Generalization error in high-dimensional perceptrons: Approaching Bayes error with convex optimization

Benjamin Aubin[†], Florent Krzakala^{*}, Yue M. Lu[°], Lenka Zdeborová[†]

[†] Université Paris-Saclay, CNRS, CEA, Institut de physique théorique, 91191, Gif-sur-Yvette, France.

* Laboratoire de Physique Statistique, CNRS & Sorbonnes Universités, École Normale Supérieure, PSL University, Paris, France.

° John A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge, MA 02138, USA

Abstract

We consider a commonly studied supervised classification of a synthetic dataset whose labels are generated by feeding a one-layer neural network with random iid inputs. We study the generalization performances of standard classifiers in the high-dimensional regime where $\alpha = n/d$ is kept finite in the limit of a high dimension d and number of samples n. Our contribution is three-fold: First, we prove a formula for the generalization error achieved by ℓ_2 regularized classifiers that minimize a convex loss. This formula was first obtained by the heuristic replica method of statistical physics. Secondly, focussing on commonly used loss functions and optimizing the ℓ_2 regularization strength, we observe that while ridge regression performance is poor, logistic and hinge regression are surprisingly able to approach the Bayes-optimal generalization error extremely closely. As $\alpha \to \infty$ they lead to Bayes-optimal rates, a fact that does not follow from predictions of margin-based generalization error bounds. Third, we design an optimal loss and regularizer that provably leads to Bayes-optimal generalization error.

1 Introduction

High-dimensional statistics, where the ratio $\alpha = n/d$ is kept finite while the dimensionality d and the number of samples n grow, often display interesting non-intuitive features. Asymptotic generalization performances for such problems in the so-called *teacher-student* setting, with synthetic data, have been the subject of intense investigations spanning many decades [1–6]. To understand the effectiveness of modern machine learning techniques, and also the limitations of the classical statistical learning approaches [7,8], it is of interest to revisit this line of research. Indeed, this direction is currently the subject to a renewal of interests, as testified by some very recent, yet already rather influential papers [9–13]. The present paper subscribes to this line of work and studies high-dimensional classification within one of the simplest models considered in statistics and machine learning: convex linear estimation with data generated by a teacher *perceptron* [14]. We will focus on the generalization abilities in this problem, and compare the performances of Bayes-optimal estimation to the more standard *empirical risk minimization*. We then compare the results with the prediction of standard generalization bounds that illustrate in particular their limitation even in this simple, yet non-trivial, setting.

Synthetic data model — We consider a supervised machine learning task, whose dataset is generated by a single layer neural network, often named a *teacher* [1–3], that belongs to the Generalized Linear Model (GLM) class. Therefore we assume the *n* samples are drawn according to

$$\mathbf{y} = \varphi_{\text{out}}^{\star} \left(\frac{1}{\sqrt{d}} \mathbf{X} \mathbf{w}^{\star} \right) \Leftrightarrow \mathbf{y} \sim P_{\text{out}}^{\star} \left(. \right) \,, \tag{1}$$

where $\mathbf{w}^{\star} \in \mathbb{R}^d$ denotes the ground truth vector drawn from a probability distribution $P_{\mathbf{w}^{\star}}$ with second moment $\rho_{\mathbf{w}^{\star}} \equiv \frac{1}{d}\mathbb{E}\left[\|\mathbf{w}^{\star}\|_2^2\right]$ and $\varphi_{\text{out}}^{\star}$ represents a deterministic or stochastic activation function equivalently associated to a distribution P_{out}^{\star} . The input data matrix $\mathbf{X} = (\mathbf{x}_{\mu})_{\mu=1}^n \in \mathbb{R}^{n \times d}$ contains iid Gaussian vectors, i.e $\forall \mu \in [1:n], \mathbf{x}_{\mu} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. Even though the framework we use and the theorems and results we derived are valid for a rather generic channel in eq. (1) –including regression problems— we will mainly focus the presentation on the commonly considered perceptron case: a binary classification task with data given by a sign activation function $\varphi_{\text{out}}^{\star}(\mathbf{z}) = \text{sign}(\mathbf{z})$, with a Gaussian weight distribution $P_{\mathbf{w}^{\star}}(\mathbf{w}^{\star}) =$ $\mathcal{N}_{\mathbf{w}^{\star}}(\mathbf{0}, \rho_{\mathbf{w}^{\star}}\mathbf{I}_d)$. The ± 1 labels are thus generated as

$$\mathbf{y} = \operatorname{sign}\left(\frac{1}{\sqrt{d}} \mathbf{X} \mathbf{w}^{\star}\right), \quad \text{with} \quad \mathbf{w}^{\star} \sim \mathcal{N}_{\mathbf{w}^{\star}}\left(\mathbf{0}, \rho_{\mathbf{w}^{\star}} \mathbf{I}_{d}\right).$$
(2)

Empirical Risk Minimization — The workhorse of machine learning is Empirical Risk Minimization (ERM), where one minimizes a *loss function* in the corresponding high-dimensional parameter space \mathbb{R}^d . To avoid overfitting of the training set one often adds a *regularization term* r. ERM then corresponds to estimating $\hat{\mathbf{w}}_{erm} = \operatorname{argmin}_{\mathbf{w}} [\mathcal{L}(\mathbf{w}; \mathbf{y}, \mathbf{X})]$ where the regularized

training loss \mathcal{L} is defined by, using the notation $z_{\mu}(\mathbf{w}, \mathbf{x}_{\mu}) \equiv \frac{1}{\sqrt{d}} \mathbf{x}_{\mu}^{\mathsf{T}} \mathbf{w}$,

$$\mathcal{L}(\mathbf{w}; \mathbf{y}, \mathbf{X}) = \sum_{\mu=1}^{n} l(y_{\mu}, z_{\mu}(\mathbf{w}, \mathbf{x}_{\mu})) + r(\mathbf{w}) .$$
(3)

The goal of the present paper is to discuss the generalization performance of these estimators for the classification task (2) in the high-dimensional limit. We focus our analysis on commonly used loss functions l, namely the square $l^{\text{square}}(y, z) = \frac{1}{2}(y - z)^2$, logistic $l^{\text{logistic}}(y, z) = \log(1 + \exp(-yz))$ and hinge losses $l^{\text{hinge}}(y, z) = \max(0, 1 - yz)$. We will mainly illustrate our results for the ℓ_2 regularization $r(\mathbf{w}) = \lambda ||\mathbf{w}||_2^2/2$, where we introduced a regularization strength hyper-parameter λ .

Related works — The above learning problem has been extensively studied in the statistical physics community using the heuristic replica method [1-3, 14, 15]. Due to the interest in high-dimensional statistics, they have experienced a resurgence in popularity in recent years. In particular, rigorous works on related problems are much more recent. The authors of [10] established rigorously the replica-theory predictions for the Bayes-optimal generalization error. Here we focus on standard ERM estimation and compare it to the results obtained in [10]. Authors of [16] analyzed rigorously M-estimators for the regression case where data are generated by a linear-activation teacher. Here we analyze classification with a more general and non-linear teacher, focusing in particular on the sign-teacher. The case of max-margin loss was studied in [17] with a technically closely related proof, but with a focus on the over-parametrized regime, thus not addressing the questions that we focus on. A range of unregularized losses was also analyzed for a sigmoid teacher (that is very similar to a sign-teacher) again in the context of the double-descent behaviour in [18, 19]. Here we focus instead on the regularized case as it drastically improves generalization performances of the ERM and that allows us to compare with the Bayes-optimal estimation as well as to standard generalization bounds. Our proof, as in the above mentioned works and [20], is based on Gordon's minimax formalism, including in particular the effect of the regularization.

Main contributions — Our first main contribution is to provide rigorously, in Sec. 2, the classification generalization performances of empirical risk minimization with the loss given by (3) in the high-dimensional limit, for any convex loss and an ℓ_2 regularization. Note that the proof is easily extended to any convex separable regularization. Additionally, we provide a proof of the equivalence between the results of our paper and the ones initially obtained by the replica method, which is of additional interest given the wide range of application of these heuristics statistical-physics technics in machine learning and computer science [21, 22]. In particular, the replica predictions in [15, 23–25] follow from our results. Another approach that originated in physics are the so-called TAP equations [26–28] that lead to the so-called Approximate Message Passing algorithm for solving linear and generalized linear problems with Gaussian matrices [29, 30]. This algorithm can be analyzed with the so-called *state evolution* method [31], and it is widely believed (and in fact proven for linear problems [4, 32]) that the fixed-point of the state evolution gives the optimal error in high-dimensional convex optimization problems.

The state evolution equations are in fact equivalent to the one given by the replica theory and therefore our results vindicate this approach as well. We also demonstrate numerically that these asymptotic results are very accurate even for moderate system sizes, and they have been performed with the scikit-learn library [33].

Secondly, and more importantly, we provide in Sec. 3 a detailed analysis of the generalization error for standard losses such as square, hinge (or equivalently support vector machine) and logistic, as a function of the regularization strength λ and the number of samples per dimension α . We observe, in particular, that while the ridge regression never closely approaches the Bayes-optimal performance, the logistic regression with optimized ℓ_2 regularization gets extremely close to optimal. And so does, to a lesser extent, the hinge regression and the max-margin estimator to which the unregularized logistic and hinge converge [34]. It is quite remarkable that these canonical losses are able to approach the error of the Bayes-optimal estimator for which, in principle, the marginals of a high-dimensional probability distribution need to be evaluated. Notably, all the later losses give —for a *good choice* of the regularization strength λ —generalization errors scaling as $\Theta(\alpha^{-1})$ for large α , just as the Bayes-optimal generalization error [10]. This is found to be at variance with the prediction of Rademacher and max-marginbased bounds that predict instead a $\Theta(\alpha^{-1/2})$ rate [35, 36], which therefore appear to be vacuous in the high-dimensional regime.

Third, in Sec. 4, we design a custom (non-convex) loss and regularizer that provably gives a plug-in estimator that efficiently achieves Bayes-optimal performances, including the optimal $\Theta(\alpha^{-1})$ rate for the generalization error. Our construction is related to the one discussed in [37–39], but is not restricted to convex losses.

2 Main technical results

In the formulas that arise for this statistical estimation problem, the correlations between the estimator $\hat{\mathbf{w}}$ and the ground truth vector \mathbf{w}^* play a fundamental role and we thus define two scalar overlap parameters to measure the statistical reconstruction:

$$m \equiv \frac{1}{d} \mathbb{E}_{\mathbf{y}, \mathbf{X}} \hat{\mathbf{w}}^{\mathsf{T}} \mathbf{w}^{\star}, \qquad \qquad q \equiv \frac{1}{d} \mathbb{E}_{\mathbf{y}, \mathbf{X}} \left[\| \hat{\mathbf{w}} \|_2 \right]^2.$$
(4)

In particular, the generalization error of the estimator $\hat{\mathbf{w}}(\alpha) \in \mathbb{R}^d$ obtained by performing Empirical Risk Minimization (ERM) on the training loss \mathcal{L} in eq. (3) with $n = \alpha d$ samples

$$e_{g}^{\text{erm}}(\alpha) \equiv \mathbb{E}_{y,\mathbf{x}} \mathbb{1}\left[y \neq \hat{y}\left(\hat{\mathbf{w}}(\alpha); \mathbf{x} \right) \right],$$
(5)

where $\hat{y}(\hat{\mathbf{w}}(\alpha); \mathbf{x})$ denotes the predicted label, has both at finite d and in the asymptotic limit an explicit expression depending only on the above overlaps m and q:

Proposition 2.1 (Generalization error of classification). *In our synthetic binary classification task, the generalization error of ERM (or equivalently the test error) is given by*

$$e_{g}^{\text{erm}}\left(\alpha\right) = \frac{1}{\pi}a\cos\left(\sqrt{\eta}\right), \text{ with } \eta \equiv \frac{m^{2}}{\rho_{w^{\star}}q} \text{ and } \rho_{w^{\star}} \equiv \frac{1}{d}\mathbb{E}\left[\|\boldsymbol{w}^{\star}\|_{2}^{2}\right].$$
(6)

Proof. The proof, shown in SM. II, is a simple computation based on a Gaussian integration. \Box

To obtain the generalization performances, it thus remains to obtain the asymptotic values of m, q (and thus of η), in the limit $d \to \infty$. With the ℓ_2 regularization, these values are characterized by a set of fixed point equations given by the next theorems ¹:

For any $\tau > 0$, let us first recall the definitions of the Moreau-Yosida regularization and the proximal operator of a convex loss function $(y, z) \mapsto \ell(y \cdot z)$:

$$\mathcal{M}_{\tau}(z) = \min_{x} \left\{ \ell(x) + \frac{(x-z)^2}{2\tau} \right\}, \qquad \mathcal{P}_{\tau}(z) = \operatorname{argmin}_{x} \left\{ \ell(x) + \frac{(x-z)^2}{2\tau} \right\}.$$
(7)

Theorem 2.2 (Gordon's min-max fixed point - Binary classification with ℓ_2 regularization). As $n, d \to \infty$ with $n/d = \alpha = \Theta(1)$, the overlap parameters m, q concentrate to

$$m \xrightarrow[d \to \infty]{} \sqrt{\rho_{\mathbf{w}^*}} \mu^*, \qquad \qquad q \xrightarrow[d \to \infty]{} (\mu^*)^2 + (\delta^*)^2, \qquad (8)$$

where parameters μ^* and δ^* are solutions of

$$(\mu^*, \delta^*) = \underset{\mu, \delta \ge 0}{\operatorname{arg\,min}} \sup_{\tau > 0} \left\{ \frac{\lambda(\mu^2 + \delta^2)}{2} - \frac{\delta^2}{2\tau} + \alpha \mathbb{E}_{g,s} \mathcal{M}_{\tau}[\delta g + \mu s \varphi_{\text{out}^*}(\sqrt{\rho_{w^*}}s)] \right\}, \quad (9)$$

and g, s are two iid standard normal random variables. The solutions (μ^*, δ^*) of (9) can be reformulated as a set of fixed point equations

$$\mu^{*} = \frac{\alpha}{\lambda\tau^{*} + \alpha} \mathbb{E}_{g,s} [s \cdot \varphi_{\text{out}^{*}}(\sqrt{\rho_{\text{w}^{*}}}s) \cdot \mathcal{P}_{\tau^{*}}(\delta^{*}g + \mu^{*}s\varphi_{\text{out}^{*}}(\sqrt{\rho_{\text{w}^{*}}}s))],$$

$$\delta^{*} = \frac{\alpha}{\lambda\tau^{*} + \alpha - 1} \mathbb{E}_{g,s} [g \cdot \mathcal{P}_{\tau^{*}}(\delta^{*}g + \mu^{*}s\varphi_{\text{out}^{*}}(\sqrt{\rho_{\text{w}^{*}}}s))],$$

$$(\delta^{*})^{2} = \alpha \mathbb{E}_{g,s} [((\delta^{*}g + \mu^{*}s\varphi_{\text{out}^{*}}(\sqrt{\rho_{\text{w}^{*}}}s)) - \mathcal{P}_{\tau^{*}}(\delta^{*}g + \mu^{*}s\varphi_{\text{out}^{*}}(\sqrt{\rho_{\text{w}^{*}}}s)))^{2}].$$
(10)

Proof. The proof, shown in SM. III.1, is an application of the Gordon minimax theory. \Box

This set of fixed point equations can be finally mapped to the ones obtained by the heuristic *replica* method from statistical physics (whose heuristic derivation is shown in SM. IV) as well as the state evolution of the approximate-message-passing algorithm [27, 30, 40]. Thus their validity for this convex estimation problem is rigorously established by the following theorem:

Corollary 2.3 (Equivalence Gordon-replicas). As $n, d \to \infty$ with $n/d = \alpha = \Theta(1)$, the overlap parameters m, q concentrate to the fixed point of the following set of equations:

$$m = \alpha \Sigma \rho_{\mathbf{w}^{\star}} \cdot \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^{\star}} \times f_{\text{out}^{\star}} \left(y, \sqrt{\rho_{\mathbf{w}^{\star}} \eta} \xi, \rho_{\mathbf{w}^{\star}} \left(1 - \eta \right) \right) \cdot f_{\text{out}} \left(y, q^{1/2} \xi, \Sigma \right) \right],$$

$$q = m^2 / \rho_{\mathbf{w}^{\star}} + \alpha \Sigma^2 \cdot \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^{\star}} \left(y, \sqrt{\rho_{\mathbf{w}^{\star}} \eta} \xi, \rho_{\mathbf{w}^{\star}} \left(1 - \eta \right) \right) \cdot f_{\text{out}} \left(y, q^{1/2} \xi, \Sigma \right)^2 \right], \quad (11)$$

$$\Sigma = \left(\lambda - \alpha \cdot \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^{\star}} \left(y, \sqrt{\rho_{\mathbf{w}^{\star}} \eta} \xi, \rho_{\mathbf{w}^{\star}} \left(1 - \eta \right) \right) \cdot \partial_{\omega} f_{\text{out}} \left(y, q^{1/2} \xi, \Sigma \right) \right] \right)^{-1},$$

¹Note that for a generic convex and non-separable regularizer (different than ℓ_2), it would contain instead six equations (see SM. III.2).

with
$$\eta \equiv \frac{m^2}{\rho_{w^*}q}$$
, $f_{out}(y,\omega,V) \equiv V^{-1}(\mathcal{P}_V[l(y,.)](\omega) - \omega)$,
 $\mathcal{Z}_{out^*}(y,\omega,V) = \mathbb{E}_z \left[P_{out^*}\left(y|\sqrt{V}z + \omega \right) \right]$, $f_{out^*}(y,\omega,V) \equiv \partial_\omega \log \left(\mathcal{Z}_{out^*}(y,\omega,V) \right)$, (12)

where ξ , z denote two iid standard normal random variables, and \mathbb{E}_y the continuous or discrete sum over all possible values y according to P_{out^*} .

Proof. For the sake of clarity, the proof is again left in SM. III.3.

Equivalent equations for the whole GLM class (classification and regression) with any separable and convex regularizer are shown in SM. III.2.

Bayes optimal baseline — Finally, we shall compare the ERM performances to the Bayesoptimal generalization error. Being the information-theoretically best possible estimator, we will use it as a reference baseline for comparison. The expression of the Bayes-optimal generalization was derived in [24] and proven in [10] and we recall here the result:

Theorem 2.4 (Bayes Asymptotic performance, from [10]). For the model (1) with $P_{w^*}(w^*) = \mathcal{N}_{w^*}(\mathbf{0}, \rho_{w^*}\mathbf{I}_d)$, the Bayes-optimal generalization error is quantified by two scalar parameters q_b and \hat{q}_b that verify the set of fixed point equations

$$q_{\rm b} = \frac{\hat{q}_{\rm b}}{1 + \hat{q}_{\rm b}}, \quad \hat{q}_{\rm b} = \alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\rm out^{\star}} \left(y, q_{\rm b}^{1/2} \xi, \rho_{\rm w^{\star}} - q_{\rm b} \right) \cdot f_{\rm out^{\star}} \left(y, q_{\rm b}^{1/2} \xi, \rho_{\rm w^{\star}} - q_{\rm b} \right)^2 \right], \quad (13)$$

and reads

$$e_{\rm g}^{\rm bayes}\left(\alpha\right) = \frac{1}{\pi}a\cos\left(\sqrt{\eta_{\rm b}}\right) \text{ with } \eta_{\rm b} = \frac{q_{\rm b}}{\rho_{\rm w}\star}.$$
 (14)

 \square

3 Generalization errors

We now move to the core of the paper and analyze the set of fixed point equations (10), or equivalently (11), leading to the generalization performances given by (6), for common classifiers on our synthetic binary classification task. As already stressed, even though the results are valid for a wide range of regularizers, we focus on estimators based on ERM with ℓ_2 regularization $r(\mathbf{w}) = \lambda ||\mathbf{w}||_2^2/2$, and with square loss (ridge regression) $l^{\text{square}}(y, z) = \frac{1}{2}(y - z)^2$, logistic loss (logistic regression) $l^{\text{logistic}}(y, z) = \log(1 + \exp(-yz))$ or hinge loss (SVM) $l^{\text{hinge}}(y, z) = \max(0, 1 - yz)$. In particular, we study the influence of the hyper-parameter λ on the generalization performances and the different large α behaviour generalization rates in the high-dimensional regime, and compare with the Bayes results. We show the solutions of the set of fixed point equations eqs. (11) in Figs. 1a, 1b, 1c respectively for ridge, hinge and logistic ℓ_2 regressions. Ridge regression is a special case, for which its quadratic loss allows to derive and fully solve the equations (see SM. V.3). However in general the set of equations has no analytical closed form and needs therefore to be solved numerically. It is in particular the case for logistic and hinge, whose Moreau-Yosida regularization eq. (12) is, however, analytical.

First, to highlight the accuracy of the theoretical predictions, we compare in Figs. 1a-1c the ERM asymptotic $(d \to \infty)$ generalization error with the performances of numerical simulations $(d = 10^3, \text{ averaged over } n_s = 20 \text{ samples})$ of ERM of the training loss eq. (3). Presented for a wide range of number of samples α and of regularization strength λ , we observe a perfect match between theoretical predictions and numerical simulations so that the error bars are barely visible and have been therefore removed. This shows that the asymptotic predictions are valid even with very moderate sizes. As an information theoretical baseline, we also show the Bayes-optimal performances (black) given by the solution of eq. (13).

Ridge estimation— As we might expect the square loss gives the worst performances. For low values of the generalization, it leads to an interpolation-peak at $\alpha = 1$. The limit of vanishing regularization $\lambda \to 0$ leads to the *least-norm* or *pseudo-inverse* estimator $\hat{\mathbf{w}}_{\text{pseudo}} = (X^TX)^{-1}X^T\mathbf{y}$. The corresponding generalization error presents the largest interpolation-peak and achieves a maximal generalization error $e_g = 0.5$. These are well known observations, discussed as early as in [23, 25], that are object of a renewal of interest under the name *double descent*, following a recent series of papers [11, 41–47]. This double descent behaviour for the pseudo-inverse is shown in Fig. 1a with a yellow line. On the contrary, larger regularization strengths do not suffer this peak at $\alpha = 1$, but their generalization error performance is significantly worse than the Bayes-optimal baseline for larger values of α . Indeed, as we might expect, for a large number of samples, a large regularization biases wrongly the training. However, even with optimized regularizations, performances of the ridge estimator remains far away from the Bayes-optimal performance.

Hinge and logistic estimation— Both these losses, which are the classical ones used in classification problems, improve drastically the generalization error. First of all, let us notice that they do not display a double-descent behavior. This is due to the fact that our results are illustrated in the noiseless case and that our synthetic dataset is always linearly separable. Optimizing the regularization, our results in Fig. 1b-1c show both hinge and logistic ERM-based classification approach very closely the Bayes error. To offset these results, note that performances of logistic regression on non-linearly separable data are however very poor, as illustrated by our analysis of a *rectangle door* teacher (see SM. V.6).

Max-margin estimation— As discussed in [34], both the logistic and hinge estimator converge, for vanishing regularization $\lambda \rightarrow 0$, to the *max-margin* solution. Taking the $\lambda \rightarrow 0$ limit in our equations, we thus obtain the *max-margin* estimator performances. While this is not what gives the best generalization error (as can be seen in Fig. 1c the logistic with an optimized λ has a lower error), the max-margin estimator gives very good results, and gets very close to the Bayes-error.

Optimal regularization— Defining the regularization value that optimizes the generalization as

$$\lambda^{\text{opt}}\left(\alpha\right) = \operatorname{argmin}_{\lambda} e_{g}^{\text{erm}}\left(\alpha,\lambda\right)\,,\tag{15}$$

we show in Figs. 1b-1c that both optimal values $\lambda^{\text{opt}}(\alpha)$ (dashed-dotted orange) for logistic and hinge regression decrease to 0 as α grows and more data are given. Somehow surprisingly, we observe in particular that the generalization performances of logistic regression with optimal regularization are *extremely close* to the Bayes performances. The difference with the optimized logistic generalization error is barely visible by eye, so that we explicitly plotted the difference, which is roughly of order 10^{-5} .

Ridge regression Fig. 1a shows a singular behaviour: there exists an optimal value (purple) which is moreover independent of α achieved for $\lambda^{\text{opt}} \simeq 0.5708$. This value was first found numerically and confirmed afterwards semi-analytically in SM. V.3.

Generalization rates at large α — Finally, we turn to the very instructive behavior at large values of α when a large amount of data is available. First, we notice that the Bayes-optimal generalization error, whose large α analysis is performed in SM. V.1, decreases as $e_{g}^{\text{bayes}} \approx 0.4417 \alpha^{-1}$. Compared to this optimal value, ridge regression gives poor performances in this regime. For any value of the regularization λ — and in particular for both the pseudo-inverse case at $\lambda = 0$ and the optimal estimator λ^{opt} — its generalization performances decrease much slower than the Bayes rate, and goes only as $e_{g}^{\text{ridge}} \approx 0.2405 \alpha^{-1/2}$ (see SM. V.3 for the derivation). Hinge and logistic regressions present a radically different, and more favorable, behaviour. Fig. 1b-1c show that keeping λ finite when α goes to ∞ , does not yield the Bayes-optimal rates. However the max-margin solution (that corresponds to the $\lambda \rightarrow 0$ limit of these estimators) gives extremely good performances $e_{g}^{\text{logistic,hinge}} \approx e_{g}^{\text{max-margin}} \approx 0.500 \alpha^{-1}$ see derivation in SM. V.4). This is the same rate as the Bayes one, only that the constant is slightly higher.

Comparison with VC and Rademacher statistical bounds— Given the fact that both the max-margin estimator and the optimized logistic achieve optimal generalization rates going as $\Theta(\alpha^{-1})$, it is of interest to compare those rates to the prediction of statistical learning theory bounds. Statistical learning analysis (see e.g. [35, 36, 48]) relies to a large extent on the Vapnik-*Chervonenkis* dimension (VC) analysis and on the so-called *Rademacher complexity*. The uniform convergence result states that if the Rademacher complexity or the Vapnik-Chervonenkis dimension $d_{\rm VC}$ is finite, then for a large enough number of samples the generalization gap will vanish uniformly over all possible values of parameters. Informally, uniform convergence tells us that with high probability, for any value of the weights w, the generalization gap satisfies $\mathcal{R}_{\text{population}}(\mathbf{w}) - \mathcal{R}_{\text{empirical}}^{n}(\mathbf{w}) = \Theta\left(\sqrt{d_{\text{VC}}/n}\right)$ where $d_{\text{VC}} = d - 1$ for our GLM hypothesis class. Therefore, given that the empirical risk can go to zero (since our data are separable), this provides a generalization error upper-bound $e_g \leq \Theta(\alpha^{-1/2})$. This is much worse that what we observe in practice, where we reach the Bayes rate $e_g = \Theta(\alpha^{-1})$. Tighter bounds can be obtained using the Rademacher complexity, and this was studied recently (using the aforementioned replica method) in [49] for the very same problem. We reproduced their results and plotted the Rademacher complexity generalization bound in Fig.1 (dashed-green) that decreases as $\Theta(\alpha^{-1/2})$ for the binary classification task eq. (2).

One may wonder if this could be somehow improved. Another statistical-physics heuristic



(a) Ridge regression: square loss with ℓ_2 regularization. Interpolation-peak, at $\alpha = 1$, is maximal for the pseudo-inverse estimator $\lambda = 0$ (yellow line) that reaches $e_{\rm g} = 0.5$.



(b) Hinge regression: hinge loss with ℓ_2 regularization. For clarity the rescaled value of $\lambda^{\rm opt}/10$ (dotted-dashed orange) is shown as well as its generalization error $e_{\rm g}^{\rm opt}$ (dotted orange) that is slightly below and almost indistinguishable of the max-margin performances (dashed black).



(c) Logistic regression: logistic loss with ℓ_2 regularization - The value of λ^{opt} (dotted-dashed orange) is shown as well as its generalization error $e_{\text{g}}^{\text{opt}}$ (dotted orange). Visually indistinguishable from the Bayes-optimal line, their difference $e_{\text{g}}^{\text{opt}} - e_{\text{g}}^{\text{bayes}}$ is shown as an inset (dashed orange).

