

The Mutual Information in Random Linear Estimation Beyond i.i.d. Matrices

Jean Barbier^{†*} and Nicolas Macris[†]

[†] Communication Theory Laboratory, EPFL, Switzerland.

* International Center for Theoretical Physics, Trieste, Italy.

Antoine Maillard and Florent Krzakala

LPS ENS, CNRS, PSL, UPMC & Sorbonne Université,

Paris, France.

Abstract—There has been definite progress recently in proving the variational single-letter formula given by the heuristic replica method for various estimation problems. In particular, the replica formula for the mutual information in the case of noisy linear estimation with random i.i.d. matrices, a problem with applications ranging from compressed sensing to statistics, has been proven rigorously. In this contribution we go beyond the restrictive i.i.d. matrix assumption and discuss the formula proposed by Takeda, Uda, Kabashima and later by Tulino, Verdu, Caire and Shamai who used the replica method. Using the recently introduced adaptive interpolation method and random matrix theory, we prove this formula for a relevant large sub-class of rotationally invariant matrices.

Few problems are as ubiquitous in computer science as the one of random linear estimation, that plays a fundamental role in machine learning [1], statistics [2] and communication [3]. Computing the information theoretic limitation for the estimation of a signal given the knowledge of its random linear projections has many applications, e.g., compressed sensing [2], code division multiple access (CDMA) [4] or error correcting codes [5, 6]. The problem is defined as follows: Consider a *signal* vector $\mathbf{X} \in \mathbb{R}^n$ with i.i.d. entries distributed according to a “prior” P_0 over \mathbb{R} with bounded support (an hypothesis that can be relaxed). One is given m measurements

$$Y_\mu = \sqrt{\frac{\lambda}{n}}(\Phi\mathbf{X})_\mu + Z_\mu, \quad \mu = 1, \dots, m, \quad (1)$$

in which $\lambda > 0$ is the *signal to noise ratio* (snr), $\mathbf{Z} = (Z_\mu)_{\mu=1}^m \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$ is a Gaussian noise and $\Phi \in \mathbb{R}^{m \times n}$ is the *measurement matrix*. We will consider the “high-dimensional limit”, namely $m, n \rightarrow \infty$ such that $\alpha \equiv m/n$ stays finite.

There has been a considerable amount of work on this model in the case where Φ is a random matrix whose elements are i.i.d. standard Gaussian. In particular a pionnering work by Tanaka [4] using a statistical physics calculation and the replica method [7] proposed a single-letter formula for the normalized mutual information between the measurements and the signal $i_n \equiv \frac{1}{n}I(\mathbf{X}; \mathbf{Y}|\Phi)$. The so-called Tanaka formula, originally written for the CDMA problem (where each element of \mathbf{X} is taken i.i.d. from ± 1), has been generalized and applied to many problems (e.g. in [8, 9]). After nearly 15 years, it has been proven with different approaches [10–13], a spectacular confirmation of the replica calculation.

In this paper, we endeavour to go beyond the very restrictive simple i.i.d. measurement matrix assumption and consider

instead a more complex situation where Φ is now taken from a non-trivial, and correlated, random ensemble. This is both more relevant to practical cases, and more realistic in terms of modeling real statistical situations.

I. RESULT AND RELATED WORKS

The random linear estimation problem with correlated, non i.i.d. matrices has been considered, again with the replica method, and a solution was proposed by Takeda, Uda and Kabashima in the case of CDMA [14] for matrices taken from a *rotationally invariant ensemble*. The replica formula was later extended by Tulino, Caire, Verdu and Shamai [9] for such similar ensembles in order to allow for more complicated priors such as Gauss-Bernoulli ones. These works point to a generic conjecture, that we now describe, giving the single-letter formula for the mutual information.

Consider a measurement matrix $\Phi = \mathbf{O}\Sigma\mathbf{N}^\top$ where \mathbf{O} and \mathbf{N} are both orthogonal matrices and Σ is diagonal (non-square if $\alpha \neq 1$). The matrices Σ , \mathbf{N} , \mathbf{O} are independent and \mathbf{N} is distributed uniformly, according to the *Haar measure* of its orthogonal group, in which case Φ is said to be *right-rotationally invariant* (if \mathbf{O} is also Haar distributed the ensemble is simply called rotationally invariant). The matrix $\mathbf{R} \equiv \frac{1}{n}\Phi^\top\Phi = \frac{1}{n}\mathbf{N}\Sigma^\top\Sigma\mathbf{N}^\top$ plays an important role. For general rotationally invariant ensembles its eigenvalues are not necessarily i.i.d. but typically are such that $\frac{1}{n}\Sigma^\top\Sigma$ has a suitable limiting eigenvalue distribution (as in fact assumed in [14]). The limit of the normalized mutual information is conjectured to be

$$\lim_{n \rightarrow \infty} i_n = \inf_{(E, r) \in \Gamma} i_{\text{RS}}(E, r; \lambda) \quad (2)$$

where the so-called replica symmetric *potential* $i_{\text{RS}}(E, r; \lambda) = i_{\text{RS}}$ is defined as

$$i_{\text{RS}} \equiv I(X; \sqrt{r}X + Z) + \frac{1}{2} \int_0^{\lambda E} \mathcal{R}_{\mathbf{R}}(-z) dz - \frac{rE}{2}. \quad (3)$$

Here we have used the *R-transform* of the matrix \mathbf{R} (see [3] for an introduction to such transforms). $I(X; \sqrt{r}X + Z)$, that we simply denote $I(r)$ later on, is the mutual information for the scalar Gaussian channel $Y = \sqrt{r}X + Z$, $X \sim P_0$ and $Z \sim \mathcal{N}(0, 1)$. Moreover, Γ is the set of critical points of the potential, or *state evolution* fixed points ($\rho \equiv \mathbb{E}_{P_0}[X^2]$)

$$\Gamma \equiv \left\{ (E, r) \in [0, \rho] \times \mathbb{R}_+ \mid \begin{aligned} E &= \text{mmse}(X|\sqrt{r}X + Z), \\ r &= \lambda \mathcal{R}_{\mathbf{R}}(-\lambda E) \end{aligned} \right\}.$$

A virtue of the formula (3) is that the details of the rotation invariant matrix ensemble only enter through the R-transform.

Note the slight difference with the potential written in [9]: In the potential (3) a factor 1/2, not present in [9], multiplies both the integrated R-transform and $-rE$. This is because they consider the complex Φ case while we consider the real case (nevertheless, we believe that our proof techniques could easily be generalized to include the complex case).

Interestingly, Manoel *et al* [15] and Reeves [16] were able to formally re-derive equivalent results independently for a subclass of rotationally invariant matrices by considering multi-layered estimation problems. This line of work combined with our present rigorous work is giving a lot of credibility to the replica conjecture for the general rotation invariant ensemble.

A. Main result

Our main result is a proof of the replica conjecture (2) for a specific, but large, set of correlated matrices. We hope it paves the way towards a completely general proof, as the non-rigorous replica calculation only assumes right-rotational invariance of Φ . We assume that the $m \times n$ matrix Φ can be decomposed as follows:

$$\Phi = \Phi' \mathbf{W} \quad (4)$$

in which all elements of the $m \times n$ matrix \mathbf{W} are i.i.d. Gaussians $\mathcal{N}(0, 1/n)$ with mean zero and variance $1/n$, Φ' is a $m \times m$ random matrix, and \mathbf{W} , Φ' are independent. Concerning Φ' , our analysis is currently complete under the assumption that it is a product of a finite number of independent matrices, each with i.i.d. matrix-elements that are either bounded or standard Gaussians. The case of a product of i.i.d. standard Gaussian matrices constitutes an interesting example that has been considered in [15].

