

How to iron out rough landscapes and get optimal performances: averaged gradient descent and its application to tensor PCA

Giulio Biroli¹ Chiara Cammarota² and
Federico Ricci-Tersenghi^{3,4} 

¹ Laboratoire de Physique de l'Ecole Normale Supérieure ENS, Université PSL, CNRS, Sorbonne Université, Université Paris-Diderot, Sorbonne Paris Cité, Paris, France

² Department of Mathematics, King's College London, Strand, London WC2R 2LS, United Kingdom

³ Dipartimento di Fisica, Sapienza Università di Roma, INFN—Sezione di Roma1, and CNR-Nanotec, unità di Roma, P.le A. Moro 5, Roma 00185 Italy

E-mail: giulio.biroli@ens.fr, chiara.cammarota@kcl.ac.uk and federico.ricci@uniroma1.it

Received 12 November 2019, revised 2 February 2020

Accepted for publication 28 February 2020

Published 8 April 2020



CrossMark

Abstract

In many high-dimensional estimation problems the main task consists in minimizing a cost function, which is often strongly non-convex when scanned in the space of parameters to be estimated. A standard solution to flatten the corresponding rough landscape consists in summing the losses associated to different data points and obtaining a smoother empirical risk. Here we propose a complementary method that works for a single data point. The main idea is that a large amount of the roughness is uncorrelated in different parts of the landscape. One can then substantially reduce the noise by evaluating an empirical average of the gradient obtained as a sum over many random independent positions in the space of parameters to be optimized. We present an algorithm, called averaged gradient descent, based on this idea and we apply it to tensor PCA, which is a very hard estimation problem. We show that averaged gradient descent over-performs physical algorithms such as gradient descent and approximate message passing and matches the best algorithmic thresholds known so far, obtained by tensor unfolding and methods based on sum-of-squares.

Keywords: gradient descent, tensor PCA, averaged gradient descent

⁴ Author to whom any correspondence should be addressed.

 Supplementary material for this article is available [online](#)

(Some figures may appear in colour only in the online journal)

1. Introduction

One recurrent central task in many modern machine learning problems is the minimization of a non-convex high-dimensional function. Gradient descent is a versatile workhorse method that is widely used in these contexts, in particular in high-dimensional estimation to optimize the likelihood function. However the performance of gradient descent can be substantially undermined in cases where the function to be optimized—or informally the landscape—is rough. One way out is to increase the signal to noise ratio by summing the losses associated to different data points and obtain a smoother empirical risk.

In this work we propose an alternative method which works for a single data point. Our main idea is that a large amount of the roughness is uncorrelated in different parts of the landscape. By evaluating an empirical average of the gradient obtained as a sum over many random independent positions in the space of parameters to be optimized, one can then substantially reduce the noise, thus effectively ironing out the landscape and letting the signal contribution emerge.

We propose an algorithm, called averaged gradient descent (AGD), based on this idea. We test it on tensor-PCA [1], a very hard high-dimensional estimation problem in which one observes a k -fold $N \times N \times \dots \times N$ tensor

$$\mathbf{T} = \mathbf{W} + \frac{\lambda}{N^{\frac{k-1}{2}}} \mathbf{v}^{\otimes k}, \quad (1)$$

where \mathbf{W} is a symmetric noise tensor with independent normally distributed elements and λ represents the signal to noise ratio (SNR). The aim is to recover the signal $\mathbf{v} \in \mathbb{R}^N$ with $\|\mathbf{v}\|_2 = \sqrt{N}$. Without loss of generality one can take \mathbf{v} pointing in a random direction on the surface $S^{N-1}(0, \sqrt{N})$ of an hyper-sphere of radius \sqrt{N} centered in the origin. The maximum likelihood estimate of \mathbf{v} is the vector \mathbf{x}^* of norm \sqrt{N} that minimizes $\sum_{i_1 \leq \dots \leq i_k} (T_{i_1, \dots, i_k} - x_{i_1} \dots x_{i_k})^2$, which leads to

$$\mathbf{x}^* = \arg \max_{\mathbf{x} \in S^{N-1}} \sum_{i_1 \leq \dots \leq i_k} T_{i_1, \dots, i_k} x_{i_1} \dots x_{i_k}. \quad (2)$$