Figure 1: Asymptotic generalization error for ℓ_2 regularization $(d \to \infty)$ as a function of α for different regularizations strengths λ , compared to numerical simulation (points) of ridge regression for $d = 10^3$ and averaged over $n_s = 20$ samples. Numerics has been performed with the default methods *Ridge*, *LinearSVC*, *LogisticRegression* of scikit-learn package [33]. Bayes optimal performances are shown with a black line and goes as $\Theta(\alpha^{-1})$, while the Rademacher complexity (dashed green) decrease as $\Theta(\alpha^{-1/2})$. Both hinge and logistic converge to maxmargin estimator (limit $\lambda = 0$) which is shown in dashed black and deceases as $\Theta(\alpha^{-1})$, while Ridge decreases as $\Theta(\alpha^{-1/2})$.

computation, however, suggests that, unfortunately, uniform bound are plagued to a slow rate $\Theta(\alpha^{-1/2})$. Indeed, the authors of [50] showed with a replica method-style computation that *there exists* some set of weights, in the binary classification task. (2), that lead to $\Theta(\alpha^{-1/2})$ rates: the uniform bound is thus tight. The gap observed between the uniform bound and the almost Bayes-optimal results observed in practice in this case is therefore not a paradox, but an illustration that the price to pay for uniform convergence is the inability to describe the optimal rates one can sometimes get in practice. Therefore, we believe, that the fact this phenomena can be observed in a such simple problem sheds an interesting light on the current debate in understanding generalization in deep learning [7].

Remarking our synthetic dataset is linearly-separable, we may try to take this fact into consideration to improve the generalization rate. In particular, it can be done using the maxmargin based generalization error for separable data:

Theorem 3.1 (Hard-margin generalization bound [35, 36, 48]). Given a $S = {x_1, \dots, x_n}$ such that $\forall \mu \in [1:n], ||x_{\mu}|| \leq r$. Let \hat{w} the hard-margin SVM estimator on S drawn with distribution D. With probability $1 - \delta$, the generalization error is bounded by

$$e_{g}(\alpha) \leq_{\alpha \to \infty} \left(4r \| \hat{\boldsymbol{w}} \| + \sqrt{\log\left(4/\delta\right) \log_{2} \| \hat{\boldsymbol{w}} \|} \right) / \sqrt{n} \,. \tag{16}$$

In our case one has $r^2 \simeq \frac{1}{d} \mathbb{E}_{\mathbf{x}} \|\mathbf{x}\|_2^2 = \frac{1}{d} \sum_{i=1}^d \mathbb{E} x_i^2 = 1$. On the other hand, in the large size limit, the norm of the estimator $\|\hat{\mathbf{w}}\|_2/\sqrt{d} \simeq \sqrt{q}$, that yields $e_{\mathrm{g}}(\alpha) \le 4\sqrt{\frac{q}{\alpha}}$. We now need to plug the values of the norm q obtained by our max-margin solution to finally obtain the results. Unfortunately, this bound turns out to be even worse than the previous one. Indeed the norm of the hard margin estimator q is found to grow with α in the solution of the fixed point equation, and therefore the margin decay rather fast, rendering the bound vacuous. For small values of α , one finds that $q \sim \alpha$ that provides a vacuous constant generalisation bound $e_{\mathrm{g}} \le \Theta(1)$, while for large α , $q \sim \alpha^2$ that yields an even worse bound $e_{\mathrm{g}} \le \Theta(\sqrt{\alpha})$. Clearly, max-margin based bounds do not perform well in this high-dimensional example.

4 Reaching Bayes optimality

Given the fact that logistic and hinge losses reach values extremely close to Bayes optimal generalization performances, one may wonder if by somehow slightly altering these losses one could actually reach the Bayesian values with a plug-in estimator obtained by ERM. This is what we achieve in this section, by constructing a (non-convex) optimization problem with a specially tuned loss and regularization, whose solution yields Bayes-optimal generalization. Recent insights have shown that indeed one can sometime re-interpret Bayesian estimation as an optimization program in inverse problems [37, 38, 51, 52]. In particular, [39] showed explicitly, on the basis of the non-rigorous replica method of statistical mechanics, that some Bayes-optimal reconstruction problems could be turned into convex M-estimation.

Matching ERM and Bayes-optimal generalization errors eqs. (6)-(14) with overlaps respectively solutions of eq. (11)-(13) and assuming that $\mathcal{Z}_{w^*}(\gamma, \Lambda) \equiv \mathbb{E}_{w \sim P_{w^*}} \exp(-1/2\Lambda w^2 + \gamma w)$ and $\mathcal{Z}_{\text{out}^{\star}}(y, \omega, V)$ are log-concave in γ and ω , we define the optimal loss and regularizer l^{opt} , r^{opt} :

$$l^{\text{opt}}(y,z) = -\min_{\omega} \left(\frac{(z-\omega)^2}{2(\rho_{w^{\star}}-q_{b})} + \log \mathcal{Z}_{\text{out}^{\star}}(y,\omega,\rho_{w^{\star}}-q_{b}) \right),$$

$$r^{\text{opt}}(w) = -\min_{\gamma} \left(\frac{1}{2}\hat{q}_{b}w^{2} - \gamma w + \log \mathcal{Z}_{w^{\star}}(\gamma,\hat{q}_{b}) \right), \text{ with } (q_{b},\hat{q}_{b}) \text{ solution of eq. (13).}$$
(17)

See SM. VI for the derivation. Following these considerations, we provide the following theorem:

Theorem 4.1. The result of empirical risk minimization eq. (3) with l^{opt} and r^{opt} in eq. (17), leads to Bayes optimal generalization error in the high-dimensional regime.

Proof. We present only the sketch of the proof here. First we note that the so called Bayesoptimal Approximate Message Passing (AMP) algorithm [30] is provably convergent, and indeed reaches Bayes-optimal performances (see [53]). Second, we remark that an AMP algorithm for the minimization of the ERM with loss and regularization given by (17) is exactly identical to the Bayes-optimal AMP. This shows that AMP applied to the ERM problem corresponding to (17) both converge to its fixed point and reach Bayes-optimal performances. The theorem finally follows by noting (see [32,54]) that the AMP fixed point corresponds to the extremization conditions of the loss.



Figure 2: Optimal loss $l^{\text{opt}}(y = 1, z)$ and regularizer $r^{\text{opt}}(w)$ for model eq. (2) as a function of α .

The optimal loss and regularizer λ^{opt} and r^{opt} for the model (2) are illustrated in Fig. (2). And numerical evidences of ERM with (17) compared to ℓ_2 logistic regression and Bayes performances are presented in SM. VI.

Acknowledgments

This work is supported by the ERC under the European Unions Horizon 2020 Research and Innovation Program 714608-SMiLe, by the French Agence Nationale de la Recherche under grant ANR-17-CE23-0023-01 PAIL and ANR-19-P3IA-0001 PRAIRIE, and by the US National Science Foundation under grants CCF-1718698 and CCF-1910410. We also acknowledge support from the chaire CFM-ENS "Science des donnees". Part of this work was done when Yue M. Lu was visiting Ecole Normale as a CFM-ENS "Laplace" invited researcher.

References

- [1] Hyunjune Sebastian Seung, Haim Sompolinsky, and Naftali Tishby. Statistical mechanics of learning from examples. *Physical review A*, 45(8):6056, 1992.
- [2] Timothy LH Watkin, Albrecht Rau, and Michael Biehl. The statistical mechanics of learning a rule. *Reviews of Modern Physics*, 65(2):499, 1993.
- [3] Andreas Engel and Christian Van den Broeck. *Statistical mechanics of learning*. Cambridge University Press, 2001.
- [4] Mohsen Bayati and Andrea Montanari. The lasso risk for gaussian matrices. *IEEE Transactions on Information Theory*, 58(4):1997–2017, 2011.
- [5] Noureddine El Karoui, Derek Bean, Peter J Bickel, Chinghway Lim, and Bin Yu. On robust regression with high-dimensional predictors. *Proceedings of the National Academy of Sciences*, 110(36):14557–14562, 2013.
- [6] David Donoho and Andrea Montanari. High dimensional robust m-estimation: Asymptotic variance via approximate message passing. *Probability Theory and Related Fields*, 166(3-4):935–969, 2016.
- [7] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- [8] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machinelearning practice and the classical bias-variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- [9] Emmanuel J. Candes and Pragya Sur. The phase transition for the existence of the maximum likelihood estimate in high-dimensional logistic regression, 2018.
- [10] Jean Barbier, Florent Krzakala, Nicolas Macris, Léo Miolane, and Lenka Zdeborová. Optimal errors and phase transitions in high-dimensional generalized linear models. *Proceedings of the National Academy of Sciences*, 116(12):5451–5460, 2019.

- [11] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J. Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation, 2019.
- [12] Mikhail Belkin, Daniel Hsu, and Ji Xu. Two models of double descent for weak features. *arXiv preprint arXiv:1903.07571*, 2019.
- [13] Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and double descent curve. *arXiv preprint arXiv:1908.05355*, 2019.
- [14] Elizabeth Gardner and Bernard Derrida. Three unfinished works on the optimal storage capacity of networks. *Journal of Physics A: Mathematical and General*, 22(12):1983, 1989.
- [15] Manfred Opper and Wolfgang Kinzel. Statistical mechanics of generalization. In Models of neural networks III, pages 151–209. Springer, 1996.
- [16] Christos Thrampoulidis, Ehsan Abbasi, and Babak Hassibi. Precise error analysis of regularized *m*-estimators in high dimensions. *IEEE Transactions on Information Theory*, 64(8):5592–5628, 2018.
- [17] Andrea Montanari, Feng Ruan, Youngtak Sohn, and Jun Yan. The generalization error of max-margin linear classifiers: High-dimensional asymptotics in the overparametrized regime, 2019.
- [18] Zeyu Deng, Abla Kammoun, and Christos Thrampoulidis. A model of double descent for high-dimensional binary linear classification. *arXiv preprint arXiv:1911.05822*, 2019.
- [19] Ganesh Kini and Christos Thrampoulidis. Analytic study of double descent in binary classification: The impact of loss. *arXiv preprint arXiv:2001.11572*, 2020.
- [20] Francesca Mignacco, Florent Krzakala, Yue M. Lu, and Lenka Zdeborová. The role of regularization in classification of high-dimensional noisy Gaussian mixture. 2(2):1–21, 2020.
- [21] Marc Mezard and Andrea Montanari. *Information, physics, and computation*. Oxford University Press, 2009.
- [22] Lenka Zdeborova. Understanding deep learning is also a job for physicists. Nature Physics, pages 1745–2481, 2020.
- [23] M. Opper, W. Kinzel, J. Kleinz, and R. Nehl. On the ability of the optimal perceptron to generalise. *Journal of Physics A: General Physics*, 23(11), 1990.
- [24] Manfred Opper and David Haussler. Generalization performance of Bayes optimal classification algorithm for learning a perceptron. *Physical Review Letters*, 66(20):2677–2680, 1991.
- [25] M. Opper and W. Kinzel. Models of neural networks III. Springer, 1996.

- [26] Marc Mézard. The space of interactions in neural networks: Gardner's computation with the cavity method. *Journal of Physics A: Mathematical and General*, 22(12):2181, 1989.
- [27] Yoshiyuki Kabashima. A cdma multiuser detection algorithm on the basis of belief propagation. *Journal of Physics A: Mathematical and General*, 36(43):11111, 2003.
- [28] Yoshiyuki Kabashima and Shinsuke Uda. A bp-based algorithm for performing bayesian inference in large perceptron-type networks. In *International Conference on Algorithmic Learning Theory*, pages 479–493. Springer, 2004.
- [29] David L Donoho, Arian Maleki, and Andrea Montanari. Message-passing algorithms for compressed sensing. *Proceedings of the National Academy of Sciences*, 106(45):18914–18919, 2009.
- [30] Sundeep Rangan. Generalized approximate message passing for estimation with random linear mixing. pages 2168–2172, 2011.
- [31] Mohsen Bayati and Andrea Montanari. The dynamics of message passing on dense graphs, with applications to compressed sensing. *IEEE Transactions on Information Theory*, 57(2):764–785, 2011.
- [32] Cédric Gerbelot, Alia Abbara, and Florent Krzakala. Asymptotic errors for convex penalized linear regression beyond gaussian matrices. *arXiv preprint arXiv:2002.04372*, 2020.
- [33] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [34] Saharon Rosset, Ji Zhu, and Trevor J. Hastie. Margin maximizing loss functions. In S. Thrun, L. K. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems 16*, pages 1237–1244. MIT Press, 2004.
- [35] Vladimir Vapnik. *Estimation of dependences based on empirical data*. Springer Science & Business Media, 2006.
- [36] Shai Shalev-Shwartz and Shai Ben-David. Understanding machine learning: From theory to algorithms. Cambridge university press, 2014.
- [37] Rémi Gribonval. Should penalized least squares regression be interpreted as maximum a posteriori estimation? *IEEE Transactions on Signal Processing*, 59(5):2405–2410, 2011.
- [38] Remi Gribonval and Pierre Machart. Reconciling "priors" & amp; "priors" without prejudice? In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 2193–2201. Curran Associates, Inc., 2013.

- [39] Madhu Advani and Surya Ganguli. An equivalence between high dimensional Bayes optimal inference and M-estimation. Advances in Neural Information Processing Systems, (1):3386–3394, 2016.
- [40] Lenka Zdeborová and Florent Krzakala. Statistical physics of inference: Thresholds and algorithms. Advances in Physics, 65(5):453–552, 2016.
- [41] Mario Geiger, Stefano Spigler, Stéphane d' Ascoli, Levent Sagun, Marco Baity-Jesi, Giulio Biroli, and Matthieu Wyart. Jamming transition as a paradigm to understand the loss landscape of deep neural networks. *Physical Review E*, 100(1), Jul 2019.
- [42] Mario Geiger, Arthur Jacot, Stefano Spigler, Franck Gabriel, Levent Sagun, Stéphane d'Ascoli, Giulio Biroli, Clément Hongler, and Matthieu Wyart. Scaling description of generalization with number of parameters in deep learning. *Journal of Statistical Mechanics: Theory and Experiment*, 2020(2):023401, 2020.
- [43] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machinelearning practice and the classical bias-variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- [44] Partha P Mitra. Understanding overfitting peaks in generalization error: Analytical risk curves for l_2 and l_1 penalized interpolation, 2019.
- [45] Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and double descent curve, 2019.
- [46] Federica Gerace, Bruno Loureiro, Florent Krzakala, Marc Mezard, and Lenka Zdeborová. Generalisation error in learning with random features and the hidden manifold model. arXiv preprint arXiv:2002.09339, 2020.
- [47] Stéphane d'Ascoli, Maria Refinetti, Giulio Biroli, and Florent Krzakala. Double trouble in double descent: Bias and variance (s) in the lazy regime. arXiv preprint arXiv:2003.01054, 2020.
- [48] Peter Bartlett and John Shawe-taylor. Generalization performance of support vector machines and other pattern classifiers, 1998.
- [49] Alia Abbara, Benjamin Aubin, Florent Krzakala, and Lenka Zdeborová. Rademacher complexity and spin glasses: A link between the replica and statistical theories of learning. arXiv preprint arXiv:1912.02729, 2019.
- [50] A Engel and W Fink. Statistical mechanics calculation of vapnik-chervonenkis bounds for perceptrons. *Journal of Physics A: Mathematical and General*, 26(23):6893, 1993.
- [51] Rémi Gribonval and Mila Nikolova. A characterization of proximity operators. *arXiv* preprint arXiv:1807.04014, 2018.
- [52] Rémi Gribonval and Mila Nikolova. On bayesian estimation and proximity operators. *Applied and Computational Harmonic Analysis*, 2019.

- [53] Jean Barbier, Florent Krzakala, Nicolas Macris, Léo Miolane, and Lenka Zdeborová. Phase Transitions, Optimal Errors and Optimality of Message-Passing in Generalized Linear Models. pages 1–59, 2017.
- [54] Andrea Montanari, YC Eldar, and G Kutyniok. Graphical models concepts in compressed sensing. *Compressed Sensing: Theory and Applications*, pages 394–438, 2012.
- [55] H Nishimori. Exact results and critical properties of the ising model with competing interactions. *Journal of Physics C: Solid State Physics*, 13(21):4071–4076, jul 1980.
- [56] Osame Kinouchi and Nestor Caticha. Learning algorithm that gives the bayes generalization limit for perceptrons. *Physical Review E - Statistical Physics, Plasmas, Fluids, and Related Interdisciplinary Topics*, 54(1):R54–R57, 1996.
- [57] Derek Bean, Peter J. Bickel, Noureddine El Karoui, and Bin Yu. Optimal M-estimation in high-dimensional regression. *Proceedings of the National Academy of Sciences of the United States of America*, 110(36):14563–14568, 2013.
- [58] David Donoho and Andrea Montanari. High dimensional robust M-estimation: asymptotic variance via approximate message passing. *Probability Theory and Related Fields*, 166(3-4):935–969, 2016.
- [59] Madhu Advani and Surya Ganguli. Statistical mechanics of optimal convex inference in high dimensions. *Physical Review X*, 6(3):1–16, 2016.

Supplementary material

In this supplementary material (SM), we provide the proofs and computation details leading to the results presented in the main manuscript. In Sec. I, we first recall the definition of the statistical model used in Sec. 1 and we give proper definitions of the denoising distributions involved in the analysis of the Bayes-optimal and Empirical Risk Minimization (ERM) estimation. In particular, we provide the analytical expressions of the denoising functions used in Sec. 3 to analyze ridge, hinge and logistic regressions. In Sec. II, we detail the computation of the binary classification generalization error leading to the expressions in Proposition. 2.1 and Thm. 2.4 respectively for ERM and Bayes-optimal estimation. In Sec. III, we present the proofs of the central theorems stated in Sec. 2. In particular, we derive the Gordon-based proof of the Thm. 2.2 in the more general regression (real-valued) version and provide as well the proof of Corollary. 2.3 which establishes the equivalence between the set of fixed-point equations of the Gordon's proof in the binary classification case and the one resulting from the heuristic replica computation. The corresponding statistical physics framework used to analyze Bayes and ERM statistical estimations and the replica computation leading to expressions in Corollary. 2.3 are detailed In Sec. IV. The section V is devoted to provide additional technical details on the results with ℓ_2 regularization addressed in Sec. 3. In particular, we present the large α expansions of the generalization error for the Bayes-optimal, ridge, pseudo-inverse and max-margin estimators, and we investigate the performances of logistic regression on non-linearly separable data. Finally in Sec. VI, we show the derivation of the fine-tuned loss and regularizer provably leading to Bayes-optimal performances, as explained and advocated in Sec. 4, and we show some numerical evidences that ERM achieves indeed Bayes-optimal error in Fig. 5.

Table of Contents

Ι	Defi	nitions and notations	20
	I.1	Statistical model	20
	I.2	Bayes-optimal and ERM estimation	20
	I.3	Denoising distributions and updates	21
	I.4	Applications	24
Π	Bina	ary classification generalization errors	28
	II.1	General case	28
	II.2	Bayes-optimal generalization error	29
	II.3	ERM generalization error	29
III Proofs of the ERM fixed points30			
	III.1	Gordon's result and proofs	30
	III.2	Replica's formulation	33
	III.3	Equivalence Gordon-Replica's formulation - ℓ_2 regularization and Gaussian	
		weights	34
IV	IV Replica computation for Bayes-optimal and ERM estimations 33		
	IV.1	Statistical inference and free entropy	38
	IV.2	Replica computation	38
	IV.3	ERM and Bayes-optimal free entropy	42
	IV.4	Sets of fixed point equations	43
	IV.5	Useful derivations	44
v	Арр	lications	50
	V.1	Bayes-optimal estimation	50
	V.2	Generalities on ERM with ℓ_2 regularization	51
	V.3	Ridge regression - Square loss with ℓ_2 regularization $\ldots \ldots \ldots \ldots \ldots$	51
	V.4	Hinge regression / SVM - Hinge loss with ℓ_2 regularization	56
	V.5	Logistic regression	58
	V.6	Logistic with non-linearly separable data - A rectangle door teacher	59
VI Reaching Bayes optimality 60			
	VI.1	Generalized Approximate Message Passing (GAMP) algorithm	60
	VI.2	Matching Bayes-optimal and ERM performances	60
	VI.3	Summary and numerical evidences	62

I Definitions and notations

I.1 Statistical model

We recall the supervised machine learning task considered in the main manuscript eq. (1), whose dataset is generated by a single layer neural network, often named a *teacher*, that belongs to the Generalized Linear Model (GLM) class. Therefore we assume the n samples are drawn according to

$$\mathbf{y} = \varphi_{\text{out}}^{\star} \left(\frac{1}{\sqrt{d}} \mathbf{X} \mathbf{w}^{\star} \right) \Leftrightarrow \mathbf{y} \sim P_{\text{out}}^{\star} \left(. \right) \,, \tag{18}$$

where $\mathbf{w}^{\star} \in \mathbb{R}^d$ denotes the ground truth vector drawn from a probability distribution $P_{\mathbf{w}^{\star}}$ with second moment $\rho_{\mathbf{w}^{\star}} \equiv \frac{1}{d}\mathbb{E}\left[\|\mathbf{w}^{\star}\|_{2}^{2}\right]$ and $\varphi_{\text{out}}^{\star}$ represents a deterministic or stochastic activation function equivalently associated to a distribution P_{out}^{\star} . The input data matrix $\mathbf{X} = (\mathbf{x}_{\mu})_{\mu=1}^{n} \in \mathbb{R}^{n \times d}$ contains iid Gaussian vectors, i.e $\forall \mu \in [1:n], \mathbf{x}_{\mu} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d})$.

I.2 Bayes-optimal and ERM estimation

Inferring the above statistical model from observations $\{y, X\}$ can be tackled in several ways. In particular, Bayesian inference provides a generic framework for statistical estimation based on the high-dimensional, often intractable, posterior distribution

$$\mathbb{P}(\mathbf{w}|\mathbf{y}, \mathbf{X}) = \frac{\mathbb{P}(\mathbf{y}|\mathbf{w}, \mathbf{X}) \mathbb{P}(\mathbf{w})}{\mathbb{P}(\mathbf{y}, \mathbf{X})}.$$
(19)

Estimating the average of the above posterior distribution in the case we have access to the ground truth prior distributions $\mathbb{P}(\mathbf{y}|\mathbf{w}, \mathbf{X}) = P_{\text{out}^{\star}}(\mathbf{y}|\mathbf{z})$ with $\mathbf{z} \equiv \frac{1}{\sqrt{d}}\mathbf{X}\mathbf{w}$ and $\mathbb{P}(\mathbf{w}) = P_{\mathbf{w}^{\star}}(\mathbf{w})$, refers to Bayes-optimal estimation and leads to the corresponding Minimal Mean-Squared Error (MMSE) estimator $\hat{\mathbf{w}}_{\text{mmse}} = \mathbb{E}_{\mathbb{P}(\mathbf{w}|\mathbf{y},\mathbf{X})}[\mathbf{w}]$. It has been rigorously analyzed in details in [10] for the whole GLM class eq. (18). Another celebrated approach and widely used in practice is the Empirical Risk Minimization (ERM) that minimizes instead a regularized loss: $\hat{\mathbf{w}}_{\text{erm}} = \operatorname{argmin}_{\mathbf{w}} [\mathcal{L}(\mathbf{w};\mathbf{y},\mathbf{X})]$ with

$$\mathcal{L}(\mathbf{w}; \mathbf{y}, \mathbf{X}) = \sum_{\mu=1}^{n} l(\mathbf{w}; y_{\mu}, \mathbf{x}_{\mu}) + r(\mathbf{w}) .$$
(20)

Interestingly analyzing the ERM estimation may be included in the above Bayesian framework. Indeed exponentiating eq. (20), we see that minimizing the loss \mathcal{L} is equivalent to maximize the posterior distribution $\mathbb{P}(\mathbf{w}|\mathbf{y}, \mathbf{X}) = e^{-\mathcal{L}(\mathbf{w};\mathbf{y},\mathbf{X})}$ if we choose carefully the prior distributions as functions of the regularizer r and the loss l:

$$-\log \mathbb{P}(\mathbf{y}|\mathbf{w}, \mathbf{X}) = l(\mathbf{w}; \mathbf{y}, \mathbf{X}), \quad -\log \mathbb{P}(\mathbf{w}) = r(\mathbf{w}).$$
(21)

Computing the maximum of the posterior $\mathbb{P}(\mathbf{y}|\mathbf{w}, X)$ refers instead to the so-called Maximum A Posteriori (MAP) estimator, and therefore analyzing the empirical minimization of (20) is equivalent to obtain the performance of the MAP estimator with prior distributions given by (21). Thus both the study of ERM (MAP) and Bayes-optimal (MMSE) estimations are simply reduced to the analysis of the posterior eq. (19).

I.3 Denoising distributions and updates

Analyzing the posterior distribution eq. (19) in the high-dimensional regime [10] will boil down to introducing the scalar denoising distributions $Q_{\rm w}, Q_{\rm out}$ and their respective normalizations $Z_{\rm w}, Z_{\rm out}$

$$Q_{\rm w}(w;\gamma,\Lambda) \equiv \frac{P_{\rm w}(w)}{\mathcal{Z}_{\rm w}(\gamma,\Lambda)} e^{-\frac{1}{2}\Lambda w^2 + \gamma w}, \quad Q_{\rm out}(z;y,\omega,V) \equiv \frac{P_{\rm out}\left(y|z\right)}{\mathcal{Z}_{\rm out}(y,\omega,V)} \frac{e^{-\frac{1}{2}V^{-1}(z-\omega)^2}}{\sqrt{2\pi V}}, \\ \mathcal{Z}_{\rm w}(\gamma,\Lambda) \equiv \mathbb{E}_{w\sim P_{\rm w}}\left[e^{-\frac{1}{2}\Lambda w^2 + \gamma w}\right], \quad \mathcal{Z}_{\rm out}(y,\omega,V) \equiv \mathbb{E}_{z\sim\mathcal{N}(0,1)}\left[P_{\rm out}\left(y|\sqrt{V}z+\omega\right)\right].$$
(22)

We define as well the denoising functions, that play a central role in Bayesian inference. Note in particular that they correspond to the *updates* of the Approximate Message Passing algorithm in [30] that we recalled in Sec. VI.1. They are defined as the derivatives of $\log Z_w$ and $\log Z_{out}$, namely

$$f_{\rm w}(\gamma,\Lambda) \equiv \partial_{\gamma} \log\left(\mathcal{Z}_{\rm w}\right) = \mathbb{E}_{Q_{\rm w}}\left[w\right] \quad \text{and} \quad \partial_{\gamma} f_{\rm w}(\gamma,\Lambda) \equiv \mathbb{E}_{Q_{\rm w}}\left[w^2\right] - f_{\rm w}^2$$
$$f_{\rm out}(y,\omega,V) \equiv \partial_{\omega} \log\left(\mathcal{Z}_{\rm out}\right) = V^{-1} \mathbb{E}_{Q_{\rm out}}\left[z-\omega\right] \quad \text{and} \quad \partial_{\omega} f_{\rm out}(y,\omega,V) \equiv \frac{\partial f_{\rm out}(y,\omega,V)}{\partial\omega}$$
(23)

I.3.1 Bayes-optimal - MMSE denoising functions

In Bayes-optimal estimation, the ground truth prior and channel distributions $P_{w^*}(w)$ and $P_{out^*}(y|z)$ of the *teacher* eq. (1) are known. Hence, replacing P_w and P_{out} in (22), we obtain the Bayes-optimal scalar denoising distributions in terms of which the Bayes-optimal free entropy eq. (95) is written

$$Q_{\mathbf{w}^{\star}}(w;\gamma,\Lambda) \equiv \frac{P_{\mathbf{w}^{\star}}(w)}{\mathcal{Z}_{\mathbf{w}^{\star}}(\gamma,\Lambda)} e^{-\frac{1}{2}\Lambda w^{2}+\gamma w}, \quad Q_{\mathrm{out}^{\star}}(z;y,\omega,V) \equiv \frac{P_{\mathrm{out}^{\star}}(y|z)}{\mathcal{Z}_{\mathrm{out}^{\star}}(y,\omega,V)} \frac{e^{-\frac{1}{2}V^{-1}(z-\omega)^{2}}}{\sqrt{2\pi V}}.$$
(24)

and the denoising updates are therefore given by eq. (23) with the corresponding distributions

$$f_{w^{\star}}(\gamma,\Lambda) \equiv \partial_{\gamma} \log \mathcal{Z}_{w^{\star}}(\gamma,\Lambda) , \qquad f_{out^{\star}}(y,\omega,V) \equiv \partial_{\omega} \log \mathcal{Z}_{out^{\star}}(y,\omega,V) .$$
(25)

I.3.2 ERM - MAP denoising functions

Before defining similar denoising functions to analyze the MAP for ERM estimation, we first recall the definition of the Moreau-Yosida regularization.

