

PAPER • OPEN ACCESS

# Fluctuations, bias, variance and ensemble of learners: exact asymptotics for convex losses in high-dimension<sup>\*</sup>

To cite this article: Bruno Loureiro *et al* *J. Stat. Mech.* (2023) 114001

View the [article online](#) for updates and enhancements.

You may also like

- [Shift-curvature, SGD, and generalization](#)  
Arwen V Bradley, Carlos A Gomez-Uribe  
and Manish Reddy Vuyyuru
- [VECM and Bayesian VECM for Overparameterization Problem](#)  
Meilina Retno Hapsari, Suci Astutik and  
Loekito Adi Soehono
- [Mean-field inference methods for neural networks](#)  
Marylou Gabrié

PAPER: ML 2023

## Fluctuations, bias, variance and ensemble of learners: exact asymptotics for convex losses in high-dimension\*

Bruno Loureiro<sup>1</sup>, Cédric Gerbelot<sup>2</sup>, Maria Refinetti<sup>2</sup>,  
Gabriele Sicuro<sup>3,\*\*</sup> and Florent Krzakala<sup>1</sup>

<sup>1</sup> École Polytechnique Fédérale de Lausanne (EPFL), Information, Learning and Physics (IdePHICS) Lab., Lausanne, CH-1015, Switzerland

<sup>2</sup> Laboratoire de Physique de l'ENS, Université PSL, CNRS, Sorbonne Université, 24 Rue Lhomond, 75005 Paris, France

<sup>3</sup> Department of Mathematics, King's College London, Strand WC2R 2LS, London, United Kingdom

E-mail: [gabriele.sicuro@unibo.it](mailto:gabriele.sicuro@unibo.it)

Received 10 May 2023

Accepted for publication 13 September 2023

Published 15 November 2023

Online at [stacks.iop.org/JSTAT/2023/114001](https://stacks.iop.org/JSTAT/2023/114001)

<https://doi.org/10.1088/1742-5468/ad0221>



CrossMark

**Abstract.** From the sampling of data to the initialisation of parameters, randomness is ubiquitous in modern Machine Learning practice. Understanding the statistical fluctuations engendered by the different sources of randomness in prediction is therefore key to understanding robust generalisation. In this manuscript we develop a quantitative and rigorous theory for the study of fluctuations in an ensemble of generalised linear models trained on different, but correlated, features in high-dimensions. In particular, we provide a complete description of the asymptotic joint distribution of the empirical risk minimiser

\*This article is an updated version of: Loureiro B, Gerbelot C, Refinetti M, Sicuro G and Krzakala F 2022 Fluctuations, bias, variance & ensemble of learners: exact asymptotics for convex losses in high-dimension *Proc. 39th Int. Conf. on Machine Learning (Proc. Machine Learning Research* vol 162) ed K Chaudhuri, S Jegelka, L Song, C Szepesvari, G Niu and S Sabato (PMLR) pp 14283–314.

\*\* Author to whom any correspondence should be addressed.



Original Content from this work may be used under the terms of the [Creative Commons Attribution 4.0 licence](https://creativecommons.org/licenses/by/4.0/). Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

for generic convex loss and regularisation in the high-dimensional limit. Our result encompasses a rich set of classification and regression tasks, such as the lazy regime of overparametrised neural networks, or equivalently the random features approximation of kernels. While allowing to study directly the mitigating effect of ensembling (or bagging) on the bias-variance decomposition of the test error, our analysis also helps disentangle the contribution of statistical fluctuations, and the singular role played by the interpolation threshold that are at the roots of the ‘double-descent’ phenomenon.

**Keywords:** learning theory, machine learning

---

## Contents

<b>1. Introduction</b> .....	<b>3</b>
1.1. Setting .....	3
1.1.1. Main contributions. ....	5
1.1.2. Related works. ....	6
<b>2. Learning with an ensemble of random features</b> .....	<b>7</b>
<b>3. Applications</b> .....	<b>9</b>
3.1. Ridge regression .....	9
3.2. Binary classification .....	10
3.3. Alignment of learners. ....	11
3.4. Ensemble predictors .....	12
<b>4. The case of general loss and regularisation</b> .....	<b>13</b>
4.1. The random feature case and the kernel limit .....	16
4.2. The kernel limit .....	17
<b>Acknowledgments</b> .....	<b>18</b>
<b>Appendix A. The replica approach</b> .....	<b>18</b>
A.1. Notation .....	18
A.2. The replica computation .....	19
A.2.1. Replica symmetric ansatz .....	21
A.2.2. Zero temperature state evolution equations .....	23
A.2.3. The case of $\ell_2$ regularisation .....	24
A.2.4. Training loss and generalisation error .....	24
A.3. Separable loss with ridge regularisation .....	25
A.3.1. The random-features model for the generative networks .....	27
A.3.2. Ridge regression .....	27
A.3.3. Binary classification problem .....	28
<b>Appendix B. Proof of the main theorem</b> .....	<b>30</b>

B.1. The learning problem .....	30
B.1.1. Assumptions .....	31
B.2. Asymptotics for the strongly convex problem .....	32
B.3. Relaxing the strong convexity constraint .....	44
B.4. A comment on non-pseudo-Lipschitz subgradients .....	44
B.5. Toolbox .....	45
B.5.1. Notations .....	45
B.5.2. Moreau envelopes and Bregman proximal operators	46
B.5.3. Gradients of Bregman envelopes .....	47
B.5.4. Gaussian concentration .....	47
B.5.5. Approximate message-passing .....	47
B.5.6. A useful result from convex analysis .....	48
<b>References</b> .....	<b>49</b>

## 1. Introduction

Randomness is ubiquitous in Machine Learning. It is present in the data (e.g. noise in acquisition and annotation), in commonly used statistical models (e.g. random features (RFs) (Rahimi and Recht 2007)), or in the algorithms used to train them (e.g. in the choice of initialisation of weights of neural networks (Narkhede *et al* 2022), or when sampling a mini-batch in Stochastic Gradient Descent (Bottou 2012)). Strikingly, fluctuations associated to different sources of randomness can have a major impact in the generalisation performance of a model. For instance, this is the case in least-squares regression with RFs, where it has been shown (D’Ascoli *et al* 2020, Geiger *et al* 2020, Jacot *et al* 2020) that the variance associated with the random projections matrix is responsible for poor generalisation near the interpolation peak (Advani and Saxe 2017, Spigler *et al* 2019, Belkin *et al* 2020). As a consequence, this *double-descent* behaviour can be mitigated by averaging over a large *ensemble* of learners, effectively suppressing this variance. Indeed, considering an ensemble (sometimes also referred to as a committee (Drucker *et al* 1994)) of independent learners provide a natural framework to study the contribution of the variance of prediction in the estimation accuracy. In this manuscript we leverage this idea to provide an exact asymptotic characterisation of the statistics of fluctuations in empirical risk minimisation with generic convex losses and penalties in high-dimensional models. We focus on the case of synthetic datasets, and we apply our results to RF learning in particular.

### 1.1. Setting

Let  $(\mathbf{x}^\mu, y^\mu) \in \mathbb{R}^d \times \mathcal{Y}$ ,  $\mu \in [n] := \{1, \dots, n\}$ , denote a labelled data set composed of  $n$  independent samples from a joint density  $p(\mathbf{x}, y)$  (e.g.  $\mathcal{Y} = \{-1, 1\}$  for a binary classification problem). In this manuscript we are interested in studying an ensemble of  $K$  parametric predictors, each of them depending on a vector of parameters  $\mathbf{w}_k \in \mathbb{R}^p$ ,  $k \in [K]$ , and independently trained on the dataset  $\{(\mathbf{x}^\mu, y^\mu)\}_{\mu \in [n]}$ . Note that even if the

vectors of parameters  $\{\mathbf{w}_k\}_{k \in [K]}$  are trained independently, they correlate through the training data. Statistical fluctuations in the learnt parameters can then arise for different reasons. For instance, a common practice is to initialise the parameters randomly during optimisation, which will induce statistical variability between the different predictors. Alternatively, each predictor could be trained on a subsample of the data, as it is commonly done in bagging (Breiman 1996). The statistical model can also be inherently stochastic, e.g. the RFs approximation for kernel methods (Rahimi and Recht 2007). Finally, the predictors could also be jointly trained, e.g. coupling them through the loss or penalty as it is done in boosting (Schapire 1990).

Our goal in this work is to provide a sharp characterisation of the statistical fluctuations of the ensemble of parameters  $\{\mathbf{w}_k\}_{k \in [K]}$  in a particular, mathematically tractable, class of predictors: *generalised linear models*,

$$\hat{y}(\mathbf{x}) = \hat{f} \left( \frac{\hat{\mathbf{w}}_1^\top \mathbf{u}_1(\mathbf{x})}{\sqrt{p}}, \dots, \frac{\hat{\mathbf{w}}_K^\top \mathbf{u}_K(\mathbf{x})}{\sqrt{p}} \right) \tag{1}$$

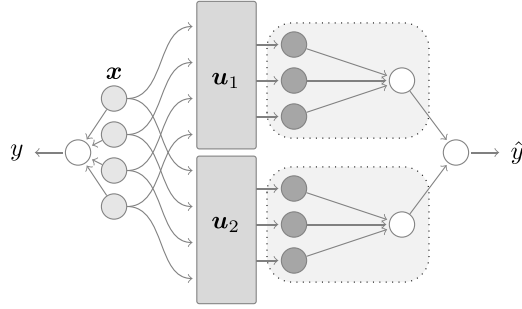
where  $\mathbf{u}_k : \mathbb{R}^d \rightarrow \mathbb{R}^p$ ,  $k \in [K]$  is an ensemble of possibly correlated features and  $\hat{f} : \mathbb{R}^K \rightarrow \mathcal{Y}$  is an activation function. For most of this work, we discuss the case in which the predictors are *independently trained* through regularised empirical risk minimisation:

$$\hat{\mathbf{w}}_k = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^p} \left[ \frac{1}{n} \sum_{\mu=1}^n \ell \left( y^\mu, \frac{\mathbf{w}^\top \mathbf{u}_k(\mathbf{x}^\mu)}{\sqrt{p}} \right) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right] \tag{2}$$

with a convex loss function  $\ell : \mathcal{Y} \times \mathbb{R} \rightarrow \mathbb{R}$  (e.g. the logistic loss) and ridge penalty whose strength is given by  $\lambda \in \mathbb{R}^+$ . However, our analysis also includes the case in which the learners are jointly trained with a generic convex penalty. This case will be further discussed in section 4. In what follows we will also concentrate in the RFs case where  $\mathbf{u}_k(\mathbf{x}) = \phi(\mathbf{F}_k \mathbf{x})$  with  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  an activation function acting component-wise and  $\mathbf{F}_k \in \mathbb{R}^{p \times d}$  a family of independently sampled random matrices. Besides being an efficient approximation for kernels (Rahimi and Recht 2007), RFs are often studied as a simple model for neural networks in the lazy and neural tangent kernel regimes of deep neural networks (Jacot *et al* 2018, Chizat *et al* 2019), in which case the matrices  $\mathbf{F}_k$  correspond to different random initialisation of hidden-layer weights. Moreover, the RFs model displays some of the exotic behaviours of high-dimensional overparametrised models, such as double-descent (Gerace *et al* 2020, Mei and Montanari 2021) and benign overfitting (Bartlett *et al* 2020), therefore providing an ideal playground to study the interplay between fluctuations and overparametrisation. A broader class of features maps is also discussed in section 4.

To provide an exact characterisation of the statistics of the estimators in equation (2), we shall assume data is generated from a target

$$y = f_0 \left( \frac{\boldsymbol{\theta}^\top \mathbf{x}}{\sqrt{d}} \right), \quad \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}_d, \rho \mathbf{I}_d), \quad \rho \in \mathbb{R}_0^+, \tag{3}$$



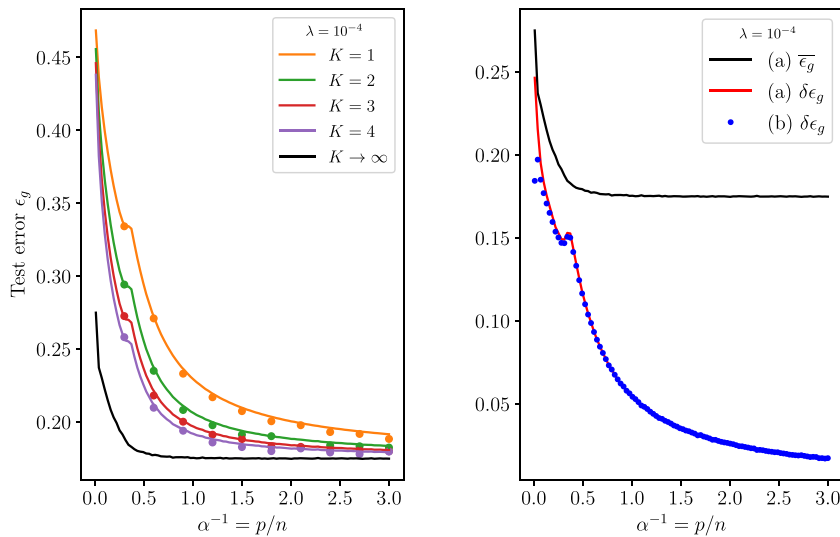
**Figure 1.** Pictorial representation of the model considered in the paper for  $K = 2$ . Two learners with the same architecture (in grey) receive a correlated input generated from the same vector  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ . The output  $\hat{y}$  is an average of their outputs. While the study of an ensemble of learners is already interesting *per se*, it is also pivotal to study the fluctuation between learners, and the error stemming from the difference in the weights in random features and lazy training.

with  $f_0 : \mathbb{R} \rightarrow \mathcal{Y}$  and  $\mathbf{I}_d$   $d$ -dimensional identity matrix. The dataset is then constructed generating *i.i.d.*  $n$  vectors  $\mathbf{x}^\mu \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ ,  $\mu \in [n]$ .

An illustration summary of the setting considered here is given in figure 1. Note that such architecture can be interpreted as a two-layer tree neural network, also known in some contexts as the *tree-committee* or *parity machine* (Schwarze and Hertz 1992).

*1.1.1. Main contributions.* The results in this manuscript can be listed as follows.

- We provide a sharp asymptotic characterisation of the joint statistics of the ensemble of empirical risk minimisers  $\{\hat{\mathbf{w}}_k\}_{k \in [K]}$  in the high-dimensional limit where  $p, n \rightarrow +\infty$  with  $n/p$  kept constant, for any convex loss and penalty. In particular, we show that the pre-activations  $\{\hat{\mathbf{w}}_k^\top \mathbf{u}_k\}_{k \in [K]}$  are jointly Gaussian, with sufficient statistics obeying a set of explicit closed-form equations. Note that the analysis of ensembling with non-square losses is out of the grasp of the most commonly adopted theoretical tools (e.g. random matrix theory). Therefore, our proof method based on recent progress on approximate message passing (AMP) techniques (Javanmard and Montanari 2013, Berthier *et al* 2020, Gerbelot and Berthier 2021) is of independent interest. Different versions of our theorem are discussed throughout the manuscript. First, in section 2 for the particular case of independently trained learners on RFs (theorem 1). Later, in section 4 for the general case of jointly trained learners on correlated Gaussian covariates (theorem 2).
- We discuss the role played by fluctuations in the non-monotonic behaviour of the generalisation performance of interpolators (a.k.a. double-descent behaviour). In particular—as discussed in Geiger *et al* (2020), d’Ascoli *et al* (2021) for the ridge case—the interpolation peak arises from the model overfitting the particular realisation of the random weights. We show the test error can be decomposed  $\epsilon_g(K = 1) = \bar{\epsilon}_g + \delta\epsilon_g$  in terms of a fluctuation-free term  $\bar{\epsilon}_g$  and a fluctuation term  $\delta\epsilon_g$  responsible for the double-descent behaviour, see figure 2 for the case of max-margin classification.



**Figure 2.** *Left.* Test error for logistic regression with  $\lambda = 10^{-4}$  and different values of  $K$  as function of  $p/n = 1/\alpha$  with  $n/d = 2$  and  $\rho = 1$ . Dots represent the average of the outcomes of  $10^3$  numerical experiments. Here we adopted  $\phi(x) = \text{erf}(x)$  and estimator  $\hat{f}(v) = \text{sign}(\sum_k v_k)$ . *Right.* Decomposition of the  $K = 1$  test error  $\epsilon_g = \bar{\epsilon}_g + \delta\epsilon_g$  for the estimator (a), with  $n/d = 2$  and  $\lambda = 10^{-4}$ . We plot also the contribution  $\delta\epsilon_g$  corresponding to the estimator (b): we numerically observed that such decomposition coincides in the two cases. Note also the presence of a kink in  $\delta\epsilon_g$  at the interpolation transition.

- In the context of classification, we discuss how *majority vote* and *score averaging*, two popular ensembling procedures, compare in terms of generalisation performance. More specifically, we show that in the setting we study score averaging consistently outperforms the majority vote predictor. However, for a large number of learners  $K \gg 1$  these two predictors agree, see figure 5(right).
- Finally, we discuss how ensembling can be used as a tool for uncertainty quantification. In particular, we connect the correlation between two learners to the probability of disagreement, and show that it decreases with overparametrisation, see figure 5(centre). We provide a full characterisation of the joint probability density of the confidence score between two independent learners, see figure 5(left).

*1.1.2. Related works.* The idea of reducing the variance of a predictor by averaging over independent learners is quite old in Machine Learning (Hansen and Salamon 1990, Perrone and Cooper 1993, Perrone 1994, Krogh and Vedelsby 1995), and an early asymptotic analysis of the regression case was given in Krogh and Sollich (1997). In particular, a variety of methods to combine an ensemble of learners appeared in the literature (Opitz and Maclin 1999). In a very inspiring work, Geiger *et al* (2020) carried out an extensive series of experiments in order to shed light on the generalisation properties of neural networks, and reported many observations and empirical arguments about the role of the variance due to the random initialisation of the weights in the double-descent curve using an ensemble of learners. This was a major motivation for the present work.



Closest to our setting is the work of Neal *et al* (2018), D'Ascoli *et al* (2020), Jacot *et al* (2020) which disentangles the various sources of variance in the process of training deep neural networks. Indeed, here we adopt the model defined by D'Ascoli *et al* (2020), and provide a rigorous justification of their results for the case of ridge regression. A slightly finer decomposition of the variance in terms of the different sources of randomness in the problem was later proposed by Adlam and Pennington (2020a). Lin and Dobriban (2021) show that such decomposition is not unique, and can be more generally understood from the point of view of the *analysis of variance* framework. Interestingly, subsequent papers were able to identify a series of triple (and more) descent, e.g. Chen *et al* (2020), Adlam and Pennington (2020b), d'Ascoli *et al* (2021).

The RFs model was introduced in the seminal work of Rahimi and Recht (2007) as an efficient approximation for kernel methods. Drawing from early ideas of Karoui (2010), Pennington and Worah (2017) showed that the empirical distribution of the Gram matrix of RF is asymptotically equivalent to a linear model with matched second statistics, and characterised in this way memorisation with RF regression. The learning problem was first analysed by Mei and Montanari (2021), who provided an exact asymptotic characterisation of the training and generalisation errors of RF regression. This analysis was later extended to generic convex losses by Gerace *et al* (2020) using the heuristic replica method, and later proved by Dhifallah and Lu (2020) using convex Gaussian inequalities.

The aforementioned asymptotic equivalence between the RF model and a Gaussian model with matched moments has been named the *Gaussian Equivalence Principle* (GEP) (Goldt *et al* 2020). Rigorous proofs in the memorisation and learning setting with square loss appeared in Pennington and Worah (2017), Mei and Montanari (2021), and for general convex penalties in Hu and Lu (2020), Goldt *et al* (2021). Goldt *et al* (2021) and Loureiro *et al* (2021b) provided extensive numerical evidence that the GEP holds for more generic feature maps, including features stemming from trained neural networks.

Most of the previously mentioned works deriving exact asymptotics for the RF model in the proportional limit use either Random Matrix Theory techniques or convex Gaussian inequalities. While these tools have been recently used in many different contexts, they ultimately fall short when considering an ensemble of predictors with generic convex loss and regularisation, along with structured design matrices. Therefore, to prove the results herein we employ an *AMP* proof technique (Bayati and Montanari 2011a, Donoho and Montanari 2016), leveraging on recently introduced progresses in Gerbelot and Berthier (2021), Loureiro *et al* (2021b) which enables to capture the full complexity of the problem and obtain the asymptotic joint distribution of the ensemble of predictors. LeJeune *et al* (2020) studies ensembles of ordinary least-squares learned from subsamples of a common data matrix, and shows its equivalence to an implicit ridge regularisation.

## 2. Learning with an ensemble of random features

In this section give a first formulation of our main result, namely the exact asymptotic characterisation of the statistics of the ensembling estimator introduced in equation (1).



We prove that, in the proportional high dimensional limit, the statistics of the arguments of the activation function in equation (1) is simply given by a multivariate Gaussian, whose covariance matrix we can completely specify. This result holds for any convex loss, any convex regularisation, and for all models of generative networks  $\mathbf{u}_k: \mathbb{R}^d \rightarrow \mathbb{R}^p$ , as we will show in full generality in section 4. However, for simplicity, in this section and in the following we focus on the setting described in section 1, in which the statistician averages over an independent ensemble of RFs, i.e.  $\mathbf{u}_k(\mathbf{x}) = \phi(\mathbf{F}_k \mathbf{x})$ . In this case, our result can be formulated as follows:

**Theorem 1 (simplified version).** *Assume that in the high-dimensional limit where  $d, p, n \rightarrow +\infty$  with  $\alpha := n/p$  and  $\gamma := d/p$  kept  $\Theta(1)$  constants, the Wishart matrix  $\mathbf{F}\mathbf{F}^\top$  has a well-defined asymptotic spectral distribution. Then in this limit, for any pseudo-Lipshitz function of order 2  $\varphi: \mathbb{R} \times \mathbb{R}^K \rightarrow \mathbb{R}$ , we have*

$$\mathbb{E}_{(x,y)} \left[ \varphi \left( y, \frac{\hat{\mathbf{w}}_1^\top \mathbf{u}_1}{\sqrt{p}}, \dots, \frac{\hat{\mathbf{w}}_K^\top \mathbf{u}_K}{\sqrt{p}} \right) \right] \xrightarrow{P} \mathbb{E}_{(\nu, \boldsymbol{\mu})} [\varphi(f_0(\nu), \boldsymbol{\mu})], \tag{4}$$

where  $(\nu, \boldsymbol{\mu}) \in \mathbb{R}^{K+1}$  is a jointly Gaussian vector  $(\nu, \boldsymbol{\mu}) \sim \mathcal{N}(\mathbf{0}_{K+1}, \boldsymbol{\Sigma})$  with covariance

$$\boldsymbol{\Sigma} = \begin{pmatrix} \rho & m \mathbf{1}_K^\top \\ m \mathbf{1}_K & \mathbf{Q} \end{pmatrix}, \quad \mathbf{Q} := (q_0 - q_1) \mathbf{I}_K + q_1 \mathbf{1}_{K,K}, \tag{5}$$

with  $\mathbf{1}_{K,K} \in \mathbb{R}^{K \times K}$  and  $\mathbf{1}_K \in \mathbb{R}^K$  are a matrix and a vector of ones respectively. The entries of  $\boldsymbol{\Sigma}$  are solutions of a set of self-consistent equations given in Corollary 4.

As discussed in the introduction, the asymptotic statistics of the *single* learner has been studied in Dhifallah and Lu (2020), Gerace *et al* (2020), Loureiro *et al* (2021b). Their result amounts to the analysis of the estimator solving the empirical risk minimisation problem in equation (2) and it is recovered imposing  $K = 1$  in the theorem above.

For  $K = 1$ ,  $(\nu, \mu) \in \mathbb{R}^2$  is jointly Gaussian with zero mean and covariance  $\boldsymbol{\Sigma} = \begin{pmatrix} \rho & m \\ m & q_0 \end{pmatrix}$ .

However, such result is not enough to quantify the correlation between different learners, induced by the training on the same dataset, which is required to compute, e.g. the test error associated with an ensembling predictor as in equation (1). For example, in the simple case where  $f_0(u) = u$  and  $\hat{f}(\mathbf{v}) = \frac{1}{K} \sum_k v_k$ , the mean-squared error on the labels is given by  $\epsilon_g = \mathbb{E}_{(x,y)} [(y - \hat{y}(\mathbf{x}))^2] = \rho + (q_0 - q_1)K^{-1} + q_1 - 2m$ , which crucially depends on the average correlation between two independent learners<sup>4</sup>  $q_1 := \frac{1}{p} \mathbb{E}[\hat{\mathbf{w}}_1^\top \hat{\mathbf{w}}_2]$ . Our main result is precisely an exact asymptotic characterisation of this correlation in the proportional limit of the previous theorem. Once  $m$ ,  $q_0$  and  $q_1$  have been determined, the generalisation error can be computed as

$$\epsilon_g := \mathbb{E}_{(x,y)} [\Delta(y, \hat{y}(\mathbf{x}))] \xrightarrow{n \rightarrow +\infty} \mathbb{E}_{(\nu, \boldsymbol{\mu})} \left[ \Delta \left( f_0(\nu), \hat{f}(\boldsymbol{\mu}) \right) \right] \tag{6}$$

for any error measure  $\Delta: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ .

<sup>4</sup> Note that since all learners are here assumed to be statistically equivalent, their pair-wise correlation is the same on average. In the general case, discussed in section 4, the correlation matrix  $\mathbf{Q} \in \mathbb{R}^{K \times K}$  can have a more complex structure.

