



On the universality of noiseless linear estimation with respect to the measurement matrix

Alia Abbbara, Antoine Baker, Florent Krzakala, Lenka Zdeborová

► To cite this version:

Alia Abbbara, Antoine Baker, Florent Krzakala, Lenka Zdeborová. On the universality of noiseless linear estimation with respect to the measurement matrix. *Journal of Physics A: Mathematical and Theoretical*, 2020, 53 (16), pp.164001. 10.1088/1751-8121/ab59ef . cea-02528193

HAL Id: cea-02528193

<https://cea.hal.science/cea-02528193>

Submitted on 1 Apr 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

On the Universality of Noiseless Linear Estimation with Respect to the Measurement Matrix

Alia Abbara¹, Antoine Baker¹, Florent Krzakala¹ and Lenka Zdeborová²

¹ Laboratoire de Physique de l'Ecole normale supérieure, Université PSL, CNRS, Sorbonne Université, Université Paris-Diderot, Sorbonne Paris Cité, Paris, France

² Institut de physique théorique, Université Paris Saclay, CNRS, CEA, 91191 Gif-sur-Yvette, France

E-mail: `alia.abbara@ens.fr`

Abstract. In a noiseless linear estimation problem, one aims to reconstruct a vector \mathbf{x}^* from the knowledge of its linear projections $\mathbf{y} = \Phi \mathbf{x}^*$. There have been many theoretical works concentrating on the case where the matrix Φ is a random i.i.d. one, but a number of heuristic evidence suggests that many of these results are universal and extend well beyond this restricted case. Here we revisit this problematic through the prism of development of message passing methods, and consider not only the universality of the ℓ_1 transition, as previously addressed, but also the one of the optimal Bayesian reconstruction. We observed that the universality extends to the Bayes-optimal minimum mean-squared (MMSE) error, and to a range of structured matrices.

1. Introduction

The problem of recovering a signal through the knowledge of its linear projections is ubiquitous in modern information theory, statistics and machine learning. In particular, many applications require to reconstruct an unknown n -dimensional signal vector \mathbf{x}^* from the linear projections

$$\mathbf{y} = \Phi \mathbf{x}^*, \quad (1)$$

where \mathbf{y} is a m -dimensional vector, and Φ is a $m \times n$ random matrix. For instance, if \mathbf{x}^* is sparse, this task of estimating the signal from its linear *random* projections is at the roots of compressed sensing [1]. A fundamental question in the field is how much the algorithmic and the information theoretic performance depends on the choice of the random matrix Φ .

In the present letter, we concentrate on the noiseless and asymptotic, large n , regime with a fixed value $\alpha = m/n$. We consider \mathbf{x}^* to be k -sparse, i.e. to have only k non-zero values, and we shall work in the limit where $n \rightarrow \infty$, $k \rightarrow \infty$, and a finite value of $\rho = k/n$. In such setting, a classical result is the following: for random matrices Φ with independent standard Gaussian entries, the (convex) reconstruction

with ℓ_1 penalty displays a precisely determined phase transition. For a certain region in the (α, ρ) -phase diagram, it typically finds back the vector \mathbf{x}^* , being the sparsest solution, whereas outside that region, it typically fails. The boundary between these two regions is called the Donoho-Tanner line [2]. It has been shown empirically that the very same phase transition location seems to hold for a wider range of random matrix ensembles, see e.g. [3, 4], suggesting a large universality of the Donoho-Tanner phase transitions. Another line of work showed that the convex ℓ_1 reconstruction problem can be treated through conic geometry, and the success probability of signal recovery only depends on a geometric number characterizing a subcone (statistical dimension or Gaussian width) [5, 6].

Here we investigate the universality of the phase transition not only for the ℓ_1 transition, but also to the performance of the optimal Bayesian reconstruction. We analyze this question through the prism of information theory, message passing methods, and random matrix theory. We shall see that the universality indeed extends to a more generic set of properties than the ℓ_1 transition, such as the minimum mean-squared (MMSE) error or the easy-hard phase transition for optimal Bayesian learning, and empirically to structured matrices such as the one appearing in [7, 8].

We note that investigation of universality are very common to physics problems, and understanding how large is the class of model for which a given result applied is a very fundamental question. The message-passing-based algorithm that we investigate in this paper to demonstrate the universality also has their origin in physics works, such as [9].

2. A short review of results for i.i.d. random matrices

A first well-understood case of universality holds for random matrices Φ where all the elements are generated i.i.d. from a well-behaved distribution -with zero mean and unit variance- which all exhibit the same transitions as Gaussian random matrices. This is known for multiple retrieval problems:

2.1. ℓ_1 recovery

Consider for instance the Donoho-Tanner line [2] that regulates the ℓ_1 recovery. Thanks to the approximate message passing solver (see below) that has been shown to be universal with respect to all i.i.d. distributions with finite moments [10, 11], we know that the Donoho-Tanner phase transition is the same for all such random matrices.