Moreau-Yosida regularization and proximal Let $\Sigma > 0$, f(z) a convex function in z. Defining the regularized functional

$$\mathcal{L}_{\Sigma}[f(,.)](z;x) = f(,z) + \frac{1}{2\Sigma} (z-x)^2 , \qquad (26)$$

the Moreau-Yosida regularization \mathcal{M}_{Σ} and the proximal map \mathcal{P}_{Σ} are defined by

$$\mathcal{P}_{\Sigma}[f(,.)](x) = \operatorname{argmin}_{z} \mathcal{L}_{\Sigma}[f(,.)](z;x) = \operatorname{argmin}_{z} \left[f(,z) + \frac{1}{2\Sigma} \left(z - x \right)^{2} \right], \qquad (27)$$

$$\mathcal{M}_{\Sigma}[f(,.)](x) = \min_{z} \mathcal{L}_{\Sigma}[f(,.)](z;x) = \min_{z} \left[f(,z) + \frac{1}{2\Sigma} \left(z - x \right)^2 \right] \,, \tag{28}$$

where (, z) denotes all the arguments of the function f, where z plays a central role. The MAP denoising functions for any convex loss l(, .) and convex separable regularizer r(.) can be written in terms of the Moreau-Yosida regularization or the proximal map as follows

$$f_{\rm w}^{\rm map,r}(\gamma,\Lambda) \equiv \mathcal{P}_{\Lambda^{-1}}\left[r(.)\right](\Lambda^{-1}\gamma) = \Lambda^{-1}\gamma - \Lambda^{-1}\partial_{\Lambda^{-1}\gamma}\mathcal{M}_{\Lambda^{-1}}\left[r(.)\right](\Lambda^{-1}\gamma),$$

$$f_{\rm out}^{\rm map,l}(y,\omega,V) \equiv -\partial_{\omega}\mathcal{M}_{V}[l(y,.)](\omega) = V^{-1}\left(\mathcal{P}_{V}[l(y,.)](\omega) - \omega\right).$$
(29)

The above updates can be considered as definitions, but it is instructive to derive them from the generic definition of the denoising distributions eq. (23) if we maximize the posterior distribution. This is done by taking, in a physics language, a *zero temperature* limit and we present it in details in the next paragraph.

Derivation of the MAP updates To have access to the maximum of the generic distributions eq. (22), we introduce a *fictive* noise/temperature Δ or inverse temperature β , $\Delta = \frac{1}{\beta}$. In particular for Bayes-optimal estimation this temperature is finite and fixed to $\Delta = \beta = 1$. Indeed with the mapping eq. (21), minimizing the loss function \mathcal{L} (20) is equivalent to maximize the posterior distribution. Therefore it can be done by taking the zero noise/temperature limit $\Delta \rightarrow 0$ of the channel and prior denoising distributions Q_{out} and Q_{w} . It is the purpose of the following paragraphs where we present the derivation leading to the result (29).

Channel Using the mapping eq. (21), we assume that the channel distribution can be expressed as $\mathbb{P}(y|z) \propto e^{-l(y,z)}$. Therefore we introduce the corresponding channel distribution P_{out} at finite temperature Δ associated to the convex loss l(y, z)

$$P_{\text{out}}^{\text{map}}\left(y|z\right) = rac{e^{-rac{1}{\Delta}l\left(y,z
ight)}}{\sqrt{2\pi\Delta}}.$$

Note that the case of the square loss $l(y, z) = \frac{1}{2} (y - z)^2$ is very specific. Its channel distribution simply reads $P_{\text{out}}(y|z) = \frac{e^{-\frac{1}{2\Delta}(y-z)^2}}{\sqrt{2\pi\Delta}}$ and is therefore equivalent to predict labels y according to a noisy Gaussian linear model $y = z + \sqrt{\Delta}\xi$, where $\xi \sim \mathcal{N}(0, 1)$ and Δ denotes therefore the *real* noise of the model.

In order to obtain a non trivial limit and a closed set of equations when $\Delta \rightarrow 0$, we must defined rescaled variables as follows:

$$V_{\dagger} \equiv \lim_{\Delta \to 0} \frac{V}{\Delta} , \qquad \qquad f_{\mathrm{out},\dagger}^{\mathrm{map}}(y,\omega,V_{\dagger}) \equiv \lim_{\Delta \to 0} \Delta \times f_{\mathrm{out}}^{\mathrm{map}}(y,\omega,V) ,$$

where we denote the rescaled quantities after taking the limit $\Delta \rightarrow 0$ by \dagger . Similarly to eq. (26), we introduce therefore the rescaled functional

$$\mathcal{L}_{V_{\dagger}}[l(y,.)](z;\omega) = l(y,z) + \frac{1}{2V_{\dagger}}(z-\omega)^2 , \qquad (30)$$

such that, injecting $P_{\text{out}}^{\text{map}}$, the channel denoising distribution $Q_{\text{out}}^{\text{map}}$ and the corresponding partition function $\mathcal{Z}_{\text{out}}^{\text{map}}$ eq. (22) simplify in the zero temperature limit as follows:

$$Q_{\text{out}}^{\text{map}}(z;y,\omega,V) \equiv \lim_{\Delta \to 0} \frac{e^{-\frac{1}{\Delta}l(y,z) + \frac{1}{2V}(z-\omega)^2}}{\sqrt{2\pi\Delta V_{\dagger}}\sqrt{2\pi\Delta}} = \lim_{\Delta \to 0} \frac{e^{-\frac{1}{\Delta}\mathcal{L}_{V_{\dagger}}[l(y,.)](z;\omega)}}{\sqrt{2\pi\Delta V_{\dagger}}\sqrt{2\pi\Delta}}, \qquad (31)$$
$$\propto \delta\left(z - \mathcal{P}_{V_{\dagger}}[l(y,.)](\omega)\right)$$

$$\mathcal{Z}_{\text{out}}^{\text{map}}\left(y,\omega,V\right) = \lim_{\Delta \to 0} \int_{\mathbb{R}} \mathrm{d}z Q_{\text{out}}^{\text{map}}(z;y,\omega,V) = \lim_{\Delta \to 0} \frac{e^{-\frac{1}{\Delta}\mathcal{M}_{V_{\dagger}}\left[l(y,.)\right](\omega)}}{\sqrt{2\pi\Delta V_{\dagger}}\sqrt{2\pi\Delta}},\qquad(32)$$

that involve the proximal map and the Moreau-Yosida regularization defined in eq. (28). Finally taking the zero temperature limit, the MAP denoising function $f_{\text{out},\dagger}^{\text{map}}$ leads to the result (29):

$$f_{\text{out},\dagger}^{\text{map}}(y,\omega,V_{\dagger}) \equiv \lim_{\Delta \to 0} \Delta \times f_{\text{out}}^{\text{map}}(y,\omega,V)$$

$$\equiv \lim_{\Delta \to 0} \Delta \times \partial_{\omega} \log \mathcal{Z}_{\text{out}}^{\text{map}} \equiv \lim_{\Delta \to 0} \Delta V^{-1} \mathbb{E}_{Q_{\text{out}}^{\text{map}}}[z-\omega] \qquad (33)$$

$$= -\partial_{\omega} \mathcal{M}_{V_{\dagger}}[l(y,.)](\omega) = V_{\dagger}^{-1} \left(\mathcal{P}_{V_{\dagger}}[l(y,.)](\omega) - \omega \right).$$

Prior Similarly as above, using the mapping eq. (21), for a convex and separable regularizer r, the corresponding prior distribution at temperature Δ can be written

$$P_{\mathbf{w}}^{\mathrm{map}}\left(w\right) = e^{-\frac{1}{\Delta}r\left(w\right)}$$

Note that at $\Delta = 1$ the classical ℓ_1 regularization with strength λ , $r^{\ell_1}(w) = -\lambda |w|$, and the ℓ_2 regularization $r^{\ell_2}(w) = -\lambda w^2/2$ are equivalent to choosing a Laplace prior $P_{\rm w}(w) \propto e^{-\lambda |w|}$ or a Gaussian prior $P_{\rm w}(w) \propto e^{-\frac{\lambda w^2}{2}}$. To obtain a meaningful limit as $\Delta \to 0$, we again introduce the following rescaled variables

$$\Lambda_{\dagger} \equiv \lim_{\Delta \to 0} \Delta \times \Lambda \,, \qquad \qquad \gamma_{\dagger} \equiv \lim_{\Delta \to 0} \Delta \times \gamma \,,$$

and the functional

$$\mathcal{L}_{\Lambda_{\dagger}^{-1}}\left[r(.)\right]\left(w;\Lambda_{\dagger}^{-1}\gamma_{\dagger}\right) = r(w) + \frac{1}{2}\Lambda_{\dagger}\left(w - \Lambda_{\dagger}^{-1}\gamma_{\dagger}\right)^{2} = \left[r(w) + \frac{1}{2}\Lambda_{\dagger}w^{2} - \gamma_{\dagger}w\right] + \frac{1}{2}\gamma_{\dagger}^{2}\Lambda_{\dagger}^{-1},$$
(34)

such that in the zero temperature limit, the prior denoising distribution $Q_{\rm w}^{\rm map}$ and the partition function $\mathcal{Z}_{\rm w}^{\rm map}$ reduce to

$$Q_{\mathbf{w}}^{\mathrm{map}}\left(w;\gamma,\Lambda\right) \equiv \lim_{\Delta \to 0} P_{\mathbf{w}}(w) e^{-\frac{1}{2}\Lambda w^{2} + \gamma w} = \lim_{\Delta \to 0} e^{-\frac{1}{\Delta}\mathcal{L}_{\Lambda_{\dagger}^{-1}}[r](w;\Lambda_{\dagger}^{-1}\gamma_{\dagger})} e^{-\frac{1}{2\Delta}\gamma_{\dagger}^{2}\Lambda_{\dagger}^{-1}}$$
$$\propto \delta\left(w - \mathcal{P}_{\Lambda_{\dagger}^{-1}}\left[r\right]\left(\Lambda_{\dagger}^{-1}\gamma_{\dagger}\right)\right) \tag{35}$$

$$\mathcal{Z}_{\mathbf{w}}^{\mathrm{map}}\left(y,\omega,V\right) = \lim_{\Delta \to 0} \int_{\mathbb{R}} \mathrm{d}w Q_{\mathbf{w}}^{\mathrm{map}}(w;\gamma,\Lambda) = \lim_{\Delta \to 0} e^{-\frac{1}{\Delta}\mathcal{M}_{\Lambda_{\dagger}^{-1}}\left[r\right]\left(\Lambda_{\dagger}^{-1}\gamma_{\dagger}\right)} e^{-\frac{1}{2\Delta}\gamma_{\dagger}^{2}\Lambda_{\dagger}^{-1}}, \quad (36)$$

that involve again the proximal map $\mathcal{P}_{\Lambda^{-1}_{\dagger}}$ and the Moreau-Yosida regularization $\mathcal{M}_{\Lambda^{-1}_{\dagger}}$ defined in eq. (28). Finally the MAP denoising update $f_{\mathrm{w},\dagger}^{\mathrm{map}}$ is simply given by:

$$\begin{split} f_{\mathrm{w},\dagger}^{\mathrm{map}}(\gamma_{\dagger},\Lambda_{\dagger}) &\equiv \lim_{\Delta \to 0} f_{\mathrm{w}}^{\mathrm{map}}(\gamma,\Lambda) = \lim_{\Delta \to 0} \partial_{\gamma} \log \mathcal{Z}_{\mathrm{w}}^{\mathrm{map}} \equiv \lim_{\Delta \to 0} \mathbb{E}_{Q_{\mathrm{w}}^{\mathrm{map}}} \left[w\right] \\ &= \lim_{\Delta \to 0} \partial_{\gamma} \left(-\frac{1}{\Delta} \mathcal{M}_{\Lambda_{\dagger}^{-1}}\left[r(.)\right] \left(\Lambda_{\dagger}^{-1}\gamma_{\dagger}\right) - \frac{1}{2\Delta} \gamma_{\dagger}^{2} \Lambda_{\dagger}^{-1}\right) \\ &= \partial_{\gamma_{\dagger}} \left(-\mathcal{M}_{\Lambda_{\dagger}^{-1}}\left[r(.)\right] \left(\Lambda_{\dagger}^{-1}\gamma_{\dagger}\right) - \frac{1}{2} \gamma_{\dagger}^{2} \Lambda_{\dagger}^{-1}\right) \\ &= \Lambda_{\dagger}^{-1} \gamma_{\dagger} - \Lambda_{\dagger}^{-1} \partial_{\Lambda_{\dagger}^{-1}\gamma_{\dagger}} \mathcal{M}_{\Lambda_{\dagger}^{-1}}\left[r(.)\right] \left(\Lambda_{\dagger}^{-1}\gamma_{\dagger}\right) = \mathcal{P}_{\Lambda_{\dagger}^{-1}}\left[r(.)\right] \left(\Lambda_{\dagger}^{-1}\gamma_{\dagger}\right) \\ &= \operatorname{argmin}_{w} \left[r(w) + \frac{1}{2} \Lambda_{\dagger} (w - \Lambda_{\dagger}^{-1} \gamma_{\dagger})^{2}\right] = \operatorname{argmin}_{w} \left[r(w) + \frac{1}{2} \Lambda_{\dagger} w^{2} - \gamma_{\dagger} w\right], \end{split}$$

and we recover the result (29).

I.4 Applications

In this section we list the explicit expressions of the Bayes-optimal eq. (25) and ERM eq. (29) denoising functions largely used to produce the examples in Sec. 3.

I.4.1 Bayes-optimal updates

The Bayes-optimal denoising functions (25) are detailed in the case of a *linear*, sign and rectangle door channel with a Gaussian noise $\xi \sim \mathcal{N}(0, 1)$ and variance $\Delta \geq 0$, and for Gaussian and sparse-binary weights.

Channel

• Linear: $y = \varphi_{\text{out}^{\star}}(z) = z + \sqrt{\Delta} \xi$

$$\mathcal{Z}_{\text{out}^{\star}}(y,\omega,V) = \mathcal{N}_{\omega}\left(y,\Delta^{\star}+V\right),$$

$$f_{\text{out}^{\star}}(y,\omega,V) = \left(\Delta^{\star}+V\right)^{-1}\left(y-\omega\right), \quad \partial_{\omega}f_{\text{out}^{\star}}(y,\omega,V) = -\left(\Delta^{\star}+V\right)^{-1}.$$
(38)

• Sign: $y = \varphi_{\text{out}^{\star}}(z) = \operatorname{sign}(z) + \sqrt{\Delta^{\star}}\xi$

$$\mathcal{Z}_{\text{out}^{\star}}(y,\omega,V) = \mathcal{N}_{y}(1,\Delta^{\star})\frac{1}{2}\left(1 + \operatorname{erf}\left(\frac{\omega}{\sqrt{2V}}\right)\right) + \mathcal{N}_{y}(-1,\Delta^{\star})\frac{1}{2}\left(1 - \operatorname{erf}\left(\frac{\omega}{\sqrt{2V}}\right)\right),$$

$$f_{\text{out}^{\star}}(y,\omega,V) = \frac{\mathcal{N}_{y}(1,\Delta^{\star}) - \mathcal{N}_{y}(-1,\Delta^{\star})}{\mathcal{Z}_{\text{out}^{\star}}(y,\omega,V)}\mathcal{N}_{\omega}(0,V).$$
(39)

• Rectangle door: $y = \varphi_{\text{out}^{\star}}(z) = \mathbb{1} (\kappa_m \leq z \leq \kappa_M) - \mathbb{1} (z \leq \kappa_m \text{ or } z \geq \kappa_M) + \sqrt{\Delta^{\star}} \xi$ For $\kappa_m < \kappa_M$, we obtain

$$\mathcal{Z}_{\text{out}^{\star}}(y,\omega,V) = \mathcal{N}_{y}(1,\Delta^{\star})\frac{1}{2}\left(\operatorname{erf}\left(\frac{\kappa_{M}-\omega}{\sqrt{2V}}\right) - \operatorname{erf}\left(\frac{\kappa_{m}-\omega}{\sqrt{2V}}\right)\right) \\ + \mathcal{N}_{y}(-1,\Delta^{\star})\frac{1}{2}\left(1 - \frac{1}{2}\left(\operatorname{erf}\left(\frac{\kappa_{M}-\omega}{\sqrt{2V}}\right) - \operatorname{erf}\left(\frac{\kappa_{m}-\omega}{\sqrt{2V}}\right)\right)\right), \quad (40)$$
$$f_{\text{out}^{\star}}(y,\omega,V) = \frac{1}{\mathcal{Z}_{\text{out}}}\left(\mathcal{N}_{y}(1,\Delta^{\star})\left(-\mathcal{N}_{\omega}(\kappa_{M},V) + \mathcal{N}_{\omega}(\kappa_{m},V)\right) \\ + \mathcal{N}_{y}(-1,\Delta^{\star})\left(\mathcal{N}_{\omega}(\kappa_{M},V) - \mathcal{N}_{\omega}(\kappa_{m},V)\right)\right).$$

Prior

• Gaussian weights: $w \sim P_{\mathrm{w}}(w) = \mathcal{N}_w(\mu, \sigma)$

$$\mathcal{Z}_{\mathbf{w}^{\star}}(\gamma,\Lambda) = \frac{e^{\frac{\gamma^{2}\sigma+2\gamma\mu-\Lambda\mu^{2}}{2(\Lambda\sigma+1)}}}{\sqrt{\Lambda\sigma+1}}, \quad f_{\mathbf{w}^{\star}}(\gamma,\Lambda) = \frac{\gamma\sigma+\mu}{1+\Lambda\sigma}, \quad \partial_{\gamma}f_{\mathbf{w}^{\star}}(\gamma,\Lambda) = \frac{\sigma}{1+\Lambda\sigma}.$$
(41)

• Sparse-binary weights: $w \sim P_w(w) = \rho \delta(w) + (\rho - 1) \frac{1}{2} \left(\delta(w - 1) + \delta(w + 1) \right)$

$$\mathcal{Z}_{\mathbf{w}^{\star}}(\gamma,\Lambda) = \rho + e^{-\frac{\Lambda}{2}}(1-\rho)\cosh(\gamma),$$

$$f_{\mathbf{w}^{\star}}(\gamma,\Lambda) = \frac{e^{-\frac{\Lambda}{2}}(1-\rho)\sinh(\gamma)}{\rho + e^{-\frac{\Lambda}{2}}(1-\rho)\cosh(\gamma)}, \quad \partial_{\gamma}f_{\mathbf{w}^{\star}}(\gamma,\Lambda) = \frac{e^{-\frac{\Lambda}{2}}(1-\rho)\cosh(\gamma)}{\rho + e^{-\frac{\Lambda}{2}}(1-\rho)\cosh(\gamma)}.$$
 (42)

I.4.2 ERM updates

The ERM denoising functions (29) have, very often, no explicit expression except for the *square* and *hinge* losses, and for ℓ_1 , ℓ_2 regularizations that are analytical. However, in the particular case of a two times differentiable convex loss the denoising functions can still be written as the solution of an implicit equation detailed below.

Convex losses

• Square loss

The proximal map for the square loss $l^{\text{square}}(y,z) = \frac{1}{2}(y-z)^2$ is easily obtained and reads

$$\mathcal{P}_{V}\left[\frac{1}{2}(y,.)^{2}\right](\omega) = \operatorname{argmin}_{z}\left[\frac{1}{2}(y-z)^{2} + \frac{1}{2V}(z-\omega)^{2}\right] = (1+V)^{-1}(\omega+yV) \ .$$

Therefore (29) yields

$$f_{\text{out}}^{\text{square}}(y,\omega,V) = V^{-1} \left(\mathcal{P}_V \left[\frac{1}{2} (y,.)^2 \right] (\omega) - \omega \right) = (1+V)^{-1} (y-\omega) , \qquad (43)$$
$$\partial_\omega f_{\text{out}}^{\text{square}}(y,\omega,V) = -(1+V)^{-1} .$$

• Hinge loss

The proximal map of the hinge loss $l^{\rm hinge}(y,z)=\max\left(0,1-yz\right)$

$$\mathcal{P}_{V}\left[l^{\text{hinge}}(y,.)\right](\omega) = \operatorname{argmin}_{z}\left[\underbrace{\max\left(0,1-yz\right) + \frac{1}{2V}\left(z-\omega\right)^{2}}_{\equiv \mathcal{L}_{0}}\right] \equiv z^{\star}(y,\omega,V).$$

can be expressed analytically by distinguishing all the possible cases:

- 1 yz < 0: $\mathcal{L}_0 = \frac{1}{2V} (z \omega)^2 \Rightarrow z^* = \omega$ if $yz^* < 1 \Leftrightarrow z^* = \omega$ if $\omega y < 1$.
- 1 yz > 0: $\mathcal{L}_0 = \frac{1}{2V} (z \omega)^2 + 1 yz \Rightarrow (z^* \omega) = yV \Leftrightarrow z^* = \omega + Vy$ if $1 yz^* > 0 \Leftrightarrow z^* = \omega + Vy$ if $\omega y < 1 y^2V = 1 V$, as $y^2 = 1$.
- Hence we have one last region to study $1-V < \omega y < 1.$ It follows $y(1-V) < \omega < y$:

$$\frac{1}{2V} (z - y)^2 \le \frac{1}{2V} (z - \omega)^2 \Rightarrow z^* = y.$$

Finally we obtain a simple analytical expression for the proximal and its derivative

$$\mathcal{P}_{V}\left[l^{\text{hinge}}(y,.)\right](\omega) = \begin{cases} \omega + Vy \text{ if } \omega y < 1 - V \\ y \text{ if } 1 - V < \omega y < 1 \\ \omega \text{ if } \omega y > 1 \end{cases}, \partial_{\omega}\mathcal{P}_{V}\left[l^{\text{hinge}}(y,.)\right](\omega) = \begin{cases} 1 \text{ if } \omega y < 1 - V \\ 0 \text{ if } 1 - V < \omega y < 1 \\ 1 \text{ if } \omega y > 1 \end{cases}$$

Hence with (29), the hinge denoising function and its derivative read

$$f_{\text{out}}^{\text{hinge}}(y,\omega,V) = \begin{cases} y \text{ if } \omega y < 1 - V \\ \frac{(y-\omega)}{V} \text{ if } 1 - V < \omega y < 1 \\ 0 \text{ otherwise} \end{cases}, \\ \partial_{\omega} f_{\text{out}}^{\text{hinge}}(y,\omega,V) = \begin{cases} -\frac{1}{V} \text{ if } 1 - V < \omega y < 1 \\ 0 \text{ otherwise} \end{cases}$$

$$(44)$$

• Generic differentiable convex loss

In general, finding the proximal map in (29) is intractable. In particular, it is the case for the logistic loss considered in Sec. V.5. However assuming the convex loss is a generic two times differentiable function $l \in D^2$, taking the derivative of the proximal map

$$\mathcal{P}_{V}\left[l(y,.)\right](\omega) = \operatorname{argmin}_{z}\left[l\left(y,z\right) + \frac{1}{2V}\left(z-\omega\right)^{2}\right] \equiv z^{\star}(y,\omega,V),$$

verifies therefore the implicit equations:

$$z^{\star}(y,\omega,V) = \omega - V\partial_z l\left(y, z^{\star}(y,\omega,V)\right), \quad \partial_\omega z^{\star}(y,\omega,V) = \left(1 + V\partial_z^2 l(y, z^{\star}(y,\omega,V))\right)^{-1}.$$
(45)

Once those equations solved, the denoising function and its derivative are simply expressed as

$$f_{\text{out}}^{\text{diff}}\left(y,\omega,V\right) = V^{-1}(z^{\star}\left(y,\omega,V\right) - \omega), \quad \partial_{\omega}f_{\text{out}}^{\text{diff}}\left(y,\omega,V\right) = V^{-1}\left(\partial_{\omega}z^{\star}\left(y,\omega,V\right) - 1\right),$$
(46)

with $z^{\star}(y, \omega, V) = \mathcal{P}_{V}[l(y, .)](\omega)$ solution of (45).

Regularizations

• ℓ_2 regularization

Using the definition of the prior update in eq. (29) for the ℓ_2 regularization $r(w) = \frac{\lambda w^2}{2}$, we obtain

$$f_{\mathbf{w}}^{\ell_{2}}(\gamma,\Lambda) = \operatorname{argmin}_{w} \left[\frac{\lambda w^{2}}{2} + \frac{1}{2}\Lambda w^{2} - \gamma w \right] = \frac{\gamma}{\lambda + \Lambda} ,$$

$$\partial_{\gamma} f_{\mathbf{w}}^{\ell_{2}}(\gamma,\Lambda) = \frac{1}{\lambda + \Lambda} \quad \text{and} \quad \mathcal{Z}_{\mathbf{w}}^{\ell_{2}}(\gamma,\Lambda) = \exp\left(\frac{\gamma^{2}\Lambda}{2(\lambda + \Lambda)^{2}}\right) .$$
(47)

• ℓ_1 regularization

Performing the same computation for the ℓ_1 regularization $r(w) = \lambda |w|$, we obtain

$$f_{\mathbf{w}}^{\ell_{1}}(\gamma,\Lambda) = \operatorname{argmin}_{w} \left[\lambda \|w\| + \frac{1}{2}\Lambda w^{2} - \gamma w \right] = \begin{cases} \frac{\gamma-\lambda}{\Lambda} & \gamma > \lambda \\ \frac{\gamma+\lambda}{\Lambda} & \gamma+\lambda < 0 \\ 0 \text{ otherwise} \end{cases},$$

$$\partial_{\gamma}f_{\mathbf{w}}^{\ell_{1}}(\gamma,\Lambda) = \begin{cases} \frac{1}{\Lambda} & \|\gamma\| > \lambda \\ 0 \text{ otherwise} \end{cases}.$$
(48)

II Binary classification generalization errors

In this section, we present the computation of the asymptotic generalization error

$$e_{g}(\alpha) \equiv \lim_{d \to \infty} \mathbb{E}_{y,\mathbf{x}} \mathbb{1} \left[y \neq \hat{y} \left(\hat{\mathbf{w}}(\alpha); \mathbf{x} \right) \right],$$
(49)

leading to expressions in Proposition. 2.1 and Thm. 2.4. The computation at finite dimension is similar if we do not consider the limit $d \to \infty$.

