymbol{\eta} - \mathbf{m}_k^*)^\top \mathbf{V}_k^{*-1}(\boldsymbol{\eta} - \mathbf{m}_k^*) - \beta \ell_k(\boldsymbol{\eta})}}{\sqrt{\det(2\pi\beta^{-1}\mathbf{V}^*)} Z(e_k, \boldsymbol{\omega}_k^*, \beta^{-1}\mathbf{V}_k^*)} d\boldsymbol{\eta}. \quad (173)$$

<sup>1</sup>To obtain the equation for  $\hat{\mathbf{V}}$  it is convenient to use Stein's lemma, so that  $\mathbb{E}[\partial_\xi \mathbf{f}_k] = \mathbb{E}[\mathbf{f}_k \xi^\top]$ .

In the  $\beta \rightarrow +\infty$  limit, the integral concentrates on the minimizer of the exponent, that is, by definition, the proximal  $\mathbf{h}_k$ . In conclusion,  $\epsilon_\ell = \sum_k \rho_k \mathbb{E}[\ell(\mathbf{h}_k)]$ . By means of the same concentration result, the training error is

$$\epsilon_t = \frac{1}{n} \sum_{\nu=1}^n \mathbb{I} \left( \varphi \left( \frac{\mathbf{W}^* \mathbf{x}^\nu}{\sqrt{d}} + \mathbf{b}^* \right) \neq \mathbf{y}^\nu \right) \xrightarrow{n \rightarrow +\infty} \sum_{k=1}^K \rho_k \mathbb{E}_\xi [\mathbb{I}(\varphi(\mathbf{h}_k) \neq \mathbf{e}_k)]. \quad (174)$$

The expressions above hold in general, but, as anticipated, important simplifications can occur in the set of saddle-point equations (169) and (171) depending on the choice of the loss  $\ell$  and of the regularization function  $r$ .

The generalisation (or test) error can be written instead as

$$\epsilon_g = \mathbb{E}_{\mathbf{y}^{\text{new}}, \mathbf{x}^{\text{new}}} \left[ \mathbb{I} \left( \varphi \left( \frac{\mathbf{W}^* \mathbf{x}^{\text{new}}}{\sqrt{d}} + \mathbf{b}^* \right) \neq \mathbf{y}^{\text{new}} \right) \right]. \quad (175)$$

This expression can be rewritten as

$$\epsilon_g = \sum_k \rho_k \int \mathbb{I}(\varphi(\boldsymbol{\eta}) = \mathbf{e}_k) \mathbb{E}_{\mathbf{x}^{\text{new}}} \left[ \delta \left( \boldsymbol{\eta} - \frac{\mathbf{W}^* \mathbf{x}^{\text{new}}}{\sqrt{d}} - \mathbf{b}^* \right) \right] d\boldsymbol{\eta} \quad (176)$$

Once again, we write

$$\mathbb{E}_{\mathbf{x}^{\text{new}}} \left[ \delta \left( \boldsymbol{\eta} - \frac{\mathbf{W}^* \mathbf{x}^{\text{new}}}{\sqrt{d}} - \mathbf{b}^* \right) \right] \xrightarrow{d \rightarrow +\infty} \mathcal{N}(\boldsymbol{\eta} | \mathbf{m}_k^* + \mathbf{b}^*, \mathbf{Q}_k^*) \quad (177)$$

so that

$$\epsilon_g = \sum_{k=1}^K \rho_k \mathbb{E}_\xi \left[ \mathbb{I} \left( \varphi \left( \mathbf{m}_k^* + \mathbf{Q}_k^{*1/2} \boldsymbol{\xi} + \mathbf{b}^* \right) \neq \mathbf{e}_k \right) \right]. \quad (178)$$

This can be easily computed numerically once that the order parameters are given.

## B.5 A note on the numerical integration of the saddle-point equations

To estimate  $\epsilon_g$ ,  $\epsilon_t$  and  $\epsilon_\ell$  we first need to find the fixed-point solutions of the saddle-point equations (169) and (171). The simplest numerical strategy consists in updating, in a self-consistent way, the order parameters until their variation according to, e.g., the Frobenius norm is smaller than a given threshold value (that we adopted to be  $10^{-5}$ ). In the simplest setting, i.e., the one discussed in Corollary 3, the update of  $(\mathbf{m}_k, \mathbf{Q}_k, \mathbf{V}_k)_{k \in [K]}$  is performed explicitly using eq. (11), where  $\mathbb{E}_{\sigma, \mu}[\bullet]$  is a shorthand for the sum over the eigenvalues and eigenvectors of the assigned covariance matrices. The update of  $(\hat{\mathbf{m}}_k, \hat{\mathbf{Q}}_k, \hat{\mathbf{V}}_k)_{k \in [K]}$  (right hand side of eq. (8)) is more involved, as it requires the computation of the proximal followed by a Gaussian average. Such average has been performed using a Monte Carlo strategy, i.e., by solving the equation for the proximal for a large number ( $10^4 - 10^5$ ) of instances of  $\boldsymbol{\xi}$  and averaging the solution. We remark that in the case of the square loss, the proximal can be computed analytically and the integration can be performed explicitly, highly simplifying the fixed-point equations (see below eq. (191)). We have found that in practice fluctuations due to the adopted Monte Carlo pool were small enough to be negligible compared with the outcomes of direct numerical experiments.

The convergence to the the correct fixed point is guaranteed (in principle) by the convexity of the problem. However, a few delicate aspects have to be taken into account in the update process described above.

1. The update requires the computation of the proximals  $\mathbf{G}$  and  $\mathbf{h}_k$ . Such computations can be performed analytically in some specific cases only (for example, in the case of ridge regression). The existence of a unique solution is guaranteed by the strong convexity of the problem defining the proximal. In our study of the cross-entropy loss function, for example, we computed the proximals  $\mathbf{h}_k$  numerically solving Eq. (194). In this problem, however, additional numerical instabilities emerged in the  $\lambda \rightarrow 0$  limit, due the fact that the discontinuity in the gradient appear, see Eq. (198). We solved this issue performing an annealing in  $\lambda$ , i.e., solving for the proximal for decreasing values of the regularization strength.

2. The numerical solution of the saddle-point equations might suffer numerical instabilities due to the operations of inversion involved, see, e.g., the equation for  $\hat{\mathbf{V}}_k$  in (169), which requires the inversion of  $\mathbf{Q}_k$ . It is convenient, in such cases, to rewrite the equation in an equivalent form which is numerically more stable. For example, in the aforementioned equation, we can observe that  $\mathbf{f}_k$  satisfies the equation  $\mathbf{f}_k + \partial_{\mathbf{x}} \ell_k(\mathbf{V}_k \mathbf{f}_k + \boldsymbol{\omega}_k) = \mathbf{0}$  so that  $\partial_{\boldsymbol{\omega}_k} \mathbf{f}_k = -(\mathbf{I}_K + \partial_{\mathbf{x}}^2 \ell_k(\mathbf{V}_k \mathbf{f}_k + \boldsymbol{\omega}_k) \mathbf{V}_k)^{-1} \partial_{\mathbf{x}}^2 \ell_k(\mathbf{V}_k \mathbf{f}_k + \boldsymbol{\omega}_k)$ . Using Stein's lemma,

$$\hat{\mathbf{V}}_k = -\alpha \rho_k \mathbb{E}_{\boldsymbol{\xi}} [\partial_{\boldsymbol{\xi}} \mathbf{f}_k] = \alpha \rho_k \mathbb{E}_{\boldsymbol{\xi}} \left[ (\mathbf{I}_K + \partial_{\mathbf{x}}^2 \ell_k(\mathbf{V}_k \mathbf{f}_k + \boldsymbol{\omega}_k) \mathbf{V}_k)^{-1} \partial_{\mathbf{x}}^2 \ell_k(\mathbf{V}_k \mathbf{f}_k + \boldsymbol{\omega}_k) \right]. \quad (179)$$

We found this equation numerically more stable than the one given in (169) when dealing with the cross-entropy loss.

Our implementation can be found at [58].

## C Some relevant particular cases

In this Appendix, we will specify the saddle-point equations for the multiclass classification problem for different choices of the loss function  $\ell$  and of the regularisation function  $r$ . From the analysis developed in the previous Appendices, it is clear that the choices of  $\ell$  and  $r$  impact separately the set of equations (169) and (171) respectively. Once the order parameters are found, it is possible to estimate the training and generalisation errors as, for example, in Section B.4.