Suppose now that

$$\hat{f}(\mathbf{v}) \equiv \hat{f}_0 \left( \frac{1}{K} \sum_k v_k \right) \tag{7}$$

for some  $\hat{f}_0: \mathbb{R} \rightarrow \mathcal{Y}$  activation function of the single learner. In this case we can introduce the random variable  $\hat{\mu} \stackrel{d}{=} \lim_{K \rightarrow +\infty} \frac{1}{K} \sum_k \mu_k$ . It is not difficult to see that the joint probability  $p(\nu, \hat{\mu}) \sim \mathcal{N}(\mathbf{0}_2, \hat{\Sigma})$  where  $\hat{\Sigma} = \begin{pmatrix} \rho & m \\ m & q_1 \end{pmatrix}$ . This formally coincides with the joint distribution for the activation fields for  $K=1$  (Gerace *et al* 2020), but with  $q_0$  replaced by  $q_1 \leq q_0$ . The smaller variance is due to the fact that the fluctuations of the activation fields are averaged out by the ensembling process. The test error in the  $K \rightarrow +\infty$  limit is then

$$\bar{\epsilon}_g := \mathbb{E}_{(\nu, \hat{\mu})} \left[ \Delta \left( f_0(\nu), \hat{f}_0(\hat{\mu}) \right) \right], \tag{8}$$

so that the fluctuation contribution to the test error for  $K=1$  can be defined as

$$\delta\epsilon_g := \mathbb{E}_{(\nu, \mu)} \left[ \Delta \left( f_0(\nu), \hat{f}_0(\mu) \right) \right] - \bar{\epsilon}_g. \tag{9}$$

The term  $\delta\epsilon_g$  is by definition the contribution suppressed by ensembling and corresponds to the *ambiguity* introduced by Krogh and Vedelsby (1995) for the square loss. This contribution expresses the variance in the ensemble and it is responsible for the non-monotonic behaviour in the test error of interpolators, also known as the double-descent behavior.

### 3. Applications

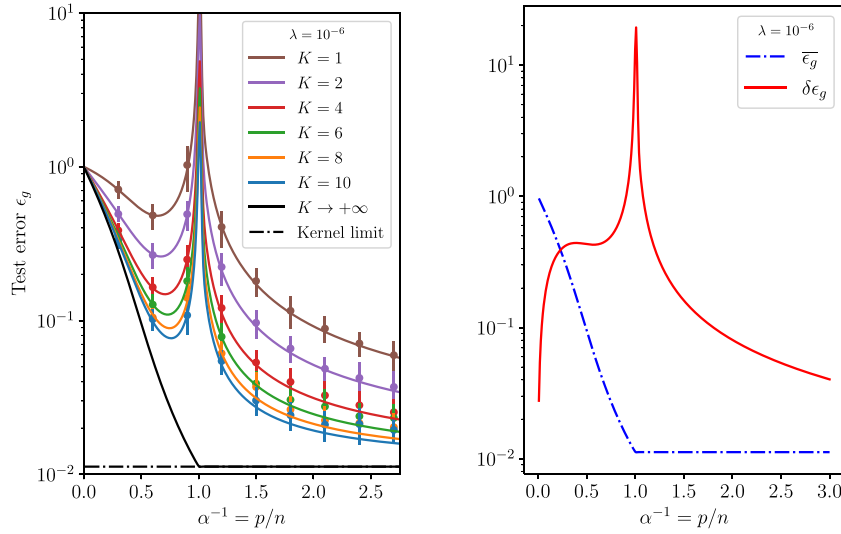
We will consider now two relevant examples of separable losses, namely a ridge loss and a logistic loss. In both cases, it is possible to derive the explicit expression of the training loss and generalisation error in terms of the elements of the correlation matrix introduced above.

#### 3.1. Ridge regression

If we assume  $f_0(x) = x$ ,  $\hat{f}(\mathbf{v}) = \frac{1}{K} \sum_k v_k$ , and a quadratic loss of the type  $\ell(y, x) = \frac{1}{2}(y - x)^2$ , it is possible to write down simple recursive equations for  $m$ ,  $q_0$  and  $q_1$  (see appendix A.3.2). Taking  $\Delta(y, \hat{y}) = (y - \hat{y})^2$ , the generalisation error is easily computed as

$$\epsilon_g = \rho + \frac{q_0 - q_1}{K} + q_1 - 2m \xrightarrow{K \rightarrow +\infty} \rho + q_1 - 2m \equiv \bar{\epsilon}_g. \tag{10}$$

Note that in this case the  $\lambda \rightarrow 0^+$  limit gives the minimum  $\ell_2$ -norm interpolator. In figure 3 we compare our theoretical prediction with numerical results for  $\lambda = 10^{-6}$  and various values of  $K$ . It is evident that the divergence of the generalisation error at  $\alpha = 1$



**Figure 3.** *Left.* Test error for ridge regression with  $\lambda = 10^{-6}$  and different values of  $K$  as function of  $p/n = 1/\alpha$  with  $n/d = 2$  and  $\rho = 1$ . Dots represent the average of the outcomes of 50 numerical experiments in which the parameters of the neurons are estimated using  $\min(d, p) = 200$ . Here we adopted  $\phi(x) = \text{erf}(x)$ . *Right.* Decomposition of  $\epsilon_g = \bar{\epsilon}_g + \delta\epsilon_g$  in the  $K = 1$  case.

is only due to the divergence of  $q_0$ , whereas the contribution  $\bar{\epsilon}_g$ , which is independent on  $q_0$ , is smooth everywhere. Alongside with the interpolation divergence,  $\delta\epsilon_g = q_0 - q_1$  has an additional bump at  $p/n = d/n$ , which corresponds to the ‘linear peak’ discussed by d’Ascoli *et al* (2021).

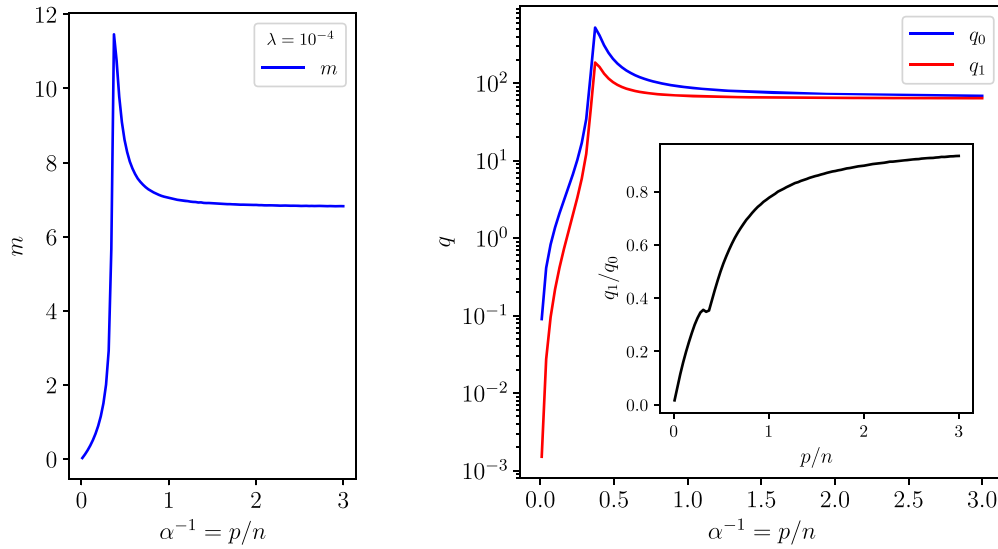
In the plot we present also the so-called kernel limit, corresponding to the limit  $n/p = \alpha \rightarrow 0$  at fixed  $n/d$ . An explicit manipulation (see appendix A.3.2) shows that  $q_1 = q_0 \equiv q$  in this limit. This implies that in the kernel limit  $\epsilon_g^k$  does not depend on  $K$ , being equal to  $\epsilon_g^k \equiv \rho + q - 2m$ . The generalisation error obtained in the kernel limit coincides with  $\bar{\epsilon}_g$  for  $p > n$ : this is expected as in  $\bar{\epsilon}_g$  the fluctuations amongst learners are averaged out, effectively recovering the cost obtained in the case of an infinite number of parameters.

### 3.2. Binary classification

Suppose now that we are considering a classification task, such that  $\mathcal{Y} = \{-1, 1\}$ . For this task we consider  $f_0(x) = \text{sign}(x)$ . A popular choice of loss in this classification task is the logistic loss,

$$\ell(y, x) = \ln(1 + e^{-yx}), \quad (11)$$

although other choices, e.g. hinge loss, can be considered. Since both the logistic and hinge losses depend only on the *margin*  $y\mathbf{w}^\top \mathbf{u}$ , the empirical risk minimiser for  $\lambda \rightarrow 0^+$  in both cases give the max-margin interpolator (Rosset *et al* 2004). We write down the explicit saddle-point equations associated to the logistic and hinge loss in



**Figure 4.** Analytical estimation of the covariance parameters characterising the correlation with the oracle  $m$  (left), the norm of the predictor in feature space  $q_0$  and the correlation between learners  $q_1$  (right) (see equation (5) for the definition) in a classification task using logistic loss with ridge penalty with  $\lambda = 10^{-4}$  at fixed  $n/d = 2$  as function of  $p/n$ . In the inset, ratio  $q_1/q_0$ , quantifying the correlation between two learners. In all parameters the interpolation kink is clearly visible.

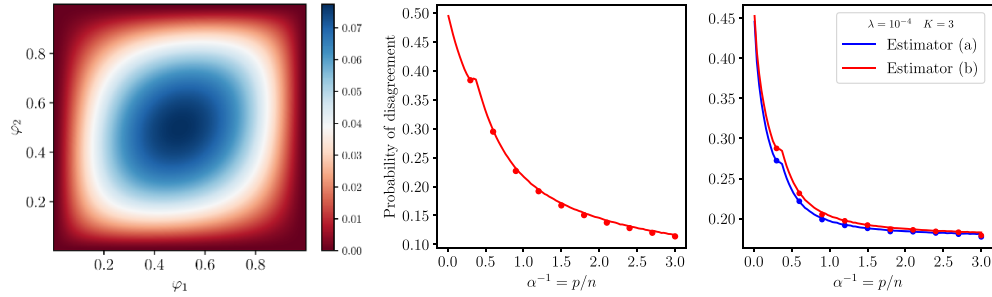
appendix A.3.3, but we will focus our attention on the logistic case for the sake of brevity. For this choice of the loss, we obtained the values of  $m$ ,  $q_0$  and  $q_1$  showed in figure 4. Using these values, a number of relevant questions can be addressed.

### 3.3. Alignment of learners

Assuming that the predictor of the learner  $k$  is  $\hat{y}_k(\mathbf{x}) = \text{sign}(\hat{\mathbf{w}}_k^\top \mathbf{u}_k(\mathbf{x}))$ , in figure 5(centre) we estimate the probability that two learners give opposite classification. This is analytically given by

$$\mathbb{P}[\hat{y}_1(\mathbf{x}) \neq \hat{y}_2(\mathbf{x})] = \mathbb{P}[\mu_1 \mu_2 < 0] = \frac{1}{\pi} \arccos\left(\frac{q_1}{q_0}\right). \quad (12)$$

Note that by definition the ratio  $q_1/q_0$  is a cosine similarity between two learners in the norm induced by the feature space. Therefore, this provides an interesting interpretation of these sufficient statistics in terms of the probability of disagreement. In particular, as illustrated in figure 5(centre) overparametrisation promotes agreement between the learners, therefore suppressing uncertainty. More generally, ensembling can be used as a technique for uncertainty estimation (Lakshminarayanan *et al* 2017). In the context of logistic regression, the pre-activation to the sign function is often interpreted as a *confidence score*. Indeed, introducing the logistic function  $\varphi_k(\mathbf{x}) = (1 + \exp(-p^{-1/2} \hat{\mathbf{w}}_k^\top \mathbf{u}_k(\mathbf{x})))^{-1}$ , it expresses the confidence of the  $k$ th classifier in associating  $\hat{y} = 1$  to the input  $\mathbf{x}$ . Therefore, it is reasonable to ask how reliable is the



**Figure 5.** *Left.* Joint probability density of the confidence score  $\varphi_i(\mathbf{x}) = (1 + \exp(-p^{-1/2} \hat{\mathbf{w}}_i^\top \mathbf{u}_i(\mathbf{x})))^{-1}$  of two learners for  $p/n \simeq 0.13$ . *Centre.* Probability that two learners give discordant predictions using logistic regression as function of  $p/n = 1/\alpha$  with  $n/d = 2$ ,  $\rho = 1$ , and  $\lambda = 10^{-4}$ . *Right.* Test error for logistic regression using the estimators in equation (13) and  $K = 3$ , with the same parameters. We adopted  $\phi(x) = \text{erf}(x)$ . We observe that the test error obtained using (a) is always smaller than the one obtained using (b). (*Centre and right*) Dots represent the average of the outcomes of  $10^3$  numerical experiments.

logistic score as a confidence measure. For instance, what is the variance of the confidence among different learners? This can be quantified by the joint probability density  $\rho(\varphi_1, \varphi_2) := \mathbb{E}_{\mathbf{x}}[\delta(\varphi_1 - \varphi_1(\mathbf{x}))\delta(\varphi_2 - \varphi_2(\mathbf{x}))]$ , which can be readily computed using our theorem 1. Figure 5(left) shows one example at fixed  $p/n$  and vanishing  $\lambda$ .

### 3.4. Ensemble predictors

In the previous two points, we discussed how ensembling can be used as a tool to quantify fluctuations. However, ensembling methods are also used in practical settings in order to mitigate fluctuations, e.g. Breiman (1996). An important question in this context is: given an ensemble of predictors  $\{\hat{\mathbf{w}}_k\}_{k \in [K]}$ , what is the best way of combining them to produce a point estimate? In our setting, this amounts to choosing the function  $\hat{f} : \mathbb{R}^K \rightarrow \mathcal{Y}$ . Let us consider two popular choices for the estimator  $\hat{f}$  in equation (1) used in practice:

$$(a) \quad \hat{f}(\mathbf{v}) = \text{sign} \left( \sum_k v_k \right), \tag{13a}$$

$$(b) \quad \hat{f}(\mathbf{v}) = \text{sign} \left( \sum_k \text{sign}(v_k) \right). \tag{13b}$$

In a sense, (a) provides an estimator based on the average of the output fields, whereas (b), which corresponds to a majority rule if  $K$  is odd (Hansen and Salamon 1990), is a function of the average of the estimators of the single learners. For both choices of the estimator we use  $\Delta(y, \hat{y}) = \delta_{\hat{y}, y}$  to measure the test error. In figure 5(right) we compare the test error obtained using (a) and (b) for  $K = 3$  with vanishing regularisation  $\lambda = 10^{-4}$ . It is observed that the estimator (a) has better performances than the estimator (b). As previously discussed, in this case logistic regression is equivalent to

max-margin estimation, and in this case the error **(a)** can be intuitively understood in terms of a robust max-margin estimation obtained by averaging the margins associated to different draws of the RFs. In the case **(a)** it is easy to show that the generalisation error takes the form

$$\epsilon_g = \frac{1}{\pi} \arccos \left( \frac{\sqrt{K} m}{\sqrt{\rho(q_0 - q_1 + K q_1)}} \right) \xrightarrow{K \rightarrow \infty} \frac{1}{\pi} \arccos \left( \frac{m}{\sqrt{\rho q_1}} \right) \equiv \bar{\epsilon}_g. \quad (14)$$

This formula is in agreement with numerical experiments, see figure 2(left). Unfortunately, we did not find a similar closed-form expression in case **(b)**. However, we can observe that in the  $K \rightarrow +\infty$  limit the generalisation error in case **(a)** coincides with the generalisation error in case **(b)**, see figure 2(right). By comparing with the results in figure 5(centre), it is evident that the benefit of ensembling in reducing the test error correlates with the tendency of learners to disagree, i.e. for small values of  $p/n$ , as stressed by Krogh and Vedelsby (1995). Finally, we observe a constant value of  $\bar{\epsilon}_g$  beyond the interpolation threshold, compatibly with the numerical results of Geiger *et al* (2020).

#### 4. The case of general loss and regularisation

In this section we generalise our results in section 2 relaxing the hypothesis on the loss, on the regularisation and on the properties of the feature maps. In the general setting we are going to consider, we denote  $P_y^0(y|x)$  the probabilistic law by which  $y$  is generated. For example, in section 2,  $P_y^0(y|x) = \delta(y - f_0(x))$ . In the treatment given here, we allow for more general cases (e.g. the presence of noise in the label generation). We make no assumptions on the generative networks  $\mathbf{u}_k$ , so that the information about the first layer is contained in the following tensors,

$$\Omega := \mathbb{E}_{\mathbf{x}} [\mathbf{U}(\mathbf{x}) \otimes \mathbf{U}(\mathbf{x})] \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}, \quad (15)$$

$$\hat{\Phi} := \mathbb{E}_{\mathbf{x}} [\mathbf{U}(\mathbf{x}) \mathbf{x}^\top \boldsymbol{\theta}] \in \mathbb{R}^{p \times K}, \quad (16)$$

$$\Theta = \hat{\Phi} \otimes \hat{\Phi} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}. \quad (17)$$

In the equations above,  $\mathbf{U}(\mathbf{x}) \in \mathbb{R}^{p \times K}$  is the matrix having as concatenated columns  $\mathbf{u}_k(\mathbf{x})$ . We aim at learning a rule as in equation (1), adopting a general convex loss  $\hat{\ell}: \mathcal{Y} \times \mathbb{R}^K \rightarrow \mathbb{R}$ , so that the weights are estimated as

$$\hat{\mathbf{W}} = \arg \min_{\mathbf{W} \in \mathbb{R}^{p \times K}} \left[ \frac{1}{n} \sum_{\mu=1}^n \hat{\ell} \left( y^\mu, \frac{\text{diag}(\mathbf{W}^\top \mathbf{U}^\mu)}{\sqrt{p}} \right) + \lambda r(\mathbf{W}) \right] \quad (18)$$

where  $r: \mathbb{R}^{p \times K} \rightarrow \mathbb{R}$  is a convex regularisation,  $\mathbf{U}^\mu \equiv \mathbf{U}(\mathbf{x}^\mu)$  and  $\hat{\mathbf{W}} \in \mathbb{R}^{p \times K}$  matrix of the concatenated columns  $\{\hat{\mathbf{w}}_k\}$ . Here, since the optimisation problem defining the estimator may be non strictly convex, the solution may not be unique. We then denote with  $\hat{\mathbf{W}}$  the unique least  $\ell_2$  norm solution of equation (18).

In the most general case, the statistical properties of  $\hat{\mathbf{W}}$  are captured by a finite set of finite-dimensional order parameters, namely  $\mathbf{V}, \hat{\mathbf{V}}, \mathbf{Q}, \hat{\mathbf{Q}} \in \mathbb{R}^{K \times K}$  and  $\mathbf{m}, \hat{\mathbf{m}} \in \mathbb{R}^K$ . These

Fluctuations, bias, variance and ensemble of learners: exact asymptotics for convex losses in high-dimension order parameters satisfy a set of fixed-point equations. To avoid a proliferation of indices in our formulas, let us introduce some notation. Let  $\mathbf{A} = (A_{kk'}^{ij})_{k,k' \in [K]}^{i,j \in [p]} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}$  be a tensor, and  $\mathbf{X} = (X_k^i)_{k \in [K]}^{i \in [p]}$ ,  $\mathbf{Y} = (Y_k^i)_{k \in [K]}^{i \in [p]}$ ,  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{p \times K}$  two matrices. We will denote

$$\langle\langle \mathbf{A} \rangle\rangle := \left( \sum_i A_{kk'}^{ii} \right)_{kk'} \in \mathbb{R}^{K \times K}, \tag{19a}$$

$$\langle\langle \mathbf{X} | \mathbf{A} | \mathbf{Y} \rangle\rangle := \left( \sum_{ij} X_k^i A_{kk'}^{ij} Y_{k'}^j \right)_{kk'} \in \mathbb{R}^{K \times K}, \tag{19b}$$

$$\langle\langle \mathbf{X} | \mathbf{Y} \rangle\rangle := \left( \sum_{ij} X_k^i Y_k^i \right)_k \in \mathbb{R}^K, \tag{19c}$$

$$\langle \mathbf{X} | \mathbf{A} | \mathbf{Y} \rangle := \sum_{ijk} X_k^i A_{kk'}^{ij} Y_k^j \in \mathbb{R} \tag{19d}$$

$$\langle \mathbf{X} | \mathbf{Y} \rangle := \sum_{ik} X_k^i Y_k^i \in \mathbb{R}. \tag{19e}$$

Given a second tensor  $\mathbf{B} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}$ , we write

$$\mathbf{A} \mathbf{B} := \left( \sum_{i'k} A_{kk'}^{ii'} B_{kk'}^{i'j} \right)_{kk'}^{ij} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}, \tag{19f}$$

$$\mathbf{A} \circ \mathbf{B} := \left( \sum_{i'} A_{kk'}^{ii'} B_{k'k}^{i'j} \right)_{kk'}^{ij} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}, \tag{19g}$$

$$\mathbf{A} \odot \mathbf{B} := \left( A_{kk'}^{ij} B_{kk'}^{ij} \right)_{kk'}^{ij} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}. \tag{19h}$$

We can now state our general result.

**Theorem 2.** *Let us consider the random quantities  $\boldsymbol{\xi} \in \mathbb{R}^K$  and  $\boldsymbol{\Xi} \in \mathbb{R}^{K \times K}$  with entries distributed as  $\mathcal{N}(0, 1)$ . Assume that in the high-dimensional limit where  $d, p, n \rightarrow +\infty$  with  $\alpha := n/p$  and  $\gamma := d/p$  kept  $\Theta(1)$  constants. Then in this limit, for any pseudo-Lispchitz functions of order 2  $\varphi: \mathbb{R} \times \mathbb{R}^K \rightarrow \mathbb{R}$  and  $\tilde{\varphi}: \mathbb{R}^{K \times p} \rightarrow \mathbb{R}$ , the estimator  $\hat{\mathbf{W}}$  verifies*

$$\begin{aligned} \mathbb{E}_{(y, \mathbf{x})} \left[ \varphi \left( y, \frac{\langle\langle \hat{\mathbf{W}} | \mathbf{U} \rangle\rangle}{\sqrt{p}} \right) \right] &\xrightarrow{P} \int_{\mathbf{y}} d\mathbf{y} \mathbb{E}_{(\nu, \boldsymbol{\mu})} [P_y^0(y|\nu) \varphi(y, \boldsymbol{\mu})], \\ \frac{1}{n} \sum_{\mu=1}^n \varphi \left( y^\mu, \frac{\langle\langle \hat{\mathbf{W}} | \mathbf{U}^\mu \rangle\rangle}{\sqrt{p}} \right) &\xrightarrow{P} \int_{\mathbf{y}} d\mathbf{y} \mathbb{E}_{\boldsymbol{\xi}} [\mathcal{Z}^0(y, \omega_0, \sigma_0) \varphi(y, \mathbf{h})], \\ \tilde{\varphi}(\hat{\mathbf{W}}) &\xrightarrow{P} \mathbb{E}_{\boldsymbol{\Xi}} [\tilde{\varphi}(\mathbf{G})], \end{aligned} \tag{20}$$



Fluctuations, bias, variance and ensemble of learners: exact asymptotics for convex losses in high-dimension where  $\mathbf{U} \equiv \mathbf{U}(\mathbf{x})$ ,  $(\nu, \boldsymbol{\mu}) \in \mathbb{R}^{1+K}$  are jointly Gaussian random variables with zero mean and covariance matrix

$$(\nu, \boldsymbol{\mu}) \sim \mathcal{N}\left(\mathbf{0}_{1+K}, \begin{pmatrix} \rho & \mathbf{m}^\top \\ \mathbf{m} & \mathbf{Q} \end{pmatrix}\right), \tag{21}$$

and we have introduced the proximals for the loss and the regularisation:

$$\begin{aligned} \mathbf{h} &:= \arg \min_{\mathbf{u}} \left[ \frac{(\mathbf{u} - \boldsymbol{\omega})^\top \mathbf{V}^{-1} (\mathbf{u} - \boldsymbol{\omega})}{2} + \hat{\ell}(y, \mathbf{u}) \right], \\ \mathbf{G} &:= \arg \min_{\mathbf{U}} \left[ \frac{\langle \mathbf{U} | (\mathbf{1}_{p,p} \otimes \hat{\mathbf{V}}) \odot \Omega | \mathbf{U} \rangle}{2} - \langle \mathbf{B} | \mathbf{U} \rangle + \lambda r(\mathbf{U}) \right], \end{aligned} \tag{22}$$

with  $\boldsymbol{\omega} := \mathbf{Q}^{1/2} \boldsymbol{\xi}$  and  $\mathbf{B} := (\mathbf{1}_p \otimes \hat{\mathbf{m}}^\top) \odot \hat{\boldsymbol{\Phi}} + ((\mathbf{1}_{p,p} \otimes \hat{\mathbf{Q}}) \odot \Omega)^{\frac{1}{2}} \boldsymbol{\Xi}$ . We have also introduced the auxiliary function

$$\mathcal{Z}^0(y, \mu, \sigma) := \int \frac{P_y^0(y|x) dx}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma}}. \tag{23}$$

and the scalar quantities  $\omega_0 := \mathbf{m}^\top \mathbf{Q}^{-1/2} \boldsymbol{\xi}$  and  $\sigma_0 := \rho - \mathbf{m}^\top \mathbf{Q}^{-1} \mathbf{m}$ . The order parameters satisfy the saddle-point equations

$$\begin{aligned} \hat{\mathbf{V}} &= -\alpha \int_y dy \mathbb{E}_\xi [\mathcal{Z}^0(y, \omega_0, \sigma_0) \partial_\omega \mathbf{f}], \\ \hat{\mathbf{Q}} &= \alpha \int_y dy \mathbb{E}_\xi [\mathcal{Z}^0(y, \omega_0, \sigma_0) \mathbf{f} \mathbf{f}^\top], \\ \hat{\mathbf{m}} &= \frac{\alpha}{\sqrt{\gamma}} \int_y dy \mathbb{E}_\xi [\partial_\mu \mathcal{Z}^0(y, \omega_0, \sigma_0) \mathbf{f}], \end{aligned} \tag{24}$$

and

$$\begin{aligned} \mathbf{V} &= \frac{2}{p} \mathbb{E}_\Xi \left\langle \mathbf{G} \left| \frac{D\left(\left(\mathbf{1}_{p,p} \otimes \hat{\mathbf{Q}}\right) \odot \Omega\right)^{1/2}}{D\hat{\mathbf{Q}}} \right| \Xi \right\rangle \\ \mathbf{Q} &= \frac{1}{p} \mathbb{E}_\Xi \langle \langle \mathbf{G} | \Omega | \mathbf{G} \rangle \rangle, \\ \mathbf{m} &= \frac{1}{\sqrt{\gamma p}} \mathbb{E}_\Xi \langle \langle \hat{\boldsymbol{\Phi}} | \mathbf{G} \rangle \rangle. \end{aligned} \tag{25}$$

In the equation above we have introduced the short-hand notation  $\mathbf{f} := \mathbf{V}^{-1}(\mathbf{h} - \boldsymbol{\omega})$ .