Our goal here is to give a rigorous proof of (2) in the setting of matrices (4) with the independence assumptions for Φ' . Some of our technical calculations are based on previous related works, and all of them will be discussed in a complete manner in a longer contribution.

Theorem 1.1 (Replica formula): Assume that the prior P_0 has compact support or, in other words, the signal is bounded. Then with $n, m, \Phi, \mathbf{R}, i_{\text{RS}}$ defined as before (in particular Φ satisfies (4) and the subsequent hypothesis), one has

$$\lim_{n \rightarrow \infty} i_n = \inf_{r \geq 0} \sup_{E \in [0, \rho]} i_{\text{RS}}(E, r; \lambda).$$

Remark 1.1 (Equivalent expressions of the replica formula): The right hand side in the theorem above can also be written as $\inf_{(E, r) \in \Gamma} i_{\text{RS}}(E, r; \lambda)$ (if the extremizers are not attained at the boundaries, which is the case when noise is present) or as $\inf_{E \in [0, \rho]} \sup_{r \geq 0} i_{\text{RS}}(E, r; \lambda)$. A proof of such equivalences, in the case of generalized linear estimation, is found in [13].

Remark 1.2 (Right-rotation invariance of the ensemble): Note that (4) implies the right-rotation invariance of Φ because of the rotation invariance of \mathbf{W} . Our result thus covers a subclass of right-rotationally invariant matrices.

Remark 1.3 (Relaxing the assumptions on Φ'): The assumptions on Φ' come from the fact that a complete rigorous analysis requires proving the concentration of the “free energy” of an *interpolating* model (see below; the free energy is equal to the mutual information up to a trivial additive constant). This involves the use of tools such as the McDiarmid bounded difference inequality and/or the Gaussian Poincaré inequality, which require some independence between degrees of freedom. If one *assumes* concentration of the free energy of the interpolating model then one can relax these assumptions to the following more general ones. It then suffices to assume that $\frac{1}{n} \Phi'^T \Phi'$ has a well-defined, positively and compactly supported, asymptotic eigenvalue distribution in the limit $n \rightarrow \infty$. We also remark that concentration proofs of the free energy of the original and interpolated models are technically similar, however we have not established a purely logical implication between the two. If such an implication holds one could also replace the independence assumption on Φ' by an assumption of concentration of the free energy.

Remark 1.4 (Relaxing the assumption on P_0): Boundedness of the signal is again used to obtain concentration results for the free energy but it can presumably be removed using a limiting argument as in [17].

B. Related works

There has been a lot of effort recently [10–13] to prove the Tanaka formula for random i.i.d. matrices. Our strategy in the present paper is to follow the *adaptive interpolation method* introduced in [18, 19]. This method, in particular, has been used by the authors of [13] to reach a rigorous demonstration of the replica formula for the mutual information for the case of i.i.d. measurement matrix Φ , in the more general situation of “generalized linear estimation” i.e., with an *arbitrary* measurement channel (instead of just a random additive noise as in (1)). Some steps of our current approach consequently follow similar ones in [13] (but with key differences) and we will refer to this work when necessary. We believe, in fact, that the approach presented in the present paper could further be generalized as well to generalized linear estimation with rotationally invariant matrices to reach the formula conjectured by Kabashima in [20]. This is left for future work.

Perhaps the most important consequence of the replica formula is that it predicts the value of the minimum mean-square error (MMSE) in the reconstruction of the signal \mathbf{X} . In fact, it is conjectured (and proved for Gaussian matrices [11, 12]) to be given by the value E that extremizes (3). While we consider here mainly the information theoretic result, a large body of work has focused on algorithmic approaches to random linear estimation, and investigated whether the MMSE is efficiently (say, in polynomial time) achievable.

For Gaussian matrices, the most successful approach, so far, again originated in statistical physics [21, 22] and is called approximate message-passing (AMP) [23]. AMP is Bayes-optimal and efficiently achieves the MMSE for a large set of parameters, as proven in [11]. There, however, might exist a region called “hard” where this is not the case, and polynomial

algorithms improving on AMP are not known. Whether there exists an efficient algorithm that is able to beat AMP in the hard region is widely considered to be a notoriously difficult problem (see e.g. [24] and reference therein).

For rotationally invariant matrices, different but related approaches were proposed [14, 25]. In particular, the general expectation-propagation (EP) [26] leads to a powerful scheme in this context [27]. Recently Ma and Ping proposed a variation of EP called OAMP [28] specially adapted to these matrices. Rangan, Schniter and Fletcher introduced a related approach called VAMP [29] and showed that it follows the fixed point equation (called state evolution) of the potential (3). Interestingly, the multi-layer AMP algorithm of Manoel *et al.* [15] also has the same fixed point. Our result thus supports that OAMP, VAMP (and multi-layer AMP) are Bayes-optimal and efficiently reach the MMSE in the “easy” region of random linear estimation with these correlated matrices, just as AMP does in the case of i.i.d. Gaussian matrices.

II. PROOF BY THE ADAPTIVE INTERPOLATION METHOD

We give here the main steps of the proof of Theorem 1.1. We will use the adaptive interpolation method, introduced in [18], and then applied in [19] and [13]. It is a powerful evolution of the interpolation method developed by Guerra and Toninelli in the context of spin glasses [30]. Many steps of the proof follow the ones of [13], and we will refer to them when necessary.

A. Interpolating estimation problem

Let us fix a sequence $s_n \in (0, 1/2]$ that goes to 0 as n goes to infinity. Let $\epsilon = (\epsilon_1, \epsilon_2) \in [s_n, 2s_n]^2$ (so that ϵ actually depends on n , but we will drop this dependency for clarity).

Let $E : [0, 1] \rightarrow [0, \rho]$ and $r : [0, 1] \rightarrow \mathbb{R}_+$ be two continuous “interpolation functions” (that will later depend on ϵ), and $R_1(t) \equiv \epsilon_1 + \int_0^t r(v)dv$, $R_2(t) \equiv \epsilon_2 + \int_0^t E(v)dv$ for $t \in [0, 1]$. Consider the following two t -dependent observation channels for $i = 1, \dots, n$ and $\mu = 1, \dots, m$:

$$\begin{cases} Y_{t,\mu} &= \sqrt{\frac{\lambda(1-t)}{n}}(\Phi\mathbf{X})_\mu + \sqrt{\frac{\lambda}{n}R_2(t)}(\Phi'\mathbf{V})_\mu + Z_{\mu}, \\ \tilde{Y}_{t,i} &= \sqrt{R_1(t)}X_i + \tilde{Z}_i. \end{cases} \quad (5)$$

In the following we assume $\lambda = 1$, as it amounts to a scaling of Φ . Here (\tilde{Z}_i) , (Z_μ) , $(V_\mu) \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$ whereas $(X_i) \stackrel{\text{i.i.d.}}{\sim} P_0$. The inference problem is to recover both $\mathbf{X} = (X_i)_{i=1}^n$ and $\mathbf{V} = (V_\mu)_{\mu=1}^m$ from the knowledge of the observations \mathbf{Y}_t , $\tilde{\mathbf{Y}}_t$ and the matrix Φ (and thus of Φ' and \mathbf{W} too as the decomposition (4) is assumed to be known).