2.2. Information theoretic optimal reconstruction

There has been a considerable amount of work in the information theory community on the computation of the mutual information and on the MMSE for problems such as (1) with Gaussian matrices. In particular, following the replica method from statistical physics (the Tanaka formula [12]), a heuristic formula has been postulated in different

situations, see e.g. [13–16]. This heuristic replica result has been recently rigorously proven in a series of papers [17, 17, 18]. In a more recent proof [19], it has been shown, again, that the formula is not specific to Gaussian i.i.d. matrices, but that any matrix with i.i.d elements of unit variance and zero mean leads to the same exact result for the mutual information and the MMSE.

2.3. Hard phase for Bayesian decoders

A third interesting point is to ask about tractable decoders that aim to perform the optimal Bayesian estimation, i.e. with a perfect prior knowledge on the distribution of \mathbf{x}^* . For simplicity, consider for instance the case where each element of \mathbf{x}^* has been sampled from a Gauss-Bernoulli distribution:

$$x_i \sim (1 - \rho)\delta(x) + \rho\mathcal{N}(0, 1).$$

In this case, the best known solver is again AMP, using a Bayesian decoder (instead of the soft thresholding function for ℓ_1 recovery) [14, 15, 20, 21]. Interestingly, it shares with the ℓ_1 recovery a similar phase transition: for a certain region in the (α, ρ) plane it typically finds back the vector \mathbf{x}^* , whereas outside that region it fails. We shall denote the limit between these regions the "Bayesian hard-phase" transition. The "Bayesian hard-phase" line, that has been precisely computed in [14, 15] is always better than the Donoho-Tanner line (as it should, since it exploits additional information). Once more, the universality of AMP shows that this phase transition is not restricted to Gaussian matrices, but extends as well to all (well normalized) i.i.d. matrices.

The fact that these three properties (the ℓ_1 , the hard-phase line, as well as the MMSE) are universal for all i.i.d. matrices makes the case for Gaussian computations, as done in theoretical computation, stronger. We shall see that this universality extends well beyond these simple cases.

3. Random rotationally invariant matrices

Moving away from the well-known i.i.d. examples, we start by considering a much larger set of random matrices defined through their singular value decomposition (SVD): any real matrix Φ can be decomposed into $\Phi = U\Sigma V$, with U and V orthogonal matrices, and Σ 's elements being Φ 's singular values. We shall look at the left rotationally invariant random matrix ensemble: these are matrices Φ that can be written as

$$\Phi = U\Sigma V$$

with an arbitrary rotation matrix U and singular values Σ , but where the matrix V has been randomly (and independently of Σ and U) generated from the Haar measure (that is, uniformly from all possible rotations).

When the singular values are different from zero, it is straightforward to justify the universality property for matrices from this subclass. We start by the definition of the

problem: we wish to find \mathbf{x} such that

$$\mathbf{y} = \Phi \mathbf{x} = U \Sigma V \mathbf{x}. \quad (2)$$

If $m \leq n$, then Σ is written as $\Sigma = \left[\begin{array}{c|c} \tilde{\Sigma} & 0 \end{array} \right]$ and we define

$$\Sigma^{inv} = \left[\begin{array}{c|c} \tilde{\Sigma}^{-1} & \\ \hline 0 & 0 \end{array} \right] \text{ such that } \Sigma^{inv} \Sigma = \left[\begin{array}{c|c} I_m & 0 \\ \hline 0 & 0 \end{array} \right].$$

Multiplying (2) on both sides by U^T , and then by Σ^{inv} ; one reaches

$$\tilde{\mathbf{y}} = \Sigma^{inv} U^T \mathbf{y} = \tilde{V} \mathbf{x} \quad (3)$$

where \tilde{V} is a $m \times n$ matrix composed of the first m lines of V . If instead $m > n$, Σ is written as

$$\Sigma^{inv} = \left[\begin{array}{c} \tilde{\Sigma} \\ \hline 0 \end{array} \right]$$

and we define $\Sigma^{inv} = \left[\begin{array}{c|c} \tilde{\Sigma}^{-1} & 0 \end{array} \right]$ such that $\Sigma^{inv} \Sigma = I_n$. Multiplying (2) by U^T then Σ^{inv} , we obtain

$$\tilde{\mathbf{y}} = \Sigma^{inv} U^T \mathbf{y} = V \mathbf{x}. \quad (4)$$

In both cases, we thus see that the problem has been transformed—in a constructive way—into a standard linear system with the sensing matrix \tilde{V} when $m \leq n$ being a (sub-sampled) random rotation one, or sensing matrix V when $m > n$. This shows that all rotationally invariant matrices, which satisfy U and Σ 's independence on V , can be transformed the same way and are in the same universality class as far as noiseless linear recovery is concerned, i.e. they will display the same phase transitions.