In the theorem above, for a tensor  $\hat{\mathbf{A}} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}$ , then  $[\frac{D\hat{\mathbf{A}}}{D\hat{\mathbf{Q}}}]_{ij}^{kk', \kappa\kappa'} \equiv \frac{\partial \hat{A}_{ij}^{kk'}}{\partial \hat{Q}_{\kappa\kappa'}}$ : in the formula, the contractions involve Latin indices only. Equations (24) are typically called *channel equations*, because depend on the form of the loss  $\hat{\ell}$ . Equations (25), instead, are usually called *prior equations*, because of their dependence on the prior,

Fluctuations, bias, variance and ensemble of learners: exact asymptotics for convex losses in high-dimension i.e.  $r$ . In the following Corollary, we specify their expression for a ridge regularisation,  $r(\mathbf{W}) = \frac{1}{2} \|\mathbf{W}\|_{\mathbb{F}}^2$ .

**Corollary 3 (ridge regularisation).** *In the hypotheses of theorem 2, if  $r(\mathbf{W}) = \frac{1}{2} \|\mathbf{W}\|_{\mathbb{F}}^2$ , then the prior equations are*

$$\begin{aligned} \mathbf{V} &= \frac{1}{p} \langle\langle \Omega \circ \mathbf{A} \rangle\rangle, \\ \mathbf{Q} &= \frac{1}{p} \langle\langle \Omega \circ \left( \mathbf{A} \left( (\mathbf{1}_{p,p} \otimes \hat{\mathbf{m}} \otimes \hat{\mathbf{m}}^\top) \odot \Theta + (\mathbf{1}_{p,p} \otimes \hat{\mathbf{Q}}) \odot \Omega \right) \mathbf{A} \right) \rangle\rangle, \\ \mathbf{m} &= \frac{1}{\sqrt{\gamma p}} \langle\langle \mathbf{A} \left( (\mathbf{1}_{p,p} \otimes \hat{\mathbf{m}} \otimes \mathbf{1}_K^\top) \odot \Theta \right) \rangle\rangle. \end{aligned} \tag{26}$$

In the equation above, we have used the auxiliary tensor  $\mathbf{A} \equiv \mathbf{A}(\hat{\mathbf{V}}; \lambda, \Omega) := (\lambda \mathbf{I}_p \otimes \mathbf{I}_K + (\mathbf{1}_{p,p} \otimes \hat{\mathbf{V}}) \odot \Omega)^{-1} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}$ .

#### 4.1. The random feature case and the kernel limit

Theorem 2 is given in a very general setting, and, in particular, no assumptions are made on the features  $\mathbf{u}_k$ . We have anticipated in section 2 that, in the case of RFs, the structure of the order parameters highly simplifies and the covariance matrix  $\Sigma$  is fully specified by only three scalar order parameters for any  $K > 1$ . Here will adapt therefore theorem 2 to the RF setting in section 2, using the notation therein. The motivation of this section is to explicitly present the self-consistent equations that are required to produce the results given in the paper.

**Corollary 4.** *Assume that in the high-dimensional limit where  $d, p, n \rightarrow +\infty$  with  $\alpha := n/p$  and  $\gamma := d/p$  kept  $\Theta(1)$  constants, the Wishart matrix  $\mathbf{F}\mathbf{F}^\top$  has a well-defined asymptotic spectral distribution. Then in this limit, for any pseudo-Lispchitz function of finite order  $\varphi: \mathbb{R} \times \mathbb{R}^K \rightarrow \mathbb{R}$ , the estimator  $\hat{\mathbf{W}}$  verifies*

$$\mathbb{E}_{(x,y)} \left[ \varphi \left( y, \frac{\langle\langle \hat{\mathbf{W}} | \mathbf{U} \rangle\rangle}{\sqrt{p}} \right) \right] \xrightarrow{\mathbb{P}} \mathbb{E}_{(\nu, \boldsymbol{\mu})} [\varphi(f_0(\nu), \boldsymbol{\mu})], \tag{27}$$

where  $(\nu, \boldsymbol{\mu}) \in \mathbb{R}^{K+1}$  is a jointly Gaussian vector with covariance

$$(\nu, \boldsymbol{\mu}) \sim \mathcal{N} \left( \mathbf{0}_{K+1}, \begin{pmatrix} \rho & m \mathbf{1}_K^\top \\ m \mathbf{1}_K & \mathbf{Q} \end{pmatrix} \right), \tag{28}$$

and  $\mathbf{Q} := (q_0 - q_1) \mathbf{I}_K + q_1 \mathbf{1}_{K,K}$ . The collection of parameters  $(q_0, q_1, m)$  is obtained solving a set of fixed point equations involving the auxiliary variables  $(\hat{q}_0, \hat{q}_1, \hat{m}, \hat{\nu}, \hat{\nu})$ , namely:

$$\hat{\nu} = -\alpha \int_y dy \mathbb{E}_\omega \left[ \mathcal{Z}^0 \left( y, \frac{m\omega}{q_0}, \rho - \frac{m^2}{q_0} \right) \partial_\omega f \right], \tag{29a}$$

$$\hat{m} = \frac{\alpha}{\sqrt{\gamma}} \int_y dy \mathbb{E}_\omega \left[ \partial_\mu \mathcal{Z}^0 \left( y, \frac{m\omega}{q_0}, \rho - \frac{m^2}{q_0} \right) f \right], \tag{29b}$$

$$\hat{q}_0 = \alpha \int_{\mathcal{Y}} dy \mathbb{E}_{\omega} \left[ \mathcal{Z}^0 \left( y, \frac{m\omega}{q_0}, \rho - \frac{m^2}{q_0} \right) f^2 \right], \quad (29c)$$

$$\hat{q}_1 = \alpha \int_{\mathcal{Y}} dy \mathbb{E}_{\omega, \omega'} \left[ \mathcal{Z}^0 \left( y, m \frac{\omega + \omega'}{q_0 + q_1}, \rho - \frac{2m^2}{q_0 + q_1} \right) f f' \right], \quad (29d)$$

$$v = \int \frac{s \varrho(s) ds}{\lambda + s \hat{v}}, \quad (29e)$$

$$m = \frac{\hat{m}}{\sqrt{\gamma}} \int \frac{s - \kappa_*^2}{\lambda + \hat{v}s} \varrho(s) ds, \quad (29f)$$

$$q_0 = \int \frac{(\hat{q}_0 + \hat{m}^2) s^2 - \hat{m}^2 \kappa_*^2 s}{(\lambda + \hat{v}s)^2} \varrho(s) ds, \quad (29g)$$

$$q_1 = \left( 1 + \frac{\hat{q}_1}{\hat{m}^2} \right) m^2. \quad (29h)$$

where  $\omega$  and  $\omega'$  are two correlated Gaussian random variables of zero mean and  $\mathbb{E}[\omega^2] = \mathbb{E}[\omega'^2] = q_0$ ,  $\mathbb{E}[\omega\omega'] = q_1$ . Moreover, we have introduced the proximals

$$f = \frac{\text{Prox}_{v\ell(y, \bullet)}(\omega) - \omega}{v}, \quad f' = \frac{\text{Prox}_{v\ell(y, \bullet)}(\omega') - \omega'}{v}, \quad (30)$$

with

$$\text{Prox}_{v\ell(y, \bullet)}(\omega) := \arg \min_x \left[ \frac{(x - \omega)^2}{2v} + \ell(y, x) \right]. \quad (31)$$

Finally,  $\varrho(s)$  is the asymptotic spectral density of the features covariance matrix

$\mathbf{\Omega} \equiv \text{Var}(\mathbf{u}) = \kappa_0^2 \mathbf{1}_{p,p} + \frac{\kappa_1^2}{d} \mathbf{F} \mathbf{F}^\top + \kappa_*^2 \mathbf{I}_p$  and the coefficients are given by  $\kappa_0 := \mathbb{E}_{\zeta}[\phi(\zeta)]$ ,  $\kappa_1 := \mathbb{E}_{\zeta}[\zeta \phi(\zeta)]$ ,  $\kappa_* := \mathbb{E}_{\zeta}[\phi^2(\zeta)] - \kappa_0^2 - \kappa_1^2$  with  $\zeta \sim \mathcal{N}(0, 1)$ .

The previous corollary recovers the results of Dhifallah and Lu (2020), Gerace *et al* (2020), and Loureiro *et al* (2021b) when restricted to the  $K = 1$  case by marginalisation.

#### 4.2. The kernel limit

The so-called kernel limit is obtained by taking the limit of infinite number of parameters so that  $\gamma \rightarrow 0$  (i.e.  $p \gg d$  and  $p \gg n$ ), but with a fixed ratio  $\alpha/\gamma = n/d$ . To balance the loss term and the regularisation it is convenient to rescale  $\lambda \mapsto \alpha\lambda$ . We can simplify the saddle-point equation in this special limit introducing  $\hat{q}_0 \mapsto \alpha\hat{q}_0$ ,  $\hat{q}_1 \mapsto \alpha\hat{q}_1$ ,  $\hat{m} \mapsto \sqrt{\alpha}\hat{m}$ ,  $\hat{v} \mapsto \alpha\hat{v}$ . The channel equations keep a simple form,

$$\hat{v} = - \int_{\mathcal{Y}} dy \mathbb{E}_{\zeta} \left[ \mathcal{Z}^0(y, \omega_0, \sigma_0) \partial_{\omega} f \right], \quad (32a)$$

$$\hat{m} = \sqrt{\delta} \int_{\mathcal{Y}} dy \mathbb{E}_{\zeta} \left[ f \partial_{\mu} \mathcal{Z}^0(y, \omega_0, \sigma_0) \right], \quad (32b)$$

$$\hat{q}_0 = \int_{\mathcal{Y}} dy \mathbb{E}_{\zeta} [\mathcal{Z}^0(y, \omega_0, \sigma_0) f^2], \quad (32c)$$

$$\hat{q}_1 = \int_{\mathcal{Y}} dy \mathbb{E}_{\omega, \omega'} \left[ \mathcal{Z}^0 \left( y, m \frac{\omega + \omega'}{q_0 + q_1}, \rho - \frac{2m^2}{q_0 + q_1} \right) f f' \right]. \quad (32d)$$

The  $p \rightarrow +\infty$  limit in the prior equations depends on the spectral density  $\varrho(s)$ . For example, if  $\mathbf{F}$  has random Gaussian entries with zero mean and unit variance, then  $\varrho(s)$  is a shifted Marchenko–Pastur distribution,

$$\varrho(s) = \nu(s - \kappa_*^2; \alpha^{-1} \delta, \kappa_1), \quad (33)$$

where, if  $[x]_+ = x\theta(x)$ ,

$$\nu(x; b, a) = \frac{\sqrt{[(a_+ - x)(x - a_-)]_+}}{2ab^2\pi^2x} + \left[1 - \frac{1}{a}\right]_+ \delta(x), \quad (34)$$

with  $a_{\pm} := b^2(1 \pm \sqrt{a})^2$ . By means of a series of algebraic manipulation, we obtain in the end at the first order in  $\alpha$

$$\begin{aligned} v &= \frac{\lambda(\kappa_1^2 + \kappa_*^2) + \delta^2 \kappa_1^2 \kappa_*^2 \hat{v}}{\lambda(\lambda + \delta \kappa_1^2 \hat{v})}, & q_0 &= \frac{\delta \kappa_1^4 (\hat{q}_0 + \delta \hat{m}^2)}{(\lambda + \delta \kappa_1^2 \hat{v})^2}, \\ m &= \frac{\sqrt{\delta} \kappa_1^2 \hat{m}}{\lambda + \delta \kappa_1^2 \hat{v}}, & q_1 &= \frac{\delta \kappa_1^4 (\hat{q}_1 + \delta \hat{m}^2)}{(\lambda + \delta \kappa_1^2 \hat{v})^2}, \end{aligned} \quad (35)$$

which complete our set of equations for the kernel limit.

## Acknowledgments

We thank Ali Beryhi, Giulio Biroli, Stéphane d’Ascoli, Justin Ko for discussions, and Francesca Mignacco for sharing code with us. We acknowledge funding from the French National Research Agency Grants ANR-17-CE23-0023-01 PAIL and ANR-19P3IA-0001 PRAIRIE.

## Appendix A. The replica approach

### A.1. Notation

We introduce here some notation that will help us to keep the expressions in this appendix more compact. Given two tensors  $\mathbf{A} = (A_{kk'}^{ij})_{kk'}^{ij} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}$  and  $\mathbf{B} = (B_{kk'}^{ij})_{kk'}^{ij} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}$ ,  $i, j \in [p]$ ,  $k, k' \in [K]$ , then

$$\hat{\mathbf{C}} = \mathbf{A}\mathbf{B} \Leftrightarrow \hat{C}_{kk'}^{ij} := \sum_{r, \kappa} A_{k\kappa}^{ir} B_{\kappa k'}^{rj} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K} \quad (36)$$

$$\mathbf{C} = \mathbf{A} \odot \mathbf{B} \Leftrightarrow C_{kk'}^{ij} := A_{kk'}^{ij} B_{kk'}^{ij} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K} \quad (37)$$

$$\tilde{\mathbf{C}} = \mathbf{A} \circ \mathbf{B} \Leftrightarrow \tilde{C}_{kk'}^{ij} := \sum_r A_{kk'}^{ir} B_{k'k}^{rj} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}. \quad (38)$$

Also, if  $\mathbf{X} = (X_k^i)_k^i \in \mathbb{R}^{p \times K}$  and  $\mathbf{Y} = (Y_k^i)_k^i \in \mathbb{R}^{p \times K}$ , we write

$$\mathbf{X} \otimes \mathbf{Y} = \left( X_k^i Y_{k'}^j \right)_{kk'}^{ij} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K} \quad (39)$$

$$\mathbf{A}|\mathbf{X} = \left( \sum_{j\kappa} A_{k\kappa}^{ij} X_k^j \right)_k^i \in \mathbb{R}^{p \times K} \quad (40)$$

$$\langle \mathbf{X}|\mathbf{Y} \rangle = \sum_{ik} X_k^i Y_k^i \in \mathbb{R} \quad (41)$$

$$\langle \mathbf{X}|\mathbf{A}|\mathbf{Y} \rangle = \sum_{ijkk'} X_k^i A_{kk'}^{ij} Y_{k'}^j \in \mathbb{R}, \quad (42)$$

$$\langle\langle \mathbf{X}|\mathbf{Y} \rangle\rangle = \left( \sum_{ij} X_k^i Y_k^j \right)_k \in \mathbb{R}^K \quad (43)$$

$$\langle\langle \mathbf{X}|\mathbf{A}|\mathbf{Y} \rangle\rangle = \left( \sum_{ij} X_k^i A_{kk'}^{ij} Y_{k'}^j \right)_{kk'} \in \mathbb{R}^{K \times K}, \quad (44)$$

$$\langle\langle \mathbf{A} \rangle\rangle = \left( \sum_i A_{kk'}^{ii} \right)_{kk'} \in \mathbb{R}^{K \times K}. \quad (45)$$

In other words, the double brackets  $\langle\langle \bullet \rangle\rangle$  express the contraction of the upper indices only. This means for example that  $\langle\langle \mathbf{X}|\mathbf{Y} \rangle\rangle = \text{diag}(\mathbf{X}^\top \mathbf{Y}) \in \mathbb{R}^K$ . Finally, if  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^K$  and  $\mathbf{A} \in \mathbb{R}^{K \times K}$ ,  $\langle \mathbf{u}|\mathbf{A}|\mathbf{v} \rangle := \mathbf{u}^\top \mathbf{A} \mathbf{v} = \sum_{kk'} u_k A_{kk'} v_{k'} \in \mathbb{R}$ . We will adopt the same notation in the simple  $K = 1$  case.

### A.2. The replica computation

The replica computation relies on the treatment of a Gibbs measure over the weights  $\mathbf{W}$  which concentrates on the weights  $\hat{\mathbf{W}}$  that minimise a certain loss  $\hat{\ell}$  when a fictitious ‘inverse temperature’ parameter is sent to infinity. Such measure reads

$$\mu_\beta(\mathbf{W}) := \frac{P_w(\mathbf{W})}{\mathcal{Z}(\beta)} \prod_{\mu=1}^n \exp \left[ -\beta \hat{\ell} \left( y^\mu, \frac{\langle\langle \mathbf{W}|\mathbf{U}^\mu \rangle\rangle}{\sqrt{p}} \right) \right], \quad (46)$$

$$\mathcal{Z}(\beta) := \int d\mathbf{W} P_w(\mathbf{W}) \prod_{\mu=1}^n \exp \left[ -\beta \hat{\ell} \left( y^\mu, \frac{\langle\langle \mathbf{W}|\mathbf{U}^\mu \rangle\rangle}{\sqrt{p}} \right) \right], \quad (47)$$

where  $P_w(\mathbf{W}) = e^{-\beta \lambda r(\mathbf{W})}$  is the prior on the weights  $\mathbf{W} = (W_k^i)_{k \in [K]}^{i \in [p]} \in \mathbb{R}^{p \times K}$ , possibly containing the regularisation. The dataset  $(y^\mu, \mathbf{U}^\mu)_\mu$  is obtained from a set of  $n$  samples  $\mathbf{x}^\mu \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ ,  $\mu \in [n]$ . For each  $\mu$ , the label  $y^\mu$  has distribution  $P_y^0(y|d^{-1/2} \langle \boldsymbol{\theta}|\mathbf{x}^\mu \rangle)$  for some fixed  $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}_d, \rho \mathbf{I}_d)$ . The array of features  $\mathbf{U}^\mu$ , instead, is obtained as function of the vector  $\mathbf{x}^\mu$  via a law  $\mathbf{U}: \mathbb{R}^d \rightarrow \mathbb{R}^{p \times K}$  such that  $\mathbf{U}^\mu := \mathbf{U}(\mathbf{x}^\mu) \in \mathbb{R}^{p \times K}$ . As we will show below, the tensors

Fluctuations, bias, variance and ensemble of learners: exact asymptotics for convex losses in high-dimension

$$\Omega := \mathbb{E}_{\mathbf{x}} [\mathbf{U}(\mathbf{x}) \otimes \mathbf{U}(\mathbf{x})] \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}, \quad (48)$$

$$\hat{\Phi} := \mathbb{E}_{\mathbf{x}} [\mathbf{U}(\mathbf{x}) \langle \mathbf{x} | \boldsymbol{\theta} \rangle] \in \mathbb{R}^{p \times K}, \quad (49)$$

will incorporate the information about the action of the law  $\mathbf{U}$ . We denote for brevity

$$P_y(y|\mathbf{u}) \propto \exp \left[ -\beta \hat{\ell}(y, \mathbf{u}) \right], \quad (50)$$

and we proceed computing the *free entropy*  $\Phi := \mathbb{E}_{(y^\mu, \mathbf{x}^\mu)_\mu} [\ln \mathcal{Z}(\beta)]$  using the replica trick, i.e. the fact that  $\mathbb{E}[\ln \mathcal{Z}(\beta)] = \lim_{s \rightarrow 0} \frac{1}{s} \ln \mathbb{E}[\mathcal{Z}^s(\beta)]$

$$\begin{aligned} \mathbb{E}_{(y^\mu, \mathbf{x}^\mu)_\mu} [\mathcal{Z}^s(\beta)] &= \prod_{a=1}^s \int d\mathbf{W}^a P_w(\mathbf{W}^a) \left( \mathbb{E}_{(y, \mathbf{x})} \left[ P_y^0 \left( y \middle| \frac{\langle \mathbf{x} | \boldsymbol{\theta} \rangle}{\sqrt{d}} \right) \right. \right. \\ &\quad \left. \left. \times \prod_{a=1}^s P_y \left( y \middle| \frac{\langle \mathbf{W}^a | \mathbf{U}(\mathbf{x}) \rangle}{\sqrt{p}} \right) \right] \right)^n. \end{aligned} \quad (51)$$

Denoting by  $\boldsymbol{\mu}^a \equiv (\mu_k^a)_{k \in [K]}$ , if we now consider

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}} \left[ P_y^0 \left( y \middle| \frac{\langle \mathbf{x} | \boldsymbol{\theta} \rangle}{\sqrt{d}} \right) \prod_{a=1}^s P_y \left( y \middle| \frac{\langle \mathbf{W}^a | \mathbf{U}(\mathbf{x}) \rangle}{\sqrt{d}} \right) \right] \\ &= \int d\nu P_y^0(y|\nu) \prod_{a=1}^s \left[ \int d\boldsymbol{\mu}^a P_y(y|\boldsymbol{\mu}^a) \right] \underbrace{\mathbb{E}_{\mathbf{x}} \left[ \delta \left( \nu - \frac{\langle \mathbf{x} | \boldsymbol{\theta} \rangle}{\sqrt{d}} \right) \prod_{a=1}^s \delta \left( \boldsymbol{\mu}^a - \frac{\langle \mathbf{W}^a | \mathbf{U}(\mathbf{x}) \rangle}{\sqrt{p}} \right) \right]}_{P(\nu, \boldsymbol{\mu})}. \end{aligned} \quad (52)$$

We apply now the *GEP* (Goldt *et al* 2021), i.e. we assume that  $P(\nu, \boldsymbol{\mu})$  is a Gaussian with covariance matrix given by

$$\Sigma(\mathbf{W}) = \begin{pmatrix} \rho & \mathbf{m}^\top \\ \mathbf{m} & \mathbf{Q} \end{pmatrix}, \quad (53)$$

where  $\mathbf{m} = (\mathbf{m}^a)_{a \in [s]} \in \mathbb{R}^{sK}$  and  $\mathbf{Q} = (\mathbf{Q}^{ab})_{a, b \in [s]} \in \mathbb{R}^{sK \times sK}$ . Here, for each  $a, b \in [n]$ ,  $\mathbf{m}^a \in \mathbb{R}^K$  and  $\mathbf{Q}^{ab} \in \mathbb{R}^{K \times K}$  and are defined as

$$\rho := \mathbb{E}[\nu^2] = \frac{\|\boldsymbol{\theta}\|_2^2}{d}, \quad (54)$$

$$\mathbf{m}^a := \mathbb{E}[\boldsymbol{\mu}^a \nu] = \frac{\langle \mathbf{W}^a | \hat{\Phi} \rangle}{\sqrt{pd}}, \quad (55)$$

$$\mathbf{Q}^{ab} := \mathbb{E}[\boldsymbol{\mu}^a \boldsymbol{\mu}^{b\top}] = \frac{\langle \mathbf{W}^a | \Omega | \mathbf{W}^b \rangle}{p}. \quad (56)$$

In the end

$$\mathbb{E}_{(y^\mu, \mathbf{x}^\mu)_\mu} [\mathcal{Z}^s(\beta)] = \mathbb{E}_{\boldsymbol{\theta}} \left[ \prod_{a=1}^s \int d\mathbf{W}^a P_w(\mathbf{W}^a) \left( \int dy \int d\nu P_y^0(y|\nu) \right. \right.$$

$$\begin{aligned} & \times \prod_{a=1}^s \left[ \int d\boldsymbol{\mu}^a P_y(y|\boldsymbol{\mu}^a) \right] P(\nu, \boldsymbol{\mu}) \Big)^n \\ & = \left( \prod_{a=1}^s \iint \frac{d\mathbf{m}^a d\hat{\mathbf{m}}^a}{(2\pi)^K} \right) \left( \prod_{a=1}^s \iint \frac{d\mathbf{Q}^{ab} d\hat{\mathbf{Q}}^{ab}}{(2\pi)^{K^2}} \right) \\ & \times \int \frac{d\rho d\hat{\rho}}{2\pi} \exp \left( p\Phi^{(s)} \left( \rho, \mathbf{m}, \mathbf{Q}, \hat{\rho}, \hat{\mathbf{m}}, \hat{\mathbf{Q}} \right) \right). \end{aligned} \tag{57}$$

Absorbing the factor  $-i$  in the integrals, and denoting by  $n/p = \alpha$  and  $d/p = \gamma$ ,

$$\begin{aligned} \Phi^{(s)} \left( \rho, \mathbf{m}, \mathbf{Q}, \hat{\rho}, \hat{\mathbf{m}}, \hat{\mathbf{Q}} \right) & = -\gamma\rho\hat{\rho} - \sqrt{\gamma} \sum_{a=1}^s \langle \hat{\mathbf{m}}^a | \mathbf{m}^a \rangle - \sum_{a \leq b} \langle \hat{\mathbf{Q}}^{ab} | \mathbf{Q}^{ab} \rangle \\ & + \alpha \Psi_y^{(s)}(\rho, \mathbf{m}, \mathbf{Q}) + \Psi_{\mathbf{w}}^{(s)}(\hat{\rho}, \hat{\mathbf{m}}, \hat{\mathbf{Q}}). \end{aligned} \tag{58}$$