From a statistical mechanics perspective the previous equation on \mathbf{x}^* can be also seen as the minimization equation of the following energy function:

$$H(\mathbf{x}) = -\frac{1}{N^{\frac{k-1}{2}}} \sum_{i_1 \leq \dots \leq i_k} T_{i_1, \dots, i_k} x_{i_1} \dots x_{i_k} \quad (3)$$

$$= -\frac{1}{N^{\frac{k-1}{2}}} \sum_{i_1 \leq \dots \leq i_k} W_{i_1, \dots, i_k} x_{i_1} \dots x_{i_k} - N \frac{\lambda}{k!} m^k \quad (4)$$

where $m = (\mathbf{v}, \mathbf{x})/N = \sum_{i=1}^N v_i x_i / N$ is the overlap with the signal or the magnetization of configuration \mathbf{x} in statistical physics language.

It is known that when $\lambda > \lambda_{\text{IT}}(k)$ with $\lambda_{\text{IT}}(k)$ of order one in the large N limit (e.g. $\lambda_{\text{IT}}(3) \simeq 2.955$) it is information theoretically possible to recover the signal [1, 2]. However, a much larger SNR, $\lambda \gg N^{\frac{k-2}{4}}$, has to be reached in order to find algorithms, such as tensor unfolding and the ones based on sum of squares [1, 3], able to recover the signal in polynomial time. Gradient descent, and other physical algorithms as approximate message passing (AMP) and Langevin dynamics, are sub-optimal and succeed only for $\lambda \gg N^{\frac{k-2}{2}}$ [1, 4]. The

inefficiency of physical algorithms is conjectured to be related to the roughness of the energy landscape, which is characterized by an exponential number of minima in the band $m < m_{tr}$, with m_{tr} shrinking to zero as an inverse power of N for any λ growing sub-exponentially with N [5, 6] (see also SM). Tensor-PCA therefore provides a very good framework to test whether our method for ironing out the landscape using multiple uncorrelated copies is efficient and able to match the performance of the best algorithms not based on the landscape (e.g. spectral methods). We show that this is indeed the case. An additional outcome of our analysis is the demonstration by systematic numerical studies that the *algorithmic threshold* of AGD is associated to a *threshold phenomenon* (a phase transition in physics jargon) that we fully characterize.

Our approach is grounded on the research axis which aims to connect the behavior of dynamics and algorithms to landscape properties, and to exploit the knowledge of the latter to improve the performance of the former. In fact, as we shall show, AGD is an algorithm fully rooted on physical intuition and aimed at optimally exploiting the information gained from the analysis of the landscape. In matching the best algorithmic performances achieved so far for tensor-PCA, AGD re-establishes the competitiveness of landscape-based algorithms originating from statistical physics and clarifies to what extent algorithmic transitions are determined by *landscape properties*. From a more general point of view, AGD inherits the versatility of gradient descent, and hence stands as a new efficient algorithm suitable to a very wide spectrum of applications.

2. Related works

Different procedures have been devised to regularize a rough landscape and improve optimization performance. One approach is based on the convolution of a rough energy function with a smoothing kernel [7]. Another procedure is based on the introduction of different copies of the system which are coupled together [8]. In both cases, the idea is to reduce the roughness by smoothing the landscape *locally* on sets of points with high overlap. Our method, instead, aims at reducing the roughness by a much more *global* average over uncorrelated copies, which have typically zero overlap.⁵

Among the many algorithms devised for tensor PCA, the one based on homotopy [9] is the closest one to AGD, although it was introduced from a very different perspective. From a general point of view, the main difference is that our method can be straightforwardly applied and extended to generic high-dimensional inference problems. We will compare in more detail later in the main text and in SM8 the two methods and their performances in the context of tensor PCA, showing the superiority of AGD in the large N limit.