Since Gaussian i.i.d. matrices belong among random rotationally invariant matrices (in this case Σ follows the Marcenko-Pastur law [22]) this means that all the information theoretic rigorous results (such as phase transitions and MMSE value) with zero noise for random Gaussian matrices applies verbatim to all rotationally invariant ensemble, as long as the SVD's matrices U and Σ are independent of V . This is a very strong universality, that applies to the three cases (1, 2, 3) from sec. 2. Note that the universality of the Donoho-Tanner line with rotationally invariant matrices was already hinted by the replica method [23].

Notice, however, that the above construction depends crucially on the fact that we consider here noiseless measurements. It would not work if an additional Gaussian noise were added in eq. (1): in this case, the transformation would make the i.i.d. Gaussian noise a correlated one. Indeed, the replica formula for noisy measurements underlines that the MMSE depends on the precise set of matrices in noisy reconstruction [13, 24] (this formula is not yet fully rigorous, but see [25] for a proof in a restricted setting). Any differences, however, must go to zero in the noiseless limit.

4. Approximate Message Passing

Having discussed the universality with respect to random rotationally invariant matrices, we now wish to discuss its effect on specific solvers, concretely the message passing algorithms.

4.1. AMP

We first consider the original approximate message passing (AMP) [26] to compute the phase transition between the phase where the algorithm reconstructs \mathbf{x}^* perfectly, and the one where reconstruction may be possible but is not achieved by the algorithm. AMP is an iterative algorithm that follows:

$$\begin{aligned}\hat{\mathbf{x}}^{t+1} &= \eta_t(\Phi^T \hat{\mathbf{x}}^t) \\ \mathbf{z}^t &= \mathbf{y} - \Phi \hat{\mathbf{x}}^t + \frac{1}{\alpha} \mathbf{z}^{t-1} \langle \eta'_{t-1}(\Phi^T \mathbf{z}^{t-1} + \hat{\mathbf{x}}^{t-1}) \rangle.\end{aligned}$$

where t is the iteration index, \mathbf{x}^t is the current estimate of \mathbf{x}^* , \mathbf{z}^t the current residual, $\langle \cdot \rangle$ is an averaged sum of components, and η_t is a prior-dependent threshold function applied component-wise (the soft thresholding for ℓ_1 , or the Bayesian decoder [14, 15]).

One of the most interesting features of AMP is that, if Φ is a Gaussian i.i.d. matrix, its mean squared error (MSE) σ_t can be tracked accurately by the state evolution formalism [10, 11, 26]. State evolution is a relatively simple recursive equation:

$$\sigma_{t+1}^2 = \Psi(\sigma_t^2), \quad \Psi(\sigma^2) = \mathbb{E} \left[\left(\eta_t(X + \frac{\sigma}{\sqrt{\alpha}} Z) - X \right)^2 \right], \quad (5)$$

where the expectation is with respect to independent random variables $Z \sim \mathcal{N}(0, 1)$ and X , whose distribution coincides with the empirical distribution of the entries of \mathbf{x}^* . Analyzing the evolution of this equation for the ℓ_1 decoder yields the Donoho-Tanner line [26], while using the Bayesian decoder it yields the hard-phase line for Bayesian decoding [14].

It would be interesting to use AMP for rotationally invariant matrices. In order to do this, we follow the construction of sec. 3: starting from equation (3) we then multiply by Σ_0 , a $m \times m$ diagonal matrix with singular values sampled from Marcenko-Pastur law (singular values of a Gaussian i.i.d. matrix \ddagger), and U_0 a $m \times m$ Haar-generated orthogonal matrix, thus ensuring that Σ_0 and U_0 are generated independently of V :

$$U_0 \Sigma_0 \tilde{\Sigma}^{-1} U^T \mathbf{y} = U_0 \Sigma_0 \tilde{V} \mathbf{x} \quad (6)$$

$$\mathbf{y}' = \Phi' \mathbf{x}. \quad (7)$$

After this transformation, $\Phi' = U_0 \Sigma_0 \tilde{V}$ is a random matrix that belongs to an ensemble very close to the Gaussian i.i.d. matrices ensemble. In fact, a recent work showed that

\ddagger The singular values of a Gaussian matrix are correlated, so in fact we may want to generate Σ_0 by first generating a random Gaussian matrix, and then calculating its singular values.

AMP applied to a Gaussian matrix follows the same state evolution as matrices such as Φ' where U_0, \tilde{V} are uniform orthogonal matrices and Σ_0 diagonal's elements are singular values sampled from the Marcenko-Pastur law [27]. Combining this result with the matrix transformation, we have thus constructively mapped the noiseless reconstruction problem back to the well-understood noiseless compressed sensing case for a Gaussian i.i.d. matrix, where we can safely apply the algorithm, and its state evolution. In the section 5.2, we apply this matrix transformation for numerical experiments using AMP.