Here we have introduced

$$\begin{aligned} \Psi_y^{(s)}(\rho, \mathbf{m}, \mathbf{Q}) & := \ln \left[ \int dy \int d\nu P_y^0(y|\nu) \prod_{a=1}^s \left[ \int d\boldsymbol{\mu}^a P_y(y|\boldsymbol{\mu}^a) \right] P(\nu, \boldsymbol{\mu}) \right] \\ \Psi_{\mathbf{w}}^{(s)}(\hat{\rho}, \hat{\mathbf{m}}, \hat{\mathbf{Q}}) & := \frac{1}{p} \ln \left[ e^{\hat{\rho} \|\boldsymbol{\theta}\|_2^2} \left( \prod_{a=1}^s \int P(\mathbf{W}^a) e^{\langle (\mathbf{1}_p \otimes \hat{\mathbf{m}}^{a\tau}) \odot \hat{\Phi} | \mathbf{W}^a \rangle} \right) \right. \\ & \left. \times \exp \left( \sum_{a \leq b} \langle \mathbf{W}^a | (\mathbf{1}_{p,p} \otimes \hat{\mathbf{Q}}^{ab}) \odot \Omega | \mathbf{W}^b \rangle \right) \right] \end{aligned} \tag{59}$$

so that in the high dimensional limit the desired average is the extremum of the functional  $1s\Phi^{(s)}$  in the  $s \rightarrow 0$  limit,

$$\mathbb{E}_{(y^\mu, \mathbf{x}^\mu)_\mu} [\ln \mathcal{Z}(\beta)] = \lim_{s \rightarrow 0} \frac{1}{s} \text{Ext} \Phi^{(s)} \left( \rho, \mathbf{m}, \mathbf{Q}, \hat{\rho}, \hat{\mathbf{m}}, \hat{\mathbf{Q}} \right). \tag{60}$$

*A.2.1. Replica symmetric ansatz.* In order to take the limit, let us assume as usual a *replica symmetric* (RS) ansatz:

$$\begin{aligned} \mathbf{m}^a & \equiv \mathbf{m} \quad a \in [s], & \hat{\mathbf{m}}^a & \equiv \hat{\mathbf{m}} \quad a \in [s], \\ \mathbf{Q}^{ab} & \equiv \begin{cases} \mathbf{R} & \text{if } a = b, \\ \mathbf{Q} & \text{if } a \neq b, \end{cases} & \hat{\mathbf{Q}}^{ab} & \equiv \begin{cases} -\frac{1}{2}\hat{\mathbf{R}} & \text{if } a = b, \\ \hat{\mathbf{Q}} & \text{if } a \neq b. \end{cases} \end{aligned} \tag{61}$$

Observe that  $\lim_{s \rightarrow 0} \Phi^{(s)} = 0$  by construction, meaning that  $\hat{\rho} = 0$  fixing  $\rho = \frac{1}{d} \mathbb{E}_{\boldsymbol{\theta}} [\|\boldsymbol{\theta}\|_2^2]$ . Before proceeding further, we note that the matrix  $\mathbf{Q}$  in the RS ansatz can be written as  $\mathbf{Q} = \mathbf{I}_s \otimes (\mathbf{R} - \mathbf{Q}) + \mathbf{1}_{s,s} \otimes \mathbf{Q}$ , where  $\mathbf{1}_{s,s}$  is the  $s \times s$  matrix of 1. Similarly,  $\mathbf{m} = \mathbf{1}_s \otimes \mathbf{m}$ , where  $\mathbf{1}_s$  is the column vector of  $s$  elements equal to 1. Following similar steps to the



Fluctuations, bias, variance and ensemble of learners: exact asymptotics for convex losses in high-dimension ones detailed, e.g. in Loureiro *et al* (2021b), we obtain

$$\begin{aligned} \Psi_y(\rho, \mathbf{m}, \mathbf{Q}, \mathbf{R}) &:= \lim_{s \rightarrow 0} \frac{1}{s} \Psi_y^{(s)}(\rho, \mathbf{m}, \mathbf{Q}) \\ &= \int_{\mathcal{Y}} dy \mathbb{E}_{\xi} \left[ \mathcal{Z}^0 \left( y, \langle \mathbf{m} | \mathbf{Q}^{-1/2} | \xi \rangle, \rho - \langle \mathbf{m} | \mathbf{Q}^{-1} | \mathbf{m} \rangle \right) \ln \mathcal{Z} \left( y, \sqrt{\mathbf{Q}} \xi, \mathbf{V} \right) \right], \end{aligned} \quad (62)$$

where  $\mathbf{V} := \mathbf{R} - \mathbf{Q}$ ,  $\xi \sim \mathcal{N}(\mathbf{0}_K, \mathbf{I}_K)$  and we have introduced

$$\mathcal{Z}(y, \mu, \Sigma) := \int \frac{P_y(y|\mu) d\mu}{\sqrt{\det(2\pi\Sigma)}} e^{-\frac{1}{2}\langle \mathbf{u} - \mu | \Sigma^{-1} | \mathbf{u} - \mu \rangle}, \quad (63)$$

$$\mathcal{Z}^0(y, \mu, \sigma) := \int \frac{P_y^0(y|x) dx}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma}}. \quad (64)$$

Similarly, defining  $\hat{\mathbf{V}} = \hat{\mathbf{R}} + \hat{\mathbf{Q}}$ , we can write down the prior channel. We can write then

$$\begin{aligned} \Psi_w(\hat{\mathbf{m}}, \hat{\mathbf{Q}}, \hat{\mathbf{R}}) &:= \lim_{s \rightarrow 0} \frac{1}{s} \Psi_w^{(s)}(0, \hat{\mathbf{m}}, \hat{\mathbf{Q}}) \\ &= \lim_{s \rightarrow 0} \frac{1}{sp} \ln \left[ \left( \prod_{a=1}^s \int d\mathbf{W}^a P_w(\mathbf{W}^a) e^{\langle (\mathbf{1}_p \otimes \hat{\mathbf{m}}^\top) \circ \hat{\Phi} | \mathbf{W}^a \rangle - \frac{1}{2} \langle \mathbf{W}^a | (\mathbf{1}_{p,p} \otimes \hat{\mathbf{R}}) \circ \Omega | \mathbf{W}^a \rangle} \right) \right. \\ &\quad \left. \times \prod_{a < b} e^{\langle \mathbf{W}^a | (\mathbf{1}_{p,p} \otimes \hat{\mathbf{Q}}) \circ \Omega | \mathbf{W}^b \rangle} \right] \\ &= \frac{1}{p} \mathbb{E}_{\Xi} \left[ \ln \left( \int d\mathbf{W} e^{-\lambda \beta r(\mathbf{W}) + \langle (\mathbf{1}_p \otimes \hat{\mathbf{m}}^\top) \circ \hat{\Phi} | \mathbf{W} \rangle + \langle \Xi | (\mathbf{1}_{p,p} \otimes \hat{\mathbf{Q}}) \circ \Omega | \mathbf{W} \rangle - \frac{1}{2} \langle \mathbf{W} | (\mathbf{1}_{p,p} \otimes \hat{\mathbf{V}}) \circ \Omega | \mathbf{W} \rangle} \right) \right] \end{aligned} \quad (65)$$

where we have performed a Hubbard–Stratonovich transformation and  $\Xi \equiv (\Xi_k^i)_k \in \mathbb{R}^{p \times K}$  has  $\Xi_k^i \sim \mathcal{N}(0, 1)$  for all  $i, k$ . The free entropy is then

$$\begin{aligned} \Phi &:= \lim_{s \rightarrow 0} \frac{1}{s} \Phi^{(s)} = -\sqrt{\gamma} \langle \hat{\mathbf{m}} | \mathbf{m} \rangle + \frac{\langle \hat{\mathbf{V}} | \mathbf{V} \rangle + \langle \hat{\mathbf{V}} | \mathbf{Q} \rangle - \langle \hat{\mathbf{Q}} | \mathbf{V} \rangle}{2} \\ &\quad + \alpha \Psi_y(\rho, \mathbf{m}, \mathbf{Q}, \mathbf{V}) + \Psi_w(\hat{\mathbf{m}}, \hat{\mathbf{Q}}, \hat{\mathbf{V}}). \end{aligned} \quad (66)$$

We are interested in the extremum of this quantity, and therefore we have to find the order parameters that maximise it by means of a set of saddle-point equations. Defining for brevity

$$\omega = \mathbf{Q}^{1/2} \xi, \quad \omega_0 = \langle \mathbf{m} | \mathbf{Q}^{-1} | \omega \rangle, \quad \sigma_0 = \rho - \langle \mathbf{m} | \mathbf{Q}^{-1} | \mathbf{m} \rangle, \quad (67)$$

Fluctuations, bias, variance and ensemble of learners: exact asymptotics for convex losses in high-dimension

a first set of saddle-point equation is

$$\hat{\mathbf{V}} = -\alpha \int_{\mathcal{Y}} dy \mathbb{E}_{\xi} [\mathcal{Z}^0(y, \omega_0, \sigma_0) \partial_{\omega} \mathbf{f}], \tag{68a}$$

$$\hat{\mathbf{Q}} = \alpha \int_{\mathcal{Y}} dy \mathbb{E}_{\xi} [\mathcal{Z}^0(y, \omega_0, \sigma_0) \mathbf{f} \mathbf{f}^{\top}], \tag{68b}$$

$$\hat{\mathbf{m}} = \frac{\alpha}{\sqrt{\gamma}} \int_{\mathcal{Y}} dy \mathbb{E}_{\xi} [\mathcal{Z}^0(y, \omega_0, \sigma_0) f^0 \mathbf{f}] \tag{68c}$$

where

$$f^0 \equiv \partial_{\omega_0} \ln \mathcal{Z}^0(y, \omega_0, \sigma_0) \tag{69}$$

$$\mathbf{f} \equiv \partial_{\omega} \ln \mathcal{Z}(y, \omega, \mathbf{V}). \tag{70}$$

*A.2.2. Zero temperature state evolution equations.* To obtain a nontrivial  $\beta \rightarrow +\infty$  limit we rescale  $\hat{\mathbf{V}} \mapsto \beta \hat{\mathbf{V}}$ ,  $\mathbf{V} \mapsto \beta^{-1} \mathbf{V}$ ,  $\hat{\mathbf{Q}} \mapsto \beta^2 \hat{\mathbf{Q}}$ ,  $\hat{\mathbf{m}} \mapsto \beta \hat{\mathbf{m}}$ . After this change of variable, equations (79) remain formally identical. To complete the set of saddle-point equations, let us observe that, defining

$$\mathcal{L}(y, \mathbf{u}) = \frac{1}{2} \langle \mathbf{u} - \omega | \mathbf{V}^{-1} | \mathbf{u} - \omega \rangle + \hat{\ell}(y, \mathbf{u}) \tag{71}$$

then after the rescaling

$$\ln \mathcal{Z}(y, \omega, \beta^{-1} \mathbf{V}) = \ln \int \frac{e^{-\beta \mathcal{L}(y, \mathbf{u})} d\mathbf{u}}{\sqrt{\det(2\pi \mathbf{V})}} \xrightarrow{\beta \gg 1} -\beta \mathcal{L}(y, \mathbf{h}) \quad \text{with } \mathbf{h} = \arg \min_{\mathbf{u}} \mathcal{L}(y, \mathbf{u}). \tag{72}$$

In this way the remaining saddle-point equations keep the form (68) but with

$$\mathbf{f} := \mathbf{V}^{-1} (\mathbf{h} - \omega). \tag{73}$$

In the  $\beta \rightarrow +\infty$  limit, we can write also

$$\begin{aligned} \Psi_{\omega}(\hat{\mathbf{m}}, \hat{\mathbf{Q}}, \hat{\mathbf{R}}) &= \frac{1}{p} \mathbb{E}_{\Xi} \left[ \ln \left( \int d\mathbf{W} e^{-\lambda \beta r(\mathbf{W}) + \beta \langle \mathbf{B} | \mathbf{W} \rangle - \frac{\beta}{2} \langle \mathbf{W} | (\mathbf{1}_{p,p} \otimes \hat{\mathbf{V}}) \odot \Omega | \mathbf{W} \rangle} \right) \right] \\ &= -\frac{\beta}{p} \mathbb{E}_{\Xi} \left[ \frac{\langle \mathbf{G} | (\mathbf{1}_{p,p} \otimes \hat{\mathbf{V}}) \odot \Omega | \mathbf{G} \rangle}{2} - \langle \mathbf{B} | \mathbf{G} \rangle + \lambda r(\mathbf{G}) \right] \end{aligned} \tag{74}$$

where

$$\mathbf{B} := (\mathbf{1}_p \otimes \hat{\mathbf{m}}^{\top}) \odot \hat{\Phi} + \left( (\mathbf{1}_{p,p} \otimes \hat{\mathbf{Q}}) \odot \Omega \right)^{1/2} \Xi \tag{75}$$

and

$$\mathbf{G} := \arg \min_{\mathbf{U}} \left[ \frac{\langle \mathbf{U} | (\mathbf{1}_{p,p} \otimes \hat{\mathbf{V}}) \odot \Omega | \mathbf{U} \rangle}{2} - \langle \mathbf{B} | \mathbf{U} \rangle + \lambda r(\mathbf{U}) \right]. \tag{76}$$

As a result, the remaining saddle point equations are

$$\mathbf{V} = \frac{2}{p} \mathbb{E}_{\Xi} \left\langle \mathbf{G} \left| \frac{\mathbb{D} \left( \left( \mathbf{1}_{p,p} \otimes \hat{\mathbf{Q}} \right) \odot \Omega \right)^{1/2}}{\mathbb{D} \hat{\mathbf{Q}}} \right| \Xi \right\rangle \quad (77a)$$

$$\mathbf{Q} = \frac{1}{p} \mathbb{E}_{\Xi} \langle \langle \mathbf{G} | \Omega | \mathbf{G} \rangle \rangle, \quad (77b)$$

$$\mathbf{m} = \frac{1}{\sqrt{\gamma p}} \mathbb{E}_{\Xi} \langle \langle \hat{\Phi} | \mathbf{G} \rangle \rangle. \quad (77c)$$

In the first equation, the derivative produce in general a 6-index tensor. Denoting  $\hat{\mathbf{A}} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K}$ , then  $[\frac{\mathbb{D} \hat{\mathbf{A}}}{\mathbb{D} \hat{\mathbf{Q}}}]_{ij}^{kk', \kappa \kappa'} \equiv \frac{\partial \hat{A}_{ij}^{kk'}}{\partial \hat{Q}_{\kappa \kappa'}}$ . In the formula, the contractions involve Latin indices only.

*A.2.3. The case of  $\ell_2$  regularisation.* If we assume an  $\hat{\ell}_2$  regularisation  $r(\mathbf{W}) = \frac{1}{2} \|\mathbf{W}\|_{\mathbb{F}}^2$ ,  $\Psi_{\mathbf{w}}$  is a Gaussian integral that can be explicitly computed before the rescaling in  $\beta$ . Denoting

$$\mathbf{A} := \left[ \lambda \mathbf{I}_K \otimes \mathbf{I}_p + \left( \mathbf{1}_{p,p} \otimes \hat{\mathbf{V}} \right) \odot \Omega \right]^{-1} \quad \text{and} \quad \Theta := \hat{\Phi} \otimes \hat{\Phi} \in \mathbb{R}^{p \times p} \otimes \mathbb{R}^{K \times K} \quad (78)$$

we obtain the following saddle-point equations for  $\mathbf{V}$ ,  $\mathbf{Q}$  and  $\mathbf{m}$ ,

$$\mathbf{V} = \frac{\langle \langle \mathbf{A} \odot \Omega \rangle \rangle}{p}, \quad (79a)$$

$$\mathbf{Q} = \frac{1}{p} \langle \langle \Omega \odot \left( \mathbf{A} \left( \left( \mathbf{1}_{p,p} \otimes \hat{\mathbf{m}} \otimes \hat{\mathbf{m}}^{\top} \right) \odot \Theta + \left( \mathbf{1}_{p,p} \otimes \hat{\mathbf{Q}} \right) \odot \Omega \right) \mathbf{A} \right) \rangle \rangle \quad (79b)$$

$$\mathbf{m} = \frac{1}{\sqrt{\gamma p}} \langle \langle \mathbf{A} \left( \left( \mathbf{1}_{p,p} \otimes \hat{\mathbf{m}} \otimes \mathbf{1}_K^{\top} \right) \odot \Theta \right) \rangle \rangle. \quad (79c)$$

*A.2.4. Training loss and generalisation error.* 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 have

$$\begin{aligned} \epsilon_{\hat{\ell}} &\equiv \frac{1}{n} \sum_{\nu=1}^n \hat{\ell} \left( y^{\nu}, \frac{\langle \langle \hat{\mathbf{W}} | \mathbf{U}^{\nu} \rangle \rangle}{\sqrt{p}} \right) \xrightarrow{n \rightarrow +\infty} -\partial_{\beta} \Psi_y \\ &= \int dy \mathbb{E}_{\xi} \left[ \frac{\mathcal{Z}^0(y, \omega_0, \sigma_0)}{\mathcal{Z}(y, \omega, \mathbf{V})} \int \frac{\hat{\ell}(y, \mathbf{u}) e^{-\frac{1}{2}(\omega - \mu | \mathbf{V}^{-1} | \omega - \mathbf{u}) - \beta \hat{\ell}(y, \mathbf{u})} d\mathbf{u}}{\sqrt{\det(2\pi \mathbf{V})}} \right] \\ &\quad \times \frac{\beta \rightarrow +\infty}{\mathbf{V} \rightarrow \beta^{-1} \mathbf{V}} \int dy \mathbb{E}_{\xi} \left[ \mathcal{Z}^0(y, \omega_0, \sigma_0) \hat{\ell}(y, \mathbf{h}) \right], \end{aligned} \quad (80)$$

where  $\mathbf{h}$  is the proximal introduced above and all overlaps have to be intended computed at the fixed point.

We can also study the generalisation error

$$\begin{aligned} \epsilon_g &= \mathbb{E}_{(y,U)} \left[ \Delta \left( y, \hat{y} \left( \frac{\langle \langle \hat{\mathbf{W}} | \mathbf{U} \rangle \rangle}{\sqrt{p}} \right) \right) \right] \\ &= \int d\boldsymbol{\mu} \int d\nu \int dy \Delta(y, \hat{y}(\boldsymbol{\mu})) P_y^0(y|\nu) \mathbb{E}_x \left[ \delta \left( \boldsymbol{\mu} - \frac{\langle \langle \hat{\mathbf{W}} | \mathbf{U} \rangle \rangle}{\sqrt{p}} \right) \delta \left( \nu - \frac{\langle \mathbf{x} | \boldsymbol{\theta} \rangle}{\sqrt{d}} \right) \right] \\ &= \int \mathbb{E}_{(\nu, \boldsymbol{\mu})} [\Delta(y, \hat{y}(\boldsymbol{\mu})) P_y^0(y|\nu)] dy \end{aligned} \tag{81}$$

where  $(\nu, \boldsymbol{\mu})$  are jointly Gaussian with covariance

$$\boldsymbol{\Sigma} = \begin{pmatrix} \rho & \mathbf{m}^\top \\ \mathbf{m} & \mathbf{Q} \end{pmatrix}. \tag{82}$$

In particular, if  $P_y^0(y|\nu) = \delta(y - f_0(\nu))$ , then  $\epsilon_g = \mathbb{E}_{(\nu, \boldsymbol{\mu})} [\Delta(f_0(\nu), \hat{y}(\boldsymbol{\mu}))]$ , which corresponds to (6).

### A.3. Separable loss with ridge regularisation

Let us focus now on the case of separable losses, i.e. losses in the form  $\hat{\ell}(y, \mathbf{u}) = \sum_k \ell(y, u_k)$ , which is a crucial special case in the analysis of our contribution. We will assume a ridge regularisation  $r(\mathbf{W}) = \frac{1}{2} \|\mathbf{W}\|_F^2$ . Let us also assume that the  $K$  generative networks are statistically equivalent. This implies a specific structure in the tensors  $\Theta$  and  $\Omega$ ,

$$\Omega_{kk'} = \Omega_{k'k}^\top \stackrel{d}{=} \begin{cases} \boldsymbol{\Omega} & \text{for } k = k', \\ \hat{\boldsymbol{\Omega}} & \text{for } k < k', \end{cases} \tag{83}$$

$$\Theta_{kk'} = \Theta_{k'k}^\top \stackrel{d}{=} \begin{cases} \boldsymbol{\Theta} & \text{for } k = k', \\ \hat{\boldsymbol{\Theta}} & \text{for } k < k'. \end{cases} \tag{84}$$

Here by  $\stackrel{d}{=}$  we mean that the equalities hold in distribution. Observe  $\Omega_{kk}$  and  $\Omega_{kk'}$  are not uncorrelated quantities. For reasons of symmetry reasons, we impose therefore the ansatz

$$\begin{aligned} \mathbf{V} &= v \mathbf{I}_K, & \hat{\mathbf{V}} &= \hat{v} \mathbf{I}_K, \\ \mathbf{m} &= m \mathbf{1}_K, & \hat{\mathbf{m}} &= \hat{m} \mathbf{1}_K, \\ \mathbf{Q} &= (q_0 - q_1) \mathbf{I}_K + q_1 \mathbf{1}_{K,K}, & \hat{\mathbf{Q}} &= (\hat{q}_0 - \hat{q}_1) \mathbf{I}_K + \hat{q}_1 \mathbf{1}_{K,K}. \end{aligned} \tag{85}$$

It is easily seen that

$$\mathbf{Q}^{1/2} = \sqrt{q_0 - q_1} \mathbf{I}_K + \frac{\sqrt{q_0 + (K-1)q_1} - \sqrt{q_0 - q_1}}{K} \mathbf{1}_{K,K}, \tag{86}$$

$$\mathbf{A}_{kk'} = (\lambda \mathbf{I}_p + \Omega_{kk})^{-1} \delta_{kk'}. \tag{87}$$

Plugging this ansatz in our equations, and introducing  $\eta = q_1/q_0$ , we obtain

$$\begin{aligned} \hat{v} &= -\alpha \int_{\mathcal{Y}} dy \mathbb{E}_{\zeta} [\mathcal{Z}^0(y, \omega_0, \sigma_0) \partial_{\omega} f], \\ \hat{m} &= \frac{\alpha}{\sqrt{\gamma}} \int_{\mathcal{Y}} dy \mathbb{E}_{\zeta} [f \partial_{\mu} \mathcal{Z}^0(y, \omega_0, \sigma_0)], \\ \hat{q}_0 &= \alpha \int_{\mathcal{Y}} dy \mathbb{E}_{\zeta} [\mathcal{Z}^0(y, \omega_0, \sigma_0) f^2], \\ \hat{q}_1 &= \alpha \mathbb{E}_{y, \zeta, \zeta'} \left[ \mathcal{Z}^0 \left( y, \frac{m}{\sqrt{q_0}} \frac{\zeta + \zeta'}{1 + \eta}, \rho - \frac{2m^2}{q_0 + q_1} \right) f f' \right], \end{aligned} \quad \text{with } \omega_0 := \frac{m\zeta}{\sqrt{q_0}}, \quad \sigma_0 := \rho - \frac{m^2}{q_0}. \quad (88)$$

The new variables  $\zeta_1$  and  $\zeta_2$  are obtained by a linear transformation from the old ones. In particular, they are distributed as two components of a vector

$$\zeta = \left( \sqrt{1 - \eta} \mathbf{I}_K + \frac{\sqrt{1 + (K - 1)\eta} - \sqrt{1 - \eta}}{K} \mathbf{1}_{K,K} \right) \boldsymbol{\xi}, \quad (89)$$

where  $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}_K, \mathbf{I}_K)$ . It follows that  $\zeta, \zeta' \sim \mathcal{N}(0, 1)$  but they are correlated as

$$\mathbb{E}[\zeta \zeta'] = \eta. \quad (90)$$

Moreover, we have introduced the proximal

$$f = \frac{h - \omega}{v} \quad \text{where} \quad h = \arg \min_x \left[ \frac{(x - \sqrt{q_0} \zeta)^2}{2v} + \ell(y, x) \right] \quad (91)$$

and the corresponding  $f'$  obtained using  $\zeta'$ . The remaining equations read

$$v = \frac{\text{tr} \left[ (\lambda \mathbf{I}_p + \hat{v} \boldsymbol{\Omega})^{-1} \boldsymbol{\Omega} \right]}{p}, \quad (92a)$$

$$m = \frac{\hat{m}}{\sqrt{\gamma}} \frac{\text{tr} \left[ \boldsymbol{\Theta} (\lambda \mathbf{I}_p + \hat{v} \boldsymbol{\Omega})^{-1} \right]}{p} \quad (92b)$$

$$q_0 = \frac{\text{tr} \left[ (\lambda \mathbf{I}_p + \hat{v} \boldsymbol{\Omega})^{-1} (\hat{m}^2 \boldsymbol{\Theta} + \hat{q}_0 \boldsymbol{\Omega}) (\lambda \mathbf{I}_p + \hat{v} \boldsymbol{\Omega})^{-1} \boldsymbol{\Omega} \right]}{p} \quad (92c)$$

$$q_1 = \frac{\text{tr} \left[ (\lambda \mathbf{I}_p + \hat{v} \boldsymbol{\Omega})^{-1} (\hat{m}^2 \hat{\boldsymbol{\Theta}} + \hat{q}_1 \hat{\boldsymbol{\Omega}}) (\lambda \mathbf{I}_p + \hat{v} \boldsymbol{\Omega}')^{-1} \hat{\boldsymbol{\Omega}}^{\top} \right]}{p}. \quad (92d)$$

A.3.1. *The random-features model for the generative networks.* To further simplify these expressions suppose now that our generative networks are such that

$$\mathbf{u}_k(\mathbf{x}) = \phi\left(\frac{\mathbf{F}_k \mathbf{x}}{\sqrt{d}}\right), \quad k \in [K], \tag{93}$$

where  $\mathbf{F}_k \in \mathbb{R}^{p \times d}$  are (fixed) random matrices extracted from some given distribution and  $\phi$  is a nonlinearity acting elementwise. As anticipated in the main text, we can use the fact that each generative network is equivalent to the following Gaussian model (Mei and Montanari 2021)

$$\mathbf{u}_k(\mathbf{x}) \mapsto \kappa_0 \mathbf{1}_p + \kappa_1 \frac{\mathbf{F}_k \mathbf{x}}{\sqrt{d}} + \kappa_* \mathbf{z}_k. \tag{94}$$

for some coefficients  $\kappa_0, \kappa_1$  and  $\kappa_*$  depending on  $\phi$  (see theorem 4), and  $\mathbf{z}^\mu \sim \mathcal{N}(\mathbf{0}_p, \mathbf{I}_p)$ . Assuming now, for the sake of simplicity, to the  $\kappa_0 = 0$  case and that  $\mathbb{E}_\theta[\boldsymbol{\theta}\boldsymbol{\theta}^\top] = \mathbf{I}_d$ , then

$$\begin{aligned} \boldsymbol{\Omega} &\stackrel{d}{=} \frac{\kappa_1^2}{d} \mathbf{F}^\top \mathbf{F} + \kappa_*^2 \mathbf{I}_p, & \hat{\boldsymbol{\Theta}} &\stackrel{d}{=} \hat{\boldsymbol{\Omega}} \stackrel{d}{=} \frac{\kappa_1^2}{d} \mathbf{F}^\top \mathbf{F}', & \mathbf{F} &\stackrel{d}{=} \mathbf{F}' \stackrel{d}{=} \mathbf{F}_k \quad \forall k \in [K]. \\ \boldsymbol{\Theta} &\stackrel{d}{=} \frac{\kappa_1^2}{d} \mathbf{F}^\top \mathbf{F}, & & & & \end{aligned} \tag{95}$$

Once the spectral density  $\varrho(s)$  of  $\boldsymbol{\Omega}$  is introduced, it is immediate to see that the equations for  $q_0, m$  and  $v$  take the forms given in the main text. The equation for  $q_1$  requires an additional step. If we introduce the symmetric random matrix

$$\hat{\mathbf{F}} := \frac{\kappa_1^2}{d} \mathbf{F} \left( (\lambda + \hat{v}\kappa_*^2) \mathbf{I}_p + \frac{\hat{v}\kappa_1^2}{d} \mathbf{F}^\top \mathbf{F} \right)^{-1} \mathbf{F}^\top \in \mathbb{R}^{p \times p} \tag{96}$$

then we can rewrite the equation as

$$q_1 = (\hat{m}^2 + \hat{q}_1) \frac{\text{tr}[\hat{\mathbf{F}}\hat{\mathbf{F}}']}{p} = \frac{\hat{m}^2 + \hat{q}_1}{\gamma} \left( \frac{\text{tr} \hat{\mathbf{F}}}{p} \right)^2 = \left( 1 + \frac{\hat{q}_1}{\hat{m}^2} \right) m^2, \tag{97}$$

where in the second equality we used the fact that  $\hat{\mathbf{F}}$  and  $\hat{\mathbf{F}}'$  are asymptotically free.