Finally, we point out that the optimality gap between algorithms not based on the landscape and statistical physics methods was very recently bridged by an extension of approximate message passing based on the Kikuchi approximation [10]. Our results show that the gap can be also closed by using an extension of gradient descent. In this way a full redemption [10] of the landscape dominated statistical physics approach against sophisticated algorithms not based on the landscape is reached.

3. Averaged gradient descent

The approach we propose here aims at being completely general. It takes advantage of physical intuition for the construction of a simple gradient-descent-based algorithm able to navigate

⁵ One may argue that Hamming distances between copies are extensive both in local and global methods, however copies in AGD are as spread as possible, while local methods must keep copies close enough.

through rough landscapes, hence, reaching very good algorithmic performances. Averaged gradient descent uses the simple idea that sampling several independent locations, called *real replicas* of the system, helps decreasing the roughness of the landscape which originates from uninformative corrupting noise. In fact, the average over the replicas leads to a relative amplification of the informative contribution produced by the signal with respect to the noise. Note that AGD can be generalized and potentially applied to the broad range of problems in high-dimensional inference (i.e. other tensor problems [11], compressed sensing [12], community detection [9, 13, 14], learning graphical models [15] just to mention a few examples) that, in certain regimes of the parameters, are characterized by a hard phase where uninformative spurious minima trap local dynamics and hamper the reconstruction of the signal. As anticipated in the introduction, in what follows we enter in the details of the application of this new algorithm to tensor-PCA, which is a notoriously hard problem in this sense, and we comment on the possibility of its generalization.

The heuristics behind this algorithm is very simple and can be discussed in full generality. For non zero signal to noise ratio λ , whenever the local information on the gradient \mathbf{g}_α contains a tiny component $\lambda \mathbf{g}_{s,\alpha}$ systematically pointing in the direction of the signal, this algorithm aims at getting it amplified with respect to the complementary component of the gradient that is originated by uninformative corrupting noise, $\mathbf{g}_{n,\alpha} = \mathbf{g}_\alpha - \lambda \mathbf{g}_{s,\alpha}$.⁶ For a given sample and considering different configurations \mathbf{x} drawn at random on the sphere, $\mathbf{g}_{n,\alpha}(\mathbf{x})$ is expected to have a strong fluctuating part and a small average $\mathbf{g}_{n,\alpha}^{\text{av}}$. The central limit theorem implies that averaging over R independent replicas of the system leads to a suppression by a factor $1/\sqrt{R}$ of the fluctuating part. By these simple arguments we conclude that the averaged algorithm will end up operating under a much higher effective signal to noise ratio. However, above a certain number of replicas, R_{opt} , the fluctuating part becomes subleading with respect to $\mathbf{g}_{n,\alpha}^{\text{av}}$ and one cannot iron out more the landscape. Thus, R_{opt} sets the optimal number of replicas that have to be used in practice. This number is evaluated for tensor PCA in the supplementary material (SM3).

The explanation above holds, and the proposed algorithm gives a net advantage in the retrieval of the signal, when the problem has only one optimal solution. This situation corresponds for instance to the case of tensor PCA with k odd. When two degenerate solutions are present, e.g. for tensor PCA with k even or any other inference problem where the global sign of the solution does not really matter, the multiple sampling of the landscape at $t = 0$ through independent different copies of the system will not be of any help. The reason is that the local gradient sampled through R different replicas will be randomly pointing toward any of the two solutions and their average will be suppressed by a factor $1/\sqrt{R}$, i.e. exactly at the same pace as the uninformative component originated by the noise. In this case we suggest to replace the averaged gradient in algorithm 1 by the eigenvector $\mathbf{w}_{\min}(t)$, with norm \sqrt{N} , corresponding to the minimum eigenvalue of the averaged Hessian $\sum_\alpha \mathcal{H}_{ij} |_{\mathbf{x}_\alpha(t)}/R$. This procedure can be seen, for $t = 0$, as a new general way to obtain spectral methods for high-dimensional inference problems [16]. An additional care is needed here to keep consistency, step by step, of the sense of the update vector. This issue is solved by asking that the scalar product is $(\mathbf{w}_{\min}(t), \mathbf{w}_{\min}(t+1)) > 0$. After a few steps t^* the symmetry between the two solutions is broken therefore it is advisable to continue with the original algorithm based on gradients, which is less computationally expensive.