4.2. Vector-AMP

While the transformation trick allows to make AMP work with random rotationally invariant matrices, another alternative is to work directly with a dedicated solver. To this means, different but related approaches were proposed [24, 28], in particular, using the general expectation-propagation (EP) [29, 30] scheme. Ma and Ping proposed a variation of EP called OAMP [31] specially adapted to rotation matrices. Rangan, Schniter and Fletcher introduced a similar approach called VAMP [32] and proved that it follows state evolution equations corresponding to the fixed point of the replica potential [13, 24, 25]. The multi-layer AMP algorithm of [33] also display the same fixed point.

We shall concentrate here on the VAMP (Vector-AMP) approach, and for a moment, put back a small additional random Gaussian i.i.d. noise of variance Δ in the measurement in eq. (1) as it is needed for stating the algorithm. VAMP then consists in the following fixed-point iteration:

$$\begin{aligned} \mathbf{u}_\ell^{t+1} &= \frac{\hat{\mathbf{x}}_\ell^t}{\langle \text{Var}_\ell^t(\mathbf{x}) \rangle} - \mathbf{u}_r^t, & \rho_\ell^{t+1} &= \frac{1}{\langle \text{Var}_\ell^t(\mathbf{x}) \rangle} - \rho_r^t, \\ \mathbf{u}_r^{t+1} &= \frac{\hat{\mathbf{x}}_r^t}{\langle \text{Var}_r^t(\mathbf{x}) \rangle} - \mathbf{u}_\ell^t, & \rho_r^{t+1} &= \frac{1}{\langle \text{Var}_r^t(\mathbf{x}) \rangle} - \rho_\ell^t, \end{aligned} \quad (8)$$

where we denote by $\mathbb{E}_{\ell,r}^t$ the expectation w.r.t. the tilted distributions $\tilde{Q}_{\ell,r}^t(\mathbf{x}) \propto P_{\ell,r}(\mathbf{x})Q_{\ell,r}^t(\mathbf{x})$, and by $\text{Var}_{\ell,r}^t(\mathbf{x})$ the variance of these distributions. Here, we have defined $Q_{l,r}(\mathbf{x}) = e^{-\frac{1}{2}\rho_{l,r}\mathbf{x}^T\mathbf{x} + \mathbf{u}_{l,r}^T\mathbf{x}}$, $P_l(\mathbf{x}) \propto e^{-\|\mathbf{y} - \Phi\mathbf{x}\|_2^2/2\Delta}$ and $P_r(\mathbf{x})$ is the prior used in the algorithm (i.e. the Laplace prior for the ℓ_1 model, or the actual distribution of the signal for Bayesian reconstruction). In particular

$$\begin{aligned} \hat{\mathbf{x}}_\ell^t &= (\Phi^T\Phi + \Delta\rho_r^t I_p)^{-1}(\Phi^T\mathbf{y} + \Delta\mathbf{u}_r^t), \\ \langle \text{Var}_\ell^t(\mathbf{x}) \rangle &= \frac{\Delta}{N} \text{Tr}(\Phi^T\Phi + \Delta\rho_r^t I_p)^{-1}, \end{aligned} \quad (9)$$

where, as for AMP, we define the denoiser that yields the estimates of x by $z(u, \rho) = \int dx P_r(x) e^{-\frac{1}{2}\rho x^2 + ux}$,

$$\begin{aligned} (\hat{x}_r)_j &= \frac{\partial}{\partial u} \log z(u, \rho) \Big|_{u_{\ell k}^t, \rho_\ell^t}, \\ \langle \text{Var}_r^t(\mathbf{x}) \rangle &= \frac{1}{n} \sum_{j=1}^p \frac{\partial^2}{\partial u^2} \log z(u, \rho) \Big|_{u_{\ell k}^t, \rho_\ell^t}. \end{aligned} \quad (10)$$

Again, the performance of the recursion can be analyzed rigorously through the state evolution [32]. For simplicity, let us concentrate on the Bayes optimal case in which case the state evolution can be closed on the variables (see [32]):

$$\sigma^t = \langle \text{Var}_r^t(\mathbf{x}) \rangle \text{ and } \epsilon^t = \langle \text{Var}_l^t(\mathbf{x}) \rangle, \quad (11)$$

by writing

$$\begin{aligned} \sigma^t(\rho_l^t) &= \Psi((\rho_l^t)^{-1}) \\ \epsilon(\rho_r^t) &= \Delta \mathbb{E} \left[\frac{1}{\Sigma^2 + \Delta \rho_r^t} \right] = \Delta S_{\Sigma^2}(-\Delta \rho_r^t) \end{aligned} \quad (12)$$

where the expectation is above the distribution of the singular values Σ of the matrix Φ , and where we recognize the Stieltjes transform $S_X(r) = \mathbb{E}[1/X - r]$.