In the Bayesian setting the posterior associated with this inference problem, written in the Gibbs-Boltzmann form, is

$$dP_{t,\epsilon}(\mathbf{x}, \mathbf{v} | \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi) = \frac{dP_0(\mathbf{x})\mathcal{D}\mathbf{v}e^{-\mathcal{H}(\epsilon,t,\mathbf{x},\mathbf{v};\mathbf{Y}_t,\tilde{\mathbf{Y}}_t,\Phi)}}{\int dP_0(\mathbf{x})\mathcal{D}\mathbf{v}e^{-\mathcal{H}(\epsilon,t,\mathbf{x},\mathbf{v};\mathbf{Y}_t,\tilde{\mathbf{Y}}_t,\Phi)}}, \quad (6)$$

where $\mathcal{D}\mathbf{v} \equiv \prod_{\mu=1}^m dv_\mu (2\pi)^{-1/2} e^{-v_\mu^2/2}$ is the standard Gaussian measure, and we have defined the interpolating *Hamiltonian* $\mathcal{H} = \mathcal{H}(\epsilon, t, \mathbf{x}, \mathbf{v}; \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi)$ as

$$\mathcal{H} \equiv \frac{1}{2} \|\mathbf{Y}_t - \sqrt{\frac{1-t}{n}}\Phi\mathbf{x} - \sqrt{\frac{R_2(t)}{n}}\Phi'\mathbf{v}\|_2^2 + \frac{1}{2} \|\tilde{\mathbf{Y}}_t - \sqrt{R_1(t)}\mathbf{x}\|_2^2,$$

It is a simple exercise (see e.g. [11]) to show that the normalized mutual information $i_n(t) \equiv \frac{1}{n}I(\mathbf{X}, \mathbf{V}; \mathbf{Y}_t, \tilde{\mathbf{Y}}_t | \Phi)$ for the interpolation estimation problem is related to the posterior normalization (or partition function) through

$$i_{n,\epsilon}(t) = -\frac{1}{n}\mathbb{E}\ln \int dP_0(\mathbf{x})\mathcal{D}\mathbf{v}e^{-\mathcal{H}(\epsilon,t,\mathbf{x},\mathbf{v};\mathbf{Y}_t,\tilde{\mathbf{Y}}_t,\Phi)} - \frac{\alpha+1}{2}. \quad (7)$$

One can verify that this *interpolating mutual information* satisfies:

$$\begin{cases} i_{n,\epsilon}(0) = i_n + o_n(1), \\ i_{n,\epsilon}(1) = I(R_1(1)) + \frac{\alpha}{2}\mathbb{E}_{X'}\ln(1+R_2(1)X') + o_n(1), \end{cases} \quad (8)$$

where $X' \sim p'_n$, with p'_n the empirical spectral distribution of the $m \times m$ matrix $\frac{1}{n}\Phi'^T\Phi'$. Here, $o_n(1) \rightarrow 0$ in the $n \rightarrow \infty$ limit, uniformly in E, r, t, ϵ . The second term in the expression of $i_{n,\epsilon}(1)$ (sometimes referred to as a Shannon transform, see e.g. [3]) is obtained using the celebrated “log-det formula” for the mutual information $\frac{1}{n}I(\mathbf{Y}_1; \mathbf{V} | \Phi')$ of an i.i.d. Gaussian input multiplied by the matrix Φ' and under additive Gaussian noise, see e.g. [3].

Now a crucial step in our proof, that is a consequence of the particular form of the measurement matrix (4), is that as n grows, the second term in $i_{n,\epsilon}(1)$ can be replaced by an integrated R -transform. Denoting $G_{\mathbf{R}}(x) \equiv \int_0^x \mathcal{R}_{\mathbf{R}}(-u)du$:

$$i_{n,\epsilon}(1) = I(R_1(1)) + \frac{1}{2}G_{\mathbf{R}}(R_2(1)) + o_n(1), \quad (9)$$

where $\mathcal{R}_{\mathbf{R}}(z)$ is the \mathbf{R} -transform associated with the asymptotic spectrum of $\mathbf{R} = \frac{1}{n}\Phi'^T\Phi'$. We give the definition of this transform as well as the proof of (9) in the next section.

A word about notations: We define the Gibbs bracket $\langle \cdot \rangle_{t,\epsilon}$ as the expectation w.r.t. the posterior (6). In contrast, we denote by \mathbb{E} the joint expectation w.r.t. all *quenched* variables (i.e. fixed by the realization of the problem), namely $(\mathbf{X}, \mathbf{V}, \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi', \mathbf{W})$, or equivalently w.r.t. $(\mathbf{X}, \mathbf{V}, \mathbf{Z}, \tilde{\mathbf{Z}}, \Phi', \mathbf{W})$.

B. Useful tools from random matrix theory

In this paragraph we show how to deduce (9) from (8). Note

$$\mathbb{E}_{X'} \ln(1 + R_2(1)X') = \mathbb{E}_{X'} \int_0^{R_2(1)} du \frac{X'}{1+uX'}.$$

The result thus follows if the following relation is true :

$$\mathcal{R}_{\mathbf{R}}(-u) = \alpha \mathbb{E}_{X'} \left[\frac{X'}{1+uX'} \right] + o_n(1). \quad (10)$$

This is a well known relation in random matrix theory. For matrices of the form $\mathbf{R} = \mathbf{W}^T\mathbf{T}\mathbf{W}$, where \mathbf{W} is a Gaussian $m \times n$ matrix with i.i.d. $\mathcal{N}(0, 1/n)$ elements, and $\mathbf{T} = \frac{1}{n}\Phi'^T\Phi'$ a non negative $m \times m$ matrix with a limiting spectral distribution, this was already shown by Marcenko and Pastur in 1967 [31] in the language of the Stieltjes transform. See also [32] for generalizations. Denoting by $g_{\mathbf{R}}(z)$ and $g_{\mathbf{T}}(z)$ (z a complex number outside the spectrum of the matrices) the limiting Stieltjes transforms of the matrices \mathbf{R} and \mathbf{T} , we have [31, 32] that the Marcenko-Pastur formula takes the form

$$z(g_{\mathbf{R}}(z))^2 + \alpha g_{\mathbf{T}}(-1/g_{\mathbf{R}}(z)) + (1 - \alpha)g_{\mathbf{R}}(z) = 0.$$

Simple algebra then implies ($g_{\mathbf{R}}^{-1}$ is the inverse function)

$$g_{\mathbf{R}}^{-1}(z) + z^{-1} = -\alpha z^{-2}(g_{\mathbf{T}}(-z^{-1}) - z).$$

Since by definition $\mathcal{R}_{\mathbf{R}}(z) \equiv g_{\mathbf{R}}^{-1}(-z) - z^{-1}$ and $g_{\mathbf{T}}(z) \equiv \int \frac{d\tau(x)}{x-z}$, $d\tau$ being the *limiting* eigenvalue distribution of \mathbf{T} , this relation is nothing else than $\mathcal{R}_{\mathbf{R}}(-z) = \alpha \int d\tau(x) \frac{x}{1+zx}$ which is equivalent to (10) when $n \rightarrow \infty$. We refer to the review [33] for a more modern discussion using free probability concepts.