⁶Note that the possibility to linearly decompose the gradient into these two terms can be considered general in the limit of small $\lambda \mathbf{g}_{s,\alpha}$, as in this limit each gradient can be expanded around the zero-signal limit and the expansion truncated to the first order. The same decomposition is straightforward in the case of tensor PCA.

Algorithm 1. Averaged gradient descent, AGD.

Input: landscape $H(\mathbf{x})$, number of replicas R , learning rate η , stopping criterion ε
Output: estimate of the location of the landscape minimum \mathbf{x}^*

- 1 $t \leftarrow 0; \mathbf{x}_{\text{CM}}(0) \leftarrow 0; r(0) \leftarrow 0$ // initialize the center of mass
- 2 Repeat
- 3 **For** $\alpha = 1, \dots, R$ **do** // given center of mass, sample R points on the sphere
- 4 $\mathbf{x}_\alpha(t) \leftarrow \mathbf{x}_{\text{CM}}(t) + \sqrt{1 - r^2(t)} \mathbf{u}_\alpha(t)$ with $\mathbf{u}_\alpha(t)$ drawn uniformly at random among vectors
such that $\|\mathbf{u}_\alpha(t)\|_2^2 = N$ and $(\mathbf{u}_\alpha(t), \mathbf{x}_{\text{CM}}(t)) = 0$
- 5 $\mathbf{g}_\alpha(t) \leftarrow \nabla H|_{\mathbf{x}_\alpha(t)}$ // evaluate the gradient on each of the R points $\mathbf{x}_\alpha(t)$
- 6 $\mathbf{x}_{\text{CM}}(t+1) \leftarrow \mathbf{x}_{\text{CM}}(t) - \eta \sum_\alpha \mathbf{g}_\alpha(t)/R$ // use the average gradient to update the position of
the center of mass
- 7 $t \leftarrow t + 1$
- 8 $r(t) \leftarrow \|\mathbf{x}_{\text{CM}}(t)\|_2 / \sqrt{N}$
- 9 **If** $r(t) > 1$ **then** // keep the center of mass inside or on the sphere
- 10 $\mathbf{x}_{\text{CM}}(t) \leftarrow \mathbf{x}_{\text{CM}}(t)/r(t)$
- 11 $r(t) \leftarrow 1$ // when $r(t) = 1$ the algorithm reduces to standard GD
- 12 **Until** $\|\mathbf{x}_{\text{CM}}(t) - \mathbf{x}_{\text{CM}}(t-1)\|_2 < \varepsilon$ // stopping condition as in standard GD
- 13 **Return** $\mathbf{x}_{\text{CM}}(t)$

In the general case a good strategy is to compute for the first steps both the average gradient and the average Hessian with its lowest eigenvector. Among the two averaged vectors the one to be used is the one leading to larger decrease in the energy function. After few steps the average gradient should become larger and should point toward the signals (even in the symmetric case of even k), one can then continue with algorithm 1.

In the next sections we are going to focus more specifically on the performances of this algorithm on tensor PCA. Interestingly in this case not only the analysis at finite R can be performed but also the study at infinite R , which turns out to be even computationally convenient. Therefore in what follows we are going to focus on the large R limit of AGD, where empirical averages are substituted by expected values on the uniform measure over the space of variables. Such algorithm involving infinite real replicas is hereafter called iAGD (infinite- R AGD). Its performance will be discussed in the section numerical results. Note that to develop an analytic understanding of these results we resort to a further simplification of the algorithm as discussed in the next section. Finally the results for finite R will be quoted and explained in the supplementary material (SM7).