### C.1 The case of $\ell_2$ regularization

In this Section we consider the relevant case of quadratic regularization,  $r(\mathbf{W}) = 1/2 \|\mathbf{W}\|_F^2$ . In this case the computation of  $\Psi_w$  can be performed explicitly via a Gaussian integration,

$$\frac{1}{\beta} \Psi_w(\hat{\mathbf{m}}, \hat{\mathbf{Q}}, \hat{\mathbf{V}}) = \frac{1}{2d} \operatorname{tr} \ln \mathbf{S} - \frac{K \ln \beta}{2\beta} + \frac{1}{2} \operatorname{tr} \left[ \mathbf{S} \odot \left( \sum_{kk'} \hat{\mathbf{m}}_k \hat{\mathbf{m}}_{k'}^\top \otimes \boldsymbol{\mu}_k \boldsymbol{\mu}_{k'}^\top + \frac{1}{d} \sum_k \hat{\mathbf{Q}}_k \otimes \boldsymbol{\Sigma}_k \right) \right]. \quad (180)$$

Here we have introduced, for notation compactness,

$$\mathbf{S} \equiv \left( \lambda \mathbf{I}_K \otimes \mathbf{I}_d + \sum_{\kappa} \hat{\mathbf{V}}_{\kappa} \otimes \boldsymbol{\Sigma}_{\kappa} \right)^{-1}. \quad (181)$$

This form of  $\Psi_w$  allows us to write in a simpler way the set of eqs. (171), that can be re-written as

$$\begin{aligned} \mathbf{Q}_k &= \operatorname{tr}_d \left[ (\mathbf{I}_K \otimes \boldsymbol{\Sigma}_k) \odot \mathbf{S} \odot \left( \sum_{kk'} \hat{\mathbf{m}}_k \hat{\mathbf{m}}_{k'}^\top \otimes \boldsymbol{\mu}_k \boldsymbol{\mu}_{k'}^\top + \frac{1}{d} \sum_{\kappa} \hat{\mathbf{Q}}_{\kappa} \otimes \boldsymbol{\Sigma}_{\kappa} \right) \odot \mathbf{S} \right] \\ \mathbf{m}_k &= \sum_{k'} \operatorname{tr}_d \left[ \mathbf{S} \odot (\hat{\mathbf{m}}_{k'} \otimes \boldsymbol{\mu}_{k'} \boldsymbol{\mu}_k^\top) \right] \\ \mathbf{V}_k &= \frac{1}{d} \operatorname{tr}_d [(\mathbf{I}_K \otimes \boldsymbol{\Sigma}_k) \odot \mathbf{S}]. \end{aligned} \quad (182)$$

In the previous equations, by  $\operatorname{tr}_d$  we denoted the trace with respect to the components living in the  $d$ -dimensional space of the dataset.

**Jointly diagonal covariances** — Suppose now that  $\boldsymbol{\Sigma}_k = \sum_i \sigma_i^k \mathbf{v}_i \mathbf{v}_i^\top$  for all  $k$ , i.e., the covariance matrices share the same basis of eigenvectors  $\{\mathbf{v}_i\}_i$ . Then, denoting  $\mu_i^k \equiv \sqrt{d} \boldsymbol{\mu}_k^\top \mathbf{v}_i$

$$\begin{aligned} \mathbf{Q}_k &= \frac{1}{d} \sum_{i=1}^d \sigma_i^k \left( \lambda \mathbf{I}_K + \sum_{\kappa} \sigma_i^{\kappa} \hat{\mathbf{V}}_{\kappa} \right)^{-1} \left( \sum_{kk'} \mu_i^k \mu_i^{k'} \hat{\mathbf{m}}_k \hat{\mathbf{m}}_{k'}^\top + \sum_{\kappa} \sigma_i^{\kappa} \hat{\mathbf{Q}}_{\kappa} \right) \left( \lambda \mathbf{I}_K + \sum_{\kappa} \sigma_i^{\kappa} \hat{\mathbf{V}}_{\kappa} \right)^{-1} \\ \mathbf{m}_k &= \frac{1}{d} \sum_{i=1}^d \sum_{k'} \mu_i^k \mu_i^{k'} \left( \lambda \mathbf{I}_K + \sum_{\kappa} \sigma_i^{\kappa} \hat{\mathbf{V}}_{\kappa} \right)^{-1} \hat{\mathbf{m}}_{k'} \\ \mathbf{V}_k &= \frac{1}{d} \sum_{i=1}^d \sigma_i^k \left( \lambda \mathbf{I}_K + \sum_{\kappa} \sigma_i^{\kappa} \hat{\mathbf{V}}_{\kappa} \right)^{-1}. \end{aligned} \quad (183)$$

Introducing the joint density

$$\frac{1}{d} \sum_{i=1}^d \prod_{\kappa=1}^K \delta(\sigma^{\kappa} - \sigma_i^{\kappa}) \delta(\mu^{\kappa} - \mu_i^{\kappa}) \xrightarrow{d \rightarrow +\infty} \rho(\boldsymbol{\sigma}, \boldsymbol{\mu}), \quad (184)$$

then we can write the saddle-point equations given in Corollary 3

$$\begin{aligned}
\mathbf{Q}_k &= \mathbb{E}_{\sigma, \mu} \left[ \sigma^k \left( \lambda \mathbf{I}_K + \sum_{\kappa} \sigma^{\kappa} \hat{\mathbf{V}}_{\kappa} \right)^{-1} \left( \sum_{kk'} \mu^k \mu^{k'} \hat{\mathbf{m}}_k \hat{\mathbf{m}}_{k'}^{\top} + \sum_{\kappa} \sigma^{\kappa} \hat{\mathbf{Q}}_{\kappa} \right) \left( \lambda \mathbf{I}_K + \sum_{\kappa} \sigma^{\kappa} \hat{\mathbf{V}}_{\kappa} \right)^{-1} \right] \\
\mathbf{m}_k &= \mathbb{E}_{\sigma, \mu} \left[ \mu^k \left( \lambda \mathbf{I}_K + \sum_{\kappa} \sigma^{\kappa} \hat{\mathbf{V}}_{\kappa} \right)^{-1} \sum_{\kappa} \mu^{\kappa} \hat{\mathbf{m}}_{\kappa} \right] \\
\mathbf{V}_k &= \mathbb{E}_{\sigma, \mu} \left[ \sigma^k \left( \lambda \mathbf{I}_K + \sum_{\kappa} \sigma^{\kappa} \hat{\mathbf{V}}_{\kappa} \right)^{-1} \right].
\end{aligned} \tag{185}$$

where the expectations  $\mathbb{E}_{\sigma, \mu}$  are taken with respect to the joint distribution  $\rho$ .

### C.1.1 Uniform covariances

Let us consider the simpler case  $\Sigma_k \equiv \Delta \mathbf{I}_d$ , with  $\Delta > 0$ . In this case, the saddle-point equations can take a more compact form that is particularly suitable for a numerical solution. Moreover, for reasons of symmetry we can write

$$\mathbf{Q}_k \equiv \mathbf{Q}, \quad \mathbf{V}_k \equiv \mathbf{V}, \quad \hat{\mathbf{Q}}_k \equiv \frac{1}{K\Delta} \hat{\mathbf{Q}}_k, \quad \hat{\mathbf{V}}_k \equiv \frac{1}{K\Delta} \hat{\mathbf{V}}, \quad \forall k. \tag{186}$$

Let us define the following  $K \times K$  matrices

- $\mathbf{M} \in \mathbb{R}^{K \times K}$  (resp.  $\hat{\mathbf{M}} \in \mathbb{R}^{K \times K}$ ) is the matrix obtained concatenating the vectors  $\mathbf{m}_k$  (resp.  $\hat{\mathbf{m}}_k$ );
- $\Theta = (\boldsymbol{\mu}_k^{\top} \boldsymbol{\mu}_{k'})_{kk'}$  is the Gram matrix of the means;
- $\mathbf{F} \in \mathbb{R}^{K \times K}$  is the matrix obtained concatenating the vectors  $\mathbf{f}_k$ ;
- $\mathbf{H} \in \mathbb{R}^{K \times K}$  is the matrix obtained concatenating the vectors  $\mathbf{h}_k$ ;
- $\mathbf{\Pi} = \text{diag}(\rho_k) \in \mathbb{R}^{K \times K}$  is a diagonal matrix with elements  $\Pi_{kk'} = \delta_{kk'} \rho_k$ .

The saddle-point equations then can be rewritten as

$$\begin{aligned}
\mathbf{Q} &= \Delta \left( \lambda \mathbf{I}_K + \hat{\mathbf{V}} \right)^{-1} \left( \hat{\mathbf{Q}} + \hat{\mathbf{M}} \Theta \hat{\mathbf{M}}^{\top} \right) \left( \lambda \mathbf{I}_K + \hat{\mathbf{V}} \right)^{-1} & \hat{\mathbf{Q}} &= \alpha \Delta \mathbb{E}_{\Xi} [\mathbf{F} \mathbf{\Pi} \mathbf{F}^{\top}] \\
\mathbf{M} &= \left( \lambda \mathbf{I}_K + \hat{\mathbf{V}} \right)^{-1} \hat{\mathbf{M}} \Theta & \hat{\mathbf{V}} &= -\alpha \Delta \mathbf{Q}^{-1/2} \mathbb{E}_{\Xi} [\mathbf{F} \mathbf{\Pi} \Xi^{\top}] \\
\mathbf{V} &= \Delta \left( \lambda \mathbf{I}_K + \hat{\mathbf{V}} \right)^{-1}, & \hat{\mathbf{M}} &= \alpha \mathbb{E}_{\Xi} [\mathbf{F} \mathbf{\Pi}] \\
& & \mathbf{b} &= \mathbb{E}_{\Xi} [(\mathbf{H} - \mathbf{M}) \mathbf{\Pi} \mathbf{1}_K].
\end{aligned} \tag{187}$$

Here and in the following  $\mathbf{1}_K$  is the vector of  $K$  components all equal to 1. These expressions are particularly suitable for a numerical implementation, because involve matrix multiplications and inversions of  $K$ -dimensional objects only.