C. Mutual information variation

In order to “compare” the potential (3) with the mutual information, we use the trivial identity

$$i_n = i_{n,\epsilon}(0) + \mathcal{O}_n(1) = i_{n,\epsilon}(1) - \int_0^1 i'_{n,\epsilon}(t) dt + \mathcal{O}_n(1)$$

which becomes, using (9),

$$i_n = I(R_1(1)) + \frac{1}{2} G_{\mathbf{R}}(R_2(1)) - \int_0^1 i'_{n,\epsilon}(t) dt + \mathcal{O}_n(1). \quad (11)$$

We now evaluate $i'_{n,\epsilon}(t)$. Define $Q \equiv \frac{1}{n} \sum_{i=1}^n X_i x_i$, called the *overlap*, and the vector $\mathbf{u}_t = (u_{t,\mu})_{\mu=1}^m$ with

$$u_{t,\mu} \equiv \sqrt{\frac{1-t}{n}} (\Phi(\mathbf{X} - \mathbf{x}))_{\mu} + \sqrt{\frac{R_2(t)}{n}} (\Phi'(\mathbf{V} - \mathbf{v}))_{\mu} + Z_{\mu}.$$

Lemma 2.1 (*Mutual information t -variation*): For $t \in (0, 1)$

$$i'_{n,\epsilon}(t) = \frac{r(t)}{2} (\rho - \mathbb{E}\langle Q \rangle_{t,\epsilon}) + \mathcal{O}_n(1) + \frac{1}{2n^2} \mathbb{E}\langle \mathbf{Z}^{\top} (\Phi' \Phi'^{\top}) \mathbf{u}_t [E(t) - (\rho - Q)] \rangle_{t,\epsilon} \quad (12)$$

The proof of this lemma is very similar to the one found in [13]. The idea is to write explicitly the derivative $i'_{n,\epsilon}(t)$ and then to integrate by parts the quenched Gaussian variables \mathbf{V} and \mathbf{W} , before using the *Nishimori identity*. This identity is a consequence of Bayes rule and the fact that we consider the *optimal Bayesian setting*, namely that all hyperparameters in the problem such as the snr and P_0 are known and used when defining the posterior, see [13, 18, 24].

D. Overlap concentration

The next lemma essentially states that the overlap concentrates around its mean, and plays a key role in our proof. The proof technique for Bayesian inference has been developed in [11, 34–36] and is akin to the analysis reviewed for example in [37]. The point however here is that in Bayesian inference overlap concentration can be proved in the whole phase diagram. We will refer to [13, 18] where the analysis has been made quite generic. We now write explicitly the dependency of $R_1(t, \epsilon)$ and $R_2(t, \epsilon)$ on ϵ .

Lemma 2.2 (*Overlap concentration*): Assume that for any $t \in (0, 1)$ the map $\epsilon = (\epsilon_1, \epsilon_2) \in [s_n, 2s_n]^2 \mapsto R(t, \epsilon) = (R_1(t, \epsilon), R_2(t, \epsilon))$ is a C^1 diffeomorphism with Jacobian determinant greater or equal to 1. Then one can find a sequence s_n going to 0 slowly enough such that there exist positive constants C and γ that only depend on the support and moments of P_0 and on α , and such that:

$$\frac{1}{s_n^2} \int_{[s_n, 2s_n]^2} d\epsilon \int_0^1 dt \mathbb{E}\langle (Q - \mathbb{E}\langle Q \rangle_{t,\epsilon})^2 \rangle_{t,\epsilon} \leq C n^{-\gamma}.$$

We refer to [13, 18] for a detailed proof (in the case where Φ' is the identity matrix). For the present model under the assumptions on Φ' , this follows from Gaussian Poincaré and McDiarmid inequalities much as in [13]. As a consequence of this result, together with Lemma. 2.1, we obtain (using continuity and boundedness properties of the functions I and $G_{\mathbf{R}}$, see again [13] for more details):

Lemma 2.3 (*Fundamental identity*): Assume $\epsilon \mapsto R(t, \epsilon)$ satisfies the hypotheses of Lemma 2.2, and choose $s_n \rightarrow 0$ according to this lemma. Assume that for all $t \in [0, 1]$ and $\epsilon \in [s_n, 2s_n]^2$ we have $E(t, \epsilon) = \rho - \mathbb{E}\langle Q \rangle_{t,\epsilon}$. Then:

$$i_n = \frac{1}{s_n^2} \int_{[s_n, 2s_n]^2} d\epsilon \left\{ I\left(\int_0^1 r(t, \epsilon) dt\right) + \frac{1}{2} G_{\mathbf{R}}\left(\int_0^1 E(t, \epsilon) dt\right) - \frac{1}{2} \int_0^1 E(t, \epsilon) r(t, \epsilon) dt \right\} + \mathcal{O}_n(1),$$

in which $\mathcal{O}_n(1)$ is uniform in the choice of the functions E, r .

E. Upper and lower bounds

Similar bounds can be found in [13, 19]. We will often refer to [13] for more details. We first prove the upper bound:

Proposition 2.1: $\limsup_{n \rightarrow \infty} i_n \leq \inf_{r \geq 0} \sup_{E \in [0, \rho]} i_{\text{RS}}(E, r; 1)$.

Proof: Choose first $r(t) = r \geq 0$ a fixed value. We then fix $R = (R_1, R_2)$ as the solution $R(t, \epsilon) = (\epsilon_1 + rt, \epsilon_2 + \int_0^t E(s, \epsilon) ds)$ to the first order differential equation: $\partial_t R_1(t) = F_1$, $\partial_t R_2(t) = F_2(t, R(t))$, and $R(0) = \epsilon$, with $F_1 = r$, $F_2(t, R(t)) = \rho - \mathbb{E}\langle Q \rangle_{t,\epsilon}$ (it is easy to show that F_2 is in $[0, \rho]$, and thus E too). One can check (see [13]) that this ODE satisfies the hypotheses of the Cauchy-Lipschitz theorem. As $F = (F_1, F_2)$ is continuous and admits continuous partial derivatives, $R(t, \epsilon)$ is C^1 (in both arguments). By the Liouville formula, the Jacobian determinant $J_{n,\epsilon}(t)$ of $\epsilon \mapsto R(t, \epsilon)$ satisfies $J_{n,\epsilon}(t) = \exp\{\int_0^t \partial_{R_2} F_2(s, R(s, \epsilon)) ds\} \geq 1$. Indeed, the partial derivative $\partial_{R_2} F_2$ is non-negative, see Prop. 6 of [13]. Also, as this Jacobian never cancels, and as $\epsilon \mapsto R(t, \epsilon)$ is injective (by unicity of $R(t, \epsilon)$), it is a diffeomorphism by the inversion theorem. Recall $G_{\mathbf{R}}(x) \equiv \int_0^x \mathcal{R}_{\mathbf{R}}(-u) du$ and (3). Then Lemma. 2.3 implies:

$$i_n = \frac{1}{s_n^2} \int_{[s_n, 2s_n]^2} d\epsilon i_{\text{RS}}\left(\int_0^1 E(t, \epsilon) dt, r; 1\right) + \mathcal{O}_n(1),$$

that directly gives the desired bound. \blacksquare

We now turn to the lower bound:

Proposition 2.2: $\liminf_{n \rightarrow \infty} i_n \geq \inf_{r \geq 0} \sup_{E \in [0, \rho]} i_{\text{RS}}(E, r; 1)$.