4. Theoretical analysis: from landscape properties to the performance of the simplest optimal algorithms

So far, we have introduced AGD and its $R \rightarrow \infty$ version called iAGD. Both algorithms are very challenging to be fully analyzed. For this reason, in this section we introduce a simplified version, SiAGD, and present its full theoretical analysis which provides several insights on the behavior of AGD and iAGD. In the next sections and in the SM, we then confirm these results and fully analyze these two algorithms by numerical experiments.

The key idea to simplify AGD and iAGD applied to tensor PCA is that both algorithms are characterized by two regimes: a first one where the norm of the center of mass increases from zero to \sqrt{N} , and a second one which corresponds to simple gradient descent (when the center of mass reaches the surface of the sphere all replicas fall on the center of mass). The simplified version that we analyze here, SiAGD, consists in modifying the first regime by

moving straight in the direction of the $t = 0$ (averaged) gradient until the center of mass hits the hyper-sphere. We shall show that SiAGD has an algorithmic threshold for the recovery of the signal, which is optimal compared to the ones of all the other algorithms known so far. Its numerical analysis and a comparison with iAGD and AGD is presented later. We will consider separately the odd and even k cases since the simplified algorithm is different, actually even simpler in the even case. Moreover, focusing on SiAGD as a simpler version of iAGD, we will work directly with averaged quantities. However it should be kept in mind that they can be estimated accurately using empirical averages over a large enough number of real replicas as discussed at the end of this section and in the SM (SM3 and SM7). Finally, we will always consider that the rate η is small enough so that the discrete updates in the algorithm can be considered a good approximation of a continuous time algorithm.

4.1. Case I: k odd

The value of the averaged gradient at $t = 0$ for iAGD reads:

$$g_i = -\frac{1}{N^{\frac{(k-1)}{2}}} \sum_{i_2 \leq \dots \leq i_k} W_{i,i_2,\dots,i_k} \mathbb{E}[x_{i_2} \dots x_{i_k}] \\ - \frac{\lambda}{(k-1)!N^{k-1}} v_i \sum_{i_2, \dots, i_k} v_{i_2} \dots v_{i_k} \mathbb{E}[x_{i_2} \dots x_{i_k}].$$

The expectation $\mathbb{E}[\cdot]$ is over the uniform measure on the sphere of radius \sqrt{N} . SiAGD consists in doing GD by using this initial averaged gradient until the norm of the center of mass reaches \sqrt{N} . Since the initial condition for the dynamics of the center of mass is the null vector, one obtains that at the end of the first regime the center of mass position equals

$$\mathbf{x}_I^{\text{CM}} = -\sqrt{N} \frac{\mathbf{g}}{\|\mathbf{g}\|_2}.$$

The second regime corresponds to gradient descent on the sphere with energy H starting from \mathbf{x}_I^{CM} .

It is easy to check that for N large the leading contribution to g_i is given by terms in which the indices i_2, \dots, i_k are grouped in $(k-1)/2$ distinct pairs of the same index. In these cases $\mathbb{E}[x_{i_2} \dots x_{i_k}]$ is simply equal to one. For example, for $k = 3$, one obtains:

$$g_i = -\frac{1}{N} \sum_j W_{i,j,j} - \frac{\lambda}{2N^2} v_i \left(\sum_j v_j^2 \right) = -\frac{1}{N} \sum_j W_{i,j,j} - \frac{\lambda}{2N} v_i.$$

The first contribution to \mathbf{g} , corresponding to \mathbf{W} , is a random Gaussian vector with norm scaling as $N^{\frac{3-k}{4}}$ and the second is a vector in the direction of the signal, v_i , of norm scaling as $\lambda N^{\frac{2-k}{2}}$. If the second term is the largest, i.e. for λ growing faster than $N^{\frac{k-1}{4}}$, a finite overlap with the signal, $m_I = (\mathbf{x}_I^{\text{CM}}, \mathbf{v})/N$ is already obtained at the end of the first regime. See the SM for a detailed derivation of the results. In the following we focus on the more challenging SNR regime, $N^{\frac{k-3}{4}} \ll \lambda \ll N^{\frac{k-1}{4}}$, where the first term has the largest norm and the overlap with the signal at the end of the first dynamical regime is approximately equal to