A.3.2. *Ridge regression.* Let us consider the simple case of ridge regression with  $f_0(x) = x$ . We will give here the channel equations that are obtained straightforwardly as

$$\begin{aligned} \hat{v} &= \frac{\alpha}{1+v} & \hat{q}_0 &= \alpha \frac{\rho - 2m + q_0}{(1+v)^2}, \\ \hat{m} &= \frac{1}{1+v} \frac{\alpha}{\sqrt{\gamma}}, & \hat{q}_1 &= \alpha \frac{\rho - 2m + q_1}{(1+v)^2}. \end{aligned} \tag{98}$$

The kernel limit is also obtained straightforwardly taking the  $\alpha \rightarrow 0$  limit and rescaling  $\hat{q}_0 \mapsto \alpha \hat{q}_0$ ,  $\hat{q}_1 \mapsto \alpha \hat{q}_1$ ,  $\hat{m} \mapsto \sqrt{\alpha} \hat{m}$ ,  $\hat{v} \mapsto \alpha \hat{v}$ .

$$\begin{aligned}
 v &= \frac{(1-\delta)\kappa_1^2 + \sqrt{(1-\delta)^2\kappa_1^4 + 2(\kappa_*^2 + \lambda)(1+\delta)\kappa_1^2 + (\kappa_*^2 + \lambda)^2} + \kappa_*^2 - \lambda}{2\lambda} \\
 m &= \frac{1}{1 + \lambda \frac{v+1}{\delta\kappa_1^2}} \\
 q_0 = q_1 &= \frac{\delta - 2m + \rho}{\left(1 + 2\lambda \frac{v+1}{\delta\kappa_1^2}\right)^2 - 1} \equiv q. \tag{99}
 \end{aligned}$$

*A.3.3. Binary classification problem.* We consider now the case  $f_0(x) = \text{sign}(x)$ , corresponding to a binary classification problem, and we write down the channel equations for this problem in the case of logistic and hinge loss. In this case we have that

$$\begin{aligned}
 \mathcal{Z}^0(y, \omega_0, \sigma_0) &= \frac{\delta(y-1) + \delta(y+1)}{2} \left( 1 + \text{erf} \left( \frac{y\omega_0}{\sqrt{2\sigma_0}} \right) \right), \\
 \partial_\mu \mathcal{Z}^0(y, \omega_0, \sigma_0) &= (\delta(y-1) - \delta(y+1)) \frac{e^{-\frac{\omega_0^2}{2\sigma_0}}}{\sqrt{2\pi\sigma_0}}. \tag{100}
 \end{aligned}$$

If we pick a logistic loss in the form  $\hat{\ell}(y, \boldsymbol{\mu}) = \sum_k \ln(1 + e^{-y\mu_k})$ , then the proximal  $h$  solves the equation

$$h = \omega + \frac{yv}{1 + e^{yh}}, \tag{101}$$

in such a way that  $f = \frac{\eta - \omega}{v}$  satisfies

$$\partial_\omega f = - \left( v + 2 \cosh \left( y \frac{vf + \omega}{2} \right) \right)^{-1}. \tag{102}$$

If we use instead a hinge loss  $\hat{\ell}(y, \boldsymbol{\mu}) = \sum_k \max(0, 1 - y\mu_k)$ , the proximal is such that

$$f = \begin{cases} y & \text{if } 1 - v > \omega y, \\ \frac{y-\omega}{v} & \text{if } 1 - v < \omega y < 1, \\ 0 & \text{otherwise,} \end{cases} \quad \partial_\omega f = \begin{cases} -\frac{1}{v} & \text{if } 1 - v < \omega y < 1, \\ 0 & \text{otherwise.} \end{cases} \tag{103}$$



In the case of the hinge loss, the simple form of the proximal allows for a more explicit expression of the channel equations. Introducing

$$\hat{\sigma} = \rho - \frac{2m^2}{q_0 + q_1}, \quad \mathcal{N}(\zeta, \zeta'; \eta) = \frac{\exp\left(-\frac{\zeta^2 + \zeta'^2 - 2\eta\zeta\zeta'}{2(1-\eta^2)}\right)}{2\pi\sqrt{1-\eta^2}} \tag{104}$$

we obtain

$$\begin{aligned} \hat{v} &= \frac{\alpha}{v} \int_{\frac{1-v}{\sqrt{q_0}}}^{1/\sqrt{q_0}} \frac{e^{-\frac{\zeta^2}{2}}}{\sqrt{2\pi}} d\zeta \left(1 + \operatorname{erf}\left(\frac{m\zeta}{\sqrt{2q_0\sigma}}\right)\right), \\ \hat{m} &= \frac{\alpha}{\sqrt{\rho\gamma}v} \left[ \frac{v + \operatorname{erf}\left(\sqrt{\frac{\rho}{2q_0\sigma}}\right) - (1-v)\operatorname{erf}\left((1-v)\sqrt{\frac{\rho}{2q_0\sigma}}\right)}{\sqrt{2\pi}} \right. \\ &\quad \left. + \sqrt{\frac{q_0\sigma}{\rho}} \frac{\exp\left(-\frac{\rho}{2q_0\sigma}\right) - \exp\left(\frac{\rho(1-v)^2}{2q_0\sigma}\right)}{2\pi} \right] \\ \hat{q}_0 &= \alpha \left[ \int_{-\infty}^{\frac{1-v}{\sqrt{q_0}}} \frac{e^{-\frac{\zeta^2}{2}}}{\sqrt{2\pi}} d\zeta \left(1 + \operatorname{erf}\left(\frac{m\zeta}{\sqrt{2q_0\sigma}}\right)\right) \right. \\ &\quad \left. + \int_{\frac{1-v}{\sqrt{q_0}}}^{1/\sqrt{q_0}} \frac{e^{-\frac{\zeta^2}{2}}}{\sqrt{2\pi}} d\zeta \left(1 + \operatorname{erf}\left(\frac{m\zeta}{\sqrt{2q_0\sigma}}\right)\right) \left(\frac{1-\sqrt{q_0}\zeta}{v}\right)^2 \right], \\ \hat{q}_1 &= \alpha \iint_{-\infty}^{\frac{1-v}{\sqrt{q_0}}} \mathcal{N}(\zeta, \zeta'; q_1/q_0) d\zeta d\zeta' \left(1 + \operatorname{erf}\left(\frac{m}{\sqrt{2q_0\hat{\sigma}}} \frac{\zeta + \zeta'}{1 + q_1/q_0}\right)\right) \\ &\quad + 2\alpha \int_{\frac{1-v}{\sqrt{q_0}}}^{1/\sqrt{q_0}} d\zeta \int_{-\infty}^{\frac{1-v}{\sqrt{q_0}}} d\zeta' \mathcal{N}(\zeta, \zeta'; q_1/q_0) \left(1 + \operatorname{erf}\left(\frac{m}{\sqrt{2q_0\hat{\sigma}}} \frac{\zeta + \zeta'}{1 + q_1/q_0}\right)\right) \left(\frac{1-\sqrt{q_0}\zeta}{v}\right) \\ &\quad + \alpha \iint_{\frac{1-v}{\sqrt{q_0}}}^{1/\sqrt{q_0}} d\zeta d\zeta' \mathcal{N}(\zeta, \zeta'; q_1/q_0) \left(1 + \operatorname{erf}\left(\frac{m}{\sqrt{2q_0\hat{\sigma}}} \frac{\zeta + \zeta'}{1 + q_1/q_0}\right)\right) \\ &\quad \times \left(\frac{1-\sqrt{q_0}\zeta}{v}\right) \left(\frac{1-\sqrt{q_0}\zeta'}{v}\right). \end{aligned} \tag{105}$$

Let us now make the change of variables  $\zeta \mapsto \frac{\sqrt{q_0+q_1}z + \sqrt{q_0-q_1}z'}{\sqrt{2q_0}}$  and  $\zeta' \mapsto \frac{\sqrt{q_0+q_1}z - \sqrt{q_0-q_1}z'}{\sqrt{2q_0}}$ . This allows us to rewrite the expression for  $q_1$  as

$$\begin{aligned}
 \hat{q}_1 &= \alpha \int_{-\infty}^{\sqrt{2}(1-v)} \mathcal{N}(z; 0, q_0 + q_1) \left( 1 + \operatorname{erf} \left( \frac{mz}{\sqrt{\hat{\sigma}}(q_0 + q_1)} \right) \right) \operatorname{erf} \left( \frac{\sqrt{2}(1-v) - z}{\sqrt{q_0 - q_1}} \right) \\
 &+ 2\alpha \int_{-\infty}^{\sqrt{2}(1-v/2)} \mathcal{N}(z; 0, q_0 + q_1) \left( 1 + \operatorname{erf} \left( \frac{mz}{\sqrt{\hat{\sigma}}(q_0 + q_1)} \right) \right) \\
 &\times \left( \frac{1}{v} - \frac{1}{v} \left[ \frac{z}{2} \operatorname{erf} \left( \frac{x}{\sqrt{q_0 - q_1}} \right) + \sqrt{\frac{q_0 - q_1}{2\pi}} e^{-\frac{x^2}{2}} \right]_{|\sqrt{2}(1-v)-z|}^{\sqrt{2}-z} \right) dz \\
 &+ \alpha \iint_{\frac{1-v}{\sqrt{q_0}}}^{1/\sqrt{q_0}} d\zeta d\zeta' \mathcal{N}(z; 0, q_0 + q_1) \mathcal{N}(z'; 0, 1) \left( 1 + \operatorname{erf} \left( \frac{mz}{\sqrt{\hat{\sigma}}(q_0 + q_1)} \right) \right) \\
 &\times \left( \frac{1 - \sqrt{q_0}\zeta}{v} \right) \left( \frac{1 - \sqrt{q_0}\zeta'}{v} \right). \tag{106}
 \end{aligned}$$

**Appendix B. Proof of the main theorem**

In this section we prove theorem 2, from which all other analytical results in the paper can be deduced. We start by reminding the learning problem defining the ensemble of estimators with a few auxiliary notations, so that this part is self contained. The sketch of proof is one pioneered in Bayati and Montanari (2011b), Donoho and Montanari (2016) and is the following: the estimator  $\mathbf{W}^*$  is expressed as the limit of a carefully chosen sequence, an *AMP iteration* (Bayati and Montanari 2011a, Zdeborová and Krzakala 2016), whose iterates can be asymptotically exactly characterised using an auxiliary, closed form iteration, the *state evolution equations*. We then show that converging trajectories of such an AMP iteration can be systematically found.

**B.1. The learning problem**

We start by reminding the definition of the problem. Consider the following generative model

$$\mathbf{y} = f_0 \left( \frac{1}{\sqrt{d}} \mathbf{X}_0 \mathbf{w}_0, \boldsymbol{\epsilon}_0 \right) \tag{107}$$

where  $\mathbf{y} \in \mathbb{R}^n$ ,  $\mathbf{X}_0 \sim \mathcal{N}(0, \boldsymbol{\Sigma}_{00}) \in \mathbb{R}^{n \times d}$ ,  $\mathbf{w}_0 \in \mathbb{R}^d$ ,  $\boldsymbol{\epsilon}_0 \in \mathbb{R}^d$  is a noise vector and  $\boldsymbol{\Sigma}_{00} \in \mathbb{R}^{d \times d}$  is a positive definite matrix. The goal is to learn this generative model using an ensemble of predictors  $\mathbf{W} = [\mathbf{w}_1 | \mathbf{w}_2 | \dots | \mathbf{w}_K] \in \mathbb{R}^{p \times K}$  where each predictor  $\mathbf{w}_k \in \mathbb{R}^p, k \in [1, K]$  is learned using a sample dataset  $\mathbf{X}_k \in \mathbb{R}^{n \times p}$ , where, for any  $i \in [1, n]$  and  $k \in [0, K]$ , we have:

$$\mathbb{E} \left[ \mathbf{x}_i^k \left( \mathbf{x}_i^{k'} \right)^\top \right] = \boldsymbol{\Sigma}_{kk'} \tag{108}$$

where each sample is Gaussian and we denote:

$$\Sigma = \begin{bmatrix} \Sigma_{00} & \Sigma_{01} & \dots & \Sigma_{0K} \\ \Sigma_{10} & \Sigma_{11} & \dots & \Sigma_{1K} \\ \dots & \dots & \dots & \dots \\ \Sigma_{K0} & \Sigma_{K1} & \dots & \Sigma_{KK} \end{bmatrix} \in \mathbb{R}^{(Kp+d) \times (Kp+d)}. \quad (109)$$

The predictors interact with each sample dataset in a linear way, i.e. we will consider a generalised linear model acting on the ensemble of products  $\{\mathbf{X}_k \mathbf{w}_k\}_{k=1}^K$ :

$$\mathbf{W}^* \in \arg \min_{\mathbf{W} \in \mathbb{R}^{p \times K}} \mathcal{L} \left( \mathbf{y}, \left\{ \frac{1}{\sqrt{p}} \mathbf{X}_k \mathbf{w}_k \right\}_{k=1}^K \right) + r_0(\mathbf{W}) \quad (110)$$

where  $\mathcal{L}, r_0$  are convex functions. We wish to determine the asymptotic properties of the estimator  $\mathbf{W}^*$  in the limit where  $n, p, d \rightarrow \infty$  with fixed ratios  $\alpha = n/p, \gamma = d/p$ . We now list the necessary assumptions for our main theorem to hold.

### B.1.1. Assumptions

- the functions  $\mathcal{L}, r_0$  are proper, closed, lower-semicontinuous, convex functions. The loss function  $\mathcal{L}$  is pseudo-lipschitz of order 2 in both its arguments and the regularisation  $r_0$  is pseudo-Lipschitz of order 2. The cost function  $\mathcal{L}(\mathbf{X}_.) + r_0(\cdot)$  is coercive.
- for any  $1 \leq k \leq K$ , the matrix  $\Sigma_k \in \mathbb{R}^{p \times p}$  is symmetric and there exist strictly positive constants  $\kappa_0, \kappa_1$  such that  $\kappa_0 \leq \lambda_{\min}(\Sigma_k) \leq \lambda_{\max}(\Sigma_k) \leq \kappa_1$ . We also assume that the matrix  $\Sigma$  is positive definite.
- There exists a positive constant  $C_{f_0}$  such that  $\left\| f_0\left(\frac{1}{\sqrt{d}} \mathbf{X}_0 \mathbf{w}_0, \epsilon_0\right) \right\|_2 \leq C_{f_0} \left( \left\| \frac{1}{\sqrt{d}} \mathbf{X}_0 \mathbf{w}_0 \right\|_2 + \|\epsilon_0\|_2 \right)$ .
- The dimensions  $n, p, d$  grow linearly with finite ratios  $\alpha = n/p$  and  $\gamma = d/p$ .
- The ground truth vector  $\mathbf{w}_0 \in \mathbb{R}^d$  and noise vector  $\epsilon_0 \in \mathbb{R}^n$  are sampled from subgaussian probability distributions independent from each other and from all other random quantities of the learning problem.

The proof method we will employ involves expressing the estimator  $\mathbf{W}^*$  as the limit of a carefully chosen sequence. In the case of non-strictly convex problems, the estimator may not be unique, making it unclear what estimator is reached by the sequence (at best we know it belongs to the set of zeroes of the subgradient of the cost function). We thus start with the following problem

$$\mathbf{W}^* \in \arg \min_{\mathbf{W} \in \mathbb{R}^{p \times K}} \mathcal{L} \left( \mathbf{y}, \{\mathbf{X}_k \mathbf{w}_k\}_{k=1}^K \right) + r_{\lambda_2}(\mathbf{W}) \quad (111)$$

$$\text{where, for any } \mathbf{W} \in \mathbb{R}^{p \times K}, r_{\lambda_2}(\mathbf{W}) = r_0(\mathbf{W}) + \frac{\lambda_2}{2} \|\mathbf{W}\|_F^2 \quad (112)$$

i.e. we add a ridge regularisation to the initial problem to make it strongly convex. We will relax this additional strong convexity constraint later on.

**B.2. Asymptotics for the strongly convex problem**

We now reformulate the minimisation problem equation (111) to make it amenable to an AMP iteration. The key feature of this ensembling problem, outside of the convexity which will be crucial to control the trajectories of the AMP iteration, is the fact that each predictor only interacts linearly with each design sample, along with the correlation structure of the overall dataset. We are effectively sampling  $n$  vectors of size  $(Kp + d)$  from the Gaussian distribution with covariance  $\Sigma$ , i.e.  $[\mathbf{x}_0|\mathbf{x}_1|\dots|\mathbf{x}_K] \sim \mathcal{N}(0, \Sigma)$ . We then write  $\{\mathbf{X}_k \mathbf{w}_k\}_{k=0}^K = [\mathbf{X}_0 \mathbf{w}_0|\dots|\mathbf{X}_K \mathbf{w}_K] \in \mathbb{R}^{n \times (K+1)}$ , such that

$$[\mathbf{X}_0 \mathbf{w}_0|\dots|\mathbf{X}_K \mathbf{w}_K] = [\mathbf{X}_0|\dots|\mathbf{X}_K] \mathbf{W} = \mathbf{Z} \Sigma^{1/2} \begin{bmatrix} \mathbf{w}_0 & 0 \\ 0 & \tilde{\mathbf{W}} \end{bmatrix} \tag{113}$$

$$\text{where } \tilde{\mathbf{W}} = \begin{bmatrix} \mathbf{w}_1 & 0 & \dots & 0 \\ 0 & \mathbf{w}_2 & \dots & 0 \\ & & \dots & \\ 0 & 0 & \dots & \mathbf{w}_K \end{bmatrix} \in \mathbb{R}^{Kp \times K} \tag{114}$$

and  $\mathbf{Z} \in \mathbb{R}^{n \times (Kp+d)}$  is a random matrix with i.i.d.  $\mathcal{N}(0, 1)$  elements. Then, any sample  $[\mathbf{x}_0|\mathbf{x}_1|\dots|\mathbf{x}_K]$  may be rewritten as

$$\mathbf{x}_0 = \Psi^{1/2} \mathbf{a} \quad \text{and} \quad [\mathbf{x}_1|\dots|\mathbf{x}_K] = \Phi^\top \Psi^{-1/2} \mathbf{a} + (\Omega - \Phi^\top \Psi^{-1} \Phi)^{1/2} \mathbf{b} \tag{115}$$

$$\mathbf{X}_0 = \mathbf{A} \Psi^{1/2} \quad \text{and} \quad [\mathbf{X}_1|\dots|\mathbf{X}_K] = \mathbf{A} \Psi^{-1/2} \Phi + \mathbf{B} (\Omega - \Phi^\top \Psi^{-1} \Phi)^{1/2} \tag{116}$$

where  $\mathbf{a} \in \mathbb{R}^d, \mathbf{b} \in \mathbb{R}^{Kp}$  are vectors with i.i.d. standard normal components,  $\mathbf{A} \in \mathbb{R}^{n \times d}, \mathbf{B} \in \mathbb{R}^{n \times Kp}$  are the corresponding design matrices, and the covariance matrices are given by  $\Psi = \Sigma_{00} \in \mathbb{R}^{d \times d}, \Phi = [\Sigma_{11}|\Sigma_{12}|\Sigma_{13}|\dots|\Sigma_{1K}] \in \mathbb{R}^{d \times Kp}$  and

$$\Omega = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} & \dots & \Sigma_{1K} \\ \Sigma_{21} & \Sigma_{22} & \dots & \Sigma_{2K} \\ \dots & & & \\ \Sigma_{K1} & \Sigma_{K2} & \dots & \Sigma_{KK} \end{bmatrix} \in \mathbb{R}^{Kp \times Kp}. \tag{117}$$

The optimisation problem may then be written, introducing the appropriate scalings

$$\begin{aligned} \tilde{\mathbf{W}}^* \in \arg \min_{\tilde{\mathbf{W}} \in \mathbb{R}^{Kp \times K}} \mathcal{L} \left( f_0 \left( \frac{1}{\sqrt{d}} \mathbf{A} \tilde{\mathbf{w}}_0 \right), \right. \\ \left. \frac{1}{\sqrt{p}} \left( \mathbf{A} \Psi^{-1/2} \Phi + \mathbf{B} (\Omega - \Phi^\top \Psi^{-1} \Phi)^{1/2} \right) \tilde{\mathbf{W}} \right) + r(\tilde{\mathbf{W}}) \end{aligned} \tag{118}$$

where we let  $\tilde{\mathbf{w}}_0 = \Psi^{1/2} \mathbf{w}_0$ , its scaled norm  $\rho_{\tilde{\mathbf{w}}_0} = \frac{1}{d} \|\tilde{\mathbf{w}}_0\|_2^2$  and we introduced the function

$$r : \mathbb{R}^{Kp \times K} \rightarrow \mathbb{R} \tag{119}$$

$$\tilde{\mathbf{W}} \rightarrow r_{\lambda_2}(\mathbf{W}). \tag{120}$$

In order to isolate the contribution correlated with the teacher, we condition the design matrix  $\mathbf{A}$  on the teacher distribution  $\mathbf{y}$ , we can write

$$\mathbf{A} = \mathbb{E}[\mathbf{A}|\mathbf{y}] + \mathbf{A} - \mathbb{E}[\mathbf{A}|\mathbf{y}] \tag{121}$$

$$= \mathbb{E}[\mathbf{A}|\mathbf{A}\tilde{\mathbf{w}}_0] + \mathbf{A} - \mathbb{E}[\mathbf{A}|\mathbf{A}\tilde{\mathbf{w}}_0] \tag{122}$$

$$= \mathbf{A}\mathbf{P}_{\tilde{\mathbf{w}}_0} + \tilde{\mathbf{A}}\mathbf{P}_{\tilde{\mathbf{w}}_0}^\perp \tag{123}$$

where  $\tilde{\mathbf{A}}$  is an independent copy of  $\mathbf{A}$ , see Bayati and Montanari (2011a) lemma 11. The cost function then becomes

$$\mathcal{L} \left( f_0(\sqrt{\rho_{\tilde{\mathbf{w}}_0}}\mathbf{s}), \frac{1}{\sqrt{p}} \left( \mathbf{s} \frac{(\Phi^\top \mathbf{w}_0)^\top}{\sqrt{d\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{A}}\mathbf{P}_{\tilde{\mathbf{w}}_0}^\perp \Psi^{-1/2}\Phi + \mathbf{B}(\Omega - \Phi^\top \Psi^{-1}\Phi)^{1/2} \right) \tilde{\mathbf{W}} \right) + r(\tilde{\mathbf{W}}) \tag{124}$$

where  $\mathbf{s} = \mathbf{A} \frac{\tilde{\mathbf{w}}_0}{\|\tilde{\mathbf{w}}_0\|_2} \in \mathbb{R}^n$  is an i.i.d. standard normal vector. The term  $\tilde{\mathbf{A}}\mathbf{P}_{\tilde{\mathbf{w}}_0}^\perp \Psi^{-1/2}\Phi + \mathbf{B}(\Omega - \Phi^\top \Psi^{-1}\Phi)^{1/2}$  can then be represented as a  $\mathbb{R}^{n \times Kp}$  Gaussian matrix with covariance

$$\Phi^\top \Psi^{-1/2} \mathbf{P}_{\tilde{\mathbf{w}}_0}^\perp \Psi^{-1/2} \Phi + \Omega - \Phi^\top \Psi^{-1} \Phi = \Omega - \Phi^\top \Psi^{-1/2} \mathbf{P}_{\tilde{\mathbf{w}}_0} \Psi^{-1/2} \Phi \tag{125}$$

$$= \Omega - \Phi^\top \Psi^{-1/2} \frac{\tilde{\mathbf{w}}_0 \tilde{\mathbf{w}}_0^\top}{\|\tilde{\mathbf{w}}_0\|_2^2} \Psi^{-1/2} \Phi = \Omega - \frac{\mathbf{c}\mathbf{c}^\top}{\|\tilde{\mathbf{w}}_0\|_2^2} \tag{126}$$

where we introduced  $\mathbf{c} = \Phi^\top \mathbf{w}_0 \in \mathbb{R}^{Kp}$  and  $\rho_{\mathbf{c}} = \frac{1}{p} \|\mathbf{c}\|_2^2$ , reaching the cost function

$$\mathcal{L} \left( f_0(\sqrt{\rho_{\tilde{\mathbf{w}}_0}}\mathbf{s}), \frac{1}{\sqrt{p}} \left( \mathbf{s} \frac{\mathbf{c}^\top}{\sqrt{d\rho_{\tilde{\mathbf{w}}_0}}} + \mathbf{Z} \left( \Omega - \frac{\mathbf{c}\mathbf{c}^\top}{\|\tilde{\mathbf{w}}_0\|_2^2} \right)^{1/2} \right) \tilde{\mathbf{W}} \right) + r(\tilde{\mathbf{W}}). \tag{127}$$