Proof: Fix R as the solution $R(t, \epsilon) = (\epsilon_1 + \int_0^t r(s, \epsilon) ds, \epsilon_2 + \int_0^t E(s, \epsilon) ds)$ to the following Cauchy problem: $\partial_t R_1(t) = F_1(t, R(t)) = \mathcal{R}_{\mathbf{R}}(\mathbb{E}\langle Q \rangle_{t,\epsilon} - \rho)$, $\partial_t R_2(t) = F_2(t, R(t)) = \rho - \mathbb{E}\langle Q \rangle_{t,\epsilon}$ and $R(0) = \epsilon$. Let us denote this equation $\partial_t R(t) = F(t, R(t))$ (F also depends on n). Note that this implies that the solutions verify $E(t, \epsilon) = \rho - \mathbb{E}\langle Q \rangle_{t,\epsilon} \in [0, \rho]$ and $r(t, \epsilon) \geq 0$. It is possible to verify (see the details in a similar case in [13]) that $F(t, R)$ is a bounded C^1 function of R , and thus the Cauchy-Lipschitz theorem implies that $R(t, \epsilon)$ is a C^1 function of both t and ϵ . The Liouville formula for the Jacobian determinant $J_{n,\epsilon}(t)$ of the map $\epsilon \mapsto R(t, \epsilon)$ yields $J_{n,\epsilon}(t) = \exp\{\int_0^t \partial_{R_1} F_1(s, R(s, \epsilon)) ds + \int_0^t \partial_{R_2} F_2(s, R(s, \epsilon)) ds\} \geq 1$. Indeed, one can show (see again

[13]) that both partial derivatives (in the exponential) are non-negative for all $s \in (0, 1)$. By the same arguments as in the previous bound, for any t , the map $\epsilon \mapsto R(t, \epsilon)$ a \mathcal{C}^1 diffeomorphism. All hypotheses of Lemma. 2.3 are verified. It leads to

$$i_n = \frac{1}{s_n^2} \int_{[s_n, 2s_n]^2} d\epsilon \left\{ I \left(\int_0^1 r(t, \epsilon) dt \right) + \frac{1}{2} G_{\mathbf{R}} \left(\int_0^1 E(t, \epsilon) dt \right) - \frac{1}{2} \int_0^1 E(t, \epsilon) r(t, \epsilon) dt \right\} + \mathcal{O}_n(1).$$

I is a concave function (see [13]), and so is $x \mapsto G_{\mathbf{R}}(x)$. Indeed, by identity (10), we have $G_{\mathbf{R}}''(x) \leq 0$. Jensen's inequality thus yields (and recalling (3))

$$i_n \geq \frac{1}{s_n^2} \int d\epsilon \int_0^1 dt i_{\text{RS}}(E(t, \epsilon), r(t, \epsilon); 1) + \mathcal{O}_n(1).$$

Notice $i_{\text{RS}}(E(t, \epsilon), r(t, \epsilon); 1) = \sup_{E \in [0, \rho]}$ $i_{\text{RS}}(E, r(t, \epsilon); 1)$. Indeed, $g_r : E \mapsto i_{\text{RS}}(E, r; 1)$ is also concave (by concavity of $G_{\mathbf{R}}$), with derivative $g_r'(E) = \frac{1}{2} \mathcal{R}_{\mathbf{R}}(-E) - \frac{\tau}{2}$. By definition of the solution $R(t, \epsilon)$, $g_r'(r(t, \epsilon)(E(t, \epsilon))) = 0$ for any (t, ϵ) , so by concavity $g_r(r(t, \epsilon))$ reaches its maximum at $E(t, \epsilon)$. Thus we finally obtain

$$i_n \geq \frac{1}{s_n^2} \int d\epsilon \int_0^1 dt \sup_{E \in [0, \rho]} i_{\text{RS}}(E, r(t, \epsilon); 1) + \mathcal{O}_n(1) \geq \inf_{r \geq 0} \sup_{E \in [0, \rho]} i_{\text{RS}}(E, r; 1) + \mathcal{O}_n(1).$$

Taking the lim inf, it ends the proof of Theorem 1.1. ■

ACKNOWLEDGMENT

We acknowledge funding from the ERC under the European Union's FP7 Grant Agreement 307087-SPARCS, the SNSF grant 200021-156672, and the ANR PAIL. We also thank Olivier Levêque and Sundeep Rangan for helpful discussions.

REFERENCES

- [1] W. B. Johnson and J. Lindenstrauss, "Extensions of lipschitz mappings into a hilbert space," *Contemporary mathematics*, 1984.
- [2] E. J. Candes and T. Tao, "Near-optimal signal recovery from random projections: Universal encoding strategies?" *IEEE Transactions on Information Theory*, vol. 52, no. 12, pp. 5406–5425, Dec 2006.
- [3] A. M. Tulino and S. Verdú, *Random matrix theory and wireless communications*. Now Publishers Inc, 2004, vol. 1.
- [4] T. Tanaka, "A statistical-mechanics approach to large-system analysis of cdma multiuser detectors," *IEEE Transactions on Information Theory*, vol. 48, no. 11, pp. 2888–2910, Nov 2002.
- [5] A. R. Barron and A. Joseph, "Toward fast reliable communication at rates near capacity with gaussian noise," in *2010 IEEE International Symposium on Information Theory*, June 2010, pp. 315–319.
- [6] J. Barbier and F. Krzakala, "Approximate message-passing decoder and capacity achieving sparse superposition codes," *IEEE Transactions on Information Theory*, vol. 63, no. 8, pp. 4894–4927, Aug 2017.
- [7] M. Mézard, G. Parisi, and M.-A. Virasoro, *Spin glass theory and beyond*. World Scientific Publishing Co., Inc., Pergamon Press, 1987.
- [8] F. Krzakala, M. Mézard, F. Sausset, Y. Sun, and L. Zdeborová, "Statistical-physics-based reconstruction in compressed sensing," *Phys. Rev. X*, vol. 2, p. 021005(18), May 2012.
- [9] A. M. Tulino, G. Caire, S. Verdú, and S. Shamai, "Support recovery with sparsely sampled free random matrices," *IEEE Transactions on Information Theory*, vol. 59, no. 7, pp. 4243–4271, July 2013.
- [10] J. Barbier, M. Dia, N. Macris, and F. Krzakala, "The mutual information in random linear estimation," in *2016 54th Annual Allerton Conference on Communication, Control, and Computing*, 2016.
- [11] J. Barbier, N. Macris, M. Dia, and F. Krzakala, "Mutual information and optimality of approximate message-passing in random linear estimation," *arXiv:1701.05823*, 2017.
- [12] G. Reeves and H. D. Pfister, "The replica-symmetric prediction for compressed sensing with gaussian matrices is exact," vol. arxiv:1607.02524.