$$m_I \approx \frac{1}{(k-1)!} \lambda N^{\frac{1-k}{4}}.$$

How large this value of m_l has to be to guarantee recovery using gradient descent in the second dynamical regime? The answer to this question comes from the analysis of the number of spurious minima of H for configurations with overlap larger or equal to m_l . The results of [5, 6] obtained by the Kac–Rice method, imply that if $\lambda m_l^{k-2} > C_k$ (C_k does not scale with N and is computed in the SM) then such number is not exponentially large in N , i.e. the initial condition for the gradient descent dynamics lies in the ‘easy’ part of the configuration space where spurious minima that can trap the dynamics do not proliferate. This is the *crucial criterion* that guarantees recovery by gradient descent dynamics.

Let us first show that this criterion allows to recover the results for the GD algorithm. The initial condition for GD is a vector drawn uniformly at random on the sphere, which has typically an overlap with the signal of the order of $1/\sqrt{N}$. Thus, the previous criterion requires λ scaling as $N^{\frac{k-2}{2}}$ for gradient descent to recover the signal, which is indeed the threshold conjectured⁷ in [4] and heuristically re-derived in more details in SM1. This is also the scaling of algorithms such as approximate message passing and Langevin dynamics [1, 4]. SiAGD instead provides for gradient descent in the second dynamical regime an initial condition which has an overlap m_l possibly larger than $1/\sqrt{N}$. Imposing that $\lambda m_l^{k-2} > C_k$ allows us to find the algorithmic threshold for SiAGD:

$$\lambda > C'_k N^{\frac{k-2}{4}}$$

where C'_k is an N -independent constant that can be straightforwardly related to C_k . Using this scaling one finds that m_l is at least of order $N^{-\frac{1}{4}}$.

We have therefore obtained two main results: we have shown that a simplified version of iAGD allows to match the performance of the best known algorithms, which is $\lambda \sim N^{\frac{k-2}{4}}$ [1, 3, 9], and we have derived such an optimal algorithmic transition directly resorting to the statistical properties of the landscape. Both results will be tested and confirmed numerically in the next section.

Finally, we notice that the second regime of SiAGD shares similarities with the homotopy-based algorithm studied in [9]: they both used the same initial condition (i.e. what is reached at the end of the first stage of dynamics of SiAGD), but the latter consists in gradient descent with $\eta = \infty$, and this is less efficient than AGD (see discussion in SM8).

4.2. Case II: k even

For even values of k , the initial value of the average gradient is exactly zero since $\mathbb{E}[x_{i_2} \cdots x_{i_k}] = 0$. In this case, as discussed previously, one has to focus on the averaged Hessian, which at the initial condition of the iAGD algorithm reads:

$$\begin{aligned} \mathcal{H}_{ij} = & -\frac{1}{N^{\frac{(k-1)}{2}}} \sum_{i_3 \leq \dots \leq i_k} W_{i,j,i_3,\dots,i_k} \mathbb{E}[x_{i_3} \cdots x_{i_k}] \\ & - \frac{\lambda v_i v_j}{(k-2)! N^{k-1}} \sum_{i_3, \dots, i_k} v_{i_3} \cdots v_{i_k} \mathbb{E}[x_{i_3} \cdots x_{i_k}] \end{aligned}$$

The leading contribution to \mathcal{H}_{ij} is given by terms in which the indices i_3, \dots, i_k are grouped in distinct pairs. In this case the average $\mathbb{E}[x_{i_3} \cdots x_{i_k}]$ is simply equal to one. For example, for

⁷ It was shown rigorously that λ scaling as $N^{\frac{k-2}{2} + \frac{1}{6}}$ is a sufficient condition for GD initialized from a random uniform initial condition to recover the signal. As argued in [4], it should be possible to obtain a tighter bound and remove the $1/6$ factor by generalizing the proof of [4].