II.1 General case

The generalization error e_g is the prediction error of the estimator $\hat{\mathbf{w}}$ on new samples {y, X}, where X is an iid Gaussian matrix and y are ± 1 labels generated according to (18):

$$\mathbf{y} = \varphi_{\text{out}^{\star}} (\mathbf{z}) \quad \text{with} \quad \mathbf{z} = \frac{1}{\sqrt{d}} \mathbf{X} \mathbf{w}^{\star} \,.$$
 (50)

As the model fitted by ERM may not lead to binary outputs, we may add a non-linearity $\varphi : \mathbb{R} \mapsto \{\pm 1\}$ (for example a sign) on top of it to insure to obtain binary outputs $\hat{\mathbf{y}} \pm 1$ according to

$$\hat{\mathbf{y}} = \varphi(\hat{\mathbf{z}}) \quad \text{with} \quad \hat{\mathbf{z}} = \frac{1}{\sqrt{d}} \mathbf{X} \hat{\mathbf{w}} \,.$$
 (51)

The classification generalization error is given by the probability that the predicted labels \hat{y} and the true labels y do not match. To compute it, first note that the vectors $(\mathbf{z}, \hat{\mathbf{z}})$ averaged over all possible ground truth vectors \mathbf{w}^* (or equivalently labels y) and input matrix X follow in the large size limit a joint Gaussian distribution with zero mean and covariance matrix

$$\sigma = \lim_{d \to \infty} \mathbb{E}_{\mathbf{w}^{\star}, \mathbf{X}} \frac{1}{d} \begin{bmatrix} \mathbf{w}^{\star \mathsf{T}} \mathbf{w}^{\star} & \mathbf{w}^{\star \mathsf{T}} \hat{\mathbf{w}} \\ \mathbf{w}^{\star \mathsf{T}} \hat{\mathbf{w}} & \hat{\mathbf{w}}^{\mathsf{T}} \hat{\mathbf{w}} \end{bmatrix} \equiv \begin{bmatrix} \sigma_{\mathbf{w}^{\star}} & \sigma_{\mathbf{w}^{\star} \hat{\mathbf{w}}} \\ \sigma_{\mathbf{w}^{\star} \hat{\mathbf{w}}} & \sigma_{\hat{\mathbf{w}}} \end{bmatrix} .$$
(52)

The asymptotic generalization error depends only on the covariance matrix σ and as the samples are iid it reads

$$e_{g}(\alpha) = \lim_{d \to \infty} \mathbb{E}_{y,\mathbf{x}} \mathbb{1} \left[y \neq \hat{y} \left(\hat{\mathbf{w}}(\alpha); \mathbf{x} \right) \right] = 1 - \mathbb{P}[y = \hat{y} \left(\hat{\mathbf{w}}(\alpha); \mathbf{x} \right)] = 1 - 2 \int_{\left(\mathbb{R}^{+}\right)^{2}} d\mathbf{x} \mathcal{N}_{\mathbf{x}} \left(\mathbf{0}, \sigma \right)$$
$$= 1 - \left(\frac{1}{2} + \frac{1}{\pi} \operatorname{atan} \left(\sqrt{\frac{\sigma_{w^{\star} \hat{w}}^{2}}{\sigma_{w^{\star}} \sigma_{\hat{w}}^{2} - \sigma_{w^{\star} \hat{w}}^{2}}} \right) \right) = \frac{1}{\pi} \operatorname{acos} \left(\frac{\sigma_{w^{\star} \hat{w}}}{\sqrt{\sigma_{w^{\star}} \sigma_{\hat{w}}^{2}}} \right),$$
(53)

where we used the fact that $\operatorname{atan}(x) = \frac{\pi}{2} - \frac{1}{2} \operatorname{acos}\left(\frac{x^2-1}{1+x^2}\right)$ and $\frac{1}{2} \operatorname{acos}(2x^2-1) = \operatorname{acos}(x)$. Finally

$$e_{g}(\alpha) \equiv \lim_{d \to \infty} \mathbb{E}_{y,\mathbf{x}} \mathbb{1}\left[y \neq \hat{y}\left(\hat{\mathbf{w}}(\alpha); \mathbf{x} \right) \right] = \frac{1}{\pi} \operatorname{acos}\left(\frac{\sigma_{\mathbf{w}^{\star} \hat{\mathbf{w}}}}{\sqrt{\rho_{\mathbf{w}^{\star}} \sigma_{\hat{\mathbf{w}}}}} \right) \,, \tag{54}$$

with

$$\sigma_{\mathbf{w}^{\star}\hat{\mathbf{w}}} \equiv \lim_{d \to \infty} \mathbb{E}_{\mathbf{w}^{\star}, \mathbf{X}} \frac{1}{d} \hat{\mathbf{w}}^{\mathsf{T}} \mathbf{w}^{\star}, \quad \rho_{\mathbf{w}^{\star}} \equiv \lim_{d \to \infty} \mathbb{E}_{\mathbf{w}^{\star}} \frac{1}{d} \|\mathbf{w}^{\star}\|_{2}^{2}, \quad \sigma_{\hat{\mathbf{w}}} \equiv \lim_{d \to \infty} \mathbb{E}_{\mathbf{w}^{\star}, \mathbf{X}} \frac{1}{d} \|\hat{\mathbf{w}}\|_{2}^{2}.$$

II.2 Bayes-optimal generalization error

The Bayes-optimal generalization error for classification is equal to eq. (54) where the Bayes estimator $\hat{\mathbf{w}}$ is the average over the posterior distribution eq. (19) denoted $\langle . \rangle$, knowing the teacher prior $P_{\mathbf{w}^*}$ and channel P_{out^*} distributions: $\hat{\mathbf{w}} = \langle \mathbf{w} \rangle_{\mathbf{w}}$. Hence the parameters $\sigma_{\hat{\mathbf{w}}}$ and $\sigma_{\mathbf{w}^*\hat{\mathbf{w}}}$ read in the Bayes-optimal case

$$\begin{split} \sigma_{\hat{\mathbf{w}}} &\equiv \lim_{d \to \infty} \mathbb{E}_{\mathbf{w}^{\star}, \mathbf{X}} \frac{1}{d} \| \hat{\mathbf{w}} \|_{2}^{2} = \lim_{d \to \infty} \mathbb{E}_{\mathbf{w}^{\star}, \mathbf{X}} \frac{1}{d} \| \langle \mathbf{w} \rangle_{\mathbf{w}} \|_{2}^{2} \equiv q_{\mathbf{b}} \,, \\ \sigma_{\mathbf{w}^{\star} \hat{\mathbf{w}}} &\equiv \lim_{d \to \infty} \mathbb{E}_{\mathbf{w}^{\star}, \mathbf{X}} \frac{1}{d} \hat{\mathbf{w}}^{\mathsf{T}} \mathbf{w}^{\star} = \lim_{d \to \infty} \mathbb{E}_{\mathbf{w}^{\star}, \mathbf{X}} \frac{1}{d} \langle \mathbf{w} \rangle_{\mathbf{w}}^{\mathsf{T}} \mathbf{w}^{\star} \equiv m_{\mathbf{b}} \,. \end{split}$$

Using Nishimori identity [55], we easily obtain $m_{\rm b} = q_{\rm b}$ which is solution of eq. (13). Therefore the generalization error simplifies

$$e_{\rm g}^{\rm bayes}(\alpha) = \frac{1}{\pi} \operatorname{acos}\left(\sqrt{\eta}_{\rm b}\right), \text{ with } \eta_{\rm b} = \frac{q_{\rm b}}{\rho_{\rm w^{\star}}}.$$
 (55)

II.3 ERM generalization error

The generalization error of the ERM estimator is given again by eq. (54) with parameters

$$\sigma_{\hat{\mathbf{w}}} \equiv \lim_{d \to \infty} \mathbb{E}_{\mathbf{w}^{\star}, \mathbf{X}} \frac{1}{d} \| \hat{\mathbf{w}} \|_{2}^{2} = \lim_{d \to \infty} \mathbb{E}_{\mathbf{w}^{\star}, \mathbf{X}} \frac{1}{d} \| \hat{\mathbf{w}}^{\text{erm}} \|_{2}^{2} \equiv q ,$$

$$\sigma_{\mathbf{w}^{\star} \hat{\mathbf{w}}} \equiv \lim_{d \to \infty} \mathbb{E}_{\mathbf{w}^{\star}, \mathbf{X}} \frac{1}{d} \hat{\mathbf{w}}^{\mathsf{T}} \mathbf{w}^{\star} = \lim_{d \to \infty} \mathbb{E}_{\mathbf{w}^{\star}, \mathbf{X}} \frac{1}{d} (\hat{\mathbf{w}}^{\text{erm}})^{\mathsf{T}} \mathbf{w}^{\star} \equiv m .$$

where the parameters m, q are the asymptotic ERM overlaps solutions of eq. (11) and that finally lead to the ERM generalization error for classification:

$$e_{\rm g}^{\rm erm}(\alpha) = \frac{1}{\pi} {\rm acos}\left(\sqrt{\eta}\right), \qquad \text{with } \eta \equiv \frac{m^2}{\rho_{\rm w} \cdot q}.$$
 (56)

III Proofs of the ERM fixed points

III.1 Gordon's result and proofs

We consider in this section that the data have been generated by a teacher (18) with Gaussian weights

$$\mathbf{w}^{\star} \sim P_{\mathbf{w}^{\star}}(\mathbf{w}^{\star}) = \mathcal{N}_{\mathbf{w}^{\star}}(\mathbf{0}, \rho_{\mathbf{w}^{\star}}\mathbf{I}_{d}) \quad \text{with} \quad \rho_{\mathbf{w}^{\star}} \equiv \mathbb{E}\left[(w^{\star})^{2}\right] \,. \tag{57}$$

III.1.1 For real outputs - Regression with ℓ_2 regularization

In what follows, we prove a theorem that characterizes the asymptotic performance of empirical risk minimization

$$\hat{\mathbf{w}}_{\text{erm}} = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^{n} l\left(y_i, \frac{1}{\sqrt{d}} \mathbf{x}_i^{\mathsf{T}} \mathbf{w}\right) + \frac{\lambda \|\mathbf{w}\|^2}{2},\tag{58}$$

where $\{y_i\}_{1 \le i \le n}$ are general real-valued outputs (that are not necessarily binary), l(y, z) is a loss function that is convex with respect to z, and $\lambda > 0$ is the strength of the ℓ_2 regularization. Note that this setting is more general than the one considered in Thm. 2.2 in the main text, which focuses on binary outputs and loss functions in the form of $l(y, z) = \ell(yz)$ for some convex function $\ell(\cdot)$.

Theorem III.1 (Regression with ℓ_2 regularization). As $n, d \to \infty$ with $n/d = \alpha = \Theta(1)$, the overlap parameters m, q concentrate to

$$m \xrightarrow[d \to \infty]{} \sqrt{\rho_{\mathbf{w}^*}} \mu^*, \qquad \qquad q \xrightarrow[d \to \infty]{} (\mu^*)^2 + (\delta^*)^2, \qquad (59)$$

where the parameters μ^*, δ^* are the solutions of

$$(\mu^*, \delta^*) = \underset{\mu, \delta \ge 0}{\operatorname{arg\,min}} \sup_{\tau > 0} \left\{ \frac{\lambda(\mu^2 + \delta^2)}{2} - \frac{\delta^2}{2\tau} + \alpha \mathbb{E}_{g,s} \mathcal{M}_{\tau}[l(\varphi_{\text{out}^*}(\sqrt{\rho_{w^*}}s), .)](\mu s + \delta g) \right\}.$$
(60)

Here, $\mathcal{M}_{\tau}[l(,.)](x)$ is the Moreau-Yosida regularization defined in (28), and g, s are two iid standard normal random variables.

Proof. Since the teacher weight vector \mathbf{w}^* is independent of the input data matrix X, we can assume without loss of generality that

$$\mathbf{w}^{\star} = \sqrt{d}\rho_d \mathbf{e}_1,$$

where \mathbf{e}_1 is the first natural basis vector of \mathbb{R}^d , and $\rho_d = \|\mathbf{w}^\star\|/\sqrt{d}$. As $d \to \infty$, $\rho_d \to \sqrt{\rho_{\mathbf{w}^\star}}$. Accordingly, it will be convenient to split the data matrix into two parts:

$$\mathbf{X} = \begin{bmatrix} \mathbf{s} & \mathbf{B} \end{bmatrix},\tag{61}$$

where $\mathbf{s} \in \mathbb{R}^{n \times 1}$ and $\mathbf{B} \in \mathbb{R}^{n \times (d-1)}$ are two sub-matrices of iid standard normal entries. The weight vector \mathbf{w} in (58) can also be written as $\mathbf{w} = [\sqrt{d\mu}, \mathbf{v}^{\mathsf{T}}]^{\mathsf{T}}$, where $\mu \in \mathbb{R}$ denotes the

projection of **w** onto the direction spanned by the teacher weight vector \mathbf{w}^{\star} , and $\mathbf{v} \in \mathbb{R}^{d-1}$ is the projection of **w** onto the complement subspace. These representations serve to simplify the notations in our subsequent derivations. For example, we can now write the output as

$$y_i = \varphi_{\text{out}^{\star}}(\rho_d s_i),\tag{62}$$

where s_i is the *i*th entry of the Gaussian vector **s** in (61).

Let Φ_d denote the cost of the ERM in (58), normalized by d. Using our new representations introduced above, we have

$$\Phi_d = \min_{\mu, \mathbf{v}} \frac{1}{d} \sum_{i=1}^n l\left(y_i, \mu s_i + \frac{1}{\sqrt{d}} \mathbf{b}_i^\mathsf{T} \mathbf{v}\right) + \frac{\lambda (d\mu^2 + \|\mathbf{v}\|^2)}{2d},\tag{63}$$

where $\mathbf{b}_i^{\mathsf{T}}$ denotes the *i*th row of B. Since the loss function $l(y_i, z)$ is convex with respect to z, we can rewrite it as

$$l(y_i, z) = \sup_{q} \{ qz - l^*(y_i, q) \},$$
(64)

where $l^*(y_i, q) = \sup_z \{qz - l(y_i, z)\}$ is its convex conjugate. Substituting (64) into (63), we have

$$\Phi_d = \min_{\mu, \mathbf{v}} \sup_{\mathbf{q}} \left\{ \frac{\mu \mathbf{q}^\mathsf{T} \mathbf{s}}{d} + \frac{1}{d^{3/2}} \mathbf{q}^\mathsf{T} \mathbf{B} \mathbf{v} - \frac{1}{d} \sum_{i=1}^n l^*(y_i, q_i) + \frac{\lambda \left(d\mu^2 + \|\mathbf{v}\|^2 \right)}{2d} \right\}.$$
(65)

Now consider a new optimization problem

$$\widetilde{\Phi}_{d} = \min_{\mu, \mathbf{v}} \sup_{\mathbf{q}} \left\{ \frac{\mu \mathbf{q}^{\mathsf{T}} \mathbf{s}}{d} + \frac{\|\mathbf{q}\|}{\sqrt{d}} \frac{\mathbf{h}^{\mathsf{T}} \mathbf{v}}{d} + \frac{\|\mathbf{v}\|}{\sqrt{d}} \frac{\mathbf{g}^{\mathsf{T}} \mathbf{q}}{d} - \frac{1}{d} \sum_{i=1}^{n} l^{*}(y_{i}, q_{i}) + \frac{\lambda \left(d\mu^{2} + \|\mathbf{v}\|^{2} \right)}{2d} \right\},$$
(66)

where $h \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d-1})$ and $g \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ are two independent standard normal vectors. It follows from Gordon's minimax comparison inequality (see, *e.g.*, [?]) that

$$\mathbb{P}(|\Phi_d - c| \ge \epsilon) \le 2\mathbb{P}\left(\left|\widetilde{\Phi}_d - c\right| \ge \epsilon\right)$$
(67)

for any constants c and $\epsilon > 0$. This implies that $\widetilde{\Phi}_d$ serves as a surrogate of Φ_d . Specifically, if $\widetilde{\Phi}_d$ concentrates around some deterministic limit c as $d \to \infty$, so does Φ_d . In what follows, we proceed to solve the surrogate problem in (66). First, let $\delta = \|\mathbf{v}\|/\sqrt{d}$. It is easy to see that (66)

can be simplified as

$$\begin{split} \widetilde{\Phi}_{d} &= \min_{\mu,\delta \ge 0} \sup_{\mathbf{q}} \left\{ \frac{\mathbf{q}^{\mathsf{T}}(\mu \mathbf{s} + \delta \mathbf{g})}{d} - \delta \frac{\|\mathbf{q}\|}{\sqrt{d}} \frac{\|\mathbf{h}\|}{\sqrt{d}} - \frac{1}{d} \sum_{i=1}^{n} l^{*}(y_{i}, q_{i}) + \frac{\lambda(\mu^{2} + \delta^{2})}{2} \right\} \\ &\stackrel{(a)}{=} \min_{\mu,\delta \ge 0} \sup_{\tau > 0} \sup_{\mathbf{q}} \left\{ -\frac{\tau \|\mathbf{q}\|^{2}}{2d} - \frac{\delta^{2} \|\mathbf{h}\|^{2}}{2\tau d} + \frac{\mathbf{q}^{\mathsf{T}}(\mu \mathbf{s} + \delta \mathbf{g})}{d} - \frac{1}{d} \sum_{i=1}^{n} l^{*}(y_{i}, q_{i}) + \frac{\lambda(\mu^{2} + \delta^{2})}{2} \right\} \\ &= \min_{\mu,\delta \ge 0} \sup_{\tau > 0} \left\{ \frac{\lambda(\mu^{2} + \delta^{2})}{2} - \frac{\delta^{2} \|\mathbf{h}\|^{2}}{2\tau d} - \frac{\alpha}{n} \inf_{\mathbf{q}} \left[\frac{\tau \|\mathbf{q}\|^{2}}{2} - \mathbf{q}^{\mathsf{T}}(\mu \mathbf{s} + \delta \mathbf{g}) + \sum_{i=1}^{n} l^{*}(y_{i}, q_{i}) \right] \right\} \\ &\stackrel{(b)}{=} \min_{\mu,\delta \ge 0} \sup_{\tau > 0} \left\{ \frac{\lambda(\mu^{2} + \delta^{2})}{2} - \frac{\delta^{2} \|\mathbf{h}\|^{2}}{2\tau d} - \frac{\alpha}{n} \sum_{i=1}^{n} \mathcal{M}_{\tau}[l(y_{i}, .)](\mu s_{i} + \delta g_{i}) \right\}. \end{split}$$

In (*a*), we have introduced an auxiliary variable τ to rewrite $-\delta \frac{\|\mathbf{q}\|}{\sqrt{d}} \frac{\|\mathbf{h}\|}{\sqrt{d}}$ as

$$-\delta \frac{\|\mathbf{q}\|}{\sqrt{d}} \frac{\|\mathbf{h}\|}{\sqrt{d}} = \sup_{\tau>0} \left\{ -\frac{\tau \|\mathbf{q}\|^2}{2d} - \frac{\delta^2 \|\mathbf{h}\|^2}{2\tau d} \right\} \,,$$

and to get (b), we use the identity

$$\inf_{q} \left\{ \frac{\tau}{2} q^2 - qz + \ell^*(q) \right\} = -\inf_{x} \left\{ \frac{(z-x)^2}{2\tau} + \ell(x) \right\}$$

that holds for any z and for any convex function $\ell(x)$ and its conjugate $\ell^*(q)$. As $d \to \infty$, standard concentration arguments give us $\frac{\|\mathbf{h}\|^2}{d} \to 1$ and $\frac{1}{n} \sum_{i=1}^n \mathcal{M}_{\tau}[l(y_i, .)](\mu s_i + \delta g_i) \to \mathbb{E}_{g,s}\mathcal{M}_{\tau}[l(y, .)](\mu s + \delta g)$ locally uniformly over τ, μ and δ . Using (67) and recalling (62), we can then conclude that the normalized cost of the ERM Φ_d converges to the optimal value of the deterministic optimization problem in (60). Finally, since $\lambda > 0$, one can show that the cost function of (60) has a unique global minima at μ^* and δ^* . It follows that the empirical values of (μ, δ) also converge to their corresponding deterministic limits (μ^*, δ^*) .

III.1.2 For binary outputs - Classification with ℓ_2 regularization

In what follows, we specialize the previous theorem to the case of binary classification, with a convex loss function in the form of $l(y, z) = \ell(yz)$ for some function $\ell(\cdot)$.

Theorem III.2 (Thm. 2.2 in the main text. Gordon's min-max fixed point - Classification with ℓ_2 regularization). As $n, d \to \infty$ with $n/d = \alpha = \Theta(1)$, the overlap parameters m, q concentrate to

$$m \xrightarrow[d \to \infty]{} \sqrt{\rho_{\mathbf{w}^*}} \mu^*, \qquad \qquad q \xrightarrow[d \to \infty]{} (\mu^*)^2 + (\delta^*)^2, \qquad (68)$$

where parameters μ^*, δ^* are solutions of

$$(\mu^*, \delta^*) = \underset{\mu, \delta \ge 0}{\operatorname{arg\,min}} \sup_{\tau > 0} \left\{ \frac{\lambda(\mu^2 + \delta^2)}{2} - \frac{\delta^2}{2\tau} + \alpha \mathbb{E}_{g,s} \mathcal{M}_{\tau} [\delta g + \mu s \varphi_{\text{out}^*}(\sqrt{\rho_{w^*}}s)] \right\}, \quad (69)$$

and g, s are two iid standard normal random variables. The solutions $(\mu^*, \delta^*, \tau^*)$ of (69) can be reformulated as a set of fixed point equations

$$\mu^{*} = \frac{\alpha}{\lambda\tau^{*} + \alpha} \mathbb{E}[s \cdot \varphi_{\text{out}^{\star}}(\sqrt{\rho_{\text{w}^{\star}}}s) \cdot \mathcal{P}_{\tau^{*}}(\delta^{*}g + \mu^{*}s\varphi_{\text{out}^{\star}}(\sqrt{\rho_{\text{w}^{\star}}}s))],$$

$$\delta^{*} = \frac{\alpha}{\lambda\tau^{*} + \alpha - 1} \mathbb{E}[g \cdot \mathcal{P}_{\tau^{*}}(\delta^{*}g + \mu^{*}s\varphi_{\text{out}^{\star}}(\sqrt{\rho_{\text{w}^{\star}}}s))],$$

$$(\delta^{*})^{2} = \alpha \mathbb{E}[(\delta^{*}g + \mu^{*}s\varphi_{\text{out}^{\star}}(\sqrt{\rho_{\text{w}^{\star}}}s) - \mathcal{P}_{\tau^{*}}(\delta^{*}g + \mu^{*}s\varphi_{\text{out}^{\star}}(\sqrt{\rho_{\text{w}^{\star}}}s)))^{2}],$$

(70)

where \mathcal{M}_{τ} and \mathcal{P}_{τ} denote the Moreau-Yosida regularization and the proximal map of a convex loss function $(y, z) \mapsto \ell(yz)$:

$$\mathcal{M}_{\tau}(z) = \min_{x} \left\{ \ell(x) + \frac{(x-z)^2}{2\tau} \right\}, \qquad \mathcal{P}_{\tau}(z) = \arg\min_{x} \left\{ \ell(x) + \frac{(x-z)^2}{2\tau} \right\}.$$

Proof. We start by deriving (69) as a special case of (60). To that end, we note that

$$\mathcal{M}_{\tau}[l(y,.)](z) = \min_{x} \left\{ l(y;x) + \frac{(x-z)^2}{2\tau} \right\}$$
$$= \min_{x} \left\{ \ell(yx) + \frac{(x-z)^2}{2\tau} \right\}$$
$$= \min_{x} \left\{ \ell(x) + \frac{(x-yz)^2}{2\tau} \right\} = \mathcal{M}_{\tau}(yz),$$

where to reach the last equality we have used the fact that $y \in \{\pm 1\}$. Substituting this special form into (60) and recalling (62), we reach (69).

Finally, to obtain the fixed point equations (70), we simply take the partial derivatives of the cost function in (69) with respect to μ , δ , τ , and use the following well-known calculus rules for the Moreau-Yosida regularization [?]:

$$\frac{\partial \mathcal{M}_{\tau}(z)}{\partial z} = \frac{z - \mathcal{P}_{\tau}(z)}{\tau},\\ \frac{\partial \mathcal{M}_{\tau}(z)}{\partial \tau} = -\frac{(z - \mathcal{P}_{\tau}(z))^2}{2\tau^2}.$$

III.2 Replica's formulation

The replica computation presented in Sec. IV boils down to the characterization of the overlaps m, q in the high-dimensional limit $n, d \to \infty$ with $\alpha = \frac{n}{d} = \Theta(1)$, given by the solution of a set of, in the most general case, six fixed point equations over $m, q, Q, \hat{m}, \hat{q}, \hat{Q}$. Introducing the natural variables $\Sigma \equiv Q - q$, $\hat{\Sigma} \equiv \hat{Q} + \hat{q}$, $\eta \equiv \frac{m^2}{\rho_w \star q}$ and $\hat{\eta} \equiv \frac{\hat{m}^2}{\hat{q}}$, the set of fixed point equations

for arbitrary P_{w^*} , P_{out^*} , convex loss l(y, z) and regularizer r(w), is finally given by

$$m = \mathbb{E}_{\xi} \left[\mathcal{Z}_{w^{\star}} \left(\sqrt{\hat{\eta}} \xi, \hat{\eta} \right) f_{w^{\star}} \left(\sqrt{\hat{\eta}} \xi, \hat{\eta} \right) f_{w} \left(\hat{q}^{1/2} \xi, \hat{\Sigma} \right) \right],$$

$$q = \mathbb{E}_{\xi} \left[\mathcal{Z}_{w^{\star}} \left(\sqrt{\hat{\eta}} \xi, \hat{\eta} \right) f_{w} \left(\hat{q}^{1/2} \xi, \hat{\Sigma} \right)^{2} \right],$$

$$\Sigma = \mathbb{E}_{\xi} \left[\mathcal{Z}_{w^{\star}} \left(\sqrt{\hat{\eta}} \xi, \hat{\eta} \right) \partial_{\gamma} f_{w} \left(\hat{q}^{1/2} \xi, \hat{\Sigma} \right) \right],$$

$$\hat{m} = \alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{out^{\star}} (.) \cdot f_{out^{\star}} (y, \sqrt{\rho_{w^{\star}} \eta} \xi, \rho_{w^{\star}} (1 - \eta)) f_{out} \left(y, q^{1/2} \xi, \Sigma \right) \right],$$

$$\hat{q} = \alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{out^{\star}} (y, \sqrt{\rho_{w^{\star}} \eta} \xi, \rho_{w^{\star}} (1 - \eta)) f_{out} \left(y, q^{1/2} \xi, \Sigma \right)^{2} \right],$$

$$\hat{\Sigma} = -\alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{out^{\star}} (y, \sqrt{\rho_{w^{\star}} \eta} \xi, \rho_{w^{\star}} (1 - \eta)) \partial_{\omega} f_{out} \left(y, q^{1/2} \xi, \Sigma \right) \right].$$
(71)

The above equations depend on the Bayes-optimal partition functions \mathcal{Z}_{w^*} , \mathcal{Z}_{out^*} defined in eq. (24), the updates f_{w^*} , f_{out^*} in eq. (25) and the ERM updates f_w , f_{out} eq. (29).

III.3 Equivalence Gordon-Replica's formulation - ℓ_2 regularization and Gaussian weights

III.3.1 Replica's formulation for ℓ_2 regularization

The proximal for the ℓ_2 penalty with strength λ can be computed explicitly in eq. (47) and the corresponding denoising function is simply given by $f_{\rm w}^{\ell_2,\lambda}(\gamma,\Lambda) = \frac{\gamma}{\lambda+\Lambda}$. Therefore, for a Gaussian teacher (57) already considered in Thm. (70) with second moment $\rho_{\rm w^{\star}}$, using the denoising function (41), the fixed point equations over m, q, Σ can be computed analytically and lead to

$$m = \frac{\rho_{\mathbf{w}^{\star}}\hat{m}}{\lambda + \hat{\Sigma}}, \qquad q = \frac{\rho_{\mathbf{w}^{\star}}\hat{m}^2 + \hat{q}}{(\lambda + \hat{\Sigma})^2}, \qquad \Sigma = \frac{1}{\lambda + \hat{\Sigma}}.$$
(72)

Hence, removing the *hat* variables in eqs. (71), the set of fixed point equations can be rewritten in a more compact way leading to the Corollary. 2.3 that we recall here:

Corollary III.3 (Corollary. 2.3 in the main text. Equivalence Gordon-Replicas). The set of fixed point equations (70) in Thm. III.2 that govern the asymptotic behaviour of the overlaps m and q is equivalent to the following set of equations, obtained from the heuristic replica computation:

$$m = \alpha \Sigma \rho_{\mathbf{w}^{\star}} \cdot \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^{\star}} \left(. \right) \cdot f_{\text{out}^{\star}} \left(y, \sqrt{\rho_{\mathbf{w}^{\star}} \eta} \xi, \rho_{\mathbf{w}^{\star}} \left(1 - \eta \right) \right) \cdot f_{\text{out}} \left(y, q^{1/2} \xi, \Sigma \right) \right]$$

$$q = m^{2} / \rho_{\mathbf{w}^{\star}} + \alpha \Sigma^{2} \cdot \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^{\star}} \left(y, \sqrt{\rho_{\mathbf{w}^{\star}} \eta} \xi, \rho_{\mathbf{w}^{\star}} \left(1 - \eta \right) \right) \cdot f_{\text{out}} \left(y, q^{1/2} \xi, \Sigma \right)^{2} \right]$$

$$\Sigma = \left(\lambda - \alpha \cdot \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^{\star}} \left(y, \sqrt{\rho_{\mathbf{w}^{\star}} \eta} \xi, \rho_{\mathbf{w}^{\star}} \left(1 - \eta \right) \right) \cdot \partial_{\omega} f_{\text{out}} \left(y, q^{1/2} \xi, \Sigma \right) \right] \right)^{-1}$$

$$T_{\mathbf{w}^{2}} = \frac{1}{2} \left[\left(\mathcal{Z}_{\mathbf{w}^{\star}} \left(y, \sqrt{\rho_{\mathbf{w}^{\star}} \eta} \xi, \rho_{\mathbf{w}^{\star}} \left(1 - \eta \right) \right) \cdot \partial_{\omega} f_{\text{out}} \left(y, q^{1/2} \xi, \Sigma \right) \right] \right]$$

$$T_{\mathbf{w}^{\star}} = \frac{1}{2} \left[\left(\mathcal{Z}_{\mathbf{w}^{\star}} \left(y, \sqrt{\rho_{\mathbf{w}^{\star}} \eta} \xi, \rho_{\mathbf{w}^{\star}} \left(1 - \eta \right) \right) \cdot \partial_{\omega} f_{\mathbf{w}^{\star}} \left(y, q^{1/2} \xi, \Sigma \right) \right] \right]$$

with $\eta \equiv \frac{m^2}{\rho_{w^*}q}$, $\xi \sim \mathcal{N}(0,1)$ and \mathbb{E}_y the continuous or discrete sum over all possible values y according to P_{out^*} .

Proof of Corollary. III.3(Corollary. 2.3). For the sake of clarity, we use the abusive notation $\mathcal{P}_V(y, \omega) = \mathcal{P}_V[l(y, .)](\omega)$, and we remove the *.