Though this transform, we see that the performance depends crucially on the distribution of eigenvalues. Let us now go back on the noiseless limit when $\Delta \rightarrow 0$ and analyze how the universality shows up. Consider again the Stieltjes transform: out of the n singular values of the $n \times n$ matrix $\Phi^T \Phi$, we shall have $(1 - \alpha)n$ of them to be zero (assuming $\alpha < 1$) while the rest are positive (since $m < n$). In this case, the limit $r \rightarrow 0$ of the Stieltjes transform will behave as $S_X(r) \approx -(1 - \alpha)/r$ so that

$$\lim_{\Delta \rightarrow 0} \epsilon(\rho_r^t) = \frac{1 - \alpha}{\rho_r^t}.$$

Again, we see that all the complicated dependence on the spectrum of the matrix Φ has been eliminated. This is a direct, alternative, proof that VAMP will also yield universal results in the zero noise limit for the Bayesian reconstruction. Given that VAMP has the same fixed point as the replica mutual information [13, 25], this argument applies to the replica prediction for the MMSE as well.

5. Structured matrices

We now move to very structured matrices, in order to test the universality as well as the quality and the prediction of the state evolution out of its comfort zone. In order to do so, we have considered different matrix ensembles:

5.1. Tested ensembled of matrices

Discrete cosine transform matrices The first ensemble we consider consists in Fourier-like matrices. A $n \times n$ discrete cosine transform (DCT) matrix Y is defined by:

$$Y_{jk} = \sqrt{\frac{2}{n}} \epsilon_k \cos \left(\frac{\pi(2j+1)k}{2n} \right), \quad (13)$$

where $j, k \in \llbracket 0, n-1 \rrbracket$, $\epsilon_0 = 1/\sqrt{2}$, $\epsilon_i = 1$ for $i = 1, \dots, n-1$. We used a sub-sampled version of these matrices in which we picked some rows randomly.

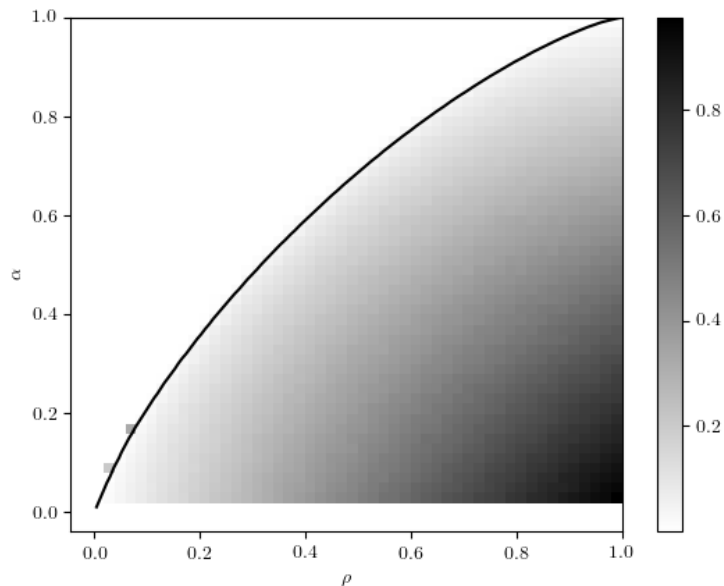


Figure 1. Phase diagram for a DCT matrix (width $n = 1000$) in the Bayes-optimal case. The averaged MSE on 50 executions of VAMP is represented by a color-code, displaying a phase transition that matches the theoretical Bayes line for Gaussian i.i.d. matrices (black line). Some finite-size effects can be seen.

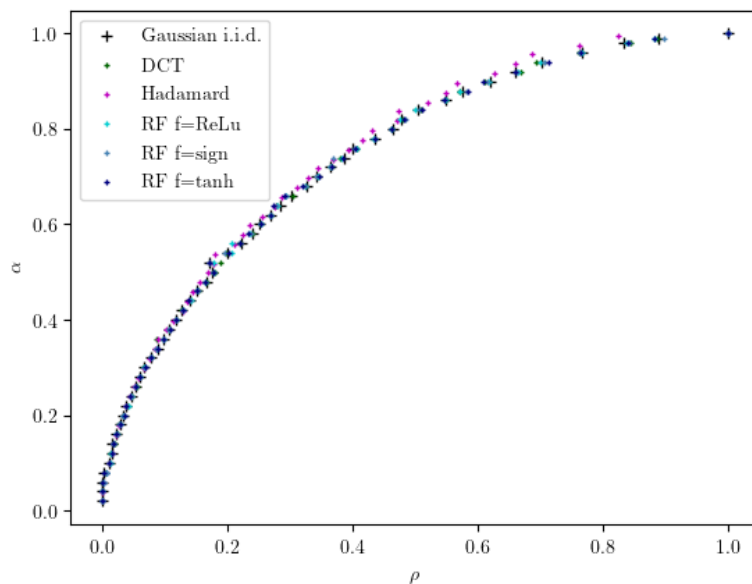


Figure 2. Phase diagram in the ℓ_1 reconstruction case obtained by averaging on 20 to 50 executions on VAMP. The dots indicate the phase transitions for Gaussian i.i.d., DCT (width $n = 2000$), Hadamard matrices ($n = 4096$), and random feature matrices $\Phi = f(WX)$ with $f = \text{ReLU}$, $f = \text{sign}$, $f = \tanh$ (W and X are Gaussian i.i.d. of size $\alpha n \times n$ and $n \times n$ with $n = 2000$). They match the theoretical Donoho-Tanner transition for Gaussian i.i.d. matrices (black line).