**Quadratic loss** — If we consider a quadratic loss  $\ell(\mathbf{y}, \mathbf{x}) = \frac{1}{2} (\mathbf{y} - \mathbf{x})^2$ , then an explicit formula for the proximal can be found, namely

$$\mathbf{f}_k = (\mathbf{I}_K + \mathbf{V})^{-1} (\mathbf{e}_K - \boldsymbol{\omega}_k) \tag{188}$$

so that the second set of saddle-point equations (187) can be written as

$$\begin{aligned}
\hat{\mathbf{Q}} &= \alpha (\mathbf{I}_K + \mathbf{V})^{-1} [(\mathbf{I}_K - \mathbf{M} - \mathbf{b} \otimes \mathbf{1}_K) \mathbf{\Pi} (\mathbf{I}_K - \mathbf{M} - \mathbf{b} \otimes \mathbf{1}_K)^{\top} + \mathbf{Q}] (\mathbf{I}_K + \mathbf{V})^{-1} \\
\hat{\mathbf{M}} &= \alpha (\mathbf{I}_K + \mathbf{V})^{-1} (\mathbf{I}_K - \mathbf{M} - \mathbf{b} \otimes \mathbf{1}_K) \mathbf{\Pi} \\
\hat{\mathbf{V}} &= \alpha \Delta (\mathbf{I}_K + \mathbf{V})^{-1}.
\end{aligned} \tag{189}$$

Observe at this point that we can explicitly solve for  $\mathbf{V}$  using the equation for it in eqs. (187). In particular,  $\mathbf{V}$  satisfies the equation  $\lambda \mathbf{V}^2 + (\alpha + \lambda - \Delta) \mathbf{V} = \Delta \mathbf{I}_K$ . Being  $\mathbf{V}$  positive definite, it follows that it is diagonal,  $\mathbf{V} = V \mathbf{I}_K$  with diagonal element

$$V = \frac{\Delta(1 - \alpha) - \lambda + \sqrt{(\Delta - \alpha\Delta - \lambda)^2 + 4\Delta\lambda}}{2\lambda}, \quad \hat{V} = \frac{\alpha\Delta}{1 + V}, \tag{190}$$

so that

$$\begin{aligned} \mathbf{Q} &= \frac{\Delta}{(\lambda + \Delta \hat{\mathbf{V}})^2} \left( \hat{\mathbf{Q}} + \hat{\mathbf{M}} \Theta \hat{\mathbf{M}}^\top \right) & \hat{\mathbf{Q}} &= \frac{\alpha [(\mathbf{I}_K - \mathbf{M} - \mathbf{b} \otimes \mathbf{1}_K) \mathbf{\Pi} (\mathbf{I}_K - \mathbf{M} - \mathbf{b} \otimes \mathbf{1}_K)^\top + \mathbf{Q}]}{(1+V)^2} \\ \mathbf{M} &= \frac{\hat{\mathbf{M}} \Theta}{\lambda + \Delta \hat{\mathbf{V}}}, & \hat{\mathbf{M}} &= -\frac{\alpha (\mathbf{I}_K - \mathbf{M} - \mathbf{b} \otimes \mathbf{1}_K) \mathbf{\Pi}}{1+V}. \\ \mathbf{b} &= (\mathbf{I}_K - \mathbf{M}) \mathbf{\Pi} \mathbf{1}_K, \end{aligned} \quad (191)$$

In the  $\lambda \rightarrow 0$  limit, for  $\alpha < 1$  it is convenient to rescale  $\hat{\mathbf{Q}} \mapsto \lambda^2 \hat{\mathbf{Q}}$  and  $\hat{\mathbf{M}} \mapsto \lambda \hat{\mathbf{M}}$ , so that

$$\begin{aligned} \mathbf{Q} &= \Delta (1-\alpha)^2 \left( \hat{\mathbf{Q}} + \hat{\mathbf{M}} \Theta \hat{\mathbf{M}}^\top \right), & \hat{\mathbf{Q}} &= \frac{\alpha [(\mathbf{I}_K - \mathbf{M} - \mathbf{b} \otimes \mathbf{1}_K) \mathbf{\Pi} (\mathbf{I}_K - \mathbf{M} - \mathbf{b} \otimes \mathbf{1}_K)^\top + \mathbf{Q}]}{\Delta^2 (1-\alpha)^2}, \\ \mathbf{M} &= (1-\alpha) \hat{\mathbf{M}} \Theta, & \hat{\mathbf{M}} &= -\frac{\alpha (\mathbf{I}_K - \mathbf{M} - \mathbf{b} \otimes \mathbf{1}_K) \mathbf{\Pi}}{\Delta (1-\alpha)}. \\ \mathbf{b} &= (\mathbf{I}_K - \mathbf{M}) \mathbf{\Pi} \mathbf{1}_K, \end{aligned} \quad (192)$$

**Cross-entropy loss** — We consider now the relevant case of the cross entropy loss

$$\ell(\mathbf{y}, \mathbf{x}) = -\sum_{k=1}^K y_k \ln \frac{e^{x_k}}{\sum_{\kappa=1}^K e^{x_\kappa}}. \quad (193)$$

If  $\mathbf{y} \in \{\mathbf{e}_k\}_{k \in [K]}$ , the loss can be written in the form  $\ell(\mathbf{y}, \mathbf{x}) = -\mathbf{y}^\top \mathbf{x} + \ln \sum_{\kappa} e^{x_\kappa}$ . If we introduce the *softmax function*  $\mathbf{soft}: \mathbb{R}^K \rightarrow \mathbb{R}^K$

$$\partial_{\mathbf{x}} \ell(\mathbf{y}, \mathbf{x}) = -\mathbf{y} + \mathbf{soft}(\mathbf{x}), \quad \mathbf{soft}_k(\mathbf{x}) \equiv \frac{\exp(x_k)}{\sum_{\kappa} \exp(x_\kappa)} \quad (194)$$

the proximal equation for the cross-entropy loss is the solution of the equations:

$$\mathbf{V}^{-1}(\mathbf{h}_k - \boldsymbol{\omega}_k) - \mathbf{e}_k + \mathbf{soft}(\mathbf{h}_k) = \mathbf{0} \iff \mathbf{f}_k = \mathbf{e}_k - \mathbf{soft}(\mathbf{V} \mathbf{f}_k + \boldsymbol{\omega}_k) \quad \forall k \in [K], \quad (195)$$

having only one solution for which, however, there is no closed-form expression. The equation can be solved numerically, and in this way we obtained the results in Section 3.2.

The saddle-point equations can be written rescaling  $\mathbf{Q} \mapsto \lambda^{-2} \mathbf{Q}$ ,  $\mathbf{V} \mapsto \lambda^{-1} \mathbf{V}$ ,  $\mathbf{M} \mapsto \lambda^{-1} \mathbf{M}$ ,  $\mathbf{b} \mapsto \lambda^{-1} \mathbf{b}$ ,  $\hat{\mathbf{V}} \mapsto \lambda \hat{\mathbf{V}}$ . They become

$$\begin{aligned} \mathbf{Q} &= \Delta \left( \mathbf{I}_K + \hat{\mathbf{V}} \right)^{-1} \left( \hat{\mathbf{Q}} + \hat{\mathbf{M}} \Theta \hat{\mathbf{M}}^\top \right) \left( \mathbf{I}_K + \hat{\mathbf{V}} \right)^{-1}, & \hat{\mathbf{Q}} &= \alpha \Delta \mathbb{E}_{\Xi} [\mathbf{F} \mathbf{\Pi} \mathbf{F}^\top], \\ \mathbf{M} &= \left( \mathbf{I}_K + \hat{\mathbf{V}} \right)^{-1} \hat{\mathbf{M}} \Theta & \hat{\mathbf{V}} &= -\alpha \Delta \mathbf{Q}^{-1/2} \mathbb{E}_{\Xi} [\mathbf{F} \mathbf{\Pi} \Xi^\top], \\ \mathbf{V} &= \Delta \left( \mathbf{I}_K + \hat{\mathbf{V}} \right)^{-1}, & \hat{\mathbf{M}} &= \alpha \mathbb{E}_{\Xi} [\mathbf{F} \mathbf{\Pi}], \\ & & \mathbf{b} &= \mathbb{E}_{\Xi} [(\mathbf{H} - \mathbf{M}) \mathbf{\Pi}], \end{aligned} \quad (196)$$

so that the dependence on  $\lambda$  disappears everywhere except in the equation for the proximal  $\mathbf{f}_k$

$$\mathbf{f}_k = \arg \min_{\mathbf{x}} \left[ \frac{1}{2} \mathbf{x}^\top \mathbf{V} \mathbf{x} + \lambda \ell \left( \mathbf{e}_k, \frac{\mathbf{V} \mathbf{x} + \boldsymbol{\omega}_k}{\lambda} \right) \right], \quad (197)$$

which, in the  $\lambda \rightarrow 0$  limit, becomes

$$\mathbf{f}_k = \arg \min_{\mathbf{x}} \left[ \frac{1}{2} \mathbf{x}^\top \mathbf{V} \mathbf{x} + \min_{\mu} \{ (\mathbf{e}_\mu - \mathbf{e}_k)^\top (\mathbf{V} \mathbf{x} + \boldsymbol{\omega}_k) \} \right]. \quad (198)$$

Note that in this limit, minimising the cross-entropy loss yields precisely the max-margin estimator [70].

## C.2 The $K = 2$ case with scalar labels

The formulas for the  $K = 2$  case can be derived directly from the general analysis given above imposing  $L = 1$ . In particular, let us assume that the two clusters are labeled with  $e_1 = +1$  and  $e_2 = -1$ . Using as classifier

$$\varphi(x) = \text{sign}(x) \quad (199)$$

the expression of the average errors is

$$\begin{aligned}
\epsilon_g &= \sum_{k \in [2]} \rho_k \mathbb{E}_\xi [\theta((-1)^k \omega_k^*)] = \sum_{k \in [2]} \frac{\rho_k}{2} \operatorname{erfc} \left( (-1)^{k-1} \frac{m_k^* + b^*}{\sqrt{2q_k^*}} \right), \\
\epsilon_t &= \sum_{k \in [2]} \rho_k \mathbb{E}_\xi [\theta((-1)^k h_k^*)], \\
\epsilon_\ell &= \sum_{k \in [2]} \rho_k \mathbb{E}_\xi [\ell((-1)^k, h_k^*)].
\end{aligned} \tag{200}$$

We will further explore this case, considering some special cases in the following.