- [13] J. Barbier, F. Krzakala, N. Macris, L. Miolane, and L. Zdeborová, "Optimal errors and phase transitions in high-dimensional generalized linear models," in *Proceedings of the 31st Conference On Learning Theory*, ser. Proceedings of Machine Learning Research, vol. 75. PMLR, July 2018, pp. 728–731. [Online]. Available: <http://arxiv.org/abs/1708.03395>
- [14] K. Takeda, S. Uda, and Y. Kabashima, "Analysis of cdma systems that are characterized by eigenvalue spectrum," *EPL (Europhysics Letters)*, vol. 76, no. 6, p. 1193, 2006.
- [15] A. Manoel, F. Krzakala, M. Mézard, and L. Zdeborová, "Multi-layer generalized linear estimation," in *2017 IEEE International Symposium on Information Theory (ISIT)*, 2017, pp. 2098–2102.
- [16] G. Reeves, "Additivity of information in multilayer networks via additive gaussian noise transforms," vol. abs/1710.04580, 2017.
- [17] M. Lelarge and L. Miolane, "Fundamental limits of symmetric low-rank matrix estimation," *ArXiv e-prints*, Nov. 2016.
- [18] J. Barbier and N. Macris, "The adaptive interpolation method: a simple scheme to prove replica formulas in bayesian inference," *Probability Theory and Related Fields*, Oct 2018. [Online]. Available: <https://doi.org/10.1007/s00440-018-0879-0>
- [19] J. Barbier, N. Macris, and L. Miolane, "The Layered Structure of Tensor Estimation and its Mutual Information," in *47th Annual Allerton Conference on Communication, Control, and Computing*, 2017.
- [20] Y. Kabashima, "Inference from correlated patterns: a unified theory for perceptron learning and linear vector channels," *Journal of Physics: Conference Series*, vol. 95, no. 1, p. 012001, 2008.
- [21] D. J. Thouless, P. W. Anderson, and R. G. Palmer, "Solution of 'solvable model of a spin glass'," *Philosophical Magazine*, vol. 35, no. 3, p. 593601, 1977.
- [22] M. Mézard, "The space of interactions in neural networks: Gardner's computation with the cavity method," *Journal of Physics A: Mathematical and General*, vol. 22, no. 12, pp. 2181–2190, 1989.
- [23] D. L. Donoho, A. Maleki, and A. Montanari, "Message-passing algorithms for compressed sensing," *Proceedings of the National Academy of Sciences*, vol. 106, no. 45, pp. 18914–18919, Nov 2009.
- [24] L. Zdeborová and F. Krzakala, "Statistical physics of inference: thresholds and algorithms," *Advances in Physics*, vol. 65, no. 5, p. 453, 2016.
- [25] B. Çakmak, O. Winther, and B. H. Fleury, "S-AMP: approximate message passing for general matrix ensembles," vol. arxiv:1405.2767.
- [26] T. P. Minka, "Expectation propagation for approximate bayesian inference," in *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, ser. UAI'01, 2001, pp. 362–369.
- [27] M. Opper and O. Winther, "Expectation consistent approximate inference," *Journal of Machine Learning Research*, vol. 6, p. 21772204, 2005.
- [28] J. Ma and L. Ping, "Orthogonal amp," *IEEE Access*, vol. 5, pp. 2020–2033, 2017.
- [29] S. Rangan, P. Schniter, and A. K. Fletcher, "Vector approximate message passing," vol. arxiv:1610.03082.
- [30] F. Guerra and F. L. Toninelli, "The thermodynamic limit in mean field spin glass models," *Communications in Mathematical Physics*, vol. 230, no. 1, pp. 71–79, 2002.
- [31] V. A. Marčenko and L. A. Pastur, "Distribution of eigenvalues for some sets of random matrices," *Mathematics of the USSR-Sbornik*, vol. 1, no. 4, p. 457, 1967.
- [32] J. W. Silverstein, "Strong convergence of the empirical distribution of eigenvalues of large dimensional random matrices," *Journal of Multivariate Analysis*, vol. 55, no. 2, pp. 331–339, 1995.
- [33] R. Speicher, "Free probability theory," *arXiv preprint arXiv:0911.0087*, 2009.
- [34] N. Macris, "Griffith-Kelly-Sherman correlation inequalities: A useful tool in the theory of error correcting codes," *IEEE Transactions on Information Theory*, vol. 53, no. 2, pp. 664–683, Feb 2007.
- [35] S. B. Korada and N. Macris, "Tight bounds on the capacity of binary input random cdma systems," *IEEE Transactions on Information Theory*, vol. 56, no. 11, pp. 5590–5613, Nov 2010.
- [36] —, "Exact solution of the gauge symmetric p-spin glass model on a complete graph," *Journal of Statistical Physics*, vol. 136, no. 2, pp. 205–230, 2009.
- [37] M. Talagrand, *Mean field models for spin glasses: Volume I: Basic examples*. Springer Science & Business Media, 2010, vol. 54.

A. Nishimori identity

Lemma A.1: Let (X, Y) be a couple of random variables on a polish space E . For a given $k \in \mathbb{N}^*$, let $(X^{(i)})_{i=1}^k$ be i.i.d random variables from the distribution (conditional over Y) $P(X = \cdot | Y)$. Denote $\langle \cdot \rangle$ the expectation with respect to this probability distribution, and \mathbb{E} the expectation with respect to the probability measure of (X, Y) . Then, for all $f : E^{k+1} \rightarrow \mathbb{R}$ continuous and bounded

$$\mathbb{E}\langle f(Y, X^{(1)}, \dots, X^{(k)}) \rangle = \mathbb{E}\langle f(Y, X^{(1)}, \dots, X^{(k-1)}, X) \rangle.$$

Proof: This is a trivial consequence of Bayes formula:

$$\mathbb{E}_{X,Y}\langle f(Y, X^{(1)}, \dots, X^{(k-1)}, X) \rangle = \mathbb{E}_Y \mathbb{E}_{X|Y}\langle f(Y, X^{(1)}, \dots, X^{(k-1)}, X) \rangle = \mathbb{E}\langle f(Y, X^{(1)}, \dots, X^{(k)}) \rangle.$$

■

B. Proof of Lemma 2.1

The proof is done in two steps. First, we show the following formula:

$$i'_{n,\epsilon}(t) = \frac{r(t)}{2}(\rho - \mathbb{E}\langle Q \rangle_{t,\epsilon}) + \frac{1}{2n^2} \sum_{\mu=1}^m \sum_{\nu=1}^m \mathbb{E}\left\langle Z_\mu(\Phi' \Phi'^\top)_{\mu\nu} u_{t,\nu} \left[E(t) - \left(\frac{1}{n} \sum_{i=1}^n X_i^2 - Q \right) \right] \right\rangle_{t,\epsilon}. \quad (13)$$

We will then conclude using the concentration of $\frac{1}{n} \sum_{i=1}^n X_i^2$ on ρ by the central limit theorem as $n \rightarrow \infty$.