Hadamard matrices A natural variant of DCT is given by the Hadamard matrices. H is a $n \times n$ Hadamard matrix if its entries are ± 1 and its rows are pairwise orthogonal, i.e. $HH^T = nI_n$. For every integer k , there exists a Hadamard matrix H_k of size 2^k . These can be created with Sylvester’s construction: Let H be a Hadamard matrix of order n . Then the partitioned matrix

$$\begin{bmatrix} H & H \\ H & -H \end{bmatrix}$$

is a Hadamard matrix of order $2n$.

Random features maps Finally, we wanted to consider here random features maps (RFM) as encountered in nonlinear regression problems. In such settings, a random features matrix $\Phi = f(WX)$ is obtained from the raw data matrix X by means of a random projection matrix W and a pointwise nonlinear activation f . Kernel regression models, nonlinear in the original data X , can then be approximately but efficiently solved by the linear estimation problem (1), with an appropriate choice for f and the W -distribution [34]. Such matrices, that can be seen as the output of a neuron with random weights, have been investigated in particular in the context of neural networks [7, 8]. Indeed, in neural networks configurations with random weights play an important role as they define the initial loss landscape. They are also fundamental in the random kitchen sinks algorithm in machine learning [34] and it is thus of interest to test our understanding of linear reconstructions with AMP and VAMP in this case.

In what follows we will test random features matrices where both W and X are random Gaussian i.i.d. matrices.

5.2. Numerical results

We provide the codes used to generate the data on github in the repo <http://sphinxteam/Universality-CS-2019>. To generate Figure 1 and 2, we ran VAMP 50 times on 50×50 points spanning the (α, ρ) -space, and computed the average mean-squared error (MSE) between the signal \mathbf{x}^* and the reconstructed configuration \mathbf{x} . The MSE is represented with a color bar (white means perfect reconstruction). For a DCT and a Hadamard matrix, we observe a phase transition in the Bayes-optimal case that matches the theoretical transition for Gaussian i.i.d. matrices. We also ran VAMP for the ℓ_1 reconstruction problem. Averaging on 20 executions (or 50 for small α where finite-size effects are more important), we recover again a phase transition matching the theoretical Donoho-Tanner line for Gaussian i.i.d. matrices [3]. Besides, we compared the MSE obtained by VAMP at each point of the phase diagram for different matrices. In figures 3 and 4, we plot the MSE averaged on 20 executions for ρ fixed and α ranging between 0 and 1. We get the same error in reconstruction for all matrices, following the MSE for Gaussian i.i.d. matrix for $\rho = 0.25, 0.5$ and 0.75 . We also checked that AMP, provided one uses the trick eq. (7), reproduce these results as well: indeed the two algorithms returned extremely similar results.

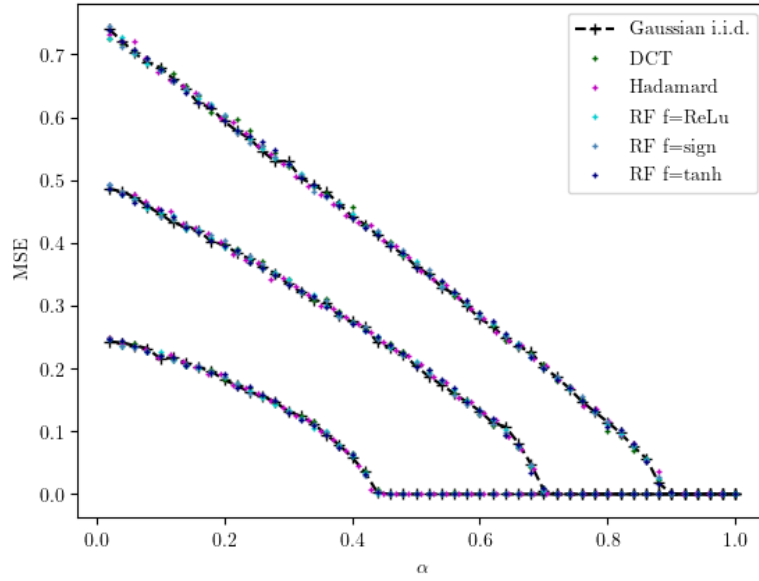


Figure 3. Mean-squared error for $\rho = 0.25, 0.5$ and 0.75 (bottom to up curves) in the Bayes-optimal case averaged on 20 executions of VAMP for Gaussian i.i.d, DCT, Hadamard, random features matrices $\Phi = f(WX)$ with $f = \text{ReLU}$, $f = \text{sign}$, $f = \text{tanh}$ (W and X are Gaussian i.i.d of size $\alpha n \times n$ and $n \times n$). The width is $n = 2000$ for all matrices.