### C.2.1 Example: $\ell_1$ regularization

In this Section we derive the saddle-point equations for the case in which the two clusters have opposite means  $\boldsymbol{\mu}_1 = -\boldsymbol{\mu}_2 \equiv \boldsymbol{\mu}$ , and the same diagonal covariance matrix,  $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 \equiv \boldsymbol{\Sigma}$ , with  $\Sigma_{ij} = \sigma_i \delta_{ij}$  and  $\sigma_i > 0$ . In this case, for symmetry reasons, the overlaps simplify and we have:

$$V_1 = V_2 \equiv V, \quad q_1 = q_2 \equiv q, \quad m_+ = -m_- \equiv m, \tag{201}$$

$$\hat{V}_+ = \hat{V}_- \equiv \frac{1}{2} \hat{V}, \quad \hat{q}_+ = \hat{q}_- \equiv \frac{1}{2} \hat{q}, \quad \hat{m}_+ = -\hat{m}_- \equiv \frac{1}{2} \hat{m}. \tag{202}$$

We define

$$\frac{1}{d} \sum_{i=1}^d \delta(\sigma - \sigma_i) \delta(\mu - \sqrt{d} \mu_i) \xrightarrow{d \rightarrow +\infty} p(\sigma, \mu) \tag{203}$$

joint distribution of the covariance diagonal elements and of the mean elements. We will denote  $\mathbb{E}_{\mu, \sigma}[\bullet]$  the average with respect to this measure. We will focus in particular on the form of the saddle-point equations obtained from the prior contribution assuming  $\ell_1$  regularization, i.e.,  $r(\mathbf{w}) = \sum_i |w_i|$ , and let us introduce the corresponding *soft-thresholding operator*:

$$\operatorname{Prox}_{\lambda, |\cdot|}(x) = \operatorname{sign}(x) \max\{|x| - \lambda, 0\}. \tag{204}$$

Observe that  $\operatorname{Prox}_{\alpha\lambda, |\cdot|}(\alpha x) = \alpha \operatorname{Prox}_{\lambda, |\cdot|}(x)$  for  $\alpha > 0$ . Its derivative given by  $\operatorname{Prox}'_{\lambda, |\cdot|}(x) = \theta(|x| > \lambda)$ . The saddle point equations from the prior part simply read:

$$V = \frac{1}{\hat{V}} \mathbb{E}_{\mu, \sigma, \xi} \left[ \operatorname{Prox}'_{\frac{\lambda}{\hat{V}}, |\cdot|} \left( \frac{\hat{m}\mu + \sqrt{\hat{q}}\sigma\xi}{\hat{V}\sigma} \right) \right], \tag{205}$$

$$q = \mathbb{E}_{\mu, \sigma, \xi} \left[ \sigma \left( \operatorname{Prox}_{\frac{\lambda}{\hat{V}}, |\cdot|} \left( \frac{\hat{m}\mu + \sqrt{\hat{q}}\sigma\xi}{\hat{V}\sigma} \right) \right)^2 \right], \tag{206}$$

$$m = \mathbb{E}_{\mu, \sigma, \xi} \left[ \mu \operatorname{Prox}_{\frac{\lambda}{\hat{V}}, |\cdot|} \left( \frac{\hat{m}\mu + \sqrt{\hat{q}}\sigma\xi}{\hat{V}\sigma} \right) \right]. \tag{207}$$

The averages over  $\xi$  can be performed explicitly using the simple expression of the proximal in this case. If we define the auxiliary functions

$$\begin{aligned}
\phi_{\pm}^0(v, u, \lambda) &\equiv \frac{1}{2} \operatorname{erfc} \left( \frac{\lambda \pm v}{\sqrt{2u}} \right) \\
\phi_{\pm}^1(u, v, \lambda) &= \sqrt{\frac{u}{2\pi}} e^{-\frac{(v \pm \lambda)^2}{2u}} - \frac{v \pm \lambda}{2} \operatorname{erfc} \left( \frac{\lambda \pm v}{\sqrt{2u}} \right), \\
\phi_{\pm}^2(v, u, \lambda) &= -\sqrt{\frac{u}{2\pi}} e^{-\frac{(\lambda \pm v)^2}{2u}} (\lambda \pm v) + \frac{u + (\lambda \pm v)^2}{2} \operatorname{erfc} \left( \frac{\lambda \pm v}{\sqrt{2u}} \right).
\end{aligned} \tag{208}$$

then

$$\begin{aligned}
V &= \frac{1}{\hat{V}} \mathbb{E}_{\mu, \sigma} [\phi_+^0(\mu \hat{m}, \sigma \hat{q}, \lambda) + \phi_-^0(\mu \hat{m}, \sigma \hat{q}, \lambda)] \\
q &= \mathbb{E}_{\mu, \sigma} \left[ \frac{\phi_+^2(\mu \hat{m}, \sigma \hat{q}, \lambda) + \phi_-^2(\mu \hat{m}, \sigma \hat{q}, \lambda)}{\sigma \hat{V}^2} \right], \\
m &= \mathbb{E}_{\mu, \sigma} \left[ \frac{\mu \phi_-^1(\mu \hat{m}, \sigma \hat{q}, \lambda) - \mu \phi_+^1(\mu \hat{m}, \sigma \hat{q}, \lambda)}{\sigma \hat{V}} \right].
\end{aligned} \tag{209}$$

**Gaussian means, homogenous covariances** — If  $p(\mu, \sigma) = \mathcal{N}(\mu|0, 1)\delta(\sigma - \Delta)$ , i.e., the means have i.i.d. Gaussian entries and  $\Sigma = \Delta \mathbf{I}_d$ , then

$$\begin{aligned} V &= \frac{1}{\hat{V}} \mathbb{E}_z \left[ \operatorname{erfc} \left( \frac{\lambda + \hat{m}z}{\sqrt{2\Delta\hat{q}}} \right) \right], \\ q &= \frac{1}{\Delta\hat{V}^2} \left\{ -\frac{e^{-\frac{1}{2}\frac{\lambda^2}{\hat{m}^2 + \Delta\hat{q}}}}{\sqrt{2\pi(\hat{m}^2 + \Delta\hat{q})}} \frac{2(\Delta\hat{q})^2\lambda}{\hat{m}^2 + \Delta\hat{q}} + \mathbb{E}_z \left[ (\lambda + \hat{m}z)^2 \operatorname{erfc} \left( \frac{\lambda + \hat{m}z}{\sqrt{2\Delta\hat{q}}} \right) \right] \right\}, \\ m &= \frac{1}{\Delta\hat{V}} \left\{ \frac{e^{-\frac{1}{2}\frac{\lambda^2}{\hat{m}^2 + \Delta\hat{q}}}}{\sqrt{2\pi(\hat{m}^2 + \Delta\hat{q})}} \frac{2\Delta\hat{q}\hat{m}\lambda}{\hat{m}^2 + \Delta\hat{q}} + \mathbb{E}_{z \sim \mathcal{N}(0,1)} \left[ (\lambda + \hat{m}z) z \operatorname{erfc} \left( \frac{\lambda + \hat{m}z}{\sqrt{2\Delta\hat{q}}} \right) \right] \right\}, \end{aligned} \quad (210)$$

with  $z \sim \mathcal{N}(0, 1)$ .

**Covariance correlated with sparse means** — In Section 3.1 we considered the case of sparse means correlated with the covariance matrices. In particular, we considered

$$p(\sigma, \mu) = p\mathcal{N}(\mu|0, 1)\delta(\sigma - \Delta_1) + (1-p)\delta(\mu)\delta(\sigma - \Delta_0). \quad (211)$$

The saddle-point equations are therefore

$$V = \frac{1}{\hat{V}} \left[ p\mathbb{E}_\mu \left[ \operatorname{erfc} \left( \frac{\lambda + \hat{m}\mu}{\sqrt{2\Delta_1\hat{q}}} \right) \right] + (1-p)\operatorname{erfc} \left( \frac{\lambda}{\sqrt{2\Delta_0\hat{q}}} \right) \right] \quad (212)$$

$$\begin{aligned} q &= \frac{p}{\Delta_1\hat{V}^2} \left\{ -\frac{e^{-\frac{1}{2}\frac{\lambda^2}{\hat{m}^2 + \Delta_1\hat{q}}}}{\sqrt{2\pi(\hat{m}^2 + \Delta_1\hat{q})}} \frac{2(\Delta_1\hat{q})^2\lambda}{\hat{m}^2 + \Delta_1\hat{q}} + \mathbb{E}_z \left[ (\lambda + \hat{m}z)^2 \operatorname{erfc} \left( \frac{\lambda + \hat{m}z}{\sqrt{2\Delta_1\hat{q}}} \right) \right] \right\} \\ &\quad - \lambda(1-p)\sqrt{\frac{\Delta_0\hat{q}}{2\pi}} e^{-\frac{\lambda^2}{2\Delta_0\hat{q}}} + \frac{1-p}{2}(\Delta_0\hat{q} + \lambda^2)\operatorname{erfc} \left( \frac{\lambda}{\sqrt{2\Delta_0\hat{q}}} \right) \end{aligned} \quad (213)$$

$$m = \frac{p}{\Delta_1\hat{V}} \left\{ \frac{e^{-\frac{1}{2}\frac{\lambda^2}{\hat{m}^2 + \Delta_1\hat{q}}}}{\sqrt{2\pi(\hat{m}^2 + \Delta_1\hat{q})}} \frac{2\Delta_1\hat{q}\hat{m}\lambda}{\hat{m}^2 + \Delta_1\hat{q}} + \mathbb{E}_z \left[ (\lambda + \hat{m}z) z \operatorname{erfc} \left( \frac{\lambda + \hat{m}z}{\sqrt{2\Delta_1\hat{q}}} \right) \right] \right\}. \quad (214)$$