Introducing  $\mathbf{m} = \frac{1}{\sqrt{dp}} \tilde{\mathbf{W}}^\top \mathbf{c} \in \mathbb{R}^K$ ,  $\mathbf{C} = \Omega - \frac{\mathbf{c}\mathbf{c}^\top}{\|\tilde{\mathbf{w}}_0\|_2^2} \in \mathbb{R}^{Kp \times Kp}$ , and the Lagrange multiplier  $\boldsymbol{\nu}$  associated to  $\mathbf{m}$ , the optimisation problem can equivalently be written

$$\inf_{\mathbf{m} \in \mathbb{R}^K, \tilde{\mathbf{W}} \in \mathbb{R}^{Kp \times K}} \sup_{\boldsymbol{\nu} \in \mathbb{R}^K} \mathcal{L} \left( f_0(\sqrt{\rho_{\tilde{\mathbf{w}}_0}}\mathbf{s}), \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \frac{1}{\sqrt{p}} \mathbf{Z}\mathbf{C}^{1/2}\tilde{\mathbf{W}} \right) + r(\tilde{\mathbf{W}}) - \boldsymbol{\nu}^\top (\tilde{\mathbf{W}}^\top \mathbf{c} - \sqrt{dp}\mathbf{m}). \tag{128}$$

We now look for an explicit expression of the matrix square root  $\mathbf{C}^{1/2}$

$$\mathbf{C} = \Omega^{1/2} \left( Id - \frac{\Omega^{-1/2} \mathbf{c} (\Omega^{-1/2} \mathbf{c})^\top}{\|\tilde{\mathbf{w}}_0\|_2^2} \right) \Omega^{1/2} \quad \text{let} \quad \tilde{\mathbf{c}} = \Omega^{-1/2} \mathbf{c} \quad (129)$$

$$= \Omega^{1/2} (\mathbf{P}_{\tilde{\mathbf{c}}}^\perp + \kappa \mathbf{P}_{\tilde{\mathbf{c}}}) \Omega^{1/2} \quad \text{where} \quad \kappa = 1 - \frac{\|\tilde{\mathbf{c}}\|_2^2}{\|\tilde{\mathbf{w}}_0\|_2^2} \quad (130)$$

$$= \Omega^{1/2} (\mathbf{P}_{\tilde{\mathbf{c}}}^\perp + \sqrt{\kappa} \mathbf{P}_{\tilde{\mathbf{c}}}) (\mathbf{P}_{\tilde{\mathbf{c}}}^\perp + \sqrt{\kappa} \mathbf{P}_{\tilde{\mathbf{c}}}) \Omega^{1/2} \quad (131)$$

where the positivity of  $\kappa$  is ensured by the positive-definiteness of  $\Sigma$ . The problem then becomes

$$\inf_{\mathbf{m}, \tilde{\mathbf{W}}} \sup_{\nu} \mathcal{L} \left( f_0(\sqrt{\rho_{\tilde{\mathbf{w}}_0}} \mathbf{s}), \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \frac{\sqrt{\kappa}}{\sqrt{p}} \mathbf{Z} \mathbf{P}_{\tilde{\mathbf{c}}} \Omega^{1/2} \tilde{\mathbf{W}} + \frac{1}{\sqrt{p}} \tilde{\mathbf{Z}} \mathbf{P}_{\tilde{\mathbf{c}}}^\perp \Omega^{1/2} \tilde{\mathbf{W}} \right) + r(\tilde{\mathbf{W}}) - \nu^\top (\tilde{\mathbf{W}}^\top \mathbf{c} - \sqrt{d p \mathbf{m}}) \quad (132)$$

where  $\tilde{\mathbf{Z}}$  is an independent copy of  $\mathbf{Z}$ , see Bayati and Montanari (2011a) lemma 11. Then

$$\frac{\sqrt{\kappa}}{\sqrt{p}} \mathbf{Z} \mathbf{P}_{\tilde{\mathbf{c}}} \Omega^{1/2} \tilde{\mathbf{W}} = \frac{\sqrt{\kappa}}{\sqrt{p}} \tilde{\mathbf{s}} \frac{\mathbf{c}^\top \tilde{\mathbf{W}}}{\|\tilde{\mathbf{c}}\|_2} \quad (133)$$

$$= \sqrt{\kappa} \tilde{\mathbf{s}} \frac{\mathbf{c}^\top \tilde{\mathbf{W}}}{p \sqrt{\rho_{\tilde{\mathbf{c}}}}} \quad (134)$$

$$= \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} \quad (135)$$

where  $\tilde{\mathbf{s}} = \mathbf{Z} \frac{\tilde{\mathbf{c}}}{\|\tilde{\mathbf{c}}\|_2}$  is an i.i.d. standard normal vector and  $\rho_{\tilde{\mathbf{c}}} = \frac{1}{p} \|\tilde{\mathbf{c}}\|_2^2$  such that the optimisation problem becomes

$$\inf_{\mathbf{m}, \tilde{\mathbf{W}}} \sup_{\nu} \mathcal{L} \left( f_0(\sqrt{\rho_{\tilde{\mathbf{w}}_0}} \mathbf{s}), \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \frac{1}{\sqrt{p}} \tilde{\mathbf{Z}} \mathbf{P}_{\tilde{\mathbf{c}}}^\perp \Omega^{1/2} \tilde{\mathbf{W}} \right) + r(\tilde{\mathbf{W}}) - \nu^\top (\tilde{\mathbf{W}}^\top \mathbf{c} - \sqrt{d p \mathbf{m}}). \quad (136)$$

Now let  $\mathbf{U} = \mathbf{P}_{\tilde{\mathbf{c}}}^\perp \Omega^{1/2} \tilde{\mathbf{W}}$ , such that  $\tilde{\mathbf{W}} = \Omega^{-1/2} \left( \frac{\sqrt{\gamma \tilde{\mathbf{c}} \mathbf{m}^\top}}{\rho_{\tilde{\mathbf{c}}}} + \mathbf{U} \right)$ . The equivalent problem in  $\mathbf{U}$  reads

$$\inf_{\mathbf{m}, \mathbf{U}} \sup_{\nu} \mathcal{L} \left( f_0(\sqrt{\rho_{\tilde{\mathbf{w}}_0}} \mathbf{s}), \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \frac{1}{\sqrt{p}} \tilde{\mathbf{Z}} \mathbf{U} \right) + r \left( \Omega^{-1/2} \left( \frac{\sqrt{\gamma \tilde{\mathbf{c}} \mathbf{m}^\top}}{\rho_{\tilde{\mathbf{c}}}} + \mathbf{U} \right) \right) - \nu^\top \mathbf{U}^\top \tilde{\mathbf{c}}. \quad (137)$$

Note that the constraint defining  $\mathbf{m}$  automatically enforces the orthogonality constraint on  $\mathbf{U}$  w.r.t.  $\tilde{\mathbf{c}}$ . The following lemma characterises properties of the feasibility sets of  $\mathbf{U}, \mathbf{m}, \nu$ .

**Lemma 5.** Consider the optimisation problem equation (137). Then there exist constants  $C_U, C_m, C_\nu$  such that

$$\frac{1}{\sqrt{p}} \|\mathbf{U}\|_F \leq C_U, \quad \|\mathbf{m}\|_2 \leq C_m, \quad \|\boldsymbol{\nu}\|_2 \leq C_\nu \quad (138)$$

with high probability as  $n, p, d \rightarrow \infty$ .

**Proof.** Consider the optimisation problem defining  $\tilde{\mathbf{W}}^*$

$$\tilde{\mathbf{W}}^* \in \underset{\tilde{\mathbf{W}} \in \mathbb{R}^{Kp \times K}}{\operatorname{argmin}} \mathcal{L}(\mathbf{y}, \mathbf{X}\tilde{\mathbf{W}}) + \tilde{r}_0(\tilde{\mathbf{W}}) + \frac{\lambda_2}{2} \|\tilde{\mathbf{W}}\|_F^2 \quad (139)$$

which, owing to the convexity of the cost function, verifies

$$\frac{1}{p} \left( \mathcal{L}(\mathbf{y}, \mathbf{X}\tilde{\mathbf{W}}^*) + \tilde{r}_0(\tilde{\mathbf{W}}^*) + \frac{\lambda_2}{2} \|\tilde{\mathbf{W}}^*\|_F^2 \right) \leq \frac{1}{p} (\mathcal{L}(\mathbf{y}, 0) + \tilde{r}_0(0)). \quad (140)$$

The functions  $\mathcal{L}$  and  $\tilde{r}_0$  are assumed to be proper, thus their sum is bounded below for any value of their arguments and we may write

$$\frac{1}{p} \frac{\lambda_2}{2} \|\tilde{\mathbf{W}}^*\|_F^2 \leq \frac{1}{p} (\mathcal{L}(\mathbf{y}, 0) + \tilde{r}_0(0)). \quad (141)$$

The pseudo-Lipschitz assumption on  $\mathcal{L}$  and  $\tilde{r}_0$  then implies that there exist positive constants  $C_{\mathcal{L}}$  and  $C_{\tilde{r}_0}$  such that

$$\frac{1}{p} \frac{\lambda_2}{2} \|\tilde{\mathbf{W}}^*\|_F^2 \leq \frac{1}{p} \left( C_{\mathcal{L}} (1 + \|\mathbf{y}\|_2^2) \right) + C_{\tilde{r}_0} \quad (142)$$

$$\leq \frac{1}{p} \left( C_{\mathcal{L}} \left( 1 + C_{f_0} \left\| \frac{1}{\sqrt{d}} \mathbf{X}_0 \mathbf{w}_0 \right\|_2^2 + C_{f_0} \|\boldsymbol{\epsilon}_0\|_2^2 \right) \right) + C_{\tilde{r}_0} \quad (143)$$

where the second line follows from the scaling assumption on the teacher function  $f_0$ . Hence

$$\frac{1}{p} \frac{\lambda_2}{2} \|\tilde{\mathbf{W}}^*\|_F^2 \leq C_{\mathcal{L}} \left( 1 + C_{f_0} \left\| \frac{1}{\sqrt{d}} \mathbf{A} \right\|_{op}^2 \|\Psi^{1/2}\|_{op}^2 \frac{\gamma}{d} \|\mathbf{w}\|_0^2 + C_{f_0} \frac{\alpha}{n} \|\boldsymbol{\epsilon}_0\|_2^2 \right) + C_{\tilde{r}_0} \quad (144)$$

where  $\|\bullet\|_{op}$  denotes the operator norm of a given matrix, and we remind that  $\mathbf{A}$  has i.i.d.  $\mathcal{N}(0,1)$  elements. By assumption the maximum singular value of  $\Psi^{1/2}$  is bounded. The maximum singular value of a random matrix with i.i.d.  $\mathcal{N}(0, \frac{1}{d})$  elements is bounded with high probability as  $n, p, d \rightarrow \infty$ , see e.g. Vershynin (2010). Finally,  $\mathbf{w}_0$  and  $\boldsymbol{\epsilon}_0$  are sampled from subgaussian probability distributions, thus their scaled norms are bounded with high probability as  $n, p, d \rightarrow \infty$  according to Bernstein's inequality, see e.g. Vershynin (2018). An application of the union bound then leads to the following statement: there exists a constant  $C_{\tilde{\mathbf{W}}}$  such that  $\frac{1}{p} \|\tilde{\mathbf{W}}\|_2^2 \leq C_{\tilde{\mathbf{W}}}$ , with high probability as  $n, p, d \rightarrow \infty$ . Now using the definition of  $\mathbf{U}$

$$\frac{1}{p} \|\mathbf{U}\|_F^2 = \frac{1}{p} \left\| \mathbf{P}_{\tilde{\mathbf{c}}}^\perp \Omega^{1/2} \tilde{\mathbf{W}} \right\|_F^2 \tag{145}$$

$$\leq \left\| \mathbf{P}_{\tilde{\mathbf{c}}}^\perp \right\|_{op}^2 \left\| \Omega^{1/2} \right\|_{op}^2 \frac{1}{p} \left\| \tilde{\mathbf{W}} \right\|_F^2 \tag{146}$$

where the singular values of  $\mathbf{P}_{\tilde{\mathbf{c}}}^\perp$  and  $\Omega^{1/2}$  are bounded with probability one. Therefore there exists a constant  $C_U$  such that  $\frac{1}{\sqrt{p}} \|\mathbf{U}\| \leq C_U$  with high probability as  $n, p, d \rightarrow \infty$ . Then, by definition of  $\mathbf{m}$  and the Cauchy–Schwarz inequality

$$\|\mathbf{m}\|_2^2 \leq \frac{1}{d} \|\mathbf{c}\|_2^2 \frac{1}{p} \left\| \tilde{\mathbf{W}} \right\|_F^2 \tag{147}$$

$$\leq \|\Phi\|_{op}^2 \frac{1}{d} \|\mathbf{w}_0\|_2^2 \frac{1}{p} \left\| \tilde{\mathbf{W}} \right\|_F^2 \tag{148}$$

combining the results previously established on  $\tilde{\mathbf{W}}$  and  $\mathbf{w}_0$  with the fact that the maximum singular value of  $\Phi$  is bounded, there exists a positive constant  $C_m$  such that  $\|\mathbf{m}\|_2 \leq C_m$  with high probability as  $n, p, d \rightarrow \infty$ . We finally turn to  $\nu$ . The optimality condition for  $\mathbf{m}$  in problem equation (128) gives

$$\nu = -\frac{1}{\sqrt{dp}} \frac{\mathbf{s}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} \partial \mathcal{L} \left( \mathbf{y}, \frac{\mathbf{m} \mathbf{s}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \frac{1}{\sqrt{p}} \mathbf{Z} \mathbf{C}^{1/2} \tilde{\mathbf{W}}^* \right). \tag{149}$$

The pseudo-Lipschitz assumption on  $\mathcal{L}$  implies that we can find a constant  $C_{\partial \mathcal{L}}$  such that

$$\|\nu\|_2^2 = \frac{1}{dp} \frac{\|\mathbf{s}\|_2^2}{\rho_{\tilde{\mathbf{w}}_0}} C_{\mathcal{L}} \left( 1 + \|\mathbf{y}\|_2^2 + \left\| \frac{\mathbf{m} \mathbf{s}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \frac{1}{\sqrt{p}} \mathbf{Z} \mathbf{C}^{1/2} \tilde{\mathbf{W}}^* \right\|_2^2 \right) \tag{150}$$

the last bound then follows from similar arguments as those employed above.  $\square$

The optimisation problem equation (137) is convex and feasible. Furthermore, we may reduce the feasibility sets of  $\mathbf{m}, \nu$  to compact spaces, and the function of  $\mathbf{U}$  is coercive and thus has bounded lower level sets. Strong duality then implies we can invert the order of minimisation to obtain the equivalent problem

$$\begin{aligned} \inf_{\mathbf{m}} \sup_{\nu} \inf_{\mathbf{U}} \mathcal{L} \left( f_0(\sqrt{\rho_{\tilde{\mathbf{w}}_0}} \mathbf{s}), \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \frac{1}{\sqrt{p}} \tilde{\mathbf{Z}} \mathbf{U} \right) \\ + r \left( \Omega^{-1/2} \left( \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} + \mathbf{U} \right) \right) - \nu^\top \mathbf{U}^\top \tilde{\mathbf{c}} \end{aligned} \tag{151}$$

and study the optimisation problem in  $\mathbf{U}$  at fixed  $\mathbf{m}, \nu$ :

$$\inf_{\mathbf{U} \in \mathbb{R}^{Kp \times K}} \tilde{\mathcal{L}} \left( \frac{1}{\sqrt{p}} \tilde{\mathbf{Z}} \mathbf{U} \right) + \tilde{r}(\mathbf{U}) \tag{152}$$



where we defined the functions

$$\tilde{\mathcal{L}} : \mathbb{R}^{n \times K} \rightarrow \mathbb{R} \tag{153}$$

$$\frac{1}{\sqrt{p}} \tilde{\mathbf{Z}} \mathbf{U} \rightarrow \mathcal{L} \left( f_0(\sqrt{\rho_{\tilde{\mathbf{w}}_0}} \mathbf{s}), \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \frac{1}{\sqrt{p}} \tilde{\mathbf{Z}} \mathbf{U} \right) \tag{154}$$

$$\tilde{r} : \mathbb{R}^{Kp \times K} \rightarrow \mathbb{R} \tag{155}$$

$$\mathbf{U} \rightarrow r \left( \Omega^{-1/2} \left( \frac{\sqrt{\gamma \tilde{\mathbf{c}}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} + \mathbf{U} \right) \right) - \boldsymbol{\nu}^\top \mathbf{U}^\top \tilde{\mathbf{c}} \tag{156}$$

and the random matrix  $\tilde{\mathbf{Z}}$  with i.i.d.  $\mathcal{N}(0,1)$  elements is independent from all other random quantities in the problem. The asymptotic properties of the unique solution to this optimisation problem can now be studied with a non-separable, matrix-valued AMP iteration. The AMP iteration solving problem equation (152) is given in the following lemma

**Lemma 6.** Consider the following AMP iteration

$$\mathbf{u}^{t+1} = \tilde{\mathbf{Z}}^\top \mathbf{h}_t(\mathbf{v}^t) - \mathbf{e}_t(\mathbf{u}^t) \langle \mathbf{h}'_t \rangle^\top \tag{157}$$

$$\mathbf{v}^t = \tilde{\mathbf{Z}} \mathbf{e}_t(\mathbf{u}^t) - \mathbf{h}_{t-1}(\mathbf{v}^{t-1}) \langle \mathbf{e}'_t \rangle^\top \tag{158}$$

where for any  $t \in \mathbb{N}$

$$\begin{aligned} \mathbf{h}_t(\mathbf{v}^t) &= \left( \mathbf{R}_{\mathcal{L}(\mathbf{y}, \cdot), \mathbf{S}^t} \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{v}^t \right) \right. \\ &\quad \left. - \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{v}^t \right) \right) (\mathbf{S}^t)^{-1} \end{aligned} \tag{159}$$

$$\mathbf{e}_t(\mathbf{u}^t) = \mathbf{R}_{r(\Omega^{-1/2}, \cdot), \hat{\mathbf{S}}^t} \left( \mathbf{u}^t \hat{\mathbf{S}}^t + \Omega^{-1/2} \mathbf{c} \boldsymbol{\nu}^\top \hat{\mathbf{S}}^t + \frac{\sqrt{\gamma \tilde{\mathbf{c}}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) - \frac{\sqrt{\gamma \tilde{\mathbf{c}}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \tag{160}$$

$$\text{and } \mathbf{S}^t = \langle (\mathbf{e}^t)' \rangle^\top, \quad \hat{\mathbf{S}}^t = - \left( \langle (\mathbf{h}^t)' \rangle^\top \right)^{-1}. \tag{161}$$

Then the fixed point  $(\mathbf{u}^\infty, \mathbf{v}^\infty)$  of this iteration verifies

$$\mathbf{R}_{r(\Omega^{-1/2}, \cdot), \hat{\mathbf{S}}^\infty} \left( \mathbf{u}^\infty \hat{\mathbf{S}}^\infty + \Omega^{-1/2} \mathbf{c} \boldsymbol{\nu}^\top \hat{\mathbf{S}}^\infty + \frac{\sqrt{\gamma \tilde{\mathbf{c}}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) - \frac{\sqrt{\gamma \tilde{\mathbf{c}}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} = \mathbf{U}^* \tag{162}$$

$$\mathbf{R}_{\mathcal{L}(\mathbf{y}, \cdot), \mathbf{S}^\infty} \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{v}^\infty \right) - \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} = \tilde{\mathbf{Z}} \mathbf{U}^* \tag{163}$$

where  $\mathbf{U}^*$  is the unique solution to the optimisation problem equation (152).

**Proof.** To find the correct form of the non-linearities in the AMP iteration, we match the optimality condition of problem equation (152) with the generic form of the fixed point of the AMP iteration equation (233). In the subsequent derivation, we absorb the scaling  $\frac{1}{\sqrt{d}}$  in the matrix  $\tilde{\mathbf{Z}}$ , such that its elements are i.i.d.  $\mathcal{N}(0,1/d)$ , and omit time

Fluctuations, bias, variance and ensemble of learners: exact asymptotics for convex losses in high-dimension indices for simplicity. Going back to problem equation (152), its optimality condition reads :

$$\tilde{\mathbf{Z}}^\top \partial \tilde{\mathcal{L}}(\tilde{\mathbf{Z}}\mathbf{U}) + \partial \tilde{r}(\mathbf{U}) = 0. \tag{164}$$

For any pair of  $K \times K$  symmetric positive definite matrices  $\mathbf{S}, \hat{\mathbf{S}}$ , this optimality condition is equivalent to

$$\tilde{\mathbf{Z}}^\top \left( \partial \tilde{\mathcal{L}}(\tilde{\mathbf{Z}}\mathbf{U}) \mathbf{S} + \tilde{\mathbf{Z}}\mathbf{U} \right) \mathbf{S}^{-1} + \left( \partial \tilde{r}(\mathbf{U}) \hat{\mathbf{S}} + \mathbf{U} \right) \hat{\mathbf{S}}^{-1} = \tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}}\mathbf{U}\mathbf{S}^{-1} + \mathbf{U}\hat{\mathbf{S}}^{-1} \tag{165}$$

where we added the same quantity on both sides of the equality. For the loss function, we can then introduce the resolvent, formally D-resolvent:

$$\hat{\mathbf{v}} = \partial \tilde{\mathcal{L}}(\tilde{\mathbf{Z}}\mathbf{U}) \mathbf{S} + \tilde{\mathbf{Z}}\mathbf{U} \iff \tilde{\mathbf{Z}}\mathbf{U} = \mathbf{R}_{\tilde{\mathcal{L}},\mathbf{S}}(\hat{\mathbf{v}}) \tag{166}$$

such that

$$\mathbf{R}_{\tilde{\mathcal{L}},\mathbf{S}}(\hat{\mathbf{v}}) = \left( \text{Id} + \partial \tilde{\mathcal{L}}(\bullet) \mathbf{S} \right)^{-1}(\hat{\mathbf{v}}) = \arg \min_{\mathbf{T} \in \mathbb{R}^{n \times K}} \left\{ \tilde{\mathcal{L}}(\mathbf{T}) + \frac{1}{2} \text{tr} \left( (\mathbf{T} - \hat{\mathbf{v}}) \mathbf{S}^{-1} (\mathbf{T} - \hat{\mathbf{v}})^\top \right) \right\}. \tag{167}$$

Similarly for the regularisation, introduce

$$\hat{\mathbf{u}} \equiv \left( \text{Id} + \partial \tilde{r}(\bullet) \hat{\mathbf{S}} \right)(\mathbf{U}) \quad \mathbf{U} = \mathbf{R}_{\tilde{r},\hat{\mathbf{S}}}(\hat{\mathbf{u}}) \tag{168}$$

where  $\mathbf{S} \in \mathbb{R}^{K \times K}$  is a positive definite matrix, and

$$\mathbf{R}_{\tilde{r},\hat{\mathbf{S}}}(\hat{\mathbf{v}}) = \left( \text{Id} + \partial \tilde{r}(\bullet) \hat{\mathbf{S}} \right)^{-1}(\hat{\mathbf{v}}) = \arg \min_{\mathbf{T} \in \mathbb{R}^{Kp \times K}} \left\{ \tilde{r}(\mathbf{T}) + \frac{1}{2} \text{tr} \left( (\mathbf{T} - \hat{\mathbf{v}}) \hat{\mathbf{S}}^{-1} (\mathbf{T} - \hat{\mathbf{v}})^\top \right) \right\} \tag{169}$$

where  $\hat{\mathbf{S}} \in \mathbb{R}^{K \times K}$  is a positive definite matrix, and  $\hat{\mathbf{v}} \in \mathbb{R}^{d \times K}$ . The optimality condition equation (165) may then be rewritten as:

$$\tilde{\mathbf{Z}}^\top \left( \mathbf{R}_{\tilde{\mathcal{L}},\mathbf{S}}(\hat{\mathbf{v}}) - \hat{\mathbf{v}} \right) \mathbf{S}^{-1} = \left( \hat{\mathbf{u}} - \mathbf{R}_{\tilde{r},\hat{\mathbf{S}}}(\hat{\mathbf{u}}) \right) \hat{\mathbf{S}}^{-1} \tag{170}$$

$$\tilde{\mathbf{Z}} \mathbf{R}_{\tilde{r},\hat{\mathbf{S}}}(\hat{\mathbf{u}}) = \mathbf{R}_{\tilde{\mathcal{L}},\mathbf{S}}(\hat{\mathbf{v}}) \tag{171}$$