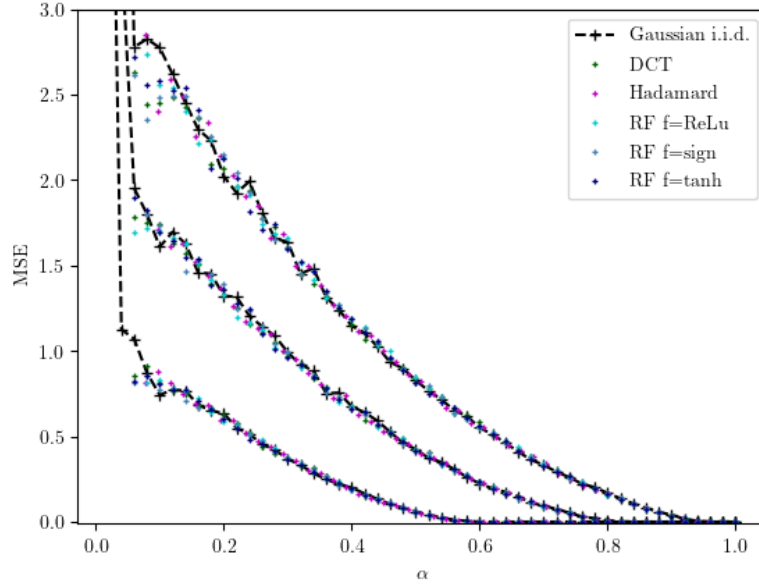


Figure 4. Mean-squared error for $\rho = 0.25, 0.5$ and 0.75 (bottom to up curves) in the ℓ_1 reconstruction case averaged on 20 executions of VAMP for Gaussian i.i.d, DCT, Hadamard, random features matrices $\Phi = f(WX)$ with $f = \text{ReLU}$, $f = \text{sign}$, $f = \text{tanh}$ (W and X are Gaussian i.i.d of size $\alpha n \times n$ and $n \times n$). The width is $n = 2000$ for all matrices.

5.3. Discussion

Figures of the previous section perfectly illustrate our main point: the universality in noiseless compressed sensing is not limited to the ℓ_1 -type reconstruction as in [3, 4], but extends to other quantities and estimators, such as the hard-phase line in Bayesian reconstruction, and the MMSE. Besides, it is not limited to random orthogonal matrices, but empirically extends to Fourier-type matrices and to the random features maps currently studied in machine learning. It is an open question to extend the proof of state evolution to these challenging matrices. It would be interesting to find a good criterion to identify which matrices satisfy this universality and which do not; this is something that we are yet unable to predict in advance. An example of structured matrices that do not seem to follow these universal phase transitions is given by Haar wavelet matrices, which can be defined recursively by:

$$W_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \text{ and } W_{2k} = \begin{bmatrix} H_k \otimes [1, -1] \\ I_k \otimes [1, 1] \end{bmatrix}$$

where I_k is the identity matrix of size k and \otimes is the Kronecker product. In fact, VAMP even fails to converge for these matrices. Investigating this behavior is an interesting direction of research.

Acknowledgment

We thank Andre Manoel and Galen Reeves for useful discussions. We acknowledge funding from the ERC under the European Union's Horizon 2020 Research and Innovation Program Grant Agreement 714608-SMiLe; and from the French National Research Agency (ANR) grant PAIL.

- [1] 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.
- [2] D. L. Donoho and J. Tanner, "Sparse nonnegative solution of underdetermined linear equations by linear programming," *Proceedings of the National Academy of Sciences*, vol. 102, no. 27, pp. 9446–9451, 2005.
- [3] D. Donoho and J. Tanner, "Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing," *Phil. Trans. A: Mathematical, Physical and Engineering Sciences*, vol. 367, no. 1906, pp. 4273–4293, 2009.
- [4] H. Monajemi, S. Jafarpour, M. Gavish, , and D. L. Donoho, "Deterministic matrices matching the compressed sensing phase transitions of gaussian random matrices," *Proceedings of the National Academy of Sciences*, vol. 110, no. 4, pp. 1181–1186, 2013.
- [5] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky, "The convex geometry of linear inverse problems," *Foundations of Computational mathematics*, vol. 12, no. 6, pp. 805–849, 2012.
- [6] D. Amelunxen, M. Lotz, M. B. McCoy, and J. A. Tropp, "Living on the edge: phase transitions in convex programs with random data," *Information and Inference: A Journal of the IMA*, vol. 3, no. 3, pp. 224–294, 2014.