In Section 3.1 we compare the performance obtained adopting an  $\ell_1$  regularization with the corresponding one obtained using  $\ell_2$ ,  $r(\mathbf{w}) = \sum_i w_i^2$ . For the sake of completeness, we give here the expression of the saddle-point equations in that case as well. In this case, the prior term  $\Psi_w$  can be written explicitly after a Gaussian integration as

$$\Psi_w(\hat{m}, \hat{Q}, \hat{V}) = -\frac{1}{2d} \operatorname{tr} \ln (\lambda \mathbf{I}_d + \hat{V}\Sigma) + \frac{1}{2} \operatorname{tr} \left[ (\lambda \mathbf{I}_d + \hat{V}\Sigma)^{-1} \left( \hat{m}_k^2 \boldsymbol{\mu}\boldsymbol{\mu}^\top + \frac{\hat{q}}{d}\Sigma \right) \right]. \quad (215)$$

In the setting given by eq. (211) the saddle point equations are then

$$q = p \frac{\hat{m}^2 \Delta_1 + \hat{q} \Delta_1^2}{(\lambda + \hat{V} \Delta_1)^2} + \frac{(1-p)\hat{q}\Delta_0^2}{(\lambda + \hat{V} \Delta_0)^2} \quad (216a)$$

$$V = p \frac{\Delta_1}{\lambda + \hat{V} \Delta_1} + \frac{(1-p)\Delta_0}{\lambda + \hat{V} \Delta_0} \quad (216b)$$

$$m = \frac{\hat{m}p}{\lambda + \hat{V} \Delta_1}. \quad (216c)$$

## D Bayes optimal error

In this Appendix, we derive a formula for the Bayes optimal classification error in the case of  $K$  clusters with the same covariance  $\Sigma_k = \Delta \mathbf{I}_d$  in the large  $d$  limit, assuming that a dataset  $\{(\mathbf{x}^\nu, \mathbf{y}^\nu)\}_{\nu \in [n]}$  of correctly labeled points is available. As usual, we will assume  $n/d = \alpha$  finite. The distribution of a pair  $(\mathbf{y}, \mathbf{x})$  is given by

$$p(\mathbf{y}, \mathbf{x} | \mathbf{M}) = \sum_k y_k \frac{\rho_k \exp\left(-\frac{1}{2\Delta} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2\right)}{(2\pi\Delta)^{\frac{d}{2}}}. \quad (217)$$

where  $\mathbf{M} \in \mathbb{R}^{d \times K}$  is the matrix of concatenated means  $\boldsymbol{\mu}_k$  estimated from the dataset, so that

$$\begin{aligned} p(\mathbf{M} | \{\mathbf{y}^\nu, \mathbf{x}^\nu\}_\nu) &\propto p(\{\mathbf{x}^\nu\}_\nu | \mathbf{M}, \{\mathbf{y}^\nu\}_\nu) P_\mu(\mathbf{M}) \\ &\propto P_\mu(\mathbf{M}) \prod_{\nu=1}^n \sum_k y_k^\nu \exp\left(-\frac{1}{2\Delta} \|\mathbf{x}^\nu - \boldsymbol{\mu}_k\|^2\right). \end{aligned} \quad (218)$$

We will assume in the following the distribution

$$P_\mu(\mathbf{M}) = \frac{\exp\left(-\frac{d}{2} \text{tr}[\mathbf{M}\boldsymbol{\Theta}^{-1}\mathbf{M}^\top]\right)}{(2\pi)^{\frac{Kd}{2}} d^{-K/2} |\boldsymbol{\Theta}|^{1/2}} \quad (219)$$

where  $\boldsymbol{\Theta} \in \mathbb{R}^{K \times K}$  is a given positive definite covariance matrix. In this way

$$\mathbb{E}[\mathbf{M}^\top \mathbf{M}] = \boldsymbol{\Theta}. \quad (220)$$

The conditional distribution for the label  $\mathbf{y}^0$  of a new point  $\mathbf{x}^0$ ,

$$\begin{aligned} p(\mathbf{y}^0 | \mathbf{x}^0, \{\mathbf{y}^\nu, \mathbf{x}^\nu\}_\nu) &\propto \mathbb{E}_{\mathbf{M} | \{\mathbf{y}^\nu, \mathbf{x}^\nu\}_\nu} [p(\mathbf{y}, \mathbf{x} | \mathbf{M})] \\ &= \int d\mathbf{M} P_\mu(\mathbf{M}) \sum_k y_k^0 \rho_k \exp\left(-\frac{\|\mathbf{x}^0 - \boldsymbol{\mu}_k\|^2}{2\Delta}\right) \prod_{\nu=1}^n \sum_k y_k^\nu \exp\left(-\frac{\|\mathbf{x}^\nu - \boldsymbol{\mu}_k\|^2}{2\Delta}\right). \end{aligned} \quad (221)$$

If  $\mathbf{n} = (n_k)_k$  is the vector of the number of examples  $n_k$  in the class  $k$ , then

$$\begin{aligned} p(\mathbf{y}^0 | \mathbf{x}^0, \{\mathbf{y}^\nu, \mathbf{x}^\nu\}_\nu) &\propto \int d\mathbf{M} P_\mu(\mathbf{M}) \prod_{k=1}^K \left[ \rho_k^{y_k^0} \exp\left(-\sum_{\nu=0}^n \frac{y_k^\nu \|\mathbf{x}^\nu - \boldsymbol{\mu}_k\|^2}{2\Delta}\right) \right] \\ &= \exp\left[\sum_k y_k^0 \left(\ln \rho_k - \frac{\|\mathbf{x}^0\|^2}{2\Delta}\right) - \frac{1}{2} \ln \det\left(1 + \frac{1}{d\Delta} \text{diag}(\mathbf{n} + \mathbf{y}^0) \boldsymbol{\Theta}\right)\right] \\ &\times \exp\left[\frac{1}{2\Delta} \text{tr}\left[\left(\sum_{\nu=0}^n \mathbf{y}^\nu \otimes \mathbf{x}^\nu\right)^\top (d\Delta \boldsymbol{\Theta}^{-1} + \text{diag}(\mathbf{n} + \mathbf{y}))^{-1} \left(\sum_{\nu=0}^n \mathbf{y}^\nu \otimes \mathbf{x}^\nu\right)\right]\right]. \end{aligned} \quad (222)$$

In the following we will denote by  $\star$  the true label of  $\mathbf{x}$ . Let  $\boldsymbol{\Pi} = \text{diag}(\rho_k)$ . Then we can write the previous expression as

$$\begin{aligned} p(\mathbf{y}^0 | \mathbf{x}^0, \{\mathbf{y}^\nu, \mathbf{x}^\nu\}_\nu) &\propto \exp\left[\sum_k y_k \left(\ln \rho_k - \frac{\|\mathbf{x}^0\|^2}{2\Delta}\right) - \frac{1}{2} \ln \det\left(1 + \frac{1}{\Delta} \alpha \boldsymbol{\Pi} \boldsymbol{\Theta}\right)\right] \\ &\times \exp\left[\frac{1}{2\Delta} \text{tr}\left[\left(\frac{1}{d} \sum_{\nu=0}^n \mathbf{y}^\nu \otimes \mathbf{x}^\nu\right)^\top (\Delta \boldsymbol{\Theta}^{-1} + \alpha \boldsymbol{\Pi})^{-1} \left(\sum_{\nu=0}^n \mathbf{y}^\nu \otimes \mathbf{x}^\nu\right)\right]\right] \end{aligned} \quad (223)$$

Observe now that

$$\frac{1}{d\Delta} \mathbf{x}^0 \sum_{\nu=1}^n y_k^\nu \mathbf{x}^\nu \xrightarrow{n, d \rightarrow +\infty} \alpha \rho_k \frac{\Theta_{\star, k} + \eta_k Z_k}{\Delta}, \quad \eta_k \equiv \sqrt{\Delta \left(1 + \frac{\Delta}{\alpha \rho_k}\right)}, \quad Z_k \sim \mathcal{N}(0, 1), \quad (224)$$

so that, defining the vector  $\mathbf{a}^* = (a_k)_{k \in [K]}$  with elements

$$a_k^* \equiv \alpha \rho_k \frac{\Theta_{*,k} + \eta_k Z_k}{\Delta}, \quad (225)$$

and neglecting the  $\mathbf{y}^0$ -independent contributions, the expression above can be rewritten as

$$p(\mathbf{y}^0 | \mathbf{x}^0, \{\mathbf{y}^\nu, \mathbf{x}^\nu\}_\nu) \propto \exp \left[ \sum_k y_k^0 \ln \rho_k + \left( \mathbf{a}^* + \frac{1}{2} \mathbf{y}^0 \right)^\top (\Delta \mathbf{\Theta}^{-1} + \alpha \mathbf{\Pi})^{-1} \mathbf{y}^0 \right] \quad (226)$$

where we have also used the fact that  $\|\mathbf{x}^0\|^2 = d\Delta + O(1)$ . This means that the Bayes optimal generalization error is

$$\varepsilon_g^{\text{BO}} = \sum_k \rho_k \mathbb{P} \left[ \arg \max_{\kappa} \left( \ln \rho_\kappa + \left( \mathbf{a}^k + \frac{1}{2} \mathbf{e}_\kappa \right)^\top (\Delta \mathbf{\Theta}^{-1} + \alpha \mathbf{\Pi})^{-1} \mathbf{e}_\kappa \right) \neq k \right]. \quad (227)$$

If  $\mathbf{\Theta} = \mathbf{I}_K$  and the clusters have same weights,  $\rho_k \equiv 1/K \Leftrightarrow \mathbf{\Pi} = 1/K \mathbf{I}_K$ , then  $\eta_k \equiv \eta$  and

$$\varepsilon_g^{\text{BO}} = \mathbb{P} \left[ \frac{1}{\eta} < \max_{\kappa \in [K-1]} Z_\kappa + Z \right], \quad (228)$$

that is the formula given in [20].