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 equation (233), omitting the time indices for simplicity, reads:

$$\mathbf{u} + \mathbf{e}(\mathbf{u}) \langle \mathbf{h}' \rangle^\top = \tilde{\mathbf{Z}}^\top \mathbf{h}(\mathbf{v}) \tag{172}$$

$$\mathbf{v} + \mathbf{h}(\mathbf{v}) \langle \mathbf{e}' \rangle^\top = \tilde{\mathbf{Z}} \mathbf{e}(\mathbf{u}). \tag{173}$$

Matching this fixed point with the optimality condition equation (170) suggests the following mapping:

$$\begin{aligned} \mathbf{h}(\mathbf{v}) &= \left( \mathbf{R}_{\tilde{\mathcal{L}},\mathbf{S}}(\mathbf{v}) - \mathbf{v} \right) \mathbf{S}^{-1}, & \mathbf{S} &= \langle \mathbf{e}' \rangle^\top, \\ \mathbf{e}(\mathbf{u}) &= \mathbf{R}_{\tilde{r},\hat{\mathbf{S}}}(\mathbf{u}\hat{\mathbf{S}}), & \hat{\mathbf{S}} &= - \left( \langle \mathbf{h}' \rangle^\top \right)^{-1}, \end{aligned} \tag{174}$$

Fluctuations, bias, variance and ensemble of learners: exact asymptotics for convex losses in high-dimension

where we redefined  $\hat{\mathbf{u}} \equiv \hat{\mathbf{u}}\hat{\mathbf{S}}$  in (168). We are now left with the task of evaluating the resolvents of  $\tilde{\mathcal{L}}, \tilde{r}$  as expressions of the original functions  $\mathcal{L}, r$ . Starting with the loss function, we get

$$\mathbf{R}_{\tilde{\mathcal{L}}, \mathbf{S}}(\mathbf{v}) = \arg \min_{\mathbf{x} \in \mathbb{R}^{n \times K}} \left\{ \mathcal{L} \left( f_0(\sqrt{\rho_{\tilde{\mathbf{w}}_0}} \mathbf{s}), \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{x} \right) + \frac{1}{2} \text{tr} \left( (\mathbf{x} - \mathbf{v}) \mathbf{S}^{-1} (\mathbf{x} - \mathbf{v})^\top \right) \right\} \quad (175)$$

letting  $\tilde{\mathbf{x}} = \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{x}$ , the problem is equivalent to

$$\begin{aligned} \mathbf{R}_{\tilde{\mathcal{L}}, \mathbf{S}}(\mathbf{v}) &= \arg \min_{\tilde{\mathbf{x}} \in \mathbb{R}^{n \times K}} \left\{ \mathcal{L}(f_0(\sqrt{\rho_{\tilde{\mathbf{w}}_0}} \mathbf{s}), \tilde{\mathbf{x}}) \right. \\ &\quad \left. + \frac{1}{2} \text{tr} \left( \left( \tilde{\mathbf{x}} - \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{v} \right) \right) \mathbf{S}^{-1} \left( \tilde{\mathbf{x}} - \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{v} \right) \right)^\top \right) \right\} \\ &\quad - \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} - \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} \end{aligned} \quad (176)$$

$$= \mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{S}} \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{v} \right) - \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} - \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} \quad (177)$$

and the corresponding non-linearity will then be

$$\mathbf{h}(\mathbf{v}) = \left( \mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{S}} \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{v} \right) - \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma \kappa} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{v} \right) \right) \mathbf{S}^{-1} \quad (178)$$

Moving to the regularisation, the resolvent reads

$$\begin{aligned} \mathbf{R}_{\tilde{r}, \hat{\mathbf{S}}}(\mathbf{u}) &= \arg \min_{\mathbf{x} \in \mathbb{R}^{Kp \times K}} \left\{ r \left( \Omega^{-1/2} \left( \frac{\sqrt{\gamma \tilde{\mathbf{c}}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} + \mathbf{x} \right) \right) - \boldsymbol{\nu}^\top \mathbf{x}^\top \Omega^{-1/2} \mathbf{c} \right. \\ &\quad \left. + \frac{1}{2} \text{tr} \left( (\mathbf{x} - \mathbf{u}) \hat{\mathbf{S}}^{-1} (\mathbf{x} - \mathbf{u})^\top \right) \right\} \end{aligned} \quad (179)$$

letting  $\tilde{\mathbf{x}} = \frac{\sqrt{\gamma \tilde{\mathbf{c}}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} + \mathbf{x}$ , we obtain

$$\mathbf{R}_{\tilde{r}, \hat{\mathbf{S}}}(\mathbf{u}) = \arg \min_{\tilde{\mathbf{x}} \in \mathbb{R}^{Kp \times K}} \left\{ r \left( \Omega^{-1/2} \tilde{\mathbf{x}} \right) - \boldsymbol{\nu}^\top \tilde{\mathbf{x}}^\top \Omega^{-1/2} \mathbf{c} \right. \quad (180)$$

$$\begin{aligned} &\quad \left. + \frac{1}{2} \text{tr} \left( \left( \tilde{\mathbf{x}} - \left( \mathbf{u} + \frac{\sqrt{\gamma \tilde{\mathbf{c}}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right) \hat{\mathbf{S}}^{-1} \right. \right. \\ &\quad \left. \left. \times \left( \tilde{\mathbf{x}} - \left( \mathbf{u} + \frac{\sqrt{\gamma \tilde{\mathbf{c}}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right)^\top \right) \right\} - \frac{\sqrt{\gamma \tilde{\mathbf{c}}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \end{aligned} \quad (181)$$

$$= \arg \min_{\tilde{\mathbf{x}} \in \mathbb{R}^{Kp \times K}} \left\{ r \left( \Omega^{-1/2} \tilde{\mathbf{x}} \right) \right. \tag{182}$$

$$+ \frac{1}{2} \text{tr} \left( \left( \tilde{\mathbf{x}} - \left( \mathbf{u} + \Omega^{-1/2} \mathbf{c} \boldsymbol{\nu}^\top \hat{\mathbf{S}} + \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right) \hat{\mathbf{S}}^{-1} \right. \\ \left. \times \left( \tilde{\mathbf{x}} - \left( \mathbf{u} + \Omega^{-1/2} \mathbf{c} \boldsymbol{\nu}^\top \hat{\mathbf{S}} + \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right)^\top \right) \left. \right\} \tag{183}$$

$$- \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \tag{184}$$

$$\mathbf{R}_{r(\Omega^{-1/2}), \hat{\mathbf{S}}} \left( \mathbf{u} + \Omega^{-1/2} \mathbf{c} \boldsymbol{\nu}^\top \hat{\mathbf{S}} + \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) - \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}}. \tag{185}$$

Which gives the following non-linearity for the AMP iteration

$$e(\mathbf{u}) = \mathbf{R}_{r(\Omega^{-1/2}), \hat{\mathbf{S}}} \left( \mathbf{u} \hat{\mathbf{S}} + \Omega^{-1/2} \mathbf{c} \boldsymbol{\nu}^\top \mathbf{V} + \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) - \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}}. \tag{186}$$

□

The following lemma then gives the exact asymptotics at each time step of the AMP iteration solving problem equation (152) : its *state evolution equations*.

**Lemma 7.** Consider the AMP iteration equations (157)–(161). Assume it is initialised with  $\mathbf{u}^0$  such that  $\lim_{d \rightarrow \infty} \frac{1}{d} \|\mathbf{e}_0(\mathbf{u}^0)^\top \mathbf{e}_0(\mathbf{u}^0)\|_F$  exists, a positive definite matrix  $\hat{\mathbf{S}}_0$ , and  $\mathbf{h}_{-1} \equiv \mathbf{0}$ . Then for any  $t \in \mathbb{N}$ , and any pair of sequences of uniformly pseudo-Lipschitz functions  $\phi_{1,n} : \mathbb{R}^{Kp \times K}$  and  $\phi_{2,n} : \mathbb{R}^{n \times K}$ , the following holds

$$\phi_{1,n}(\mathbf{u}^t) \stackrel{P}{\simeq} \mathbb{E} \left[ \phi_{1,n} \left( \mathbf{G} \left( \hat{\mathbf{Q}}^t \right)^{1/2} \right) \right] \tag{187}$$

$$\phi_{2,n}(\mathbf{v}^t) \stackrel{P}{\simeq} \mathbb{E} \left[ \phi_{2,n} \left( \mathbf{H} \left( \mathbf{Q}^t \right)^{1/2} \right) \right] \tag{188}$$

where  $\mathbf{G} \in \mathbb{R}^{Kp \times K}$  and  $\mathbf{H} \in \mathbb{R}^{n \times K}$  are independent random matrices with i.i.d. standard normal elements, and  $\mathbf{Q}^t, \hat{\mathbf{Q}}^t, \mathbf{V}^t, \hat{\mathbf{V}}^t$  are given by the equations

$$\mathbf{Q}^t = \frac{1}{p} \mathbb{E} \left[ \left( \mathbf{R}_{r(\Omega^{-1/2}), (\hat{\mathbf{V}}^t)^{-1}} \left( \mathbf{G} \left( \hat{\mathbf{Q}}^t \right)^{1/2} \left( \hat{\mathbf{V}}^t \right)^{-1} + \Omega^{-1/2} \mathbf{c} \boldsymbol{\nu}^\top \left( \hat{\mathbf{V}}^t \right)^{-1} + \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) - \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right)^\top \right. \\ \left. \times \left( \mathbf{R}_{r(\Omega^{-1/2}), (\hat{\mathbf{V}}^t)^{-1}} \left( \mathbf{G} \left( \hat{\mathbf{Q}}^t \right)^{1/2} \left( \hat{\mathbf{V}}^t \right)^{-1} + \Omega^{-1/2} \mathbf{c} \boldsymbol{\nu}^\top \left( \hat{\mathbf{V}}^t \right)^{-1} + \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) - \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right] \tag{189}$$

$$\hat{\mathbf{Q}}^t = \frac{1}{p} \mathbb{E} \left[ \left( \left( \mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}^{t-1}}(\cdot) - Id \right) \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H} \left( \mathbf{Q}^{t-1} \right)^{1/2} \right) \left( \mathbf{V}^{t-1} \right)^{-1} \right)^\top \right. \tag{190}$$

$$\left. \times \left( \mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}^{t-1}}(\cdot) - Id \right) \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H} \left( \mathbf{Q}^{t-1} \right)^{1/2} \right) \left( \mathbf{V}^{t-1} \right)^{-1} \right] \tag{191}$$

$$\mathbf{V}^t = \frac{1}{p} \mathbb{E} \left[ (\hat{\mathbf{Q}}^t)^{-1/2} \mathbf{G}^\top R_{r(\Omega^{-1/2}, \hat{\mathbf{V}}^t)^{-1}} \left( \mathbf{G}(\hat{\mathbf{Q}}^t)^{1/2}(\hat{\mathbf{V}}^t)^{-1} + \Omega^{-1/2} \mathbf{c}\boldsymbol{\nu}^\top (\hat{\mathbf{V}}^t)^{-1} + \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right] \quad (192)$$

$$\hat{\mathbf{V}}^t = -\frac{1}{p} \mathbb{E} \left[ (\mathbf{Q}^{t-1})^{-1/2} \mathbf{H}^\top \left( \mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}^{t-1}}(\cdot) - Id \right) \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma} \kappa \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H}(\mathbf{Q}^{t-1})^{1/2} \right) (\mathbf{V}^{t-1})^{-1} \right]. \quad (193)$$

**Proof.** Owing to the properties of Bregman proximity operators (Bauschke *et al* 2003, 2006), the update functions in the AMP iteration equations (157)–(161) are Lipschitz continuous. Thus under the assumptions made on the initialisation, the assumptions of theorem 11 are verified, which gives the desired result.  $\square$

**Lemma 8.** Consider iteration equations (157)–(161), where the parameters  $\mathbf{Q}, \hat{\mathbf{Q}}, \mathbf{V}, \hat{\mathbf{V}}$  are initialised at any fixed point of the state evolution equations of lemma 7. For any sequence initialised with  $\hat{\mathbf{V}}_0 = \hat{\mathbf{V}}$  and  $\mathbf{u}_0$  such that

$$\lim_{d \rightarrow \infty} \frac{1}{d} \mathbf{e}_0(\mathbf{u}_0)^\top \mathbf{e}_0(\mathbf{u}_0) = \mathbf{Q} \quad (194)$$

the following holds

$$\lim_{t \rightarrow \infty} \lim_{p \rightarrow \infty} \frac{1}{\sqrt{p}} \|\mathbf{u}^t - \mathbf{u}^*\|_F = 0 \quad \lim_{t \rightarrow \infty} \lim_{d \rightarrow \infty} \frac{1}{\sqrt{p}} \|\mathbf{v}^t - \mathbf{v}^*\|_F = 0. \quad (195)$$

**Proof.** The proof of this lemma is identical to that of lemma 7 from Loureiro *et al* (2021b).  $\square$

Combining these results, we obtain the following asymptotic characterisation of  $\mathbf{U}^*$ .

**Lemma 9.** For any fixed  $\mathbf{m}$  and  $\boldsymbol{\nu}$  in their feasibility sets, let  $\mathbf{U}^*$  be the unique solution to the optimisation problem equation (152). Then, for any sequences (in the problem dimension) of pseudo-Lipschitz functions of order 2  $\phi_{1,n} : \mathbb{R}^{n \times K} \rightarrow \mathbb{R}$  and  $\phi_{2,n} : \mathbb{R}^{Kp \times K} \rightarrow \mathbb{R}$ , the following holds

$$\phi_{1,n}(\mathbf{U}^*) \stackrel{P}{\simeq} \mathbb{E} \left[ \phi_{1,n} \left( \mathbf{R}_{r(\Omega^{-1/2}, \hat{\mathbf{V}}^{-1})} \left( \mathbf{G} \hat{\mathbf{Q}}^{1/2} \hat{\mathbf{V}}^{-1} + \Omega^{-1/2} \mathbf{c}\boldsymbol{\nu}^\top \hat{\mathbf{V}}^{-1} + \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) - \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right] \quad (196)$$

$$\phi_{2,n} \left( \frac{1}{\sqrt{p}} \tilde{\mathbf{Z}} \mathbf{U}^* \right) \stackrel{P}{\simeq} \mathbb{E} \left[ \phi_{2,n} \left( \mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}} \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma} \kappa \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H} \hat{\mathbf{Q}}^{1/2} \right) - \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} - \tilde{\mathbf{s}} \frac{\sqrt{\gamma} \kappa \mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} \right) \right] \quad (197)$$

where  $\mathbf{G} \in \mathbb{R}^{Kp \times K}$  and  $\mathbf{H} \in \mathbb{R}^{n \times K}$  are independent random matrices with i.i.d. standard normal elements, and  $\mathbf{Q}, \hat{\mathbf{Q}}, \mathbf{V}, \hat{\mathbf{V}}$  are given by the fixed point (assumed to be unique) of the following set of self consistent equations

$$\mathbf{Q} = \frac{1}{p} \mathbb{E} \left[ \left( \mathbf{R}_{r(\Omega^{-1/2}, \hat{\mathbf{V}}^{-1})} \left( \mathbf{G} \hat{\mathbf{Q}}^{1/2} \hat{\mathbf{V}}^{-1} + \Omega^{-1/2} \mathbf{c}\boldsymbol{\nu}^\top \hat{\mathbf{V}}^{-1} + \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) - \frac{\sqrt{\gamma} \tilde{\mathbf{c}} \mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right)^\top \right] \quad (198)$$

$$\left( \mathbf{R}_{r(\Omega^{-1/2}), \hat{\mathbf{V}}^{-1}} \left( \mathbf{G}\hat{\mathbf{Q}}^{1/2}\hat{\mathbf{V}}^{-1} + \Omega^{-1/2}\mathbf{c}\boldsymbol{\nu}^\top\hat{\mathbf{V}}^{-1} + \frac{\sqrt{\gamma}\tilde{\mathbf{c}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) - \frac{\sqrt{\gamma}\tilde{\mathbf{c}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \quad (199)$$

$$\begin{aligned} \hat{\mathbf{Q}} = & \frac{1}{p} \mathbb{E} \left[ \left( (\mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}}(\cdot) - Id) \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma\kappa}\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H}\mathbf{Q}^{1/2} \right) \mathbf{V}^{-1} \right)^\top \right. \\ & \left. \times \left( (\mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}}(\cdot) - Id) \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma\kappa}\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H}\mathbf{Q}^{1/2} \right) \mathbf{V}^{-1} \right) \right] \quad (200) \end{aligned}$$

$$\mathbf{V} = \frac{1}{p} \mathbb{E} \left[ \hat{\mathbf{Q}}^{-1/2} \mathbf{G}^\top \mathbf{R}_{r(\Omega^{-1/2}), \hat{\mathbf{V}}^{-1}} \left( \mathbf{G}\hat{\mathbf{Q}}^{1/2}\hat{\mathbf{V}}^{-1} + \Omega^{-1/2}\mathbf{c}\boldsymbol{\nu}^\top\hat{\mathbf{V}}^{-1} + \frac{\sqrt{\gamma}\tilde{\mathbf{c}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right] \quad (201)$$

$$\hat{\mathbf{V}} = -\frac{1}{p} \mathbb{E} \left[ \mathbf{Q}^{-1/2} \mathbf{H}^\top \left( (\mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}}(\cdot) - Id) \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma\kappa}\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H}\mathbf{Q}^{1/2} \right) \mathbf{V}^{-1} \right) \right]. \quad (202)$$

**Proof.** Combining the results of the previous lemmas, this proof is close to that of theorem 1.5 in Bayati and Montanari (2011b).  $\square$

Returning to the optimisation problem on  $\mathbf{m}, \boldsymbol{\nu}$  in equation (151), the solution  $\mathbf{U}^*$ , at any dimension, verifies the zero gradient conditions on  $\mathbf{m}, \boldsymbol{\nu}$ :

$$\partial \boldsymbol{\nu} = 0 \iff (\mathbf{U}^*)^\top \tilde{\mathbf{c}} = 0 \quad (203)$$

$$\begin{aligned} \partial \mathbf{m} = 0 \iff & \left( \frac{\mathbf{s}}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \frac{\tilde{\mathbf{s}}\sqrt{\gamma\kappa}}{\rho_{\tilde{\mathbf{c}}}} \right)^\top \mathcal{L} \left( f_0(\sqrt{\rho_{\tilde{\mathbf{w}}_0}}\mathbf{s}), \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma\kappa}\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \frac{1}{\sqrt{p}} \tilde{\mathbf{Z}}\mathbf{U} \right) \\ & + \frac{\sqrt{\gamma}\tilde{\mathbf{V}}^\top}{\rho_{\tilde{\mathbf{c}}}} \Omega^{-1/2} \partial r \left( \Omega^{-1/2} \left( \frac{\sqrt{\gamma}\tilde{\mathbf{c}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} + \mathbf{U} \right) \right) = 0. \quad (204) \end{aligned}$$

Using lemma 9 while assuming the subgradients of  $\mathcal{L}, r$  are pseudo-Lipschitz (we discuss this assumption in section B.4), we obtain for  $\mathbf{m}$

$$\frac{1}{p} \mathbb{E} \left[ \left( \mathbf{R}_{r(\Omega^{-1/2}), \hat{\mathbf{V}}^{-1}} \left( \mathbf{G}\hat{\mathbf{Q}}^{1/2}\hat{\mathbf{V}}^{-1} + \Omega^{-1/2}\mathbf{c}\boldsymbol{\nu}^\top\hat{\mathbf{V}}^{-1} + \frac{\sqrt{\gamma}\tilde{\mathbf{c}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) - \frac{\sqrt{\gamma}\tilde{\mathbf{c}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right)^\top \tilde{\mathbf{c}} \right] = 0 \quad (205)$$

$$\iff \mathbf{m} = \frac{1}{\sqrt{dp}} \mathbb{E} \left[ \tilde{\mathbf{c}}^\top \mathbf{R}_{r(\Omega^{-1/2}), \hat{\mathbf{V}}^{-1}} \left( \mathbf{G}\hat{\mathbf{Q}}^{1/2}\hat{\mathbf{V}}^{-1} + \Omega^{-1/2}\mathbf{c}\boldsymbol{\nu}^\top\hat{\mathbf{V}}^{-1} + \frac{\sqrt{\gamma}\tilde{\mathbf{c}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right] \quad (206)$$

and for  $\nu$

$$\frac{1}{p} \mathbb{E} \left[ \left( \frac{\mathbf{s}}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \frac{\tilde{\mathbf{s}}\sqrt{\gamma\kappa}}{\rho_{\tilde{\mathbf{c}}}} \right)^\top \partial \mathcal{L} \left( f_0(\sqrt{\rho_{\tilde{\mathbf{w}}_0}}\mathbf{s}), \mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}} \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma\kappa}\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H}\hat{\mathbf{Q}}^{1/2} \right) \right) \right] \quad (207)$$

$$+ \frac{\sqrt{\gamma\tilde{\mathbf{c}}}^\top}{\rho_{\tilde{\mathbf{c}}}} \Omega^{-1/2} \partial r \left( \Omega^{-1/2} \left( \mathbf{R}_{r(\Omega^{-1/2}), \hat{\mathbf{V}}^{-1}} \left( \mathbf{G}\hat{\mathbf{Q}}^{1/2}\hat{\mathbf{V}}^{-1} + \Omega^{-1/2}\mathbf{c}\nu^\top\hat{\mathbf{V}}^{-1} + \frac{\sqrt{\gamma\tilde{\mathbf{c}}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right) \right) = 0. \quad (208)$$

Using the definition of D-resolvents, this is equivalent to

$$\frac{1}{p} \mathbb{E} \left[ \left( \frac{\mathbf{s}}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \frac{\tilde{\mathbf{s}}\sqrt{\gamma\kappa}}{\rho_{\tilde{\mathbf{c}}}} \right)^\top (Id - \mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}}(\cdot)) \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma\kappa}\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H}\hat{\mathbf{Q}}^{1/2} \right) \mathbf{V}^{-1} \right] \quad (209)$$

$$+ \frac{\sqrt{\gamma\tilde{\mathbf{c}}}^\top}{\rho_{\tilde{\mathbf{c}}}} (Id - \mathbf{R}_{r(\Omega^{-1/2}), \hat{\mathbf{V}}^{-1}}(\cdot)) \left( \mathbf{G}\hat{\mathbf{Q}}^{1/2}\hat{\mathbf{V}}^{-1} + \Omega^{-1/2}\mathbf{c}\nu^\top\hat{\mathbf{V}}^{-1} + \frac{\sqrt{\gamma\tilde{\mathbf{c}}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \hat{\mathbf{V}} = 0 \quad (210)$$

which simplifies to

$$\begin{aligned} \nu^\top = & -\frac{1}{\sqrt{\gamma p}} \mathbb{E} \left[ \left( \frac{\mathbf{s}}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \frac{\tilde{\mathbf{s}}\sqrt{\gamma\kappa}}{\rho_{\tilde{\mathbf{c}}}} \right)^\top (Id - \mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}}(\cdot)) \right. \\ & \left. \times \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma\kappa}\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H}\hat{\mathbf{Q}}^{1/2} \right) \mathbf{V}^{-1} \right] \end{aligned} \quad (211)$$

which brings us to the following set of six self consistent equations

$$\mathbf{Q} = \frac{1}{p} \mathbb{E} \left[ \left( \mathbf{R}_{r(\Omega^{-1/2}), \hat{\mathbf{V}}^{-1}} \left( \mathbf{G}\hat{\mathbf{Q}}^{1/2}\hat{\mathbf{V}}^{-1} + \Omega^{-1/2}\mathbf{c}\nu^\top\hat{\mathbf{V}}^{-1} + \frac{\sqrt{\gamma\tilde{\mathbf{c}}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) - \frac{\sqrt{\gamma\tilde{\mathbf{c}}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right)^\top \right] \quad (212)$$

$$\left( \mathbf{R}_{r(\Omega^{-1/2}), \hat{\mathbf{V}}^{-1}} \left( \mathbf{G}\hat{\mathbf{Q}}^{1/2}\hat{\mathbf{V}}^{-1} + \Omega^{-1/2}\mathbf{c}\nu^\top\hat{\mathbf{V}}^{-1} + \frac{\sqrt{\gamma\tilde{\mathbf{c}}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) - \frac{\sqrt{\gamma\tilde{\mathbf{c}}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right] \quad (213)$$

$$\begin{aligned} \hat{\mathbf{Q}} = & \frac{1}{p} \mathbb{E} \left[ \left( (\mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}}(\cdot) - Id) \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma\kappa}\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H}\hat{\mathbf{Q}}^{1/2} \right) \mathbf{V}^{-1} \right)^\top \right. \\ & \left. \times \left( (\mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}}(\cdot) - Id) \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma\kappa}\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H}\hat{\mathbf{Q}}^{1/2} \right) \mathbf{V}^{-1} \right) \right] \end{aligned} \quad (214)$$

$$\mathbf{V} = \frac{1}{p} \mathbb{E} \left[ \hat{\mathbf{Q}}^{-1/2} \mathbf{G}^\top \mathbf{R}_{r(\Omega^{-1/2}), \hat{\mathbf{V}}^{-1}} \left( \mathbf{G}\hat{\mathbf{Q}}^{1/2}\hat{\mathbf{V}}^{-1} + \Omega^{-1/2}\mathbf{c}\nu^\top\hat{\mathbf{V}}^{-1} + \frac{\sqrt{\gamma\tilde{\mathbf{c}}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right] \quad (215)$$

$$\hat{\mathbf{V}} = -\frac{1}{p} \mathbb{E} \left[ \mathbf{Q}^{-1/2} \mathbf{H}^\top \left( (\mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}}(\cdot) - Id) \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma\kappa}\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H}\hat{\mathbf{Q}}^{1/2} \right) \mathbf{V}^{-1} \right) \right] \quad (216)$$

$$\mathbf{m} = \frac{1}{\sqrt{dp}} \mathbb{E} \left[ \tilde{\mathbf{c}}^\top \mathbf{R}_{r(\Omega^{-1/2}), \hat{\mathbf{V}}^{-1}} \left( \mathbf{G}\hat{\mathbf{Q}}^{1/2}\hat{\mathbf{V}}^{-1} + \Omega^{-1/2}\mathbf{c}\nu^\top\hat{\mathbf{V}}^{-1} + \frac{\sqrt{\gamma\tilde{\mathbf{c}}}\mathbf{m}^\top}{\rho_{\tilde{\mathbf{c}}}} \right) \right] \quad (217)$$

$$\begin{aligned} \boldsymbol{\nu}^\top &= -\frac{1}{\sqrt{\gamma p}} \mathbb{E} \left[ \left( \frac{\mathbf{s}}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \frac{\tilde{\mathbf{s}}\sqrt{\gamma\kappa}}{\rho_{\tilde{\mathbf{c}}}} \right)^\top (Id - \mathbf{R}_{\mathcal{L}(y, \cdot), \mathbf{V}}(\cdot)) \right. \\ &\quad \left. \times \left( \mathbf{s} \frac{\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{w}}_0}}} + \tilde{\mathbf{s}} \frac{\sqrt{\gamma\kappa}\mathbf{m}^\top}{\sqrt{\rho_{\tilde{\mathbf{c}}}}} + \mathbf{H}\hat{\mathbf{Q}}^{1/2} \right) \mathbf{V}^{-1} \right]. \end{aligned} \tag{218}$$

This set of equations then characterises the asymptotic distribution of the estimator  $\mathbf{U}^*$  in the sense of lemma 9, with the optimal values of  $\mathbf{m}$  and  $\boldsymbol{\nu}$ . Using the definition of  $\mathbf{U}^*$  and  $\tilde{\mathbf{Z}}\mathbf{U}^*$ , along with the definition of the function  $r$  w.r.t. the original regularisation function, a tedious but straightforward calculation allows reconstruct the asymptotic properties of  $\mathbf{W}^*$  and of the set  $\{\mathbf{X}_k \mathbf{w}_k^*\}_{k=1}^K$  given in the main text.