Recall that we defined the Gibbs bracket

$$\langle A(\mathbf{x}, \mathbf{v}) \rangle_{t,\epsilon} = \frac{\int dP_0(\mathbf{x}) \mathcal{D}\mathbf{v} e^{-\mathcal{H}(\epsilon, t, \mathbf{x}, \mathbf{v}; \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi)} A(\mathbf{x}, \mathbf{v})}{\int dP_0(\mathbf{x}) \mathcal{D}\mathbf{v} e^{-\mathcal{H}(\epsilon, t, \mathbf{x}, \mathbf{v}; \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi)}}. \quad (14)$$

From this and the definition of $i_{n,\epsilon}(t)$ (7), one gets

$$i'_{n,\epsilon}(t) = \frac{1}{n} \mathbb{E}[\mathcal{H}'(\epsilon, t, \mathbf{X}, \mathbf{V}; \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi) \ln \mathcal{Z}] + \frac{1}{n} \mathbb{E}\langle \mathcal{H}'(\epsilon, t, \mathbf{x}, \mathbf{v}; \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi) \rangle_{t,\epsilon}, \quad (15)$$

where the partition function \mathcal{Z} and Hamiltonian derivative with respect to t read

$$\mathcal{H}'(\epsilon, t, \mathbf{X}, \mathbf{V}; \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi) = \frac{1}{2} \sum_{\mu=1}^m Z_\mu \left(\sqrt{\frac{1}{n(1-t)}} (\Phi \mathbf{X})_\mu - \frac{E(t)}{\sqrt{nR_2(t)}} (\Phi' \mathbf{V})_\mu \right) - \frac{1}{2} \sum_{i=1}^n \tilde{Z}_i \frac{r(t)}{\sqrt{R_1(t)}} X_i, \quad (16)$$

$$\mathcal{Z} = \mathcal{Z}(\epsilon, t, \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi) \equiv \int dP_0(\mathbf{x}) \mathcal{D}\mathbf{v} e^{-\mathcal{H}(\epsilon, t, \mathbf{x}, \mathbf{v}; \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi)}. \quad (17)$$

The Nishimori identity (Lemma A.1) directly implies

$$\mathbb{E}\langle \mathcal{H}'(\epsilon, t, \mathbf{x}, \mathbf{v}; \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi) \rangle_t = \mathbb{E} \mathcal{H}'(\epsilon, t, \mathbf{X}, \mathbf{V}; \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi) = 0. \quad (18)$$

We now compute $\mathbb{E}[\tilde{Z}_i X_i \ln \mathcal{Z}]$. Using a Gaussian integration by parts, which reads for any real function f with continuous derivative $\mathbb{E}[\tilde{Z}_i f(\tilde{Z}_i)] = \mathbb{E}[f'(\tilde{Z}_i)]$ for $\tilde{Z}_i \sim \mathcal{N}(0, 1)$, we obtain the first term of (15) as

$$\begin{aligned} & -\frac{1}{2n} \frac{r(t)}{\sqrt{R_1(t)}} \sum_{i=1}^n \mathbb{E}\left[X_i \tilde{Z}_i \ln \int dP_0(\mathbf{x}) \mathcal{D}\mathbf{v} e^{-\mathcal{H}(\epsilon, t, \mathbf{x}, \mathbf{v}; \mathbf{Y}_t, \tilde{\mathbf{Y}}_t, \Phi)} \right] \\ &= -\frac{1}{2n} \frac{r(t)}{\sqrt{R_1(t)}} \sum_{i=1}^n \mathbb{E}\left[X_i \tilde{Z}_i \ln \int dP_0(\mathbf{x}) \mathcal{D}\mathbf{v} e^{\text{term}_1(\mathbf{v}, \mathbf{x}) - \frac{1}{2} \sum_i (\sqrt{R_1(t)}(X_i - x_i) + \tilde{Z}_i)^2} \right] \\ &= \frac{1}{2n} \frac{r(t)}{\sqrt{R_1(t)}} \sum_{i=1}^n \mathbb{E}\left[X_i \langle \sqrt{R_1(t)}(X_i - x_i) + \tilde{Z}_i \rangle_{t,\epsilon} \right] \\ &= \frac{r(t)\rho}{2} - \frac{r(t)}{2} \mathbb{E}\left\langle \frac{1}{n} \sum_{i=1}^n X_i x_i \right\rangle_{t,\epsilon} \\ &= \frac{r(t)}{2} (\rho - \mathbb{E}\langle Q \rangle_{t,\epsilon}). \end{aligned} \quad (19)$$

In the same way, an integration by parts with respect to $V_i \sim \mathcal{N}(0, 1)$ yields

$$\begin{aligned}
& -\frac{1}{2n} \frac{E(t)}{\sqrt{nR_2(t)}} \sum_{\mu=1}^m \mathbb{E} \left[Z_\mu (\Phi' \mathbf{V})_\mu \ln \mathcal{Z} \right] \\
&= -\frac{1}{2n} \frac{E(t)}{\sqrt{nR_2(t)}} \sum_{\mu=1}^m \sum_{i=1}^n \mathbb{E} \left[Z_\mu \Phi'_{\mu i} V_i \ln \int dP_0(\mathbf{x}) \mathcal{D}\mathbf{v} e^{-\frac{1}{2} \sum_\nu \left(\sqrt{\frac{1-t}{n}} (\Phi(\mathbf{X}-\mathbf{x}))_\nu + \sqrt{\frac{R_2(t)}{n}} (\Phi'(\mathbf{V}-\mathbf{v}))_\nu + Z_\nu \right)^2 + \text{term}_2(\mathbf{x})} \right] \\
&= \frac{E(t)}{2n^2} \sum_{\mu=1}^m \sum_{\nu=1}^m \sum_{i=1}^n \mathbb{E} \left[Z_\mu \Phi'_{\mu i} \Phi'_{\nu i} \left\langle \sqrt{\frac{1-t}{n}} (\Phi(\mathbf{X}-\mathbf{x}))_\nu + \sqrt{\frac{R_2(t)}{n}} (\Phi'(\mathbf{V}-\mathbf{v}))_\nu + Z_\nu \right\rangle_{t,\epsilon} \right] \\
&= \frac{E(t)}{2n^2} \sum_{\mu=1}^m \sum_{\nu=1}^m \mathbb{E} \left[Z_\mu (\Phi' \Phi'^\top)_{\mu\nu} \langle u_{t,\nu} \rangle_{t,\epsilon} \right]. \tag{20}
\end{aligned}$$