## E Experiments with real data

In this Appendix we discuss the experiments of Section 3.3 with real data sets.

**Numerical details** — Consider a real data set  $\{(\mathbf{x}^\nu, y^\nu)\}_{\nu=1}^{n_{\text{tot}}}$  with  $n_{\text{tot}}$  samples which we assume are independent. As a pre-processing step we center, normalise and flatten the inputs  $\mathbf{x}^\nu$  into  $d$ -dimensional vectors. For both the MNIST [61] and Fashion-MNIST [62] data sets used in the experiments we have normalised the inputs by 255, such that components  $x_i^\nu \in [0, 1]$ . In what follows we focus on binary classification tasks and encode the labels as  $y^\nu \in \{-1, 1\}$ . For example, for the MNIST and Fashion-MNIST data sets we have  $d = 784$  and  $n_{\text{tot}} = 7 \times 10^4$ , and we split the inputs into two classes depending on the task of interest, e.g. odd vs. even digits and clothes vs. accessories items, respectively. Define the empirical distribution over the data set:

$$\hat{P}(\mathbf{x}, y) = \frac{1}{n_{\text{tot}}} \sum_{\nu=1}^{n_{\text{tot}}} \delta(\mathbf{x} - \mathbf{x}^\nu) \delta(y - y^\nu) \quad (229)$$

The question we want to answer is: how well can we approximate the learning curves  $(\epsilon_g, \epsilon_t)$  on a given ERM classification task by approximating  $\hat{P}$  with a Gaussian mixture distribution? To answer this question, we consider a Gaussian mixture distribution  $P_2$  as defined in Eq. (1) with the same means and covariances as  $\hat{P}$ :

$$\hat{\boldsymbol{\mu}}_k = \frac{1}{n_{\text{tot}}} \sum_{\nu=1}^{n_{\text{tot}}} \mathbf{x}^\nu \mathbb{I}(\mathbf{x}^\nu \in \mathcal{C}_k), \quad \hat{\boldsymbol{\Sigma}}_k = \frac{1}{n_{\text{tot}}} \sum_{\nu=1}^{n_{\text{tot}}} (\mathbf{x}^\nu - \boldsymbol{\mu}_k)(\mathbf{x}^\nu - \boldsymbol{\mu}_k)^\top \mathbb{I}(\mathbf{x}^\nu \in \mathcal{C}_k) \quad (230)$$

for  $k \in \{+, -\}$  labelling the two clusters. Similarly, the class probabilities  $\rho_k$  are also estimated from the full data set:

$$\hat{\rho}_k = \frac{1}{n_{\text{tot}}} \sum_{\nu=1}^{n_{\text{tot}}} \mathbb{I}(\mathbf{x}^\nu \in \mathcal{C}_k). \quad (231)$$

The parameters  $(\hat{\boldsymbol{\mu}}_k, \hat{\boldsymbol{\Sigma}}_k, \hat{\rho}_k)$  completely characterise the approximating Gaussian mixture distribution  $P_2$ , and together with Theorem 1 can be used to compute the theoretical learning curves  $(\epsilon_g, \epsilon_t)$  as in Fig. 5 of the main. Note that this discussion can be easily generalised to the case in which a non-linear feature map  $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^p$  is applied to the data prior to fitting. The only difference is that the empirical distribution  $\hat{P}$  is defined over the features  $\{(\mathbf{v}^\nu, y^\nu)\}_{\nu=1}^{n_{\text{tot}}}$  where  $\mathbf{v}^\nu = \varphi(\mathbf{x}^\nu)$ , and the Gaussian mixture approximation  $P_2$  is defined with respect to the empirical features distribution. Figure 6 of the main manuscript shows an example where a random feature map  $\mathbf{v} = \text{erf}(\mathbf{F}\mathbf{x})$  with  $\mathbf{F} \in \mathbb{R}^{p \times d}$  a random Gaussian projection applied to MNIST and fashion MNIST before the fitting with different ratios  $\gamma = p/d$ .

The theoretical learning curves are then compared with two sets of finite instance simulations. First, we simulate the learning problem on synthetic data sampled from the approximating Gaussian mixture distribution  $P_2$ , and the learning curves are computed by averaging over 10 instances of the problem. Second, we simulate the learning problem on the real data set. The real data set is split into training and test sets, and for a given sample complexity  $\alpha = n/d$  we sub-sample  $n = \alpha d$  points from the training set. The averaged learning curves are computed over different instances of the sub-sampling, with replacement.

**Discussion** — As expected, we find good agreement between theory and simulations with synthetic data drawn from the approximating Gaussian mixture distribution  $P_2$ , even for relatively small input dimensions (e.g.  $d = 784$  for MNIST). Surprisingly, we have found that in many cases the Gaussian mixture is a good approximation to the real data curves, see Figs. 5 and 6 for examples of logistic regression on input space and with random features. Figure 7 shows an example where the feature map  $\varphi$  is given by removing the last layer of the following fully-connected 2-layer neural network pre-trained on the full MNIST odd vs. even data set:

```
Sequential(
  (0): Linear(in_features=784, out_features=784, bias=False)
  (1): ReLU()
  (2): Linear(in_features=784, out_features=1, bias=False)
  (3): Tanh()
)
```

with the training performed by minimising the square loss with the Adam optimiser and random initialisation. However, we have also found cases in which the approximation is not as sharp, see blue curves in Fig. 10. Understanding the factors determining the quality of the approximation in real data sets is an interesting question we expect to address in future work.

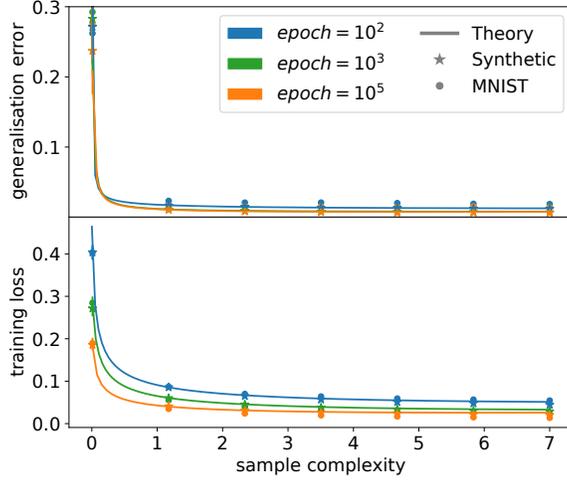


Figure 7: Generalisation error and training loss for logistic regression on MNIST with a feature map  $\varphi$  obtained by training 2-layer fully connected neural network, with  $\ell_2$  penalty and fixed  $\lambda = 0.05$ . The different curves show the performance at different stages of training.

**Multiclass vs. binary approximation** – In the cases previously discussed, we have considered a  $K = 2$  cluster approximation  $P_2$  to the empirical data distribution  $\hat{P}$ . However, the data sets considered here (MNIST and Fashion-MNIST) are originally composed of 10 classes, and therefore we should ask the question of whether a  $K = 10$  cluster approximation  $P_{10}$  where we fit the means and covariances of each original class is any different from the approximation studied above. In principle, these two approximations can have very different statistical properties. For instance, from Theorem 2 it follows that the generalisation and training errors of Gaussian mixtures only depend on the statistics of the local field  $\lambda = \mathbf{W}\mathbf{x}$  conditioned on the labels, which in the binary setting considered here is  $y \in \{+, -\}$ . Conditioned on  $y = \pm$ , this local field is simply a Gaussian random variable under  $P_2$ , while it is a multi-modal random variable under  $P_{10}$ . Therefore, there is *a priori* no reason for these two approximations to give the same learning curves.

As an example, consider a  $K = 4$  Gaussian mixture distribution with a common variance  $\Sigma_k = \Delta \mathbf{I}_d$  and with means:

$$\boldsymbol{\mu}_1 = \mathbf{e}_1 + \mathbf{e}_2, \quad \boldsymbol{\mu}_2 = \mathbf{e}_1 - \mathbf{e}_2, \quad \boldsymbol{\mu}_3 = -\mathbf{e}_1 + \mathbf{e}_2, \quad \boldsymbol{\mu}_4 = -\mathbf{e}_1 - \mathbf{e}_2 \quad (232)$$

where  $\mathbf{e}_i \in \mathbb{R}^d$  is the canonical basis vector of  $\mathbb{R}^d$ , with entries  $e_{ij} = \delta_{ij}$ . We consider two label assignments: a) a realisable case in which clusters 1 and 2 are assigned label +1, and clusters 3 and 4 are assigned -1 and b) a non-realisable case in which clusters 1 and 4 are assigned +1 and clusters 2 and 3 are assigned -1 (XOR function), see Fig. 8 (top) for an illustration. Now consider a dual  $K = 2$  Gaussian mixture model with means and covariances  $(\boldsymbol{\mu}_\pm, \Sigma_\pm)$  chosen to match the class means and covariances of the  $K = 4$  mixture, see Fig. 8 (bottom) for an illustration. In Fig. 9 we compare the learning curves of the  $K = 4$  model with the  $K = 2$  counterpart with matched class means and covariances. While in the realisable case a) both have identical performance under the error bars, in the non-realisable case b) the performance in are significantly different.