Dictionary We first map the Gordon's parameters (μ, δ, τ) in eq. (70) to (m, q, Σ) in eq. (73):

$$\sqrt{\rho_{\mathbf{w}^{\star}}}\mu \leftrightarrow m$$
, $\mu^2 + \delta^2 \leftrightarrow q$, $\tau \leftrightarrow \Sigma$.

so that

$$\eta = \frac{m^2}{\rho_{\mathbf{w}^\star}q} = \frac{\mu^2}{\mu^2 + \delta^2} \,, \qquad \qquad 1 - \eta = \frac{\delta^2}{\mu^2 + \delta^2} \,,$$

From eq. (24), we can rewrite the channel partition function $\mathcal{Z}_{\mathrm{out}^\star}$ and its derivative

$$\mathcal{Z}_{\text{out}^{\star}}(y,\omega,V) = \mathbb{E}_{z} \left[P_{\text{out}^{\star}}\left(y|\sqrt{V}z+\omega \right) \right],$$

$$\partial_{\omega}\mathcal{Z}_{\text{out}^{\star}}(y,\omega,V) = \frac{1}{\sqrt{V}} \mathbb{E}_{z} \left[zP_{\text{out}^{\star}}\left(y|\sqrt{V}z+\omega \right) \right],$$
(74)

.

where z denotes a standard normal random variable.

Equation over m Let us start with the equation over m in eq. (73):

$$= \frac{\alpha}{\lambda \tau + \alpha} \mathbb{E}_{s,g} \left[g \cdot \mathcal{P}_{\tau} \left(\varphi_{\text{out}^{\star}} \left(\sqrt{\rho_{\text{w}^{\star}}} s \right), \delta g + \mu s \right) \right]$$
$$= \frac{\alpha}{\lambda \tau + \alpha} \mathbb{E}_{s,g} \left[s \cdot \varphi_{\text{out}^{\star}} \left(\sqrt{\rho_{\text{w}^{\star}}} s \right) \left(\mathcal{P}_{\tau} \left(\delta g + \mu s \right) \varphi_{\text{out}^{\star}} \left(\sqrt{\rho_{\text{w}^{\star}}} s \right) \right) \right],$$

(Second fixed point equation)

where we used the fact that $P_{\mathrm{out}^\star}\left(y|z\right)=\delta(y-\varphi_{\mathrm{out}^\star}(z)),$ the change of variables

$$\begin{cases} s = \frac{\mu\xi + \delta z}{\sqrt{\mu^2 + \delta^2}} \\ g = \frac{\delta\xi - \mu z}{\sqrt{\mu^2 + \delta^2}} \end{cases} \Leftrightarrow \begin{cases} \xi = \frac{\delta g + \mu s}{\sqrt{\mu^2 + \delta^2}} \\ z = \frac{\delta s - \mu g}{\sqrt{\mu^2 + \delta^2}} \end{cases}, \tag{75}$$

and finally in the last equality the definition of the second fixed point equation in eqs. (70):

$$\delta = \alpha \frac{\mathbb{E}_{s,g} \left[g \cdot \mathcal{P}_{\tau} \left(\varphi_{\text{out}^{\star}} \left(\sqrt{\rho_{\text{w}^{\star}} s} \right), \delta g + \mu s \right) \right]}{\lambda \tau + \alpha - 1} \,. \tag{76}$$

Equation over q Let us now compute the equation over q in eq. (73):

$$q - m^{2}/\rho_{w^{\star}} = \Sigma^{2} \alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^{\star}} \left(y, \sqrt{\rho_{w^{\star}} \eta} \xi, \rho_{w^{\star}} \left(1 - \eta \right) \right) f_{\text{out}} \left(y, q^{1/2} \xi, \Sigma \right)^{2} \right]$$

$$= \Sigma^{2} \alpha \mathbb{E}_{y,\xi,z} \left[P_{\text{out}^{\star}} \left(y | \sqrt{\rho_{w^{\star}}} \left(\sqrt{1 - \eta} z + \sqrt{\eta} \xi \right) \right) \frac{1}{\Sigma^{2}} \left(p_{\Sigma} \left(y, \sqrt{q} \xi \right) - \sqrt{q} \xi \right)^{2} \right]$$

$$(Using eq. (74))$$

$$\Leftrightarrow \delta^{2} = \alpha \mathbb{E}_{y,\xi,z} \left[P_{\text{out}^{\star}} \left(y | \sqrt{\rho_{w^{\star}}} \frac{\delta z + \mu \xi}{\sqrt{\mu^{2} + \delta^{2}}} \right) \left(p_{\tau} \left(y, \sqrt{\mu^{2} + \delta^{2}} \xi \right) - \sqrt{\mu^{2} + \delta^{2}} \xi \right)^{2} \right]$$

$$(Dictionary)$$

$$= \alpha \mathbb{E}_{\xi,z} \left[\left(p_{\tau} \left(\varphi_{\text{out}^{\star}} \left(\sqrt{\rho_{w^{\star}}} \frac{\delta z + \mu \xi}{\sqrt{\mu^{2} + \delta^{2}}} \right), \sqrt{\mu^{2} + \delta^{2}} \xi \right) - \sqrt{\mu^{2} + \delta^{2}} \xi \right)^{2} \right]$$

$$(Integration over y)$$

$$= \alpha \mathbb{E}_{g,s} \left[\left(p_{\tau} \left(\varphi_{\text{out}^{\star}} \left(\sqrt{\rho_{\text{w}^{\star}}} s \right), \delta g + \mu s \right) - \left(\delta g + \mu s \right) \right)^2 \right]$$
(Change of variables $(\xi, z) \to (g, s)$)

Equation over Σ Let us conclude with the equation over Σ in eq. (73) that we encountered in eq. (76). Let us first compute

$$\alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^{\star}} \left(y, \sqrt{\rho_{\text{w}^{\star}} \eta} \xi, \rho_{\text{w}^{\star}} \left(1 - \eta \right) \right) \partial_{\omega} f_{\text{out}} \left(y, q^{1/2} \xi, \Sigma \right) \right]$$

$$= \alpha \mathbb{E}_{y,\xi,z} \left[P_{\text{out}^{\star}} \left(y | \sqrt{\rho_{\text{w}^{\star}}} \left(\sqrt{1 - \eta} z + \sqrt{\eta} \xi \right) \right) \frac{1}{\Sigma} \left(\partial_{\omega} p_{\Sigma} \left(y, \sqrt{q} \xi \right) - 1 \right) \right] \quad \text{(Using eq. (74))}$$

$$= \frac{\alpha}{\mathbb{E}_{y,\xi,z}} \left[P_{\text{out}^{\star}} \left(y | \sqrt{\rho_{\text{w}^{\star}}} \frac{\delta z + \mu \xi}{\sqrt{\rho_{\text{w}^{\star}}}} \right) \left(\partial_{\omega} \mathcal{P}_{\tau} \left(y, \sqrt{\mu^{2} + \delta^{2}} \xi \right) - 1 \right) \right] \quad \text{(Dictionary)}$$

$$= \frac{\alpha}{\tau} \mathbb{E}_{y,\xi,z} \left[P_{\text{out}^{\star}} \left(y | \sqrt{\rho_{\text{w}^{\star}}} \frac{\delta x + \mu_{\xi}}{\sqrt{\mu^2 + \delta^2}} \right) \left(\partial_{\omega} \mathcal{P}_{\tau} \left(y, \sqrt{\mu^2 + \delta^2 \xi} \right) - 1 \right) \right] \qquad \text{(Dictionary)}$$

$$= \frac{\alpha}{\tau} \mathbb{E}_{\xi,z} \left[\partial_{\omega} \mathcal{P}_{\tau} \left(\varphi_{\text{out}^{\star}} \left(\sqrt{\rho_{\text{w}^{\star}}} \frac{\delta z + \mu \xi}{\sqrt{\mu^2 + \delta^2}} \right), \sqrt{\mu^2 + \delta^2} \xi \right) \right] - \frac{\alpha}{\tau} \qquad \text{(Integration over } y\text{)}$$

$$= \frac{1}{\tau} \alpha \left(\mathbb{E}_{g,s} \left[\partial_{\omega} \mathcal{P}_{\tau} \left(\varphi_{\text{out}^{\star}} \left(\sqrt{\rho_{\text{w}^{\star}}} s \right), \delta g + \mu s \right) \right] - 1 \right) \quad \text{(Change of variables } (\xi, z) \to (g, s) \text{)}$$

therefore, the last equation over Σ in eq. (73) reads

$$\begin{split} \Sigma &= \left(\lambda - \alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^{\star}} \left(y, \sqrt{\rho_{\text{w}^{\star}} \eta} \xi, \rho_{\text{w}^{\star}} \left(1 - \eta \right) \right) \partial_{\omega} f_{\text{out}} \left(y, q^{1/2} \xi, \Sigma \right) \right] \right)^{-1} \\ &\Leftrightarrow \\ \tau &= \left(\lambda - \frac{1}{\tau} \alpha \left(\mathbb{E}_{g,s} \left[\partial_{\omega} \mathcal{P}_{\tau} \left(\varphi_{\text{out}^{\star}} \left(\sqrt{\rho_{\text{w}^{\star}} s} \right), \delta g + \mu s \right) \right] - 1 \right) \right)^{-1} \\ &\Leftrightarrow \\ \alpha \mathbb{E}_{g,s} \left[\partial_{\omega} \mathcal{P}_{\tau} \left(\varphi_{\text{out}^{\star}} \left(\sqrt{\rho_{\text{w}^{\star}} s} \right), \delta g + \mu s \right) \right] = \tau \lambda + \alpha - 1 \,. \end{split}$$

Noting that

where we used the Stein's lemma in the last equality, we finally obtain

$$\alpha \mathbb{E}_{g,s} \left[\partial_{\omega} \mathcal{P}_{\tau} \left(\varphi_{\text{out}^{\star}} \left(\sqrt{\rho_{\text{w}^{\star}}} s \right), \delta g + \mu s \right) \right] = \tau \lambda + \alpha - 1$$

$$\Leftrightarrow \delta = \frac{\alpha}{\tau \lambda + \alpha - 1} \mathbb{E}_{g,s} \left[g \cdot \mathcal{P}_{\tau} \left(\varphi_{\text{out}^{\star}} \left(\sqrt{\rho_{\text{w}^{\star}}} s \right), \delta g + \mu s \right) \right] \,.$$

Gauge transformation We still remain to prove that

$$\mathbb{E}_{s,g}\left[g \cdot \mathcal{P}_{\tau}\left(\varphi_{\text{out}^{\star}}\left(\sqrt{\rho_{\text{w}^{\star}}s}\right), \delta g + \mu s\right)\right] = \mathbb{E}_{s,g}\left[g \cdot \mathcal{P}_{\tau}\left(\delta g + \mu s \varphi_{\text{out}^{\star}}\left(\sqrt{\rho_{\text{w}^{\star}}s}\right)\right)\right]$$
$$\mathbb{E}_{s,g}\left[s \cdot \mathcal{P}_{\tau}\left(\varphi_{\text{out}^{\star}}\left(\sqrt{\rho_{\text{w}^{\star}}s}\right), \delta g + \mu s\right)\right] = \mathbb{E}_{s,g}\left[s \cdot \mathcal{P}_{\tau}\left(\delta g + \mu s \varphi_{\text{out}^{\star}}\left(\sqrt{\rho_{\text{w}^{\star}}s}\right)\right)\right]$$
$$\mathbb{E}_{g,s}\left[\left(p_{\tau}\left(\varphi_{\text{out}^{\star}}\left(\sqrt{\rho_{\text{w}^{\star}}s}\right), \delta g + \mu s\right) - \left(\delta g + \mu s\right)\right)^{2}\right] = \mathbb{E}_{g,s}\left[\left(\left(p_{\tau} - 1\right)\left(\delta g + \mu s \varphi_{\text{out}^{\star}}\left(\sqrt{\rho_{\text{w}^{\star}}s}\right)\right)\right)^{2}\right]$$
(77)

As $\varphi_{\text{out}^{\star}}(\sqrt{\rho_{\text{w}^{\star}}}s) = \pm 1$, we can transform $s \to s\varphi_{\text{out}^{\star}}(\sqrt{\rho_{\text{w}^{\star}}}s) = \tilde{s}$. It does not change the distribution of the random variable \tilde{s} that is still a normal random variable. Finally denoting $\mathcal{P}_{\tau}(1, \delta g + \mu s \varphi_{\text{out}^{\star}}(\sqrt{\rho_{\text{w}^{\star}}}s)) = \mathcal{P}_{\tau}(\delta g + \mu s \varphi_{\text{out}^{\star}}(\sqrt{\rho_{\text{w}^{\star}}}s))$, we obtain the equivalence with eq. (70), which concludes the proof.

IV Replica computation for Bayes-optimal and ERM estimations

In this section, we present the statistical physics framework and the replica computation leading to the general set of fixed point equations (11) and to the Bayes-optimal fixed point equations (13).

IV.1 Statistical inference and free entropy

As stressed in Sec. I, both ERM and Bayes-optimal estimations can be analyzed in a unified framework that consists in studying the joint distribution $\mathbb{P}(\mathbf{y}, X)$ in the following posterior distribution

$$\mathbb{P}(\mathbf{w}|\mathbf{y}, \mathbf{X}) = \frac{\mathbb{P}(\mathbf{y}|\mathbf{w}, \mathbf{X}) \mathbb{P}(\mathbf{w})}{\mathbb{P}(\mathbf{y}, \mathbf{X})},$$
(78)

known as the so-called *partition function* in the physics literature. It is the generating function of many useful statistical quantities and is defined by

$$\mathcal{Z}(\mathbf{y}, \mathbf{X}) \equiv P(\mathbf{y}, \mathbf{X}) = \int_{\mathbb{R}^d} \mathrm{d}\mathbf{w} P_{\mathrm{out}}(\mathbf{y} | \mathbf{w}, \mathbf{X}) P_{\mathrm{w}}(\mathbf{w})$$

=
$$\int_{\mathbb{R}^n} \mathrm{d}\mathbf{z} P_{\mathrm{out}}(\mathbf{y} | \mathbf{z}) \int_{\mathbb{R}^d} \mathrm{d}\mathbf{w} P_{\mathrm{w}}(\mathbf{w}) \,\delta\left(\mathbf{z} - \frac{1}{\sqrt{d}} \mathbf{X} \mathbf{w}\right) \,,$$
(79)

where we introduced the variable $\mathbf{z} = \frac{1}{\sqrt{d}} X \mathbf{w}$. However in the considered high-dimensional regime $(d \to \infty, n \to \infty, \alpha = \Theta(1))$, we are interested instead in the *averaged* (over instances of input data X and teacher weights \mathbf{w}^* or equivalently over the output labels \mathbf{y}) free entropy Φ defined as

$$\Phi(\alpha) \equiv \mathbb{E}_{\mathbf{y}, \mathbf{X}} \left[\lim_{d \to \infty} \frac{1}{d} \log \mathcal{Z}(\mathbf{y}, \mathbf{X}) \right].$$
(80)

The replica method is an heuristic method of statistical mechanics that allows to compute the above average over the random dataset $\{y, X\}$. We show in the next section the classical computation for the Generalized Linear Model hypothesis class and iid data X.

IV.2 Replica computation

IV.2.1 Derivation

We present here the replica computation of the averaged free entropy $\Phi(\alpha)$ in eq. (80) for general prior distributions $P_{\rm w}$, $P_{\rm w^{\star}}$ and channel distributions $P_{\rm out}$, $P_{\rm out^{\star}}$, so that the computation remain valid for both Bayes-optimal and ERM estimation (with any convex loss l and regularizer r).

Replica trick The average in eq. (80) is intractable in general, and the computation relies on the so called *replica trick* that consists in applying the identity

$$\mathbb{E}_{\mathbf{y},\mathbf{X}}\left[\lim_{d\to\infty}\frac{1}{d}\log\mathcal{Z}\left(\mathbf{y},\mathbf{X}\right)\right] = \lim_{r\to0}\left[\lim_{d\to\infty}\frac{1}{d}\frac{\partial\log\mathbb{E}_{\mathbf{y},\mathbf{X}}\left[\mathcal{Z}\left(\mathbf{y},\mathbf{X}\right)^{r}\right]}{\partial r}\right].$$
(81)

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

Thus the replicated partition function in eq. (81) can be written as

$$\mathbb{E}_{\mathbf{y},\mathbf{X}}\left[\mathcal{Z}\left(\mathbf{y},\mathbf{X}\right)^{r}\right] = \mathbb{E}_{\mathbf{w}^{\star},\mathbf{X}}\left[\prod_{a=1}^{r}\int_{\mathbb{R}^{n}} \mathrm{d}\mathbf{z}^{a}P_{\mathrm{out}^{a}}\left(\mathbf{y}|\mathbf{z}^{a}\right)\int_{\mathbb{R}^{d}}\mathrm{d}\mathbf{w}^{a}P_{\mathrm{w}^{a}}\left(\mathbf{w}^{a}\right)\delta\left(\mathbf{z}^{a}-\frac{1}{\sqrt{d}}\mathbf{X}\mathbf{w}^{a}\right)\right]$$

$$= \mathbb{E}_{\mathbf{X}}\int_{\mathbb{R}^{n}}\mathrm{d}\mathbf{y}\int_{\mathbb{R}^{n}}\mathrm{d}\mathbf{z}^{\star}P_{\mathrm{out}^{\star}}\left(\mathbf{y}|\mathbf{z}^{\star}\right)\int_{\mathbb{R}^{d}}\mathrm{d}\mathbf{w}^{\star}P_{\mathrm{w}^{\star}}\left(\mathbf{w}^{\star}\right)\delta\left(\mathbf{z}^{\star}-\frac{1}{\sqrt{d}}\mathbf{X}\mathbf{w}^{\star}\right)$$

$$\times\left[\prod_{a=1}^{r}\int_{\mathbb{R}^{n}}\mathrm{d}\mathbf{z}^{a}P_{\mathrm{out}^{a}}\left(\mathbf{y}|\mathbf{z}^{a}\right)\int_{\mathbb{R}^{d}}\mathrm{d}\mathbf{w}^{a}P_{\mathrm{w}^{a}}\left(\mathbf{w}^{a}\right)\delta\left(\mathbf{z}^{a}-\frac{1}{\sqrt{d}}\mathbf{X}\mathbf{w}^{a}\right)\right]$$

$$= \mathbb{E}_{\mathbf{X}}\int_{\mathbb{R}^{n}}\mathrm{d}\mathbf{y}\prod_{a=0}^{r}\int_{\mathbb{R}^{n}}\mathrm{d}\mathbf{z}^{a}P_{\mathrm{out}^{a}}\left(\mathbf{y}|\mathbf{z}^{a}\right)\int_{\mathbb{R}^{d}}\mathrm{d}\mathbf{w}^{a}P_{\mathrm{w}^{a}}\left(\mathbf{w}^{a}\right)\delta\left(\mathbf{z}^{a}-\frac{1}{\sqrt{d}}\mathbf{X}\mathbf{w}^{a}\right)$$

$$(82)$$

with the decoupled channel $P_{\text{out}}(\mathbf{y}|\mathbf{z}) = \prod_{\mu=1}^{n} P_{\text{out}}(y_{\mu}|z_{\mu})$. Note that the average over \mathbf{y} is equivalent to the one over the ground truth vector \mathbf{w}^{\star} , which can be considered as a new replica \mathbf{w}^{0} with index a = 0 leading to a total of r + 1 replicas.

We suppose that inputs are drawn from an iid distribution, for example a Gaussian $\mathcal{N}(0, 1)$. More precisely, for $i, j \in [1 : d]$, $\mu, \nu \in [1 : n]$, $\mathbb{E}_{\mathbf{X}} \left[x_i^{(\mu)} x_j^{(\nu)} \right] = \delta_{\mu\nu} \delta_{ij}$. Hence $z_{\mu}^a = \frac{1}{\sqrt{d}} \sum_{i=1}^d x_i^{(\mu)} w_i^a$ is the sum of iid random variables. The central limit theorem insures that $z_{\mu}^a \sim \mathcal{N} \left(\mathbb{E}_{\mathbf{X}}[z_{\mu}^a], \mathbb{E}_{\mathbf{X}}[z_{\mu}^a z_{\mu}^b] \right)$, with the two first moments given by:

$$\begin{cases} \mathbb{E}_{\mathbf{X}}[z_{\mu}^{a}] = \frac{1}{\sqrt{d}} \sum_{i=1}^{d} \mathbb{E}_{\mathbf{X}} \left[x_{i}^{(\mu)} \right] w_{i}^{a} = 0 \\ \mathbb{E}_{\mathbf{X}}[z_{\mu}^{a} z_{\mu}^{b}] = \frac{1}{d} \sum_{ij} \mathbb{E}_{\mathbf{X}} \left[x_{i}^{(\mu)} x_{j}^{(\mu)} \right] w_{i}^{a} w_{j}^{b} = \frac{1}{d} \sum_{ij} \delta_{ij} w_{i}^{a} w_{j}^{b} = \frac{1}{d} \mathbf{w}^{a} \cdot \mathbf{w}^{b} . \end{cases}$$
(83)

In the following we introduce the symmetric *overlap* matrix $Q(\{\mathbf{w}^a\}) \equiv \left(\frac{1}{d}\mathbf{w}^a \cdot \mathbf{w}^b\right)_{a,b=0..r}$. Let us define $\tilde{\mathbf{z}}_{\mu} \equiv (z^a_{\mu})_{a=0..r}$ and $\tilde{\mathbf{w}}_i \equiv (w^a_i)_{a=0..r}$. The vector $\tilde{\mathbf{z}}_{\mu}$ follows a multivariate Gaussian distribution $\tilde{\mathbf{z}}_{\mu} \sim P_{\tilde{\mathbf{z}}}(\tilde{\mathbf{z}}; Q) = \mathcal{N}_{\tilde{\mathbf{z}}}(\mathbf{0}_{r+1}, Q)$ and as $P_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}}) = \prod_{a=0}^{r} P_{\mathbf{w}}(\tilde{w}^{a})$ it follows

$$\begin{split} \mathbb{E}_{\mathbf{y},\mathbf{X}}\left[\mathcal{Z}\left(\mathbf{y},\mathbf{X}\right)^{r}\right] &= \mathbb{E}_{\mathbf{X}} \int_{\mathbb{R}^{n}} \mathrm{d}\mathbf{y} \prod_{a=0}^{r} \int_{\mathbb{R}^{n}} \mathrm{d}\mathbf{z}^{a} P_{\mathrm{out}^{a}}\left(\mathbf{y}|\mathbf{z}^{a}\right) \int_{\mathbb{R}^{d}} \mathrm{d}\mathbf{w}^{a} P_{\mathrm{w}^{a}}\left(\mathbf{w}^{a}\right) \delta\left(\mathbf{z}^{a} - \frac{1}{\sqrt{d}}\mathbf{X}\mathbf{w}^{a}\right) \\ &= \left[\int_{\mathbb{R}} \mathrm{d}y \int_{\mathbb{R}^{r+1}} \mathrm{d}\tilde{\mathbf{z}} P_{\mathrm{out}}\left(y|\tilde{\mathbf{z}}\right) P_{\tilde{\mathbf{z}}}(\tilde{\mathbf{z}};Q(\tilde{\mathbf{w}}))\right]^{n} \left[\int_{\mathbb{R}^{r+1}} \mathrm{d}\tilde{\mathbf{w}} P_{\tilde{\mathbf{w}}}\left(\tilde{\mathbf{w}}\right)\right]^{d} \,, \end{split}$$

because the channel and the prior distributions factorize. Introducing the change of variable and the Fourier representation of the δ -Dirac function, which involves a new ad-hoc parameter \hat{Q} :

$$\begin{split} 1 &= \int_{\mathbb{R}^{r+1\times r+1}} \mathrm{d}Q \prod_{a \leq b} \delta\left(dQ_{ab} - \sum_{i=1}^{d} w_i^a w_i^b \right) \\ &\propto \int_{\mathbb{R}^{r+1\times r+1}} \mathrm{d}Q \int_{\mathbb{R}^{r+1\times r+1}} \mathrm{d}\hat{Q} \exp\left(-d\mathrm{Tr}\left[Q\hat{Q} \right] \right) \exp\left(\frac{1}{2} \sum_{i=1}^{d} \tilde{\mathbf{w}}_i^{\mathsf{T}} \hat{Q} \tilde{\mathbf{w}}_i \right) \,, \end{split}$$

the replicated partition function becomes an integral over the symmetric matrices $Q \in \mathbb{R}^{r+1 \times r+1}$ and $\hat{Q} \in \mathbb{R}^{r+1 \times r+1}$, that can be evaluated using a Laplace method in the $d \to \infty$ limit,

$$\mathbb{E}_{\mathbf{y},\mathbf{X}}\left[\mathcal{Z}\left(\mathbf{y},\mathbf{X}\right)^{r}\right] = \int_{\mathbb{R}^{r+1\times r+1}} \mathrm{d}Q \int_{\mathbb{R}^{r+1\times r+1}} \mathrm{d}\hat{Q}e^{d\Phi^{(r)}(Q,\hat{Q})}$$
(84)

$$\underset{d \to \infty}{\simeq} \exp\left(d \cdot \mathbf{extr}_{Q,\hat{Q}}\left\{\Phi^{(r)}(Q,\hat{Q})\right\}\right),\tag{85}$$

where we defined

$$\begin{cases} \Phi^{(r)}(Q,\hat{Q}) = -\operatorname{Tr}\left[Q\hat{Q}\right] + \log \Psi_{w}^{(r)}(\hat{Q}) + \alpha \log \Psi_{out}^{(r)}(Q) \\ \Psi_{w}^{(r)}(\hat{Q}) = \int_{\mathbb{R}^{r+1}} d\mathbf{\tilde{w}} P_{\tilde{w}}(\mathbf{\tilde{w}}) e^{\frac{1}{2}\mathbf{\tilde{w}}^{\intercal}\hat{Q}\mathbf{\tilde{w}}} \\ \Psi_{out}^{(r)}(Q) = \int dy \int_{\mathbb{R}^{r+1}} d\mathbf{\tilde{z}} P_{\tilde{z}}(\mathbf{\tilde{z}};Q) P_{out}(y|\mathbf{\tilde{z}}), \end{cases}$$

$$(86)$$

and $P_{\tilde{z}}(\tilde{\mathbf{z}};Q) = \frac{e^{-\frac{1}{2}\tilde{\mathbf{z}}^{T}Q^{-1}\tilde{\mathbf{z}}}}{\det(2\pi Q)^{1/2}}.$

Finally switching the two limits $r \to 0$ and $d \to \infty$, the quenched free entropy Φ simplifies as a saddle point equation

$$\Phi(\alpha) = \mathbf{extr}_{Q,\hat{Q}} \left\{ \lim_{r \to 0} \frac{\partial \Phi^{(r)}(Q,\hat{Q})}{\partial r} \right\},\tag{87}$$

over symmetric matrices $Q \in \mathbb{R}^{r+1 \times r+1}$ and $\hat{Q} \in \mathbb{R}^{r+1 \times r+1}$. In the following we will assume a simple ansatz for these matrices in order to first obtain an analytic expression in r before taking the derivative with respect to r.

RS free entropy Let's compute the functional $\Phi^{(r)}(Q, \hat{Q})$ appearing in the free entropy eq. (87) in the simplest ansatz: the Replica Symmetric ansatz. This later assumes that all replica remain equivalent with a common overlap $q = \frac{1}{d} \mathbf{w}^a \cdot \mathbf{w}^b$ for $a \neq b$, a norm $Q = \frac{1}{d} ||\mathbf{w}^a||_2^2$, and an overlap with the ground truth $m = \frac{1}{d} \mathbf{w}^a \cdot \mathbf{w}^*$, leading to the following expressions of the replica symmetric matrices $Q_{rs} \in \mathbb{R}^{r+1 \times r+1}$ and $\hat{Q}_{rs} \in \mathbb{R}^{r+1 \times r+1}$:

$$Q_{\rm rs} = \begin{pmatrix} Q^0 & m & \dots & m \\ m & Q & \dots & \dots \\ \dots & \dots & \dots & q \\ m & \dots & q & Q \end{pmatrix} \quad \text{and} \quad \hat{Q}_{\rm rs} = \begin{pmatrix} \hat{Q}^0 & \hat{m} & \dots & \hat{m} \\ \hat{m} & -\frac{1}{2}\hat{Q} & \dots & \dots \\ \dots & \dots & \dots & \hat{q} \\ \hat{m} & \dots & \hat{q} & -\frac{1}{2}\hat{Q} \end{pmatrix}, \quad (88)$$

with $Q^0 = \rho_{w^*} = \frac{1}{d} \|\mathbf{w}^*\|_2^2$. Let's compute separately the terms involved in the functional $\Phi^{(r)}(Q, \hat{Q})$ eq. (86) with this ansatz: the first is a trace term, the second a term $\Psi_w^{(r)}$ depending on the prior distributions P_{w} , P_{w^*} and finally the third a term $\Psi_{out}^{(r)}$ that depends on the channel distributions P_{out^*}, P_{out} .