### B.3. Relaxing the strong convexity constraint

Assuming the set of self consistent equations (212) have a unique fixed point regardless of the strong convexity assumption, this solution defines a unique set of six order parameters for the  $\lambda_2 = 0$  case. Furthermore, using proposition 12, the unique estimator  $\mathbf{W}^*(\lambda_2)$  solving problem equation (111) for strictly positive  $\lambda_2$  converges to the least-norm solution to the convex (but not strongly) equation (110). Thus, for any pseudo-Lipschitz observable of  $\mathbf{U}^*(\lambda_2)$ , we have, on the one side a continuous function of  $\lambda_2$  with a unique continuous extension at  $\lambda_2 = 0$ , and on the other side a function of  $\lambda_2$  prescribed by the expectation taken w.r.t. the asymptotic Gaussian model parametrised by the state evolution parameters which is defined for all positive values of  $\lambda_2$ . Since both functions match for any strictly positive  $\lambda_2$ , continuity implies they also match for  $\lambda_2 = 0$  and we obtain the exact asymptotics of the least  $\ell_2$  norm solution of problem equation (110). Regarding the uniqueness of the solution to the fixed point equations (212), it is shown in Loureiro *et al* (2021a) that a similar set of equations, although for a vector valued variable, i.e. no ensembling, the solution is unique even if the original problem is not strictly convex. This is proven by showing that the fixed point equations are the solution of a strictly convex problem. We expect this to be true here as well, and leave this part for a longer version of this paper.

### B.4. A comment on non-pseudo-Lipschitz subgradients

Provided the subgradients in equation (203) are pseudo-Lipschitz continuous, the proof goes through. However some convex functions commonly used in machine learning, such as the hinge loss or the  $\ell_1$  norm for the penalty, have non-pseudo-Lipschitz gradient. To circumvent this issue, one can consider the optimisation problem where both loss and regularisation are replaced by their Moreau envelopes with strictly positive parameters  $\tau_1, \tau_2$ , as is done in Celentano *et al* (2020) for the LASSO. Moreau envelopes are everywhere differentiable and have Lipschitz gradient for strictly positive values of their parameter (Bauschke *et al* 2011), thus the asymptotic characterisation holds. One can then take the parameters to zero, using the fact that the limit at zero in the parameters of Moreau envelopes is well defined (Bauschke *et al* 2011), recovering the original function. Since proximity operators are defined as strongly convex problems, the sequence of problems defined by the proximal operator of a Moreau envelope with



decreasing parameter converges to the proximal operator of the original function when the parameter is taken to zero. Finally, inverting the expectations on random quantities with the limit taking the parameters of the Moreau envelopes to zero can be done by verifying the dominated convergence theorem using the firm-nonexpansiveness of proximity operators and the corresponding bounds on their norms, see Bauschke *et al* (2011) chapter 4, section 1. We leave the details of this part to a longer version of this paper.

### B.5. Toolbox

In this section, we reproduce part of the appendix of Loureiro *et al* (2021b) for completeness, in order to give an overview of the main concepts and tools on AMP algorithms which will be required for the proof.

*B.5.1. Notations.* For a given function  $\phi: \mathbb{R}^{d \times K} \rightarrow \mathbb{R}^{d \times K}$ , we write:

$$\phi(\mathbf{X}) = \begin{bmatrix} \phi^1(\mathbf{X}) \\ \vdots \\ \phi^d(\mathbf{X}) \end{bmatrix} \in \mathbb{R}^{d \times K} \tag{219}$$

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

$$\frac{\partial \phi^i}{\partial \mathbf{X}_j}(\mathbf{X}) = \begin{bmatrix} \frac{\partial \phi^i_1(\mathbf{X})}{\partial X_{j1}} & \dots & \frac{\partial \phi^i_1(\mathbf{X})}{\partial X_{jK}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \phi^i_K(\mathbf{X})}{\partial X_{j1}} & \dots & \frac{\partial \phi^i_K(\mathbf{X})}{\partial X_{jK}} \end{bmatrix} \in \mathbb{R}^{K \times K}. \tag{220}$$

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  $\overset{\text{P}}{\simeq}$  denotes convergence in probability. We start with some definitions that commonly appear in the AMP literature, see e.g. Bayati and Montanari (2011a), Javanmard and Montanari (2013), Berthier *et al* (2020). 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 10 (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  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times K}$ ,

$$\frac{\|\phi(\mathbf{X}) - \phi(\mathbf{Y})\|_F}{\sqrt{m}} \leq L(k, K) \left( 1 + \left( \frac{\|\mathbf{X}\|_F}{\sqrt{n}} \right)^{k-1} + \left( \frac{\|\mathbf{Y}\|_F}{\sqrt{n}} \right)^{k-1} \right) \frac{\|\mathbf{X} - \mathbf{Y}\|_F}{\sqrt{n}} \tag{221}$$

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  $f : \mathbb{R}^{n \times K} \rightarrow \mathbb{R}, \mathbf{X} \mapsto \frac{1}{n} \|\mathbf{X}\|_F^2$  is pseudo-Lipshitz of order 2.

*B.5.2. 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. Bauschke *et al* (2011), Parikh and Boyd (2014). 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 Thrampoulidis *et al* (2018), Loureiro *et al* (2021a). 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

$$\operatorname{arg\,min}_{\mathbf{T} \in \mathbb{R}^{d \times K}} \left\{ f(\mathbf{T}) + \frac{1}{2} \operatorname{tr} \left( (\mathbf{T} - \mathbf{X}) \mathbf{V}^{-1} (\mathbf{T} - \mathbf{X})^\top \right) \right\}. \tag{222}$$

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

$$\operatorname{prox}_{f(\bullet \mathbf{V}^{1/2})} \left( \mathbf{X} \mathbf{V}^{-1/2} \right) \mathbf{V}^{1/2} \in \mathbb{R}^{d \times K} \tag{223}$$

or as a Bregman proximal operator (Bauschke *et al* 2003) defined with the Bregman distance induced by the strictly convex, coercive function (for positive definite  $\mathbf{V}$ )

$$\mathbf{X} \mapsto \frac{1}{2} \operatorname{tr} \left( \mathbf{X} \mathbf{V}^{-1} \mathbf{X}^\top \right) \tag{224}$$

which justifies the use of the Bregman resolvent

$$\operatorname{arg\,min}_{\mathbf{T} \in \mathbb{R}^{d \times K}} \left\{ f(\mathbf{T}) + \frac{1}{2} \operatorname{tr} \left( (\mathbf{T} - \mathbf{X}) \mathbf{V}^{-1} (\mathbf{T} - \mathbf{X})^\top \right) \right\} = (\operatorname{Id} + \partial f(\bullet) \mathbf{V})^{-1}(\mathbf{X}). \tag{225}$$

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 (224), see e.g. Bauschke *et al* (2003, 2006). 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  $\mathbf{x} \in \mathcal{X}$

$$T(\mathbf{x}) = \operatorname{arg\,min}_{\mathbf{y} \in \mathcal{X}} \{ f(\mathbf{x}) + D_g(\mathbf{x}, \mathbf{y}) \}. \tag{226}$$

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

$$\langle T\mathbf{x} - T\mathbf{y}, \nabla g(T\mathbf{x}) - \nabla g(T\mathbf{y}) \rangle \leq \langle T\mathbf{x} - T\mathbf{y}, \nabla g(\mathbf{x}) - \nabla g(\mathbf{y}) \rangle \quad (227)$$

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

*B.5.3. Gradients of Bregman envelopes.* Consider, for any  $\mathbf{X} \in \mathbb{R}^{d \times K}$  the Bregman envelope

$$\mathcal{M}_{f, \mathbf{V}}(\mathbf{X}) = \inf_{\mathbf{T} \in \mathbb{R}^{d \times K}} \left\{ f(\mathbf{T}) + \frac{1}{2} \text{tr} \left( (\mathbf{T} - \mathbf{X}) \mathbf{V}^{-1} (\mathbf{T} - \mathbf{X})^\top \right) \right\} \quad (228)$$

then

$$\nabla_{\mathbf{X}} \mathcal{M}_{f, \mathbf{V}}(\mathbf{X}) = \left( \mathbf{X} - (\text{Id} + \partial f(\bullet) \mathbf{V})^{-1}(\mathbf{X}) \right) \mathbf{V}^{-1} \quad (229)$$

and

$$\nabla_{\mathbf{V}} \mathcal{M}_{f, \mathbf{V}}(\mathbf{X}) = -\frac{1}{2} \left\| \left( \mathbf{X} - (\text{Id} + \partial f(\bullet) \mathbf{V})^{-1}(\mathbf{X}) \right) \mathbf{V}^{-1} \right\|_F^2. \quad (230)$$

*B.5.4. Gaussian concentration.* Gaussian concentration properties are at the root of this proof. Such properties are reviewed in more detail, for example, in Vershynin (2018). We refer the interested reader to related works for a more detailed discussion.

*B.5.5. Approximate message-passing.* Approximate message-passing algorithms (Donoho *et al* 2009, Rangan 2011, Donoho and Montanari 2016) are a statistical physics inspired (Mézard *et al* 1987, Zdeborová and Krzakala 2016) family of iterations which can be used to solve high dimensional inference problems. 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 AMP 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 Javanmard and Montanari (2013), Berthier *et al* (2020) included in Gerbelot and Berthier (2021). 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

$$\left\{ \mathbf{h}_t : \mathbb{R}^{n \times K} \rightarrow \mathbb{R}^{n \times K} \right\}_{t \in \mathbb{N}} \quad \left\{ \mathbf{e}_t : \mathbb{R}^{d \times K} \rightarrow \mathbb{R}^{d \times K} \right\}_{t \in \mathbb{N}} \quad (231)$$

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

$$\lim_{d \rightarrow \infty} \frac{1}{d} \left\| \mathbf{e}_0(\mathbf{u}^0)^\top \mathbf{e}_0(\mathbf{u}^0) \right\|_F \quad (232)$$

exists, and recursively define:

Fluctuations, bias, variance and ensemble of learners: exact asymptotics for convex losses in high-dimension

$$\mathbf{u}^{t+1} = \mathbf{A}^\top \mathbf{h}_t(\mathbf{v}^t) - \mathbf{e}_t(\mathbf{u}^t) \langle \mathbf{h}'_t \rangle^\top \quad (233)$$

$$\mathbf{v}^t = \mathbf{A} \mathbf{e}_t(\mathbf{u}^t) - \mathbf{h}_{t-1}(\mathbf{v}^{t-1}) \langle \mathbf{e}'_t \rangle^\top \quad (234)$$

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:

$$\langle \mathbf{h}'_t \rangle = \frac{1}{d} \sum_{i=1}^n \frac{\partial \mathbf{h}_t^i}{\partial \mathbf{v}_i}(\mathbf{v}^t) \in \mathbb{R}^{K \times K} \quad \langle \mathbf{e}'_t \rangle = \frac{1}{d} \sum_{i=1}^d \frac{\partial \mathbf{e}_t^i}{\partial \mathbf{u}_i}(\mathbf{u}^t) \in \mathbb{R}^{K \times K}. \quad (235)$$

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

$$\mathbf{Q}_{t+1,s} = \mathbf{Q}_{s,t+1} = \lim_{d \rightarrow \infty} \frac{1}{d} \mathbb{E} \left[ \mathbf{e}_s(\hat{\mathbf{Z}}^s)^\top \mathbf{e}_{t+1}(\hat{\mathbf{Z}}^{t+1}) \right] \in \mathbb{R}^{K \times K} \quad (236)$$

$$\hat{\mathbf{Q}}_{t+1,s+1} = \hat{\mathbf{Q}}_{s+1,t+1} = \lim_{d \rightarrow \infty} \frac{1}{d} \mathbb{E} \left[ \mathbf{h}_s(\mathbf{Z}^s)^\top \mathbf{h}_t(\mathbf{Z}^t) \right] \in \mathbb{R}^{K \times K} \quad (237)$$

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 11.** *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  $n, d \rightarrow +\infty$ :*

$$\phi_n(\mathbf{u}^0, \mathbf{v}^0, \mathbf{u}^1, \mathbf{v}^1, \dots, \mathbf{v}^{t-1}, \mathbf{u}^t) \stackrel{\text{P}}{\simeq} \mathbb{E} \left[ \phi_n(\mathbf{u}^0, \mathbf{Z}^0, \hat{\mathbf{Z}}^1, \mathbf{Z}^1, \dots, \mathbf{Z}^{t-1}, \hat{\mathbf{Z}}^t) \right] \quad (238)$$

where

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

**B.5.6. A useful result from convex analysis** Here we remind a result from Bauschke *et al* (2011) describing the limiting behaviour of regularised estimators for vanishing regularisation.

**Proposition 12 (theorem 26.20 from Bauschke *et al* (2011)).** *Let  $f$  and  $h$  be proper, lower semi-continuous, convex functions. Suppose that  $\text{argmin } f \cap \text{dom}(g) \neq \emptyset$  and that  $h$  is coercive and strictly convex. Then  $g$  admits a unique minimiser  $\mathbf{x}_0$  over  $\text{argmin } f$  and, for every  $\epsilon \in ]0, 1[$ , the regularised problem*

$$\text{argmin}_{\mathbf{x}} f(\mathbf{x}) + \epsilon h(\mathbf{x}) \quad (239)$$

*admits a unique solution  $\mathbf{x}_\epsilon$ . If we assume further that  $h$  is uniformly convex on any closed ball of the input space, then  $\lim_{\epsilon \rightarrow 0} \mathbf{x}_\epsilon = \mathbf{x}_0$ .*

## References

- Adlam B and Pennington J 2020b The neural tangent kernel in high dimensions: triple descent and a multi-scale theory of generalization *Int. Conf. on Machine Learning* (PMLR) pp 74–84
- Adlam B and Pennington J 2020a Understanding double descent requires a fine-grained bias-variance decomposition *Advances in Neural Information Processing Systems* vol 33, ed H Larochelle, M Ranzato, R Hadsell, M F Balcan and H Lin (Curran Associates, Inc.) pp 11022–32
- Advani M S and Saxe A M 2017 High-dimensional dynamics of generalization error in neural networks (arXiv:1710.03667)
- Bartlett P L, Long P M, Lugosi G and Tsigler A 2020 Benign overfitting in linear regression *Proc. Natl Acad. Sci.* **117** 30063–70
- Bauschke H H *et al* 2011 *Convex Analysis and Monotone Operator Theory in Hilbert Spaces* vol 408 (Springer)
- Bauschke H H, Borwein J M and Combettes P L 2003 Bregman monotone optimization algorithms *SIAM J. Control Optim.* **42** 596–636
- Bauschke H, Combettes P and Noll D 2006 Joint minimization with alternating Bregman proximity operators *Pac. J. Optim.* **2** 3
- Bayati M and Montanari A 2011a The dynamics of message passing on dense graphs, with applications to compressed sensing *IEEE Trans. Inf. Theory* **57** 764–85
- Bayati M and Montanari A 2011b The LASSO risk for Gaussian matrices *IEEE Trans. Inf. Theory* **58** 1997–2017
- Belkin M, Hsu D, Ma S and Mandal S 2020 Reply to Loog *et al.*: looking beyond the peaking phenomenon *Proc. Natl Acad. Sci.* **117** 10627
- Berthier R, Montanari A and Nguyen P-M 2020 State evolution for approximate message passing with non-separable functions *Inf. Inference* **9** 33–79
- Bottou L 2012 Stochastic gradient descent tricks *Neural Networks: Tricks of the Trade* (Springer) pp 421–36
- Breiman L 1996 Bagging predictors *Mach. Learn.* **24** 123–40
- Celentano M, Montanari A and Wei Y 2020 The Lasso with general Gaussian designs with applications to hypothesis testing (arXiv:2007.13716)
- Chen L, Min Y, Belkin M and Karbasi A 2020 Multiple descent: design your own generalization curve (arXiv:2008.01036)
- Chizat L, Oyallon E and Bach F 2019 On lazy training in differentiable programming *Advances in Neural Information Processing Systems* vol 32, ed H Wallach, H Larochelle, A Beygelzimer, F d'Alché-Buc, E Fox and R Garnett (Curran Associates, Inc.)
- D'Ascoli S, Refinetti M, Biroli G and Krzakala F 2020 Double trouble in double descent: bias and variance(s) in the lazy regime vol 119 *Proc. 37th Int. Conf. on Machine Learning* ed H Daumé III and A Singh (PMLR) pp 2280–90
- d'Ascoli S, Sagun L and Biroli G 2021 Triple descent and the two kinds of overfitting: where and why do they appear? *J. Stat. Mech.* **124002**
- Dhifallah O and Lu Y M 2020 A precise performance analysis of learning with random features (arXiv:2008.11904)
- Donoho D L, Maleki A and Montanari A 2009 Message-passing algorithms for compressed sensing *Proc. Natl Acad. Sci.* **106** 18914–9
- Donoho D and Montanari A 2016 High dimensional robust M-estimation: asymptotic variance via approximate message passing *Probab. Theory Relat. Fields* **166** 935–69
- Drucker H, Cortes C, Jackel L D, LeCun Y and Vapnik V 1994 Boosting and other ensemble methods *Neural Comput.* **6** 1289–301
- El Karoui N 2010 The spectrum of kernel random matrices *Ann. Stat.* **38** 1–50
- Geiger M, Jacot A, Spigler S, Gabriel F, Sagun L, d'Ascoli S, Biroli G, Hongler C and Wyart M 2020 Scaling description of generalization with number of parameters in deep learning *J. Stat. Mech.* **023401**
- Gerace F, Loureiro B, Krzakala F, Mezard M and Zdeborova L 2020 Generalisation error in learning with random features and the hidden manifold model vol 1119 *Proc. 37th Int. Conf. on Machine Learning* ed H Daumé III and A Singh (PMLR) pp 3452–62
- Gerbelot C and Berthier R 2021 Graph-based approximate message passing iterations (arXiv:2109.11905)
- Goldt S, Loureiro B, Reeves G, Krzakala F, Mézard M and Zdeborová L 2021 The Gaussian equivalence of generative models for learning with shallow neural networks *Proc. 2nd Mathematical and Scientific Machine Learning Conf.* vol 145 pp 1–46
- Goldt S, Mézard M, Krzakala F and Zdeborová L 2020 Modeling the influence of data structure on learning in neural networks: the hidden manifold model *Phys. Rev. X* **10** 041044
- Hu H and Lu Y M 2020 Universality laws for high-dimensional learning with random features (arXiv:2009.07669)
- Jacot A, Gabriel F and Hongler C 2018 Neural tangent kernel: convergence and generalization in neural networks *Advances in Neural Information Processing Systems* vol 31, ed S Bengio, H Wallach, H Larochelle, K Grauman, N Cesa-Bianchi and R Garnett (Curran Associates, Inc.)
- Jacot A, Simsek B, Spadaro F, Hongler C and Gabriel F 2020 Implicit regularization of random feature models *Int. Conf. on Machine Learning* (PMLR) pp 4631–40

- Javanmard A and Montanari A 2013 State evolution for general approximate message passing algorithms, with applications to spatial coupling *Inf. Inference* **2** 115–44
- Kai Hansen L and Salamon P 1990 Neural network ensembles *IEEE Trans. Pattern Anal. Mach. Intell.* **12** 993–1001
- Krogh A and Sollich P 1997 Statistical mechanics of ensemble learning *Phys. Rev. E* **55** 811–25
- Krogh A and Vedelsby J 1995 Neural network ensembles, cross validation and active learning *Advances in Neural Information Processing Systems* vol 7, ed G Tesauro, D Touretzky and T Leen (MIT Press)
- Lakshminarayanan B, Pritzel A and Blundell C 2017 Simple and scalable predictive uncertainty estimation using deep ensembles (arXiv:1612.01474)
- LeJeune D, Javadi H and Baraniuk R 2020 The implicit regularization of ordinary least squares ensembles *Int. Conf. on Artificial Intelligence and Statistics* (PMLR) pp 3525–35
- Lin L and Dobriban E 2021 What causes the test error? Going beyond bias-variance via ANOVA *J. Mach. Learn. Res.* **22** 1–82
- Loureiro B, Gerbelot C, Cui H, Goldt S, Krzakala F, Mézard M and Zdeborová L 2021a Capturing the learning curves of generic features maps for realistic data sets with a teacher-student model (arXiv:2102.08127)
- Loureiro B, Sicuro G, Gerbelot C, Pocco A, Krzakala F and Zdeborová L 2021b Learning Gaussian mixtures with generalised linear models: precise asymptotics in high-dimensions (arXiv:2106.03791)
- Mei S and Montanari A 2021 The generalization error of random features regression: precise asymptotics and the double descent curve *Commun. Pure Appl. Math.* **75** 667–766
- Mézard M, Parisi G and Virasoro M 1987 *Spin Glass Theory and Beyond: An Introduction to the Replica Method and Its Applications* vol 9 (World Scientific Publishing Company)
- Narkhede M V, Bartakke P P and Sutaone M S 2022 A review on weight initialization strategies for neural networks *Artif. Intell. Rev.* **55** 291–322
- Neal B, Mittal S, Baratin A, Tantia V, Scicluna M, Lacoste-Julien S and Mitliagkas I 2018 A modern take on the bias-variance tradeoff in neural networks (arXiv:1810.08591)
- Opitz D and Maclin R 1999 Popular ensemble methods: an empirical study *J. Artif. Intell. Res.* **11** 169–98
- Parikh N and Boyd S 2014 Proximal algorithms *Found. Trends Optim.* **1** 127–239
- Pennington J and Worah P 2017 Nonlinear random matrix theory for deep learning *Advances in Neural Information Processing Systems* vol 30, ed I Guyon, U V Luxburg, S Bengio, H Wallach, R Fergus, S Vishwanathan and R Garnett (Curran Associates, Inc.)
- Perrone M P and Cooper L N 1993 When networks disagree: ensemble methods for hybrid neural networks *How We Learn; How We Remember: Toward an Understanding of Brain and Neural Systems* (Chapman and Hall) pp 126–42
- Perrone M 1994 Putting it all together: methods for combining neural networks *Advances in Neural Information Processing Systems* vol 6, ed J Cowan, G Tesauro and J Alspector (Morgan-Kaufmann)
- Rahimi A and Recht B 2007 Random features for large-scale kernel machines *Advances in Neural Information Processing Systems (NIPS)* pp 1177–84
- Rangan S 2011 Generalized approximate message passing for estimation with random linear mixing *2011 IEEE Int. Symp. on Information Theory Proc.* (IEEE) pp 2168–72
- Rosset S, Zhu J and Hastie T 2004 Margin maximizing loss functions *Advances in Neural Information Processing Systems* vol 16, ed S Thrun, L Saul and B Schölkopf (MIT Press)
- Schapire R E 1990 The strength of weak learnability *Mach. Learn.* **5** 197–227
- Schwarze H and Hertz J 1992 Generalization in a large committee machine *Europhys. Lett.* **20** 375–80
- Spigler S, Geiger M, d’Ascoli S, Sagun L, Biroli G and Wyart M 2019 A jamming transition from under- to over-parametrization affects generalization in deep learning *J. Phys. A: Math. Theor.* **52** 474001
- Thrapoulidis C, Abbasi E and Hassibi B 2018 Precise error analysis of regularized  $M$ -estimators in high dimensions *IEEE Trans. Inf. Theory* **64** 5592–628
- Vershynin R 2010 Introduction to the non-asymptotic analysis of random matrices (arXiv:1011.3027)
- Vershynin R 2018 *High-Dimensional Probability: An Introduction with Applications in Data Science* vol 47 (Cambridge University Press)
- Zdeborová L and Krzakala F 2016 Statistical physics of inference: thresholds and algorithms *Adv. Phys.* **65** 453–552