Indeed, a similar behaviour can be observed in the real data experiments. Fig. 10 compares the real learning curves of a MNIST 5v5 binary classification task (classifying five first digits vs. five last) with the two different Gaussian mixture approximations:  $P_{10}$  where we fit the means and covariances of each individual cluster and  $P_2$ , where we fit only the class-wise means and covariances. While both approximations capture the high-level behaviour of the learning curves,  $P_{10}$  is closer to the real learning curve than  $P_2$ .

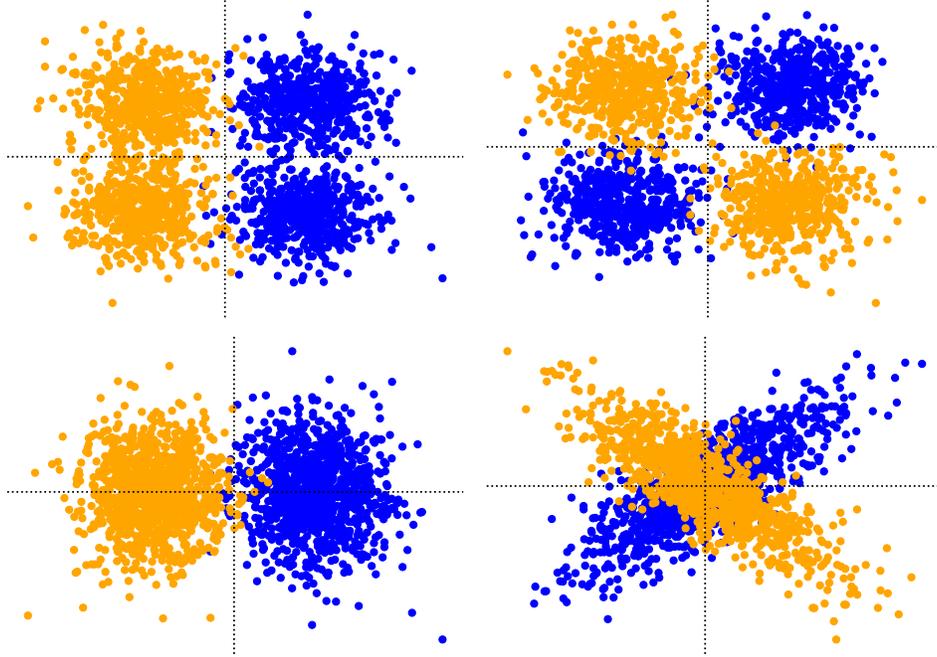


Figure 8: Two dimensional projection of the setting described in eq. (232). **(Left)** Realisable case, **(Right)** Non-realisable case (XOR function).

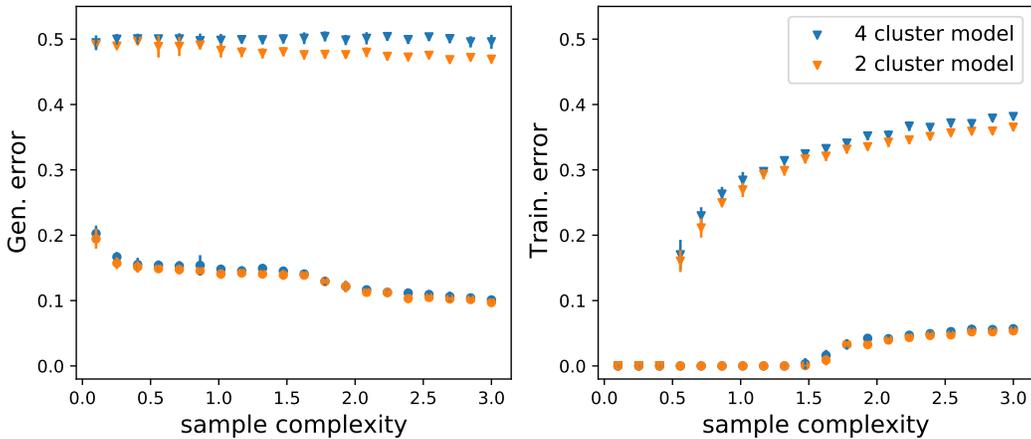


Figure 9: **(Left)** Generalisation and **(right)** training errors as a function of the sample complexity for logistic regression with  $\ell_2$  penalty and  $\lambda = 10^{-4}$  for the four models pictured in Fig. 8. Points denote the separable model (bottom curve), and triangles denote the non-realizable xor model (top curves). We have chosen a balanced scenario with  $\Delta = 0.5$ .

**Note on numerical instabilities** — When dealing with means and covariance matrices estimated from real data sets, we have observed that for small regularisation strength  $\lambda \ll 1$  the self-consistent equations from Theorem 1 can develop spurious fixed points corresponding to negative values of the overlap parameters  $q_{\pm} = \mathbf{W}^{\top} \Sigma_{\pm} \mathbf{W}$  — which is clearly not possible since  $\Sigma_{\pm}$  is a positive-definite matrix. This is observed across different scenarios, and is independent of the choice of loss or the particular way the equations are solved. In fact, the minimum value of  $\lambda$  below which the spurious fixed point develop seems to depend only on the conditioning number of the covariance matrices.

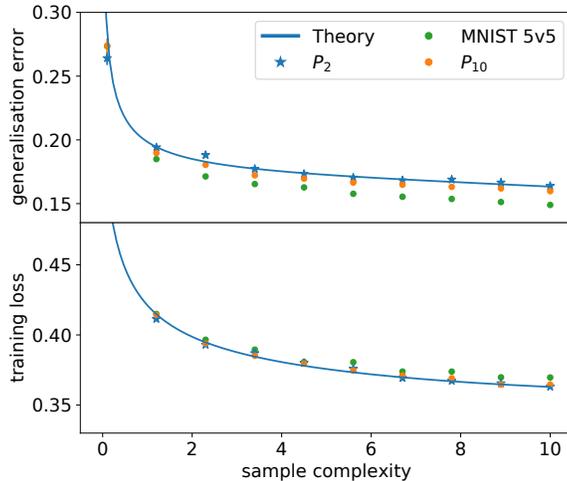


Figure 10: Generalisation error and training loss for logistic regression on the task of classifying  $\{0, 1, 2, 3, 4\}$  vs  $\{5, 6, 7, 8, 9\}$  digits of MNIST, as a function of the sample complexity for fixed  $\ell_2$  penalty  $\lambda = 0.1$ . The blue curves show the 2-Gaussian cluster approximation  $P_2$  (solid for theory, points for finite size simulations), while the orange points show the 10-Gaussian cluster approximation  $P_{10}$ , which lies systematically below. The green points denote simulations on the true data set.

## References

- [1] Stuart Geman, Elie Bienenstock, and René Doursat. Neural networks and the bias/variance dilemma. *Neural Computation*, 4(1):1–58, 1992.
- [2] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. Springer New York Inc., New York, NY, USA, 2001.
- [3] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- [4] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J. Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *Preprint arXiv:1903.08560*, 2020.
- [5] Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand kernel learning. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 541–549. PMLR, 10–15 Jul 2018.
- [6] Peter L. Bartlett, Philip M. Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- [7] Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and double descent curve. *Communications on Pure and Applied Mathematics*, 2019. To appear, preprint arXiv:1908.05355.
- [8] Federica Gerace, Bruno Loureiro, Flornet Krzakala, Marc Mézard, and Lenka Zdeborová. Generalisation error in learning with random features and the hidden manifold model. In *37th International Conference on Machine Learning*, 2020.
- [9] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Limitations of lazy training of two-layers neural network. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32, pages 9111–9121, 2019.