Trace term The trace term can be easily computed and takes the following form:

$$\operatorname{Tr}\left(Q\hat{Q}\right)\Big|_{\mathrm{rs}} = Q^{0}\hat{Q}^{0} + rm\hat{m} - \frac{1}{2}rQ\hat{Q} + \frac{r(r-1)}{2}q\hat{q}.$$
(89)

Prior integral Evaluated at the RS fixed point, and using a Gaussian identity also known as a Hubbard-Stratonovich transformation $\mathbb{E}_{\xi} \exp(\sqrt{a\xi}) = e^{\frac{a}{2}}$, the prior integral can be further simplified

$$\begin{split} \Psi_{\mathbf{w}}^{(r)}(\hat{Q})\Big|_{\mathrm{rs}} &= \int_{\mathbb{R}^{r+1}} d\tilde{\mathbf{w}} P_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}}) e^{\frac{1}{2}\tilde{\mathbf{w}}^{\mathsf{T}}\hat{Q}_{\mathrm{rs}}\tilde{\mathbf{w}}} \\ &= \mathbb{E}_{w^{\star}} e^{\frac{1}{2}\hat{Q}^{0}(w^{\star})^{2}} \int_{\mathbb{R}^{r}} d\tilde{\mathbf{w}} P_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}}) e^{w^{\star}\hat{m}\sum_{a=1}^{r}\tilde{w}^{a} - \frac{1}{2}(\hat{Q} + \hat{q})\sum_{a=1}^{r}(\tilde{w}^{a})^{2} + \frac{1}{2}\hat{q}(\sum_{a=1}^{r}\tilde{w}^{a})^{2}} \\ &= \mathbb{E}_{\xi,w^{\star}} e^{\frac{1}{2}\hat{Q}^{0}(w^{\star})^{2}} \left[\mathbb{E}_{w} \exp\left(\left[\hat{m}w^{\star}w - \frac{1}{2}(\hat{Q} + \hat{q})w^{2} + \hat{q}^{1/2}\xi w\right]\right)\right]^{r} \,. \end{split}$$

$$(90)$$

Channel integral Let's focus on the inverse matrix

$$Q_{\rm rs}^{-1} = \begin{bmatrix} Q_{00}^{-1} & Q_{01}^{-1} & Q_{01}^{-1} & Q_{01}^{-1} \\ Q_{01}^{-1} & Q_{11}^{-1} & Q_{12}^{-1} & Q_{12}^{-1} \\ Q_{01}^{-1} & Q_{12}^{-1} & Q_{11}^{-1} & Q_{12}^{-1} \\ Q_{01}^{-1} & Q_{12}^{-1} & Q_{12}^{-1} & Q_{11}^{-1} \end{bmatrix}$$
(91)

with

$$\begin{cases} Q_{00}^{-1} &= \left(Q^0 - rm(Q + (r-1)q)^{-1}m\right)^{-1} \\ Q_{01}^{-1} &= -\left(Q^0 - rm(Q + (r-1)q)^{-1}m\right)^{-1}m(q + (r-1)q)^{-1} \\ Q_{11}^{-1} &= (Q-q)^{-1} - (Q + (r-1)q)^{-1}q(Q-q)^{-1} \\ &+ (Q + (r-1)q)^{-1}m\left(Q^0 - rm(Q + (r-1)q)^{-1}m\right)^{-1}m(Q + (r-1)q)^{-1} \\ Q_{12}^{-1} &= -(Q + (r-1)q)^{-1}q(Q-q)^{-1} \\ &+ (Q + (r-1)q)^{-1}m\left(Q - rm(Q + (r-1)q)^{-1}m\right)^{-1}m(Q + (r-1)q)^{-1} \end{cases}$$

and its determinant:

det
$$Q_{\rm rs} = (Q-q)^{r-1} (Q+(r-1)q) (Q^0 - rm(Q+(r-1)q)^{-1}m)$$

Using the same kind of Gaussian transformation, we obtain

$$\begin{split} \Psi_{\text{out}}^{(r)}(Q)\Big|_{\text{rs}} &= \int \mathrm{d}y \int_{\mathbb{R}^{r+1}} d\tilde{\mathbf{z}} e^{-\frac{1}{2}\tilde{\mathbf{z}}^{\intercal}Q_{\text{rs}}^{-1}\tilde{\mathbf{z}}-\frac{1}{2}\log(\det(2\pi Q_{\text{rs}}))} P_{\text{out}}(y|\tilde{\mathbf{z}}) \\ &= \mathbb{E}_{y,\xi} \mathbf{e}^{-\frac{1}{2}\log(\det(2\pi Q_{\text{rs}}))} \\ &\times \int \mathrm{d}z^{\star}P_{\text{out}^{\star}}\left(y|z^{\star}\right) e^{-\frac{1}{2}Q_{00}^{-1}(z^{\star})^{2}} \left[\int dz P_{\text{out}}\left(y|z\right) e^{-Q_{01}^{-1}z^{\star}z - \frac{1}{2}\left(Q_{11}^{-1} - Q_{12}^{-1}\right)z^{2} - Q_{12}^{-1/2}\xi z}\right]^{r} \end{split}$$

IV.3 ERM and Bayes-optimal free entropy

Taking carefully the derivative and the $r \to 0$ limit imposes $\hat{Q}^0 = 0$ and we finally obtain the replica symmetric free entropy Φ_{rs} :

$$\Phi_{\rm rs}(\alpha) \equiv \mathbb{E}_{\mathbf{y},\mathbf{X}} \left[\lim_{d \to \infty} \frac{1}{d} \log \left(\mathcal{Z}\left(\mathbf{y},\mathbf{X}\right) \right) \right]$$

$$= \mathbf{extr}_{Q,\hat{Q},q,\hat{q},m,\hat{m}} \left\{ -m\hat{m} + \frac{1}{2}Q\hat{Q} + \frac{1}{2}q\hat{q} + \Psi_{\rm w}\left(\hat{Q},\hat{m},\hat{q}\right) + \alpha\Psi_{\rm out}\left(Q,m,q;\rho_{\rm w^{\star}}\right) \right\},$$
(92)

where $\rho_{w^{\star}} = \lim_{d \to \infty} \mathbb{E}_{w^{\star}} \frac{1}{d} \|w^{\star}\|_2^2$ and the channel and prior integrals are defined by

$$\Psi_{\mathrm{w}}\left(\hat{Q},\hat{m},\hat{q}\right) \equiv \mathbb{E}_{\xi}\left[\mathcal{Z}_{\mathrm{w}^{\star}}\left(\hat{m}\hat{q}^{-1/2}\xi,\hat{m}\hat{q}^{-1}\hat{m}\right)\log\mathcal{Z}_{\mathrm{w}}\left(\hat{q}^{1/2}\xi,\hat{Q}+\hat{q}\right)\right],$$

$$\Psi_{\mathrm{out}}\left(Q,m,q;\rho_{\mathrm{w}^{\star}}\right) \equiv \mathbb{E}_{y,\xi}\left[\mathcal{Z}_{\mathrm{out}^{\star}}\left(y,mq^{-1/2}\xi,\rho_{\mathrm{w}^{\star}}-mq^{-1}m\right)\log\mathcal{Z}_{\mathrm{out}}\left(y,q^{1/2}\xi,Q-q\right)\right],$$
(93)

where again Z_{out^*} and Z_{w^*} are defined in eq. (24) and depend on the *teacher*, while the denoising functions Z_{out} and Z_w depend on the inference model. In particular, we explicit in the next sections the above free entropy in the case of ERM and Bayes-optimal estimation.

IV.3.1 ERM estimation

As described in eq. (21), the free entropy for ERM estimation is therefore given by eq. (92) if we take $-\log \mathbb{P}(\mathbf{y}|\mathbf{z}) = l(\mathbf{y}, \mathbf{z})$ and $-\log \mathbb{P}(\mathbf{w}) = r(\mathbf{w})$. As described in Sec. I.3.2 they lead to the following partition functions:

$$\mathcal{Z}_{w}^{\lambda}(\gamma,\Lambda) = \lim_{\Delta \to 0} e^{-\frac{1}{\Delta}\mathcal{M}_{\Lambda^{-1}}[r(\lambda,.)](\Lambda^{-1}\gamma)} e^{-\frac{1}{2\Delta}\gamma^{2}\Lambda^{-1}},$$

$$\mathcal{Z}_{out}(y,\omega,V) = \lim_{\Delta \to 0} \frac{e^{-\frac{1}{\Delta}\mathcal{M}_{\underline{X}}[l(y,.)](\omega)}}{\sqrt{2\pi V}\sqrt{2\pi\Delta}},$$
(94)

with the Moreau-Yosida regularization (28).

IV.3.2 Bayes-optimal estimation

In the Bayes-optimal case, we have access to the ground truth distributions $\mathbb{P}(\mathbf{y}|\mathbf{z}) = P_{\text{out}^*}(\mathbf{y}|\mathbf{z})$ and $\mathbb{P}(\mathbf{w}) = P_{\mathbf{w}^*}(\mathbf{w})$, and therefore $\mathcal{Z}_{\text{out}} = \mathcal{Z}_{\text{out}^*}$, $\mathcal{Z}_{\mathbf{w}} = \mathcal{Z}_{\mathbf{w}^*}$. Nishimori conditions in the Bayes-optimal case [55] imply that $Q = \rho_{\mathbf{w}^*}$, $m = q = q_{\text{b}}$, $\hat{Q} = 0$, $\hat{m} = \hat{q} = \hat{q}_{\text{b}}$. Therefore the free entropy eq. (92) simplifies as an optimization problem over two scalar *overlaps* q_{b} , \hat{q}_{b} :

$$\Phi^{\rm b}(\alpha) = \mathbf{extr}_{q_{\rm b},\hat{q}_{\rm b}} \left\{ -\frac{1}{2} q_{\rm b} \hat{q}_{\rm b} + \Psi^{\rm b}_{\rm w}\left(\hat{q}_{\rm b}\right) + \alpha \Psi^{\rm b}_{\rm out}\left(q_{\rm b};\rho_{\rm w^{\star}}\right) \right\},\tag{95}$$

with free entropy terms $\Psi^{\rm b}_{\rm w}$ and $\Psi^{\rm b}_{\rm out}$ given by

$$\begin{split} \Psi^{\mathrm{b}}_{\mathrm{w}}\left(\hat{q}\right) &= \mathbb{E}_{\xi}\left[\mathcal{Z}_{\mathrm{w}^{\star}}\left(\hat{q}^{1/2}\xi,\hat{q}\right)\log\mathcal{Z}_{\mathrm{w}^{\star}}\left(\hat{q}^{1/2}\xi,\hat{q}\right)\right]\,,\\ \Psi^{\mathrm{b}}_{\mathrm{out}}\left(q;\rho_{\mathrm{w}^{\star}}\right) &= \mathbb{E}_{y,\xi}\left[\mathcal{Z}_{\mathrm{out}^{\star}}\left(y,q^{1/2}\xi,\rho_{\mathrm{w}^{\star}}-q\right)\log\mathcal{Z}_{\mathrm{out}^{\star}}\left(y,q^{1/2}\xi,\rho_{\mathrm{w}^{\star}}-q\right)\right]\,. \end{split}$$

and again Z_{out^*} and Z_{w^*} are defined in eq. (24). The above replica symmetric free entropy in the Bayes-optimal case has been rigorously proven in [10].

IV.4 Sets of fixed point equations

As highlighted in Sec. II, the asymptotic overlaps m, q measure the performances of the ERM or Bayes-optimal statistical estimators, whose behaviours are respectively characterized by extremizing the free entropy (92) and (95). This section is devoted to derive the corresponding sets of fixed point equations.

IV.4.1 ERM estimation

Extremizing the free entropy eq. (92), we easily obtain the set of six fixed point equations

These equations can be formulated as functions of the partition functions Z_{out^*} , Z_{w^*} and the denoising functions f_{out^*} , f_{w^*} , f_{out} , f_{w} defined in eq. (25) and eq. (29). The derivation is shown in Appendix. IV.5.3 and defining the natural variables $\Sigma = Q - q$, $\hat{\Sigma} = \hat{Q} + \hat{q}$, $\eta \equiv \frac{m^2}{\rho_{\text{w}^*}q}$ and $\hat{\eta} \equiv \frac{\hat{m}^2}{\hat{q}}$, it can be written as

$$m = \mathbb{E}_{\xi} \left[\mathcal{Z}_{w^{\star}} \left(\sqrt{\hat{\eta}} \xi, \hat{\eta} \right) f_{w^{\star}} \left(\sqrt{\hat{\eta}} \xi, \hat{\eta} \right) f_{w} \left(\hat{q}^{1/2} \xi, \hat{\Sigma} \right) \right],$$

$$q = \mathbb{E}_{\xi} \left[\mathcal{Z}_{w^{\star}} \left(\sqrt{\hat{\eta}} \xi, \hat{\eta} \right) f_{w} \left(\hat{q}^{1/2} \xi, \hat{\Sigma} \right)^{2} \right],$$

$$\Sigma = \mathbb{E}_{\xi} \left[\mathcal{Z}_{w^{\star}} \left(\sqrt{\hat{\eta}} \xi, \hat{\eta} \right) \partial_{\gamma} f_{w} \left(\hat{q}^{1/2} \xi, \hat{\Sigma} \right) \right],$$

$$\hat{m} = \alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{out^{\star}} (.) \cdot f_{out^{\star}} (y, \sqrt{\rho_{w^{\star}} \eta} \xi, \rho_{w^{\star}} (1 - \eta)) f_{out} \left(y, q^{1/2} \xi, \Sigma \right) \right],$$

$$\hat{q} = \alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{out^{\star}} (y, \sqrt{\rho_{w^{\star}} \eta} \xi, \rho_{w^{\star}} (1 - \eta)) f_{out} \left(y, q^{1/2} \xi, \Sigma \right)^{2} \right],$$

$$\hat{\Sigma} = -\alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{out^{\star}} (y, \sqrt{\rho_{w^{\star}} \eta} \xi, \rho_{w^{\star}} (1 - \eta)) \partial_{\omega} f_{out} \left(y, q^{1/2} \xi, \Sigma \right) \right],$$
(97)

and we finally obtain the set of equations eqs. (71).

IV.4.2 Bayes-optimal estimation

Extremizing the Bayes-optimal free entropy eq. (95), we easily obtain the set of 2 fixed point equations over the scalar parameters $q_{\rm b}$, $\hat{q}_{\rm b}$. In fact, it can also be deduced from eq. (97) using the Nishimori conditions $f_{\rm w} = f_{\rm w^{\star}}$, $f_{\rm out} = f_{\rm out^{\star}}$, $m = q = q_{\rm b}$, $\Sigma = \rho_{\rm w^{\star}} - q$, $\hat{m} = \hat{q} = \hat{q}_{\rm b}$ and $\hat{Q} = 0$ that lead to the result (13) in Thm. 2.4, from [10]

$$\hat{q}_{\rm b} = \alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\rm out^{\star}} \left(y, q_{\rm b}^{1/2} \xi, \rho_{\rm w^{\star}} - q_{\rm b} \right) f_{\rm out^{\star}} \left(y, q_{\rm b}^{1/2} \xi, \rho_{\rm w^{\star}} - q_{\rm b} \right)^2 \right] ,$$

$$q_{\rm b} = \mathbb{E}_{\xi} \left[\mathcal{Z}_{\rm w^{\star}} \left(\hat{q}_{\rm b}^{1/2} \xi, \hat{q}_{\rm b} \right) f_{\rm w^{\star}} \left(\hat{q}_{\rm b}^{1/2} \xi, \hat{q}_{\rm b} \right)^2 \right] .$$
(98)

IV.5 Useful derivations

In this section, we give useful computation steps that we used to transform the sets of fixed point equations (96).

IV.5.1 Prior free entropy term

In specific simple cases, the prior free entropy term

$$\Psi_{\mathbf{w}}\left(\hat{Q},\hat{m},\hat{q}\right) \equiv \mathbb{E}_{\xi}\left[\mathcal{Z}_{\mathbf{w}^{\star}}\left(\hat{m}\hat{q}^{-1/2}\xi,\hat{m}\hat{q}^{-1}\hat{m}\right)\log\mathcal{Z}_{\mathbf{w}}\left(\hat{q}^{1/2}\xi,\hat{Q}+\hat{q}\right)\right]$$

in (93) can be computed explicitly. This is the case of Gaussian and binary priors P_{w^*} with ℓ_2 regularization. In particular, they lead surprisingly to the same expression meaning that choosing a binary or Gaussian teacher distribution does not affect the ERM performances with ℓ_2 regularization.

Gaussian prior Let us compute the corresponding free entropy term with partition functions $\mathcal{Z}_{w^{\star}}$ for a Gaussian prior $P_{w^{\star}}(w^{\star}) = \mathcal{N}_{w^{\star}}(0, \rho_{w^{\star}})$ and $\mathcal{Z}_{w}^{\ell_{2},\lambda}$ for a ℓ_{2} regularization respectively given by eq. (41) and eq. (47):

$$\mathcal{Z}_{\mathbf{w}^{\star}}\left(\boldsymbol{\gamma},\boldsymbol{\Lambda}\right) = \frac{e^{\frac{\boldsymbol{\gamma}^{2}\boldsymbol{\rho}_{\mathbf{w}^{\star}}}{2\left(\boldsymbol{\Lambda}\boldsymbol{\rho}_{\mathbf{w}^{\star}}+1\right)}}}{\sqrt{\boldsymbol{\Lambda}\boldsymbol{\rho}_{\mathbf{w}^{\star}}+1}}\,,\quad \mathcal{Z}_{\mathbf{w}}^{\ell_{2},\boldsymbol{\lambda}}\left(\boldsymbol{\gamma},\boldsymbol{\Lambda}\right) = \frac{e^{\frac{\boldsymbol{\gamma}^{2}}{2\left(\boldsymbol{\Lambda}+\boldsymbol{\lambda}\right)}}}{\sqrt{\boldsymbol{\Lambda}+\boldsymbol{\lambda}}}\,.$$

The prior free entropy term reads

$$\begin{split} \Psi_{\rm w}\left(\hat{Q},\hat{m},\hat{q}\right) &= \mathbb{E}_{\xi}\left[\mathcal{Z}_{\rm w^{\star}}\left(\hat{m}\hat{q}^{-1/2}\xi,\hat{m}\hat{q}^{-1}\hat{m}\right)\log\mathcal{Z}_{\rm w}^{\ell_{2},\lambda}\left(\hat{q}^{1/2}\xi,\hat{q}+\hat{Q}\right)\right] \\ &= \mathbb{E}_{\xi}\left[\mathcal{Z}_{\rm w^{\star}}\left(\hat{m}\hat{q}^{-1/2}\xi,\hat{m}\hat{q}^{-1}\hat{m}\right)\left(\frac{\hat{q}\xi^{2}}{2\left(\lambda+\hat{Q}+\hat{q}\right)}-\frac{1}{2}\log\left(\lambda+\hat{Q}+\hat{q}\right)\right)\right] \\ &= \int\mathrm{d}\xi\mathcal{N}_{\xi}\left(0,1+\rho_{\rm w^{\star}}\hat{m}^{2}\hat{q}^{-1}\right)\left(\frac{\hat{q}\xi^{2}}{2\left(\lambda+\hat{Q}+\hat{q}\right)}-\frac{1}{2}\log\left(\lambda+\hat{Q}+\hat{q}\right)\right) \\ &= \frac{1}{2}\left(\frac{\hat{q}+\rho_{\rm w^{\star}}\hat{m}^{2}}{\lambda+\hat{Q}+\hat{q}}-\log\left(\lambda+\hat{Q}+\hat{q}\right)\right) \end{split}$$
(99)

In the Bayes-optimal case for $\rho_{w^{\star}} = 1$, the computation is similar and is given by the above expression with $\lambda = 1$, $\hat{Q} = 0$, $\hat{m} = \hat{q}$:

$$\Psi_{\rm w}^{\rm bayes}(\hat{q}) = = \frac{1}{2} \left(\hat{q} - \log\left(1 + \hat{q}\right) \right)$$
(100)

Binary prior Let us compute the corresponding free entropy term with partition functions $\mathcal{Z}_{w^{\star}}$ for a binary prior $P_{w^{\star}}(w^{\star}) = \frac{1}{2} \left(\delta(w^{\star} - 1) + \delta(w^{\star} + 1) \right)$ and $\mathcal{Z}_{w}^{\ell_{2},\lambda}$ for a ℓ_{2} regularization respectively given by eq. (42) and eq. (47):

$$\mathcal{Z}_{\mathbf{w}^{\star}}(\gamma,\Lambda) = e^{-\frac{\Lambda}{2}} \cosh(\gamma), \quad \mathcal{Z}_{\mathbf{w}}^{\ell_{2},\lambda}(\gamma,\Lambda) = \frac{e^{\frac{\gamma^{2}}{2(\Lambda+\lambda)}}}{\sqrt{\Lambda+\lambda}}.$$

The entropy term Ψ_w reads

$$\begin{split} \Psi_{\rm w}\left(\hat{Q},\hat{m},\hat{q}\right) &= \mathbb{E}_{\xi}\left[\mathcal{Z}_{\rm w^{\star}}\left(\hat{m}\hat{q}^{-1/2}\xi,\hat{m}\hat{q}^{-1}\hat{m}\right)\log\mathcal{Z}_{\rm w}^{\ell_{2},\lambda}\left(\hat{q}^{1/2}\xi,\hat{q}+\hat{Q}\right)\right] \\ &= \mathbb{E}_{\xi}\left[\mathcal{Z}_{\rm w^{\star}}\left(\hat{m}\hat{q}^{-1/2}\xi,\hat{m}\hat{q}^{-1}\hat{m}\right)\left(\frac{\hat{q}\xi^{2}}{2\left(\lambda+\hat{Q}+\hat{q}\right)}-\frac{1}{2}\log\left(\lambda+\hat{Q}+\hat{q}\right)\right)\right] \\ &= \int d\xi\frac{e^{-\frac{\xi^{2}}{2}}}{\sqrt{2\pi}}e^{-\frac{\hat{m}\hat{q}^{-1}\hat{m}}{2}}\cosh\left(\hat{m}\hat{q}^{-1/2}\xi\right)\left(\frac{\hat{q}\xi^{2}}{2\left(\lambda+\hat{Q}+\hat{q}\right)}-\frac{1}{2}\log\left(\lambda+\hat{Q}+\hat{q}\right)\right) \\ &= \frac{1}{2}\left(\frac{\hat{q}+\hat{m}^{2}}{\lambda+\hat{Q}+\hat{q}}-\log\left(\lambda+\hat{Q}+\hat{q}\right)\right) \end{split}$$
(101)

We recover exactly the same free entropy term than for Gaussian prior teacher eq. (99) for $\rho_{w^{\star}} = 1$.

IV.5.2 Updates derivatives

Let's compute, in full generality, the derivative of the partition functions defined in Sec. 22 and that will be useful to simplify the set (96).

$$\partial_{\gamma} \mathcal{Z}_{w} (\gamma, \Lambda) = \mathcal{Z}_{w} (\gamma, \Lambda) \times \mathbb{E}_{Q_{w}} [w] = \mathcal{Z}_{w} (\gamma, \Lambda) f_{w} (\gamma, \Lambda)$$

$$\partial_{\Lambda} \mathcal{Z}_{w} (\gamma, \Lambda) = -\frac{1}{2} \mathcal{Z}_{w} (\gamma, \Lambda) \times \mathbb{E}_{Q_{w}} [w^{2}] = -\frac{1}{2} \left(\partial_{\gamma} f_{w} (\gamma, \Lambda) + f_{w}^{2} (\gamma, \Lambda) \right)$$

$$\partial_{\omega} \mathcal{Z}_{out} (y, \omega, V) = \mathcal{Z}_{out} (y, \omega, V) \times V^{-1} \mathbb{E}_{Q_{out}} [z - \omega]$$

$$= \mathcal{Z}_{out} (y, \omega, V) f_{out} (y, \omega, V)$$

$$\partial_{V} \mathcal{Z}_{out} (y, \omega, V) = \frac{1}{2} \mathcal{Z}_{out} (y, \omega, V) \times \left(\mathbb{E}_{Q_{out}} \left[V^{-2} (z - \omega)^{2} \right] - V^{-1} \right)$$

$$= \frac{1}{2} \mathcal{Z}_{out} (y, \omega, V) \left(\partial_{\omega} f_{out} (y, \omega, V) + f_{out}^{2} (y, \omega, V) \right)$$
(102)

IV.5.3 Simplifications of the fixed point equations

We recall the set of fixed point equations eq. (96)

$$\hat{Q} = -2\alpha \partial_Q \Psi_{\text{out}}, \qquad \qquad Q = -2\partial_{\hat{Q}} \Psi_{\text{w}}
\hat{q} = -2\alpha \partial_q \Psi_{\text{out}}, \qquad \qquad q = -2\partial_{\hat{q}} \Psi_{\text{w}}, \qquad (103)
\hat{m} = \alpha \partial_m \Psi_{\text{out}}, \qquad \qquad m = \partial_{\hat{m}} \Psi_{\text{w}},$$

that can be simplified and formulated as functions of \mathcal{Z}_{out^*} , \mathcal{Z}_{w^*} , f_{out^*} , f_{w^*} , f_{out} , and f_w defined in eq. (25) and eq. (29), using the derivatives in (102).