Let us now look at the final term we need to compute. By our hypothesis (4), this term reads, using again a Gaussian integration by part but this time with respect to $W_{ji} \sim \mathcal{N}(0, 1/n)$,

$$\begin{aligned}
& \frac{1}{2n} \sqrt{\frac{1}{n(1-t)}} \sum_{\mu=1}^m \mathbb{E} \left[Z_\mu (\Phi' \mathbf{W} \mathbf{X})_\mu \ln \mathcal{Z} \right] \\
&= \frac{1}{2n} \sqrt{\frac{1}{n(1-t)}} \sum_{\mu=1}^m \sum_{i,j=1}^n \mathbb{E} \left[Z_\mu \Phi'_{\mu j} W_{ji} X_i \ln \int dP_0(\mathbf{x}) \mathcal{D}\mathbf{v} e^{-\frac{1}{2} \sum_\nu \left(\sqrt{\frac{1-t}{n}} (\Phi' \mathbf{W}(\mathbf{X}-\mathbf{x}))_\nu + \sqrt{\frac{R_2(t)}{n}} (\Phi'(\mathbf{V}-\mathbf{v}))_\nu + Z_\nu \right)^2 + \text{term}_2(\mathbf{x})} \right] \\
&= -\frac{1}{2n^3} \sum_{\mu,\nu=1}^m \sum_{i,j=1}^n \mathbb{E} \left[Z_\mu \Phi'_{\mu j} \Phi'_{\nu j} X_i \left\langle (X_i - x_i) \left(\sqrt{\frac{1-t}{n}} (\Phi(\mathbf{X}-\mathbf{x}))_\nu + \sqrt{\frac{R_2(t)}{n}} (\Phi'(\mathbf{V}-\mathbf{v}))_\nu + Z_\nu \right) \right\rangle_{t,\epsilon} \right] \\
&= -\frac{1}{2n^2} \sum_{\mu,\nu=1}^m \mathbb{E} \left[Z_\mu (\Phi' \Phi'^\top)_{\mu\nu} \left\langle u_{t,\nu} \left(\frac{1}{n} \sum_{i=1}^n X_i^2 - \frac{1}{n} \sum_{i=1}^n X_i x_i \right) \right\rangle_{t,\epsilon} \right].
\end{aligned}$$

Combining all three terms leads to (13).

We now go to the last step. By adding and subtracting a term to (13) we reach

$$\begin{aligned}
i'_{n,\epsilon}(t) &= \frac{r(t)}{2} (\rho - \mathbb{E}\langle Q \rangle_{t,\epsilon}) + \frac{1}{2n^2} \sum_{\mu=1}^m \sum_{\nu=1}^m \mathbb{E} \langle Z_\mu (\Phi' \Phi'^\top)_{\mu\nu} u_{t,\nu} [E(t) - (\rho - Q)] \rangle_{t,\epsilon} \\
&\quad + \frac{1}{2n^2} \sum_{\mu=1}^m \sum_{\nu=1}^m \mathbb{E} \langle Z_\mu (\Phi' \Phi'^\top)_{\mu\nu} u_{t,\nu} \left(\rho - \frac{1}{n} \sum_{i=1}^n X_i^2 \right) \rangle_{t,\epsilon}. \tag{21}
\end{aligned}$$

Using the Cauchy-Schwarz inequality we obtain that the last term can be bounded as

$$\frac{1}{n^2} \left| \mathbb{E} \langle \mathbf{Z}^\top (\Phi' \Phi'^\top) \mathbf{u}_t \left(\rho - \frac{1}{n} \sum_{i=1}^n X_i^2 \right) \rangle_{t,\epsilon} \right| \leq \left\{ \frac{1}{n^4} \mathbb{E} \langle (\mathbf{Z}^\top (\Phi' \Phi'^\top) \mathbf{u}_t)^2 \rangle_{t,\epsilon} \right\}^{1/2} \mathbb{E} \left[\left(\rho - \frac{1}{n} \sum_{i=1}^n X_i^2 \right)^2 \right]^{1/2}. \tag{22}$$

As the X_i are independent the central limit theorem implies that $\mathbb{E}[(\rho - \frac{1}{n} \sum_{i=1}^n X_i^2)^2] = \mathcal{O}(1/n)$. Thus it remains to show that the multiplicative term in front is bounded:

$$\begin{aligned}
& \frac{1}{n^4} \mathbb{E} \langle (\mathbf{Z}^\top (\Phi' \Phi'^\top) \mathbf{u}_t)^2 \rangle_{t,\epsilon} \leq \frac{1}{n^4} \mathbb{E} \langle \|\mathbf{Z}\|^2 \|\mathbf{u}_t\|^2 \|\Phi' \Phi'^\top\|_{\mathbb{F}}^2 \rangle_{t,\epsilon} \\
& \leq \frac{1}{n^4} \sqrt{\mathbb{E} \langle \|\mathbf{Z}\|^4 \|\mathbf{u}_t\|^4 \rangle_{t,\epsilon}} \mathbb{E} [\|\Phi' \Phi'^\top\|_{\mathbb{F}}^4] \leq \frac{1}{n^4} \sqrt{\mathbb{E} [\|\mathbf{Z}\|^8] \mathbb{E} \langle \|\mathbf{u}_t\|^8 \rangle_{t,\epsilon}} \mathbb{E} [\|\Phi' \Phi'^\top\|_{\mathbb{F}}^4] = \mathcal{O}(1). \tag{23}
\end{aligned}$$

The last equality follows from the following observations. By construction of Φ' , $\mathbb{E}[\|\Phi' \Phi'^\top\|_{\mathbb{F}}^4]^{1/2} = \mathcal{O}(n^2)$. Moreover, as \mathbf{Z} is a m -dimensional Gaussian vector with i.i.d. components $\mathbb{E}[\|\mathbf{Z}\|^8]^{1/4} = \mathcal{O}(n)$. Finally, the Nishimori identity leads to $\mathbb{E}[\langle \|\mathbf{u}_t\|^8 \rangle_{t,\epsilon}]^{1/4} = \mathcal{O}(n)$. This claim is proven using a consequence of the triangle inequality:

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n \quad \|\mathbf{x} + \mathbf{y}\|^8 \leq 2^7 (\|\mathbf{x}\|^8 + \|\mathbf{y}\|^8),$$

which is combined with the Nishimori identity:

$$\begin{aligned}
\mathbb{E}[\langle \|\mathbf{u}_t\|^8 \rangle_{t,\epsilon}] &= \mathbb{E} \left\langle \left\| \sqrt{\frac{1-t}{n}} \Phi(\mathbf{X}-\mathbf{x}) + \sqrt{\frac{R_2(t)}{n}} \Phi'(\mathbf{V}-\mathbf{v}) + \mathbf{Z} \right\|^8 \right\rangle_{t,\epsilon} \\
&\leq 2^7 \mathbb{E}[\|\mathbf{Z}\|^8] + 2^{22} (1-t)^4 \mathbb{E} \left[\left\| \frac{1}{\sqrt{n}} \Phi \mathbf{X} \right\|^8 \right] + 2^{22} R_2(t)^4 \mathbb{E} \left[\left\| \frac{1}{\sqrt{n}} \Phi' \mathbf{V} \right\|^8 \right]. \tag{24}
\end{aligned}$$

One can now use that both $\frac{1}{n}\Phi^\top\Phi$ and $\frac{1}{n}\Phi'^\top\Phi'$ have almost surely bounded Euclidian (or Frobenius) norm when $n \rightarrow \infty$. This implies that there exists $C > 0$ such that

$$\mathbb{E}\left[\left\|\frac{1}{\sqrt{n}}\Phi\mathbf{X}\right\|^8\right] \leq C\mathbb{E}\left[\|\mathbf{X}\|^8\right].$$

Moreover $\mathbb{E}\left[\|\mathbf{X}\|^8\right] = \mathcal{O}(n^4)$ because we assumed the prior distribution P_0 to be compactly supported. The same argument can be conducted for bounding $\mathbb{E}\left[\left\|\frac{1}{\sqrt{n}}\Phi'\mathbf{V}\right\|^8\right]$ since $\mathbb{E}\left[\|\mathbf{V}\|^8\right] = \mathcal{O}(n^4)$ as \mathbf{V} is a standard Gaussian vector.