- [10] Sebastian Goldt, Marc Mézard, Florent Krzakala, and Lenka Zdeborová. Modeling the influence of data structure on learning in neural networks: The hidden manifold model. *Physical Review X*, 10(4):041044, 2020.
- [11] Sebastian Goldt, Bruno Loureiro, Galen Reeves, Florent Krzakala, Marc Mézard, and Lenka Zdeborová. The Gaussian equivalence of generative models for learning with shallow neural networks. *Preprint arXiv:2006.14709*, 2020.
- [12] Bruno Loureiro, Cédric Gerbelot, Hugo Cui, Sebastian Goldt, Florent Krzakala, Marc Mézard, and Lenka Zdeborová. Capturing the learning curves of generic features maps for realistic data sets with a teacher-student model. *Preprint arXiv:2102.08127*, 2021.
- [13] Tengyuan Liang and Pragma Sur. A precise high-dimensional asymptotic theory for Boosting and minimum- $\ell_1$ -norm interpolated classifiers. *Preprint arXiv:2002.01586*, 2020.
- [14] Francesca Mignacco, Florent Krzakala, Yue Lu, Pierfrancesco Urbani, and Lenka Zdeborová. The role of regularization in classification of high-dimensional noisy Gaussian mixture. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 6874–6883. PMLR, 13–18 Jul 2020.
- [15] Maria Refinetti, Sebastian Goldt, Florent Krzakala, and Lenka Zdeborová. Classifying high-dimensional gaussian mixtures: Where kernel methods fail and neural networks succeed. *Preprint arXiv:2102.11742*, 2021.
- [16] Emmanuel J Candès, Pragma Sur, et al. The phase transition for the existence of the maximum likelihood estimate in high-dimensional logistic regression. *The Annals of Statistics*, 48(1):27–42, 2020.
- [17] Vardan Papyan, X. Y. Han, and David L. Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):24652–24663, 2020.
- [18] Mohamed El Amine Seddik, Cosme Louart, Mohamed Tamaazousti, and Romain Couillet. Random matrix theory proves that deep learning representations of GAN-data behave as Gaussian mixtures. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 8573–8582. PMLR, 13–18 Jul 2020.
- [19] David Donoho and Jiashun Jin. Higher criticism thresholding: Optimal feature selection when useful features are rare and weak. *Proceedings of the National Academy of Sciences*, 105(39):14790–14795, 2008.
- [20] Christos Thrampoulidis, Samet Oymak, and Mahdi Soltanolkotabi. Theoretical insights into multiclass classification: A high-dimensional asymptotic view. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 8907–8920. Curran Associates, Inc., 2020.
- [21] Zeyu Deng, Abla Kammoun, and Christos Thrampoulidis. A model of double descent for high-dimensional binary linear classification. *Preprint arXiv:1911.05822*, 2020.
- [22] Xiaoyi Mai, Zhenyu Liao, and Romain Couillet. A large scale analysis of logistic regression: Asymptotic performance and new insights. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3357–3361, 2019.
- [23] Xiaoyi Mai and Zhenyu Liao. High dimensional classification via regularized and unregularized empirical risk minimization: Precise error and optimal loss. *Preprint arXiv:1905.13742*, 2020.
- [24] Edgar Dobriban and Stefan Wager. High-dimensional asymptotics of prediction: Ridge regression and classification. *The Annals of Statistics*, 46(1):247 – 279, 2018.
- [25] Ganesh Kini and Christos Thrampoulidis. Analytic study of double descent in binary classification: The impact of loss. *Preprint arXiv:2001.11572*, 2020.

- [26] Housseem Sifaou, Abla Kammoun, and Mohamed-Slim Alouini. Phase transition in the hard-margin support vector machines. In *2019 IEEE 8th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, pages 415–419, 2019.
- [27] Ke Wang and Christos Thrampoulidis. Binary classification of gaussian mixtures: Abundance of support vectors, benign overfitting and regularization. 2021.
- [28] Niladri S. Chatterji and Philip M. Long. Finite-sample analysis of interpolating linear classifiers in the overparameterized regime. *Preprint arXiv:2004.12019*, 2021.
- [29] Yuan Cao, Quanquan Gu, and Mikhail Belkin. Risk bounds for over-parameterized maximum margin classification on sub-gaussian mixtures. *Preprint arXiv:2104.13628*, 2021.
- [30] Fariborz Salehi, Ehsan Abbasi, and Babak Hassibi. The impact of regularization on high-dimensional logistic regression. *Preprint arXiv:1906.03761*, 2019.
- [31] Christos Thrampoulidis, Ehsan Abbasi, and Babak Hassibi. Precise error analysis of regularized  $m$ -estimators in high dimensions. *IEEE Transactions on Information Theory*, 64(8):5592–5628, 2018.
- [32] Mihailo Stojnic. A framework to characterize performance of lasso algorithms. *Preprint arXiv:1303.7291*, 2013.
- [33] 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.
- [34] 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.
- [35] David L Donoho, Adel Javanmard, and Andrea Montanari. Information-theoretically optimal compressed sensing via spatial coupling and approximate message passing. *IEEE transactions on information theory*, 59(11):7434–7464, 2013.
- [36] 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.
- [37] Raphael Berthier, Andrea Montanari, and Phan-Minh Nguyen. State evolution for approximate message passing with non-separable functions. *Information and Inference: A Journal of the IMA*, 9(1):33–79, 2020.
- [38] Andre Manoel, Florent Krzakala, Marc Mézard, and Lenka Zdeborová. Multi-layer generalized linear estimation. In *2017 IEEE International Symposium on Information Theory (ISIT)*, pages 2098–2102. IEEE, 2017.
- [39] Jiashun Jin. Impossibility of successful classification when useful features are rare and weak. *Proceedings of the National Academy of Sciences*, 106(22):8859–8864, 2009.
- [40] Jun Shao, Yazhen Wang, Xinwei Deng, and Sijian Wang. Sparse linear discriminant analysis by thresholding for high dimensional data. *The Annals of Statistics*, 39(2):1241 – 1265, 2011.
- [41] Qing Mai, Hui Zou, and Ming Yuan. A direct approach to sparse discriminant analysis in ultra-high dimensions. *Biometrika*, 99(1):29–42, 12 2012.
- [42] Yanfang Li and Jinzhu Jia. L1 least squares for sparse high-dimensional LDA. *Electronic Journal of Statistics*, 11(1):2499 – 2518, 2017.
- [43] Thomas M Cover. Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition. *IEEE transactions on electronic computers*, (3):326–334, 1965.
- [44] Elizabeth Gardner. The space of interactions in neural network models. *Journal of physics A: Mathematical and general*, 21(1):257, 1988.

- [45] Arthur Jacot, Berfin Şimşek, Francesco Spadaro, Clément Hongler, and Franck Gabriel. Kernel alignment risk estimator: Risk prediction from training data. *Preprint arXiv:2006.09796*, 2020.
- [46] Blake Bordelon, Abdulkadir Canatar, and Cengiz Pehlevan. Spectrum dependent learning curves in kernel regression and wide neural networks. In *International Conference on Machine Learning*, pages 1024–1034. PMLR, 2020.
- [47] Neal Parikh and Stephen Boyd. Proximal algorithms. *Foundations and Trends in optimization*, 1(3):127–239, 2014.
- [48] Heinz H Bauschke, Patrick L Combettes, et al. *Convex analysis and monotone operator theory in Hilbert spaces*, volume 408. Springer, 2011.
- [49] Michael Celentano, Andrea Montanari, and Yuting Wei. The Lasso with general Gaussian designs with applications to hypothesis testing. *Preprint arXiv:2007.13716*, 2020.
- [50] Erwin Bolthausen. An iterative construction of solutions of the TAP equations for the Sherrington–Kirkpatrick model. *Communications in Mathematical Physics*, 325(1):333–366, 2014.
- [51] Mohsen Bayati and Andrea Montanari. The LASSO risk for Gaussian matrices. *IEEE Transactions on Information Theory*, 58(4):1997–2017, 2011.
- [52] Cedric Gerbelot, Alia Abbara, and Florent Krzakala. Asymptotic Errors for Teacher-Student Convex Generalized Linear Models (or: How to Prove Kabashima’s Replica Formula). *Preprint arXiv:2006.06581*, 2020.
- [53] Benjamin Aubin, Antoine Maillard, Jean Barbier, Florent Krzakala, Nicolas Macris, and Lenka Zdeborová. The committee machine: Computational to statistical gaps in learning a two-layers neural network. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124023, 2019.
- [54] Florent Krzakala, Marc Mézard, François Sausset, YF Sun, and Lenka Zdeborová. Statistical-physics-based reconstruction in compressed sensing. *Physical Review X*, 2(2):021005, 2012.
- [55] Cynthia Rush and Ramji Venkataramanan. Finite sample analysis of approximate message passing algorithms. *IEEE Transactions on Information Theory*, 64(11):7264–7286, 2018.
- [56] Mohsen Bayati, Marc Lelarge, Andrea Montanari, et al. Universality in polytope phase transitions and message passing algorithms. *Annals of Applied Probability*, 25(2):753–822, 2015.
- [57] Wei-Kuo Chen and Wai-Kit Lam. Universality of approximate message passing algorithms. *Electronic Journal of Probability*, 26:1–44, 2021.
- [58] Bruno Loureiro, Gabriele Sicuro, Cédric Gerbelot, Alessandro Pocco, Florent Krzakala, and Lenka Zdeborová. GaussMixtureProject, October 2021. <https://github.com/IdePHICS/GaussMixtureProject>.
- [59] Scott Shaobing Chen, David L. Donoho, and Michael A. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20(1):33–61, 1998.
- [60] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *The Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [61] Yann LeCun and Corinna Cortes. *ATT Labs [Online]*, 2010. Database released under CC BY-SA 3.0 license at <http://yann.lecun.com/exdb/mnist/>.
- [62] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *Preprint arXiv:1708.07747*, 2017. Database released under MIT licence at <https://github.com/zalandoresearch/fashion-mnist>.

- [63] Ali Rahimi and Benjamin Recht. Random Features for Large-Scale Kernel Machines. In *NIPS*, pages 1177–1184, 2007.
- [64] Heinz H Bauschke, Jonathan M Borwein, and Patrick L Combettes. Bregman monotone optimization algorithms. *SIAM Journal on control and optimization*, 42(2):596–636, 2003.
- [65] Heinz H Bauschke, Minh N Dao, and Scott B Lindstrom. Regularizing with bregman–moreau envelopes. *SIAM Journal on Optimization*, 28(4):3208–3228, 2018.
- [66] Lenka Zdeborová and Florent Krzakala. Statistical physics of inference: Thresholds and algorithms. *Advances in Physics*, 65(5):453–552, 2016.
- [67] Cédric Gerbelot and Raphaël Berthier. Graph-based approximate message passing iterations. *arXiv preprint arXiv:2109.11905*, 2021.
- [68] Ryan J Tibshirani. The lasso problem and uniqueness. *Electronic Journal of statistics*, 7:1456–1490, 2013.
- [69] 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.
- [70] Saharon Rosset, Ji Zhu, and Trevor Hastie. Margin maximizing loss functions. In *NIPS*, pages 1237–1244, 2003.