Equation over \hat{q}

$$\begin{split} \partial_{q}\Psi_{\text{out}} &= \partial_{q}\mathbb{E}_{y,\xi}\left[\mathcal{Z}_{\text{out}^{\star}}\left(y,mq^{-1/2}\xi,\rho_{\text{w}^{\star}}-mq^{-1}m\right)\log\mathcal{Z}_{\text{out}}\left(y,q^{1/2}\xi,Q-q\right)\right]\right] \\ &= \mathbb{E}_{y,\xi}\left[\partial_{q}\omega^{\star}\partial_{\omega}\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}^{\star}}\right] \\ &+ \frac{\mathcal{Z}_{\text{out}^{\star}}}{\mathcal{Z}_{\text{out}}}\left(\frac{1}{2}q^{-3/2}\xi f_{\text{out}^{\star}}\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}} + \frac{m^{2}q^{-2}}{2}\left(\partial_{\omega}f_{\text{out}^{\star}} + f_{\text{out}^{\star}}^{2}\right)\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}} \\ &+ \frac{\mathcal{Z}_{\text{out}^{\star}}}{\mathcal{Z}_{\text{out}}}\left(\frac{1}{2}q^{-1/2}\xi f_{\text{out}^{\star}}\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}} + f_{\text{out}}^{2}\right)\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}} \\ &+ \frac{\mathcal{Z}_{\text{out}^{\star}}}{\mathcal{Z}_{\text{out}^{\star}}}\left(\frac{1}{2}q^{-1/2}\xi f_{\text{out}^{\star}}\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}}\right) + m^{2}q^{-2}\left(\partial_{\omega}f_{\text{out}^{\star}} + f_{\text{out}^{\star}}^{2}\right)\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}} \\ &+ \frac{\mathcal{Z}_{\text{out}^{\star}}}{\mathcal{Z}_{\text{out}^{\star}}}\left(\frac{1}{2}q^{-1/2}\xi f_{\text{out}^{\star}}\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}}\right) + m^{2}q^{-2}\left(\partial_{\omega}f_{\text{out}^{\star}} + f_{\text{out}^{\star}}^{2}\right)\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}} \\ &+ \frac{\mathcal{Z}_{\text{out}^{\star}}}{\mathcal{Z}_{\text{out}^{\star}}}\left(\frac{1}{2}q^{-1/2}\xi f_{\text{out}^{\star}}\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}}\right) + m^{2}q^{-2}\left(\partial_{\omega}f_{\text{out}^{\star}} + f_{\text{out}^{\star}}^{2}\right)\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}} \\ &+ \frac{\mathcal{Z}_{\text{out}^{\star}}}{\left(\frac{1}{2}q^{-1/2}\xi f_{\text{out}^{\star}}\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}}\right) + m^{2}q^{-2}\left(\partial_{\omega}f_{\text{out}^{\star}} + f_{\text{out}^{\star}^{2}\right)\mathcal{Z}_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}} \\ &+ \left(\partial_{\xi}\left(f_{\text{out}\mathcal{Z}_{\text{out}^{\star}}\right) - \left(\partial_{\omega}f_{\text{out}^{\star}}\mathcal{Z}_{\text{out}^{\star}}\right)\right] \qquad (\text{Stein lemma)} \\ &= \frac{1}{2}\mathbb{E}_{y,\xi}\left[-m^{2}q^{-2}\left(\partial_{\omega}f_{\text{out}^{\star}}\log\mathcal{Z}_{\text{out}^{\star}}\right] + \frac{1}{2}\mathbb{E}_{y,\xi}\left[-mq^{-1}\mathcal{Z}_{\text{out}^{\star}}f_{\text{out}^{\star}}\int_{\mathcal{U}^{\star}}\mathcal{Z}_{\text{out}^{\star}}\int_{\mathcal{U}^{\star}}\mathcal{Z}_{\text{out}^{\star}}\int_{\mathcal{U}^{\star}}\mathcal{Z}_{\text{out}^{\star}}\mathcal{Z}_{\text{out}^{\star}} \log\mathcal{Z}_{\text{out}^{\star}}\right] \\ &\quad \left(\partial_{\omega}f_{\text{out}^{\star}} + f_{\text{out}^{\star}}^{2}\right)\mathcal{Z}_{\text{out}^{\star}}f_{\text{out}^{\star}}f_{\text{out}^{\star}}f_{\text{out}^{\star}}f_{\text{out}^{\star}}\int_{\mathcal{U}^{\star}}\mathcal{Z}_{\text{out}^{\star}}\mathcal{Z}_{$$

that leads to

$$\hat{q} = -2\alpha \partial_q \Psi_{\text{out}} = \alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^\star} \left(y, mq^{-1/2}\xi, \rho_{\text{w}^\star} - mq^{-1}m \right) f_{\text{out}} \left(y, q^{1/2}\xi, Q - q \right)^2 \right].$$
(104)

Equation over \hat{m}

that leads to

$$\hat{m} = \alpha \partial_m \Psi_{\text{out}}$$

$$= \alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^\star} \left(., ., . \right) f_{\text{out}^\star} \left(y, mq^{-1/2}\xi, \rho_{\text{w}^\star} - mq^{-1}m \right) f_{\text{out}} \left(y, q^{1/2}\xi, Q - q \right) \right].$$
(105)

Equation over \hat{Q}

$$\begin{aligned} \partial_Q \Psi_{\text{out}} &= \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^\star} \left(y, mq^{-1/2}\xi, \rho_{\text{w}^\star} - mq^{-1}m \right) \partial_Q \log \mathcal{Z}_{\text{out}} \left(y, q^{1/2}\xi, Q - q \right) \right] \\ &= \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^\star} \left(y, mq^{-1/2}\xi, \rho_{\text{w}^\star} - mq^{-1}m \right) \partial_Q V \partial_V \log \mathcal{Z}_{\text{out}} \left(y, q^{1/2}\xi, Q - q \right) \right] \\ &= \frac{1}{2} \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^\star} \left(y, mq^{-1/2}\xi, \rho_{\text{w}^\star} - mq^{-1}m \right) \left(\partial_\omega f_{\text{out}} + f_{\text{out}}^2 \right) \left(y, q^{1/2}\xi, Q - q \right) \right] \end{aligned}$$

leading to

$$\hat{Q} = -2\alpha \partial_Q \Psi_{\text{out}}$$

= $-\alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\text{out}^{\star}} \left(y, mq^{-1/2}\xi, \rho_{\text{w}^{\star}} - mq^{-1}m \right) \partial_{\omega} f_{\text{out}} \left(y, q^{1/2}\xi, Q - q \right) \right] - \hat{q}.$ (106)

Equation over q

$$\begin{split} \partial_{\hat{q}}\Psi_{w} &= \partial_{\hat{q}}\mathbb{E}_{\xi}\left[\mathcal{Z}_{w^{*}}\left(\hat{m}\hat{q}^{-1/2}\xi,\hat{m}\hat{q}^{-1}\hat{m}\right)\log\mathcal{Z}_{w}\left(\hat{q}^{1/2}\xi,\hat{Q}+\hat{q}\right)\right] \\ &= \mathbb{E}_{\xi}\left[\partial_{\hat{q}}\omega^{*}\partial_{\omega}\mathcal{Z}_{w^{*}}\log\mathcal{Z}_{w}+\partial_{\hat{q}}V^{*}\partial_{V}\mathcal{Z}_{w^{*}}\log\mathcal{Z}_{w}+\frac{\mathcal{Z}_{w^{*}}}{\mathcal{Z}_{w}}\left(\partial_{\hat{q}}\omega\partial_{\omega}\mathcal{Z}_{w}+\partial_{\hat{q}}V\partial_{V}\mathcal{Z}_{w}\right)\right] \\ &= \mathbb{E}_{\xi}\left[-\frac{\hat{m}}{2}\hat{q}^{-3/2}\xi f_{w^{*}}\mathcal{Z}_{w^{*}}\log\mathcal{Z}_{w}+\frac{\hat{m}^{2}\hat{q}^{-2}}{2}\left(\partial_{\omega}f_{w^{*}}+f_{w^{*}}^{2}\right)\mathcal{Z}_{w^{*}}\log\mathcal{Z}_{w} \\ &+\frac{\mathcal{Z}_{w^{*}}}{\mathcal{Z}_{w}}\left(\frac{1}{2}\hat{q}^{-1/2}\xi f_{w}\mathcal{Z}_{w}-\frac{1}{2}\left(\partial_{\omega}f_{w}+f_{w}^{2}\right)\mathcal{Z}_{w}\right)\right] \\ &= \mathbb{E}_{\xi}\left[-\frac{\hat{m}}{2}\hat{q}^{-3/2}\partial_{\xi}\left(f_{w^{*}}\mathcal{Z}_{w^{*}}\log\mathcal{Z}_{w}\right)+\frac{\hat{m}^{2}\hat{q}^{-2}}{2}\left(\partial_{\omega}f_{w^{*}}+f_{w^{*}}^{2}\right)\mathcal{Z}_{w^{*}}\log\mathcal{Z}_{w} \\ &+\left(\frac{1}{2}\hat{q}^{-1/2}\partial_{\xi}\left(f_{w}\mathcal{Z}_{w^{*}}\right)-\frac{1}{2}\left(\partial_{\omega}f_{w}+f_{w}^{2}\right)\mathcal{Z}_{w^{*}}\right)\right] \\ &= \frac{1}{2}\mathbb{E}_{\xi}\left[-\hat{m}^{2}\hat{q}^{-2}\left(\partial_{\omega}f_{w^{*}}\mathcal{Z}_{w^{*}}\log\mathcal{Z}_{w}+\mathcal{Z}_{w^{*}}f_{w^{*}}^{2}\log\mathcal{Z}_{w}-\left(\partial_{\omega}f_{w^{*}}+f_{w^{*}}^{2}\right)\mathcal{Z}_{w^{*}}\log\mathcal{Z}_{w}\right) \\ &-\hat{m}\hat{q}^{-1}\mathcal{Z}_{w^{*}}f_{w^{*}}f_{w}+\left(\hat{m}\hat{q}^{-1}\mathcal{Z}_{w^{*}}f_{w}f_{w^{*}}+\mathcal{Z}_{w^{*}}\partial_{\omega}f_{w}-\left(\partial_{\omega}f_{w}+f_{w}^{2}\right)\mathcal{Z}_{w^{*}}\right)\right] \\ &= -\frac{1}{2}\mathbb{E}_{\xi}\left[\mathcal{Z}_{w^{*}}\left(\hat{m}\hat{q}^{-1/2}\xi,\hat{m}\hat{q}^{-1}\hat{m}\right)f_{w}\left(\hat{q}^{1/2}\xi,\hat{Q}+\hat{q}\right)^{2}\right] \text{ (Simplifications with (102))} \end{split}$$

leading to

$$q = -2\partial_{\hat{q}}\Psi_{w} = \mathbb{E}_{\xi} \left[\mathcal{Z}_{w^{\star}} \left(\hat{m}\hat{q}^{-1/2}\xi, \hat{m}\hat{q}^{-1}\hat{m} \right) f_{w} \left(\hat{q}^{1/2}\xi, \hat{q} + \hat{Q} \right)^{2} \right]$$
(107)

Equation over m

$$\begin{split} \partial_{\hat{m}} \Psi_{w} &= \partial_{m} \mathbb{E}_{\xi} \left[\mathcal{Z}_{w^{\star}} \left(\hat{m} \hat{q}^{-1/2} \xi, \hat{m} \hat{q}^{-1} \hat{m} \right) \log \mathcal{Z}_{w} \left(\hat{q}^{1/2} \xi, \hat{Q} + \hat{q} \right) \right] \\ &= \mathbb{E}_{\xi} \left[(\partial_{\hat{m}} \omega^{\star} \partial_{\omega} \mathcal{Z}_{w^{\star}} + \partial_{\hat{m}} V^{\star} \partial_{V} \mathcal{Z}_{w^{\star}}) \log \mathcal{Z}_{w} \right] \\ &= \mathbb{E}_{\xi} \left[\left(\hat{q}^{-1/2} \xi f_{w^{\star}} \mathcal{Z}_{w^{\star}} - \hat{m} \hat{q}^{-1} \left(\partial_{\omega} f_{w^{\star}} + f_{w^{\star}}^{2} \right) \mathcal{Z}_{w^{\star}} \right) \log \mathcal{Z}_{w} \right] \\ &= \mathbb{E}_{\xi} \left[\hat{m} \hat{q}^{-1} \partial_{\xi} \left(f_{w^{\star}} \mathcal{Z}_{w^{\star}} \log \mathcal{Z}_{w} \right) - \left(\partial_{\omega} f_{w^{\star}} + f_{w^{\star}}^{2} \right) \mathcal{Z}_{w^{\star}} \log \mathcal{Z}_{w} \right] \quad \text{(Stein Lemma)} \\ &= \mathbb{E}_{\xi} \left[\hat{m} \hat{q}^{-1} \left(\partial_{\omega} f_{w^{\star}} \mathcal{Z}_{w^{\star}} \log \mathcal{Z}_{w} + \mathcal{Z}_{w^{\star}} f_{w^{\star}}^{2} \log \mathcal{Z}_{w} - \left(\partial_{\omega} f_{w^{\star}} + f_{w^{\star}}^{2} \right) \mathcal{Z}_{w^{\star}} \log \mathcal{Z}_{w} \right) \\ &\quad + \mathcal{Z}_{w^{\star}} f_{w^{\star}} f_{w} \right] \\ &= \mathbb{E}_{\xi} \left[\mathcal{Z}_{w^{\star}} \left(\hat{m} \hat{q}^{-1/2} \xi, \hat{m} \hat{q}^{-1} \hat{m} \right) f_{w^{\star}} \left(\hat{m} \hat{q}^{-1/2} \xi, \hat{m} \hat{q}^{-1} \hat{m} \right) f_{w} \left(\hat{q}^{1/2} \xi, \hat{Q} + \hat{q} \right) \right] \\ \quad \text{(Simplifications with (102))} \end{split}$$

leading to

$$m = 2\partial_{\hat{m}}\Psi_{w} = \mathbb{E}_{\xi} \left[\mathcal{Z}_{w^{\star}} \left(\hat{m}\hat{q}^{-1/2}\xi, \hat{m}\hat{q}^{-1}\hat{m} \right) f_{w^{\star}} \left(\hat{m}\hat{q}^{-1/2}\xi, \hat{m}\hat{q}^{-1}\hat{m} \right) f_{w} \left(\hat{q}^{1/2}\xi, \hat{q} + \hat{Q} \right) \right]$$
(108)

Equation over ${\cal Q}$

$$\begin{split} \partial_{\hat{Q}}\Psi_{w}\left(\hat{Q},\hat{m},\hat{q}\right) &= \partial_{\hat{Q}}\mathbb{E}_{\xi}\left[\mathcal{Z}_{w^{\star}}\left(\hat{m}\hat{q}^{-1/2}\xi,\hat{m}\hat{q}^{-1}\hat{m}\right)\log\mathcal{Z}_{w}\left(\hat{q}^{1/2}\xi,\hat{Q}+\hat{q}\right)\right] \\ &= \mathbb{E}_{\xi}\left[\mathcal{Z}_{w^{\star}}\left(\hat{m}\hat{q}^{-1/2}\xi,\hat{m}\hat{q}^{-1}\hat{m}\right)\frac{1}{\mathcal{Z}_{w}}\partial_{\hat{Q}}\Lambda\partial_{\Lambda}\mathcal{Z}_{w}\left(\hat{q}^{1/2}\xi,\hat{Q}+\hat{q}\right)\right] \\ &= -\frac{1}{2}\mathbb{E}_{\xi}\left[\mathcal{Z}_{w^{\star}}\left(\hat{m}\hat{q}^{-1/2}\xi,\hat{m}\hat{q}^{-1}\hat{m}\right)\left(\partial_{\gamma}f_{w}+f_{w}^{2}\right)\right] \qquad (\text{with (102)})$$

hence

$$Q = -2\partial_{\hat{Q}}\Psi_{w} = \mathbb{E}_{\xi} \left[\mathcal{Z}_{w^{\star}} \left(\hat{m}\hat{q}^{-1/2}\xi, \hat{m}\hat{q}^{-1}\hat{m} \right) \partial_{\gamma}f_{w} \left(\hat{q}^{1/2}\xi, \hat{q} + \hat{Q} \right) \right] + q.$$
(109)

V Applications

In this section, we provide details of the results presented in Sec. 3. In particular as an illustration, we consider a Gaussian teacher ($\rho_{w^*} = 1$) with a noiseless sign activation:

$$P_{\text{out}^{\star}}(y|z) = \delta\left(y - \text{sign}(z)\right), \qquad P_{\text{w}^{\star}}(w^{\star}) = \mathcal{N}_{w^{\star}}\left(0, \rho_{\text{w}^{\star}}\right), \qquad (110)$$

whose corresponding denoising functions are derived in eq. (39) and eq. (41).

Remark V.1. Note that performances of ERM with ℓ_2 regularization for a teacher with Gaussian weights $P_{w^*}(w) = \mathcal{N}_w(0,1)$ or binary weights $P_{w^*}(w) = \frac{1}{2}(\delta(w-1) + \delta(w+1))$, will be similar. Indeed free entropy terms Ψ_w eq. (93) for a Gaussian prior (99) and for binary weights (101) are equal in this setting, so do the set of fixed point equations.

V.1 Bayes-optimal estimation

Using expressions eq. (39) and eq. (41), corresponding to the *teacher* model eq. (110), the prior equation eq. (98) can be simplified while the channel one has no analytical expression. Hence the set of fixed point equations eqs. (100) for the model eq. (110) read

$$q_{\rm b} = \frac{\hat{q}_{\rm b}}{1 + \hat{q}_{\rm b}}, \quad \hat{q}_{\rm b} = \alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\rm out^{\star}} \left(y, q_{\rm b}^{1/2} \xi, \rho_{\rm w^{\star}} - q_{\rm b} \right) f_{\rm out^{\star}} \left(y, q_{\rm b}^{1/2} \xi, \rho_{\rm w^{\star}} - q_{\rm b} \right)^2 \right].$$
(111)

Large α **behaviour** Let us derive the large α behaviour of the Bayes-optimal generalization error eq. (55) that depends only on the overlap $q_{\rm b}$ solution of eq. (111). $q_{\rm b}$ measures the correlation with the ground truth, so we expect that in the limit $\alpha \to \infty$, $q_{\rm b} \to 1$. Therefore, we need to extract the behaviour of $\hat{q}_{\rm b}$ in eq. (111). Injecting expressions $\mathcal{Z}_{\rm out^{\star}}$ and $f_{\rm out^{\star}}$ from eq. (39), we obtain

$$\begin{split} \hat{q}_{\rm b} &= \alpha \mathbb{E}_{y,\xi} \left[\mathcal{Z}_{\rm out^{\star}} \left(y, q_{\rm b}^{1/2}\xi, 1 - q_{\rm b} \right) f_{\rm out^{\star}} \left(y, q_{\rm b}^{1/2}\xi, 1 - q_{\rm b} \right)^2 \right] \\ &= 2\alpha \int D\xi y^2 \frac{\mathcal{N}_{\sqrt{q}\xi}(0, 1 - q_{\rm b})^2}{\frac{1}{2} \left(1 + \operatorname{erf} \left(\frac{\sqrt{q_{\rm b}}\xi}{\sqrt{2(1 - q_{\rm b})}} \right) \right)} = \frac{2}{\pi} \frac{\alpha}{1 - q_{\rm b}} \int D\xi \frac{e^{-\frac{q_{\rm b}\xi^2}{1 - q_{\rm b}}}}{\left(1 + \operatorname{erf} \left(\frac{\sqrt{q_{\rm b}}\xi}{\sqrt{2(1 - q_{\rm b})}} \right) \right)} \,, \end{split}$$

where the last integral can be computed in the limit $q_{\rm b} \rightarrow 1$:

$$\begin{split} \int D\xi \frac{e^{-\frac{q_{\rm b}\xi^2}{1-q_{\rm b}}}}{\left(1+\operatorname{erf}\left(\frac{\sqrt{q_{\rm b}}\xi}{\sqrt{2(1-q_{\rm b})}}\right)\right)} &= \int \mathrm{d}\xi \frac{\frac{-e^{\frac{\xi^2(q_{\rm b}+1)}{2(1-q_{\rm b})}}}{\sqrt{2\pi}}}{\left(1+\operatorname{erf}\left(\frac{\sqrt{q_{\rm b}}\xi}{\sqrt{2(1-q_{\rm b})}}\right)\right)} \\ &\simeq \int \mathrm{d}\xi \frac{\frac{-e^{\frac{\xi^2}{1-q_{\rm b}}}}{\sqrt{2\pi}}}{\left(1+\operatorname{erf}\left(\frac{\xi}{\sqrt{2(1-q_{\rm b})}}\right)\right)} &= \frac{\sqrt{1-q_{\rm b}}}{\sqrt{2\pi}} \int \mathrm{d}\eta \frac{e^{-\eta^2}}{1+\operatorname{erf}\left(\frac{\eta}{\sqrt{2}}\right)} = \frac{c_0}{\sqrt{2\pi}}\sqrt{1-q_{\rm b}}\,, \end{split}$$

with $c_0 \equiv \int d\eta \frac{e^{-\eta^2}}{1 + erf(\frac{\eta}{\sqrt{2}})} \simeq 2.83748$. Finally, we obtain in the large α limit:

$$\hat{q}_{\rm b} = k \frac{\alpha}{\sqrt{1-q_{\rm b}}} \,, \qquad \qquad q_{\rm b} = \frac{\hat{q}_{\rm b}}{1+\hat{q}_{\rm b}} \,,$$

with $k \equiv \frac{2c_0}{\pi\sqrt{2\pi}} \simeq 0.720647$. The above equations can be solved analytically and lead to:

$$q_{\rm b} = \frac{1}{2} \left(\alpha k \sqrt{\alpha^2 k^2 + 4} - \alpha^2 k^2 \right) \underset{\alpha \to \infty}{\simeq} 1 - \frac{1}{\alpha^2 k^2}, \qquad \hat{q}_{\rm b} = k^2 \alpha^2,$$

and therefore the Bayes-optimal asymptotic generalization error is given by

$$e_{\rm g}^{\rm bayes}(\alpha) = \frac{1}{\pi} \arccos\left(\sqrt{q_{\rm b}}\right) \underset{\alpha \to \infty}{\simeq} \frac{1}{k\pi} \frac{1}{\alpha} \simeq \frac{0.4417}{\alpha} \,.$$
 (112)

V.2 Generalities on ERM with ℓ_2 regularization

Combining the teacher update for Gaussian weights eq. (41) with the update associated to the ℓ_2 regularization eq. (41), the free entropy term can be explicitly derived in (99). Taking the corresponding derivatives, the fixed point equations for m, q, Σ eq. (96) are thus explicit and simply read

$$\Sigma = \frac{1}{\lambda + \hat{\Sigma}}, \qquad q = \frac{\rho_{\mathbf{w}^{\star}} \hat{m}^2 + \hat{q}}{(\lambda + \hat{\Sigma})^2}, \qquad m = \frac{\rho_{\mathbf{w}^{\star}} \hat{m}}{\lambda + \hat{\Sigma}}.$$
(113)

All the following examples have been performed with a ℓ_2 regularization, so that the above equations (113) remain valid for the different losses considered in Sec. 3. In the next subsections, we provide some details on the asymptotic performances of ERM with various losses with ℓ_2 regularization and $\rho_{w^*} = 1$.

In general for a generic loss, the proximal eq. (29) has no analytical expression, just as the fixed point equations (97). The square loss is particular in the sense eqs. (97) have a closed form solution. Also the Hinge loss has an analytical proximal. Apart from that, eqs. (97) must be solved numerically. However it is useful to notice that the proximal can be easily found for a two times differentiable loss using eq. (46). This is for example the case of the logistic loss.

V.3 Ridge regression - Square loss with ℓ_2 regularization

The prior equations over m, q, Σ are already derived in eq. (113) and remain valid. Combining eq. (39) for the considered sign channel with a potential additional Gaussian noise Δ^* in (110) and the square loss eq. (43), the channel fixed point equations for $\hat{q}, \hat{m}, \hat{\Sigma}$ eqs. (97) lead to

V.3.1 Pseudo-inverse estimator

We analyze the fixed point equations eqs. (114) for the *pseudo-inverse* estimator, that is in the limit $\lambda \rightarrow 0$.

Solving Σ Combining the two first equations over Σ and $\hat{\Sigma}$ in (114), we obtain

$$\Sigma = \frac{\sqrt{(\alpha + \lambda - 1)^2 + 4\lambda} - \alpha - \lambda + 1}{2\lambda} \underset{\lambda \to 0}{\simeq} \frac{1 - \alpha + |\alpha - 1|}{2\lambda} + \frac{1}{2} \left(\frac{\alpha + 1}{|\alpha - 1|} - 1 \right), \quad (115)$$

that exhibits two different behaviour depending if $\alpha < 1$ or $\alpha > 1$.

Regime $\alpha < 1$ In this regime $\alpha < 1$, eq. (115) becomes

$$\Sigma = \frac{1-\alpha}{\lambda} + \frac{\alpha}{1-\alpha} \,,$$

that leads to the closed set of equations in the limit $\lambda \to 0$

$$\Sigma = \frac{(1-\alpha)^2 + \alpha\lambda}{\lambda(1-\alpha)} \underset{\lambda \to 0}{\simeq} \frac{1-\alpha}{\lambda}, \qquad \qquad \hat{\Sigma} = \frac{(1-\alpha)\alpha\lambda}{(\alpha-1)^2 + \lambda} \underset{\lambda \to 0}{\simeq} \frac{\lambda\alpha}{1-\alpha}, \\ m = \frac{\alpha(1-\alpha)}{\lambda + (1-\alpha)} \sqrt{\frac{2}{\pi}} \underset{\lambda \to 0}{\simeq} \alpha \sqrt{\frac{2}{\pi}}, \qquad \qquad \hat{m} = \frac{\lambda\alpha\sqrt{\frac{2}{\pi}}}{\lambda + (1-\alpha)} \underset{\lambda \to 0}{\simeq} \frac{\lambda\alpha\sqrt{\frac{2}{\pi}}}{1-\alpha}, \\ q \underset{\lambda \to 0}{\simeq} \frac{\alpha(\pi(1+\Delta^*) - 2\alpha)}{\pi(1-\alpha)}, \qquad \qquad \hat{q} \underset{\lambda \to 0}{\simeq} \frac{\alpha\lambda^2(2(\alpha-2)\alpha + \pi(\Delta^* + 1))}{\pi(1-\alpha)(1-\alpha+\lambda)^2}.$$
(116)

Hence we obtain for $\alpha < 1$:

$$m^{\text{pseudo}} = \alpha \sqrt{\frac{2}{\pi}}$$
 $q^{\text{pseudo}} = \frac{\alpha(\pi(1 + \Delta^*) - 2\alpha)}{\pi(1 - \alpha)}$ (117)

and the corresponding generalization error

$$e_{\rm g}^{\rm pseudo}\left(\alpha\right) = \frac{1}{\pi} \, {\rm acos}\left(\sqrt{\frac{2\alpha(1-\alpha)}{\pi\left(1+\Delta^{\star}\right)-2\alpha}}\right) \, {\rm if} \, \alpha < 1 \,. \tag{118}$$

Note in particular that $e_{\rm g}^{\rm pseudo}(\alpha) \xrightarrow[\alpha \to 1]{\alpha \to 1} 0.5$, meaning that the interpolation peak at $\alpha = 1$ reaches the maximum generalization error.

Regime $\alpha > 1$ Eq. (115) becomes

$$\Sigma = \frac{1}{2} \left(\frac{\alpha + 1}{\alpha - 1} - 1 \right) = \frac{1}{2} \left(\frac{\alpha + 1}{\alpha - 1} - 1 \right) = \frac{1}{\alpha - 1}.$$

In the limit $\lambda \rightarrow 0$, the fixed point equations eqs. (114) reduce to

$$\Sigma + 1 = \frac{\alpha}{\alpha - 1}, \qquad \qquad \hat{\Sigma} = \alpha - 1,$$

$$q = \frac{(\alpha - 1)^2 \frac{2}{\pi} + \hat{q}}{(\alpha - 1)^2}, \qquad \qquad \hat{q} = \frac{(\alpha - 1)^2}{\alpha} \left((1 + q + \Delta^*) - \frac{4}{\pi} \right), \qquad (119)$$

$$m = \sqrt{\frac{2}{\pi}}, \qquad \qquad \hat{m} = (\alpha - 1) \sqrt{\frac{2}{\pi}}.$$

In particular we obtain for $\alpha > 1$:

$$m^{\text{pseudo}} = \sqrt{\frac{2}{\pi}}, \qquad q^{\text{pseudo}} = \frac{1}{\alpha - 1} \left(1 + \Delta^* + \frac{2}{\pi} \left(\alpha - 2 \right) \right), \qquad (120)$$

and the corresponding generalization error

$$e_{\rm g}^{\rm pseudo}\left(\alpha\right) = \frac{1}{\pi} \, \mathrm{acos}\left(\sqrt{\frac{\alpha - 1}{\frac{\pi}{2}\left(1 + \Delta^{\star}\right) + \left(\alpha - 2\right)}}\right) \, \mathrm{if} \, \alpha > 1 \,. \tag{121}$$

Large α **behaviour** From this expression we easily obtain the large α behaviour of the pseudo-inverse estimator:

$$e_{\rm g}^{\rm pseudo}(\alpha) = \frac{1}{\pi} \operatorname{acos}\left(\sqrt{\frac{\alpha - 1}{\frac{\pi}{2}\left(1 + \Delta^{\star}\right) + \left(\alpha - 2\right)}}\right) = \frac{1}{\pi} \operatorname{acos}\left(\left(1 + \frac{C}{\alpha - 1}\right)^{1/2}\right) \underset{\alpha \to \infty}{\simeq} \frac{c}{\sqrt{\alpha}}$$

where $C = \frac{\pi}{2}(1 + \Delta^*) - 1$ and $c = \frac{\sqrt{C}}{\pi}$. In particular for a noiseless teacher $\Delta^* = 0$, $c = \sqrt{\frac{\pi-2}{2\pi^2}} \simeq 0.240487$, leading to

$$e_{\rm g}^{\rm pseudo}(\alpha) \simeq \frac{0.2405}{\sqrt{\alpha}}$$
 (122)

V.3.2 Ridge at finite λ

Let us now consider the set of fixed point equation eq. (114) for finite $\lambda \neq 0$. Defining

$$t_0 \equiv \sqrt{(\alpha + \lambda - 1)^2 + 4\lambda}$$

$$t_1 \equiv (t_0 + \alpha + \lambda + 1)^{-1}$$

$$t_2 \equiv \sqrt{2(\alpha + 1)\lambda + (\alpha - 1)^2 + \lambda^2}$$

$$t_3 \equiv (t_2 + \alpha + \lambda + 1)^{-1}$$

$$t_4 \equiv \sqrt{\alpha^2 + 2\alpha(\lambda - 1) + (\lambda + 1)^2}$$

the equations can be in fact fully solved analytically and read

$$\begin{split} \Sigma &= \frac{1}{2} \frac{t_0 - \alpha - \lambda + 1}{\lambda} \\ \hat{\Sigma} &= \frac{1}{2} (t_0 + \alpha - \lambda - 1) \\ q &= \frac{2\alpha \left(-8\alpha^2 t_1 + 2\alpha + \pi\Delta^* + \pi \right)}{\pi \left(\alpha^2 + \alpha \left(t_2 + 2\lambda - 2 \right) + (\lambda + 1) \left(t_2 + \lambda + 1 \right) \right)}, \\ \hat{q} &= \left(4\alpha\lambda^2 \left(\pi (\Delta^* + 1) \left(t_4 + (\alpha + \lambda) \left(t_2 + \alpha + \lambda \right) + 2\lambda + 1 \right) \right) \\ &- 8\alpha t_3 (t_4 + (\alpha + \lambda) \left(\sqrt{2(\alpha + 1)\lambda} + (\alpha - 1)^2 + \lambda^2 + \alpha + \lambda \right) + 2\lambda \right) - 8\alpha t_3 + 4\alpha^2)), \\ m &= \frac{2\sqrt{\frac{2}{\pi}\alpha}}{t_2 + \alpha + \lambda + 1}, \\ \hat{m} &= \frac{2\sqrt{\frac{2}{\pi}\alpha\lambda}}{t_0 - \alpha + \lambda + 1}. \end{split}$$

Generalization error behaviour at large α Expanding the ratio $\frac{m}{\sqrt{q}}$ in the large α limit, we obtain

$$\frac{m}{\sqrt{q}} \simeq 1 - \frac{C}{2\alpha}$$
 with $C = \frac{\pi}{2} \left(1 + \Delta^{\star} \right) - 1$

leading to

$$e_{\rm g}^{{\rm ridge},\lambda}\left(\alpha\right) = \frac{1}{\pi} \cos\left(\frac{m}{\sqrt{q}}\right) \underset{\alpha \to \infty}{\simeq} \frac{c}{\sqrt{\alpha}} \text{ with } c = \frac{\sqrt{C}}{\pi}.$$
 (123)

Thus, the asymptotic generalization error for ridge regression with any regularization strength $\lambda \ge 0$ decrease as $\frac{0.2405}{\sqrt{\alpha}}$, similarly to the pseudo-inverse result.