- [7] J. Pennington and P. Worah, “Nonlinear random matrix theory for deep learning,” in *Advances in Neural Information Processing Systems*, 2017, pp. 2637–2646.
- [8] Z. Liao and R. Couillet, “On the spectrum of random features maps of high dimensional data,” in *International Conference on Machine Learning*, 2018, pp. 3069–3077.
- [9] D. J. Thouless, P. W. Anderson, and R. G. Palmer, “Solution of ‘solvable model of a spin glass,’” *Philosophical Magazine*, vol. 35, no. 3, pp. 593–601, 1977.
- [10] M. Bayati and A. Montanari, “The dynamics of message passing on dense graphs, with applications to compressed sensing,” *IEEE Transactions on Information Theory*, vol. 57, no. 2, pp. 764–785, 2011.
- [11] M. Bayati, M. Lelarge, A. Montanari *et al.*, “Universality in polytope phase transitions and message passing algorithms,” *The Annals of Applied Probability*, vol. 25, no. 2, pp. 753–822, 2015.
- [12] T. Tanaka, “Statistical mechanics of CDMA multiuser demodulation,” *EPL (Europhysics Letters)*, vol. 54, no. 4, p. 540, 2001.
- [13] 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, 2013.
- [14] F. Krzakala, M. Mézard, F. Sausset, Y. Sun, and L. Zdeborová, “Probabilistic reconstruction in compressed sensing: algorithms, phase diagrams, and threshold achieving matrices,” *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2012, no. 08, p. P08009, 2012.
- [15] —, “Statistical-physics-based reconstruction in compressed sensing,” *Physical Review X*, vol. 2, no. 2, p. 021005, 2012.
- [16] J. Zhu and D. Baron, “Performance regions in compressed sensing from noisy measurements,” in *2013 47th Annual Conference on Information Sciences and Systems (CISS)*. IEEE, 2013, pp. 1–6.
- [17] J. Barbier, M. Dia, N. Macris, and F. Krzakala, “The mutual information in random linear estimation,” in *Communication, Control, and Computing (Allerton), 2016 54th Annual Allerton Conference on*, 2016, pp. 625–632.
- [18] G. Reeves and H. D. Pfister, “The replica-symmetric prediction for compressed sensing with gaussian matrices is exact,” in *Information Theory (ISIT), 2016 IEEE International Symposium on*, 2016, pp. 665–669.
- [19] J. Barbier, F. Krzakala, N. Macris, L. Miolane, and L. Zdeborová, “Optimal errors and phase transitions in high-dimensional generalized linear models,” *Proceedings of the National Academy of Sciences*, vol. 116, no. 12, pp. 5451–5460, 2019.
- [20] J. Vila and P. Schniter, “Expectation-maximization bernoulli-gaussian approximate message passing,” in *Signals, Systems and Computers (ASILOMAR), 2011 Conference Record of the Forty Fifth Asilomar Conference on*. IEEE, 2011, pp. 799–803.
- [21] A. Montanari, “Graphical models concepts in compressed sensing,” *Compressed Sensing: Theory and Applications*, pp. 394–438, 2012.
- [22] A. M. Tulino, S. Verdú *et al.*, “Random matrix theory and wireless communications,” *Foundations and Trends® in Communications and Information Theory*, vol. 1, no. 1, pp. 1–182, 2004.
- [23] Y. Kabashima, T. Wadayama, and T. Tanaka, “A typical reconstruction limit for compressed sensing based on lp-norm minimization,” *J. of Stat. Mech.: Theory and Experiment*, vol. 2009, no. 09, p. L09003, 2009.
- [24] 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.
- [25] J. Barbier, N. Macris, A. Maillard, and F. Krzakala, “The mutual information in random linear estimation beyond i.i.d. matrices,” in *2018 IEEE International Symposium on Information Theory (ISIT)*, 2018.
- [26] 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. 18 914–18 919, 2009.
- [27] K. Takeuchi, “A unified framework of state evolution for message-passing algorithms,” *arXiv*

- preprint arXiv:1901.03041*, 2019.
- [28] B. Cakmak, O. Winther, and B. H. Fleury, “S-amp: Approximate message passing for general matrix ensembles,” in *2014 IEEE Information Theory Workshop (ITW 2014)*. IEEE, 2014, pp. 192–196.
 - [29] 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.
 - [30] M. Oppor and O. Winther, “Expectation consistent approximate inference,” *Journal of Machine Learning Research*, vol. 6, p. 2177, 2005.
 - [31] J. Ma and L. Ping, “Orthogonal amp,” *IEEE Access*, vol. 5, pp. 2020–2033, 2017.
 - [32] S. Rangan, P. Schniter, and A. K. Fletcher, “Vector approximate message passing,” in *2017 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2017, pp. 1588–1592.
 - [33] 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.
 - [34] A. Rahimi and B. Recht, “Random features for large-scale kernel machines,” in *Advances in neural information processing systems*, 2008, pp. 1177–1184.