Optimal regularization The optimal value $\lambda^{\text{opt}}(\alpha)$, introduced in Sec. 3, which minimizes the generalization error at a given α can be found taking the derivative of $\frac{m}{\sqrt{q}}$ and is written as the root of the following functional

$$\begin{split} F[\alpha, \lambda, \Delta^{\star}] &= \partial_{\lambda} \left(\frac{m}{\sqrt{q}} \right) = \frac{a_1 a_2}{a_3 a_4^2} \,, \\ \text{with} \\ a_1 &= -4\alpha \sqrt{\frac{a_4}{\alpha^2 + \alpha \left(t_2 + 2\lambda - 2 \right) + \left(\lambda + 1 \right) \left(t_2 + \lambda + 1 \right)}} \,, \\ a_2 &= 2 \left(\alpha^2 t_3 + \alpha \left(2\lambda t_3 + \left(t_2 + 2 \right) t_3 - 1 \right) + \left(\lambda + 1 \right) \left(t_2 + \lambda + 1 \right) t_3 \right) - \pi (1 + \Delta^{\star}) \,, \\ a_3 &= \frac{t_0}{t_1} \,, \\ a_4 &= \alpha \left(2 - 8t_1 \right) + \pi \left(1 + \Delta^{\star} \right) \,. \end{split}$$



Figure 3: (Left) Absolute value of the derivative of m/\sqrt{q} with respect to λ plotted in a logarithmic scale. λ^{opt} is reached at the root of the functional $F[\alpha, \lambda]$ that corresponds to the divergence in the logarithmic scale. Plotted for a wide range of α , the optimal value is clearly constant and independent of α . Its value is approximately $\lambda^{\text{opt}} \simeq 0.570796$. (Right) Bayes-optimal (black) vs ridge regression (dashed red) generalization errors with optimal ℓ_2 regularization $\lambda^{\text{opt}} \simeq 0.570796$.

Unfortunately, this functional cannot be analyzed analytically. Instead we plot its value for a wide range of α as a function of λ (for $\Delta^* = 0$) and we observe in particular that there exists a unique value $\lambda^{\text{opt}} \simeq 0.570796$ as illustrated in Fig. 3 (left) that is independent of α . As an illustration, we show the generalization error of ridge regression with the optimal regularization $\lambda^{\text{opt}} = 0.5708$ compared to the Bayes-optimal performances in Fig. 3 (right).

V.4 Hinge regression / SVM - Hinge loss with ℓ_2 regularization

The hinge loss $l^{\text{hinge}}(y, z) = \max(0, 1 - yz)$ is linear by part and is therefore another simple example of analytical loss to analyze. In particular its proximal map can computed in eq. (44) and the corresponding denoising functions read:

$$f_{\text{out}}\left(y,q^{1/2}\xi,\Sigma\right) = \begin{cases} y \text{ if } \xi y < \frac{1-\Sigma}{\sqrt{q}} \\ \frac{y-\sqrt{q}\xi}{\Sigma} \text{ if } \frac{1-\Sigma}{\sqrt{q}} < \xi y < \frac{1}{\sqrt{q}} \\ 0 \text{ otherwise} \end{cases}$$
(124)
$$\partial_{\omega}f_{\text{out}}\left(y,q^{1/2}\xi,\Sigma\right) = \begin{cases} -\frac{1}{\Sigma} \text{ if } \frac{1-\Sigma}{\sqrt{q}} < \xi y < \frac{1}{\sqrt{q}} \\ 0 \text{ otherwise} \end{cases}$$
.

The fixed point equations eq. (97) have unfortunately no closed form and need to be solved numerically.

V.4.1 Max-margin estimator

As proven in [34] both the hinge and logistic estimators converge to the *max-margin* solution in the limit $\lambda \to 0$ as soon as the data are linearly separable. We will start with the fixed point equations for hinge, whose denoising functions (124) are analytical. Taking the $\lambda \to 0$ limit is non-trivial and we need therefore to introduce some rescaled variables to obtain a closed set of equations. Numerical evidences at finite α show that we shall use the following rescaled variables:

$$\hat{m} = \Theta(\lambda), \quad \hat{q} = \Theta(\lambda^2), \quad \hat{\Sigma} = \Theta(\lambda), \quad m = \Theta(1), \quad q = \Theta(1), \quad \Sigma = \Theta(\lambda^{-1}).$$

The fixed point equations eq. (97) simplify and become

$$m = \frac{\hat{m}}{1+\hat{\Sigma}}, \qquad q = \frac{\hat{m}^2 + \hat{q}}{(1+\hat{\Sigma})^2}, \qquad \Sigma = \frac{1}{1+\hat{\Sigma}},$$

$$\hat{m} = \frac{2\alpha}{\Sigma} \mathcal{I}_{\hat{m}}(q,\eta), \quad \hat{q} = \frac{2\alpha}{\Sigma^2} \mathcal{I}_{\hat{q}}(q,\eta), \quad \hat{\Sigma} = \frac{2\alpha}{\Sigma} \mathcal{I}_{\hat{\Sigma}}(q,\eta),$$

(125)

with

$$\mathcal{I}_{\hat{m}}(q,\eta) \equiv \int_{-\infty}^{\frac{1}{\sqrt{q}}} \mathrm{d}\xi \mathcal{N}_{\xi}(0,1) \mathcal{N}_{\xi}\left(0,\frac{1-\eta}{\sqrt{\eta}}\right) (1-\sqrt{q}\xi) , \\
= \frac{\sqrt{2\pi} \left(\operatorname{erf}\left(\frac{1}{\sqrt{2}\sqrt{q(1-\eta)}}\right) + 1\right) + 2e^{-\frac{1}{2q(1-\eta)}}\sqrt{q(1-\eta)}}{4\pi} \\
\mathcal{I}_{\hat{q}}(q,\eta) \equiv \int_{-\infty}^{\frac{1}{\sqrt{q}}} \mathrm{d}\xi \mathcal{N}_{\xi}(0,1) \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{\sqrt{\eta}\xi}{\sqrt{2(1-\eta)}}\right)\right) (1-\sqrt{q}\xi)^{2} , \\
\mathcal{I}_{\hat{\Sigma}}(q,\eta) \equiv \int_{-\infty}^{\frac{1}{\sqrt{q}}} \mathrm{d}\xi \mathcal{N}_{\xi}(0,1) \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{\sqrt{\eta}\xi}{\sqrt{2(1-\eta)}}\right)\right) .$$
(126)

Large α **expansion** Numerically at large α (and $\lambda \rightarrow 0$), we obtain the following scalings

$$q = \Theta(\alpha^2), \quad m = \Theta(\alpha), \quad \Sigma = \Theta(1), \quad \hat{q} = \Theta(1), \quad \hat{m} = \Theta(\alpha), \quad \hat{\Sigma} = \Theta(1).$$
(127)

Therefore, in order to close the equations, we introduce new variables (c_q,c_η) such that

$$q \underset{\alpha \to \infty}{=} c_q \alpha^2, \qquad \eta = 1 - \frac{c_\eta}{\alpha^2}.$$
 (128)

In this limit, we can extract the large α behaviours of integrals $\mathcal{I}_{\hat{m}}, \mathcal{I}_{\hat{q}}, \mathcal{I}_{\hat{\Sigma}}$:

$$\mathcal{I}_{\hat{m}}(q,\eta) = \mathcal{I}_{\hat{m}}^{\infty}(c_q, c_\eta) , \quad \mathcal{I}_{\hat{q}}(q,\eta) = \frac{\mathcal{I}_{\hat{q}}^{\infty}(c_q, c_\eta)}{\alpha} , \quad \mathcal{I}_{\hat{\Sigma}}(q,\eta) = \frac{\mathcal{I}_{\hat{\Sigma}}^{\infty}(c_q, c_\eta)}{\alpha} , \quad (129)$$

where $\mathcal{I}^\infty_{\hat{m}}, \mathcal{I}^\infty_{\hat{q}}, \mathcal{I}^\infty_{\hat{\Sigma}}$ are $\Theta(1)$ and read

$$\mathcal{I}_{\hat{m}}^{\infty}(c_{q},c_{\eta}) \equiv \frac{\sqrt{2\pi} \left(\operatorname{erf} \left(\frac{1}{\sqrt{2}\sqrt{c_{\eta}c_{q}}} \right) + 1 \right) + 2e^{-\frac{1}{2c_{\eta}c_{q}}} \sqrt{c_{\eta}c_{q}}}{4\pi}, \\
\mathcal{I}_{\hat{q}}^{\infty}(c_{q},c_{\eta}) \equiv \frac{e^{-\frac{1}{2c_{\eta}c_{q}}} \left(\sqrt{2\pi} (3c_{\eta}c_{q}+1)e^{\frac{1}{2c_{\eta}c_{q}}} \left(\operatorname{erf} \left(\frac{1}{\sqrt{2}\sqrt{c_{\eta}c_{q}}} \right) + 1 \right) + 4(c_{\eta}c_{q})^{3/2} + 2\sqrt{c_{\eta}c_{q}}} \right)}{12\pi\sqrt{c_{q}}}, \\
\mathcal{I}_{\hat{\Sigma}}^{\infty}(c_{q},c_{\eta}) \equiv \frac{\sqrt{2\pi} \left(\operatorname{erf} \left(\frac{1}{\sqrt{2}\sqrt{c_{\eta}c_{q}}} \right) + 1 \right) + 2e^{-\frac{1}{2c_{\eta}c_{q}}} \sqrt{c_{\eta}c_{q}}}{4\pi\sqrt{c_{q}}}. \tag{130}$$

Hence the set of fixed-point equations eq. (125) simplifies to:

$$\hat{\Sigma} = \frac{2\mathcal{I}_{\hat{\Sigma}}^{\infty}(c_{q}, c_{\eta})}{1 - 2\mathcal{I}_{\hat{\Sigma}}^{\infty}(c_{q}, c_{\eta})}, \qquad \Sigma = 1 - 2\mathcal{I}_{\hat{\Sigma}}^{\infty}(c_{q}, c_{\eta})$$

$$\hat{m} = \frac{2\alpha \mathcal{I}_{\hat{m}}^{\infty}(c_{q}, c_{\eta})}{1 - 2\mathcal{I}_{\hat{\Sigma}}^{\infty}(c_{q}, c_{\eta})}, \qquad m = 2\alpha \mathcal{I}_{\hat{m}}^{\infty}(c_{q}, c_{\eta})$$

$$\hat{q} = \frac{2\mathcal{I}_{\hat{q}}^{\infty}(c_{q}, c_{\eta})}{\left(1 - 2\mathcal{I}_{\hat{\Sigma}}^{\infty}(c_{q}, c_{\eta})\right)^{2}}, \quad q = 4\alpha^{2} \left(\mathcal{I}_{\hat{m}}^{\infty}(c_{q}, c_{\eta})\right)^{2} + 2\mathcal{I}_{\hat{q}}^{\infty}(c_{q}, c_{\eta}),$$
(131)

which can be closed by rewriting the equations eqs. (128):

$$\eta = \frac{m^2}{q} \equiv 1 - \frac{c_{\eta}}{\alpha^2} = 1 - \frac{\mathcal{I}_{\hat{q}}^{\infty}(c_q, c_{\eta})}{2\left(\mathcal{I}_{\hat{m}}^{\infty}(c_q, c_{\eta})\right)^2} \frac{1}{\alpha^2},$$

$$q = c_q \alpha^2 \simeq 4\alpha^2 \left(\mathcal{I}_{\hat{m}}^{\infty}(c_q, c_{\eta})\right)^2.$$
(132)

Equivalently $(c_q^{\star}, c_{\eta}^{\star})$ is the root of the set of non-linear fixed point equations $(F_{\eta}(c_q, c_{\eta}), F_q(c_q, c_{\eta}))$:

$$F_{\eta}(c_q, c_{\eta}) \equiv \frac{\mathcal{I}_{\hat{q}}^{\infty}(c_q, c_{\eta})}{2\left(\mathcal{I}_{\hat{m}}^{\infty}(c_q, c_{\eta})\right)^2} - c_{\eta}, \qquad F_q(c_q, c_{\eta}) \equiv 4\left(\mathcal{I}_{\hat{m}}^{\infty}(c_q, c_{\eta})\right)^2 - c_q, \qquad (133)$$

that cannot be solved analytically. However a unique numerical solution is found and lead to $(c_q^{\star}, c_{\eta}^{\star}) = (0.9911, 2.4722)$. Therefore the generalization error of the max-margin estimator in the large α regime is given by

$$e_{\rm g}^{\rm max-margin}(\alpha) = \frac{1}{\pi}\arccos\left(\frac{m}{\sqrt{q}}\right) \underset{\alpha \to \infty}{\simeq} \frac{1}{\pi}\arccos\left(1 - \frac{c_{\eta}^{\star}}{\alpha^2}\right) \underset{\alpha \to \infty}{\simeq} \frac{K}{\alpha}, \quad (134)$$

with $K=\frac{\sqrt{c_{\eta}^{\star}}}{\pi}\simeq 0.5005,$ leading to

$$e_{\rm g}^{\rm max-margin}(\alpha) \simeq \frac{0.5005}{\alpha}$$
 (135)

V.5 Logistic regression

The logistic loss is a combination of the cross entropy loss $l(y, z) = -y \log(\sigma(z)) - (1 - y) \log(1 - \sigma(z))$ with as sigmoid activation function σ , that simplifies for binary labels $y \pm 1$ to $l^{\text{logistic}}(y, z) = \log(1 + \exp(-yz))$ with the two first derivatives given by

$$\partial_z l^{\text{logistic}}(y,z) = -\frac{y}{e^{zy}+1}\,,\qquad \partial_z^2 l^{\text{logistic}}(y,z) = \frac{y^2}{2(1+\cosh{(zy)})} = \frac{y^2}{4\cosh{\left(\frac{yz}{2}\right)}}\,.$$

Its proximal is not analytical, but it can be written as the solution of the implicit equation (45) providing the corresponding denoising functions (46). Solving the fixed point equations (97), we obtain performances that approach closely the Bayes-optimal baseline as illustrated in Fig. 4 (left).



Figure 4: (Left) Logistic regression - Generalization error as a function of α for different regularizations strength λ . Decreasing λ , the generalization error approaches very closely the Bayes-optimal error (black line). The difference with the Bayes error is shown as an inset. Logistic flirts with Bayes error but never achieves it exactly. The asymptotic behaviour is compared to numerical logistic regression with $d = 10^3$ and averaged over $n_s = 20$ samples, performed with the default method *LogisticRegression* of the scikit-learn package [33]. (Right) Rectangle door teacher with $\kappa = 0.6745$ - Bayes-optimal generalization error (black) compared to asymptotic generalization performances of ℓ_2 logistic regression (dashed yellow line) and numerical ERM (crosses).

V.6 Logistic with non-linearly separable data - A rectangle door teacher

The analysis of ERM for the linearly separable dataset generated by (110) reveals that logistic regression with ℓ_2 regularization was able to approach very closely Bayes-optimal error. Therefore it seems us very interesting to investigate if logistic regression could perform as well on a more complicated non-linearly separable dataset obtained by a *rectangle door* channel

$$\mathbf{y} = \operatorname{sign}\left(\frac{1}{\sqrt{d}}\mathbf{X}\mathbf{w}^{\star} - \kappa\right). \tag{136}$$

This channel has been already considered in [10] and we fix the width of the door to $\kappa = 0.6745$ to obtain labels ± 1 with probability 0.5. We then compare the ERM performances of logistic regression with ℓ_2 regularization to the Bayes-optimal performances given by (111) with denoising functions derived in eq. (40). We show in Fig. 4 (**right**) the comparison only for an arbitrary hyper-parameter $\lambda = 1.10^{-2}$, as results are similar for any regularization. As we might expect, the logistic regression is not able to reach the Bayes-optimal generalization error. Both Bayes-optimal and ERM performances are stuck in the symmetric fixed point m = 0 up to $\alpha_{\rm it} \simeq 1.393$. Above this threshold it becomes unstable and Bayes error decreases to zero in the $\alpha \rightarrow 0$ limit, while the logistic regression with arbitrary λ remains stuck to its maximal generalization error, meaning that in this non-linearly separable case, the logistic regression largely underperforms Bayes-optimal performances.

VI Reaching Bayes optimality

In this section, we propose a derivation inspired by [24, 37-39, 51, 52, 56-59] of the fine-tuned loss and regularizer (17) discussed in Sec. 4. We assume that the dataset is generated by a teacher (18) such that $\mathcal{Z}_{out^*}(., \omega, .)$ and $\mathcal{Z}_{w^*}(\gamma, .)$ are respectively log-concave in ω and γ . The derivation is based on the GAMP algorithm introduced in [30] for the model eq. (1), that we start by recalling.

VI.1 Generalized Approximate Message Passing (GAMP) algorithm

The GAMP algorithm can be written as the following set of iterative equations that depend on the update functions (23):

$$\begin{cases} \hat{\mathbf{w}}^{t+1} = f_{\mathbf{w}}(\boldsymbol{\gamma}^{t}, \Lambda^{t}) \\ \hat{\mathbf{c}}_{\mathbf{w}}^{t+1} = \partial_{\boldsymbol{\gamma}} f_{\mathbf{w}}(\boldsymbol{\gamma}^{t}, \Lambda^{t}) \\ \mathbf{f}_{\mathbf{out}}^{t} = f_{\mathbf{out}}\left(y, \boldsymbol{\omega}^{t}, V^{t}\right) \end{cases} \text{ and } \begin{cases} \Lambda_{i}^{t} = -\frac{1}{d} \sum_{\mu=1}^{n} \mathbf{X}_{\mu i}^{2} \partial_{\boldsymbol{\omega}} f_{\mathbf{out},\mu}^{t} \\ \boldsymbol{\gamma}_{i}^{t} = \frac{1}{\sqrt{d}} \sum_{\mu=1}^{n} \mathbf{X}_{\mu i} f_{\mathbf{out},\mu}^{t} + \Lambda_{i}^{t} \hat{w}_{i}^{t} \\ V_{\mu}^{t} = \frac{1}{d} \sum_{i=1}^{d} \mathbf{X}_{\mu i}^{2} \hat{c}_{w,i}^{t} \\ \omega_{\mu}^{t} = \frac{1}{\sqrt{d}} \sum_{i=1}^{d} \mathbf{X}_{\mu i} \hat{w}_{i}^{t} - V_{\mu}^{t} f_{\mathbf{out},\mu}^{t-1} \end{cases}$$
(137)

It has been proven in [53] that the GAMP algorithm with Bayes-optimal update functions $f_{\rm w} = f_{\rm w^{\star}}$ and $f_{\rm out} = f_{\rm out^{\star}}$ (25) converges to the Bayes-optimal performances in the large size limit. Yet the GAMP denoising functions are generic and can be chosen as will depending on the statistical estimation method. In particular we may choose the denoising functions for Bayes-optimal estimation (25) or the ones corresponding to ERM estimation (29)

$$f_{w}^{\text{bayes}}(\gamma, \Lambda) = \partial_{\gamma} \log \left(\mathcal{Z}_{w^{\star}} \right) ,$$

$$f_{\text{out}}^{\text{bayes}}(y, \omega, V) = \partial_{\omega} \log \left(\mathcal{Z}_{\text{out}^{\star}} \right) ,$$

$$f_{w}^{\text{erm}, r}(\gamma, \Lambda) = \Lambda^{-1} \gamma - \Lambda^{-1} \partial_{\Lambda^{-1} \gamma} \mathcal{M}_{\Lambda^{-1}} \left[r(.) \right] \left(\Lambda^{-1} \gamma \right) ,$$

$$f_{\text{out}}^{\text{erm}, l}(y, \omega, V) = -\partial_{\omega} \mathcal{M}_{V} [l(y, .)](\omega) ,$$
(138)

whose corresponding GAMP algorithms (137) will achieve potentially different fixed points and thus different performances. As it is proven that GAMP with Bayes-optimal updates lead to the optimal generalization error, so that ERM matches the same performances it is sufficient to enforce that at each time step t the Bayes-optimal and ERM denoising functions are equal $f^{\text{bayes}} = f^{\text{erm}}$. Enforcing these two constraints will lead to the expressions for the optimal loss l^{opt} and regularizer r^{opt} , so that ERM matches Bayes-optimal performances.

VI.2 Matching Bayes-optimal and ERM performances

Imposing the equality on the channel updates we obtain

$$f_{\text{out}}^{\text{bayes}}\left(y,\omega,V\right) = f_{\text{out}}^{\text{erm},l}\left(y,\omega,V\right) \Leftrightarrow \partial_{\omega}\log\left(\mathcal{Z}_{\text{out}^{\star}}\right)\left(y,\omega,V\right) = -\partial_{\omega}\mathcal{M}_{V}\left[l^{\text{opt}}\left(y,.\right)\right]\left(\omega\right)$$

Integrating, leaving aside the constant that will not influence the final result, and taking the Moreau-Yosida regularization on both sides, we obtain:

$$\mathcal{M}_{V}\left[\log \mathcal{Z}_{\text{out}^{\star}}\left(y,.,V\right)\right]\left(\omega\right) = \mathcal{M}_{V}\left[-\mathcal{M}_{V}\left[l^{\text{opt}}\left(y,.\right)\right]\left(\omega\right)\right] = -l^{\text{opt}}\left(y,\omega\right),$$

where we invert the Moreau-Yosida regularization in the last equality that is valid as long as $\mathcal{Z}_{\text{out}^{\star}}(y, \omega, V)$ is assumed to be log-concave in ω , (see [39] for a derivation). We finally obtain

$$l^{\text{opt}}(y,z) = -\mathcal{M}_{V}\left[\log\left(\mathcal{Z}_{\text{out}^{\star}}\right)(y,.,V)\right](z) = -\min_{\omega}\left(\frac{(z-\omega)^{2}}{2V} + \log\mathcal{Z}_{\text{out}^{\star}}(y,\omega,V)\right).$$
(139)

Let us perform the same computation for the prior updates. First we introduce a rescaled denoising distribution:

$$\tilde{Q}_{\mathbf{w}^{\star}}(w;\gamma,\Lambda) \equiv \frac{1}{\tilde{\mathcal{Z}}_{\mathbf{w}^{\star}}(\gamma,\Lambda)} P_{\mathbf{w}^{\star}}(w) e^{-\frac{1}{2}\Lambda \left(w-\Lambda^{-1}\gamma\right)^{2}},
\log\left(\tilde{\mathcal{Z}}_{\mathbf{w}^{\star}}(\gamma,\Lambda)\right) = \log\left(\mathcal{Z}_{\mathbf{w}^{\star}}(\gamma,\Lambda)\right) - \frac{1}{2}\Lambda^{-1}\gamma^{2},$$
(140)

so that the the prior updates read

$$f_{w}^{\text{bayes}}(\gamma,\Lambda) = \partial_{\gamma}\log\left(\mathcal{Z}_{w^{\star}}\right) = \Lambda^{-1}\gamma + \Lambda^{-1}\partial_{\Lambda^{-1}\gamma}\log\left(\tilde{\mathcal{Z}}_{w^{\star}}\right),$$

$$f_{w}^{\text{erm},r}(\gamma,\Lambda) = \mathcal{P}_{\Lambda^{-1}}\left[r\right]\left(\Lambda^{-1}\gamma\right) = \Lambda^{-1}\gamma - \Lambda^{-1}\partial_{\Lambda^{-1}\gamma}\mathcal{M}_{\Lambda^{-1}}\left[r\right]\left(\Lambda^{-1}\gamma\right).$$
(141)

Imposing the equivalence of the Bayes-optimal and ERM prior update,

$$f_{\mathbf{w}}^{\mathrm{bayes}}\left(\gamma,\Lambda\right) = f_{\mathbf{w}}^{\mathrm{erm},r}\left(\gamma,\Lambda\right) \Leftrightarrow \partial_{\Lambda^{-1}\gamma}\log\left(\tilde{\mathcal{Z}}_{\mathbf{w}^{\star}}\right) = -\partial_{\Lambda^{-1}\gamma}\mathcal{M}_{\Lambda^{-1}}\left[r^{\mathrm{opt}}\right]\left(\Lambda^{-1}\gamma\right), \quad (142)$$

and assuming that $\mathcal{Z}_w(\gamma, \Lambda)$ is log-concave in γ , we may invert the Moreau-Yosida regularization, that leads to:

$$r^{\text{opt}}\left(\Lambda^{-1}\gamma\right) = -\mathcal{M}_{\Lambda^{-1}}\left[\log\left(\tilde{\mathcal{Z}}_{w^{\star}}\right)\left(.,\Lambda^{-1}\right)\right](w)$$

$$= -\min_{\Lambda^{-1}\gamma}\left(\frac{(w-\Lambda^{-1}\gamma)^{2}}{2\Lambda^{-1}} + \log\tilde{\mathcal{Z}}_{w^{\star}}\left(\gamma,\Lambda\right)\right) = -\min_{\gamma}\left(\frac{1}{2}\Lambda w^{2} - \gamma w + \log\mathcal{Z}_{w^{\star}}\left(\gamma,\Lambda\right)\right).$$
(143)

The last step, is to characterize the variances V and Λ involved in (139) and (143) that are so far undetermined. To achieve the Bayes-optimal performances, we therefore need to used the variances V and Λ solutions of the Bayes-optimal GAMP algorithm (137). In the large size limit, these quantities concentrate and are giveen by the State Evolution (SE) of the GAMP algorithm, that we recall herein.

State evolution of GAMP In the large size limit, the expectation of the parameter V and Λ over the ground truth \mathbf{w}^* and the input data X lead to [53]:

$$\mathbb{E}_{\mathbf{w}^{\star},\mathbf{X}}\left[V\right] = \rho_{\mathbf{w}^{\star}} - q_{\mathbf{b}}, \qquad \qquad \mathbb{E}_{\mathbf{w}^{\star},\mathbf{X}}\left[\Lambda\right] = \hat{q}_{\mathbf{b}}, \qquad (144)$$

where $q_{\rm b}$ and $\hat{q}_{\rm b}$ are solutions of the Bayes-optimal set of fixed point equations eq. (13).

VI.3 Summary and numerical evidences

Choosing the fine-tuned (potentially non-convex depending on Z_{out^*} and Z_{w^*}) loss and regularizer

$$l^{\text{opt}}(y,z) = -\min_{\omega} \left(\frac{(z-\omega)^2}{2(\rho_{w^{\star}}-q_{b})} + \log \mathcal{Z}_{\text{out}^{\star}}(y,\omega,\rho_{w^{\star}}-q_{b}) \right)$$

$$r^{\text{opt}}(w) = -\min_{\gamma} \left(\frac{1}{2} \hat{q}_{b} w^2 - \gamma w + \log \mathcal{Z}_{w^{\star}}(\gamma,\hat{q}_{b}) \right)$$
(145)

with $q_{\rm b}$ and $\hat{q}_{\rm b}$ are solutions of the Bayes-optimal set of fixed point equations eq. (13), we showed that ERM can provably match the Bayes-optimal performances. In particular we illustrated the behaviour of the optimal loss and regularizer $\lambda^{\rm opt}$ and $r^{\rm opt}$ for the model (2) in Fig. 2 of the main text. Note in particular that even though the loss $l^{\rm opt}$ is not convex (but seems quasi-convex), numerical simulations of ERM with (145) (black dots) presented in Fig. 5 show that ERM achieves indeed the Bayes-optimal performances (black line) even at finite dimension.



Figure 5: Generalization error obtained by optimization of the optimal loss l^{opt} and r^{opt} for the model (2), compared to ℓ_2 logistic regression and Bayes-optimal performances. Numerics has been performed with scipy.optimize.minimize with the L-BFGS-B solver for $d = 10^3$ and averaged over $n_s = 10$ instances. The error bars are barely visible.