$k = 4$, one obtains:

$$\mathcal{H}_{ij} = -\frac{1}{N^{3/2}} \sum_k W_{i,j,k,k} - \frac{\lambda}{2N^3} v_i v_j \left(\sum_k v_k^2 \right) = -\frac{1}{N^{3/2}} \sum_k W_{i,j,k,k} - \frac{\lambda}{2N^2} v_i v_j$$

The first term of \mathcal{H} is a random matrix belonging to the Gaussian orthogonal ensemble [17], whereas the second term is a rank one perturbation proportional to the projector in the direction of the signal. Such random matrices display an interesting phenomenon called BBP transition (Baik, Ben Arous and Peché [18, 19]): given a symmetric matrix with random elements extracted from a normal distribution $\mathcal{N}(0, 1/N)$ perturbed by a rank one matrix $-\alpha v_i v_j / N$ with $\|\mathbf{v}\|_2 = N$, in the large N limit there exists a finite $\alpha_{\text{BBP}} = 1$ such that for $\alpha > \alpha_{\text{BBP}}$ the eigenvector associated to the smallest eigenvalue of the matrix has a finite overlap with \mathbf{v} . By taking into account the specific scaling with N of the two terms in the Hessian we get that at large N for $\lambda > (k-2)! \alpha_{\text{BBP}} N^{\frac{k-2}{4}}$ the eigenvector corresponding to the smallest eigenvalue of the Hessian has a finite overlap with the signal. In consequence, for SNR above $N^{(k-2)/4}$, in the even k case, the information about the signal is present in the initial averaged Hessian: already at the beginning of the dynamics, by averaging over different replicas, a downward direction toward the signal emerges. At variance with the k odd case, a simplified iAGD algorithm that consists in moving the center of mass in the direction of the eigenvector associated to the smallest eigenvalue of the initial averaged Hessian until hitting the sphere with radius \sqrt{N} is already enough to obtain the best algorithmic performance. For even values of k , the second regime of SiAGD, corresponding to gradient descent on the sphere, is not even needed to obtain a finite overlap.

It is interesting to contrast the result above with the one for the Hessian obtained for a random vector drawn uniformly on the sphere, which is a typical initial condition for the GD algorithm. Repeating the previous analysis, one finds a similar result—a GOE matrix perturbed by a rank one perturbation in the direction of the signal—but now the BBP transition takes place for $\lambda > (k-2)! N^{\frac{k-2}{2}}$, which is indeed the conjectured scaling to recover the signal by gradient descent [4].

The analysis performed above can be repeated for a finite number of replicas, hence bridging the gap between the performance of the GD and SiAGD algorithm. For finite R one finds that the algorithmic transition is at $\lambda(R) \sim N^{(k-2)/2} R^{-0.5(k-2)/(k-1)}$ (see SM3). The use of $R > 1$ different initial configurations helps reducing the algorithmic gap: the larger is R the smaller the algorithmic threshold is. As explained in the SM, the smoothing of the landscape using different replicas becomes ineffective when $R \gg R_{\text{opt}} \sim N^{(k-1)/2}$. However, for these values of R one has already reached the regime studied above.

In summary, in the odd and even k cases, we find that the analysis of the ‘bare’ landscape naturally leads to the scaling of the algorithmic threshold as $N^{\frac{k-2}{2}}$ whereas the analysis performed using many replicas allow to substantially averaging out the noise and to match the best scaling currently known, which is $N^{\frac{k-2}{4}}$.

We have found that the k -even case is simpler than the k -odd one; this finding emerges also from the previous literature (more involved methods were used to obtain the scaling $N^{\frac{k-2}{4}}$ for odd values of k), but was not explained. Our landscape based analysis offers a simple reason for it.

5. Numerical results

In this section we present the results of our numerical tests, which are limited to the $k = 3$ case because the memory requirements scale like N^k and thus for larger values of k one is limited to very small values of N . The aim of this section is twofold: on the one hand we want to

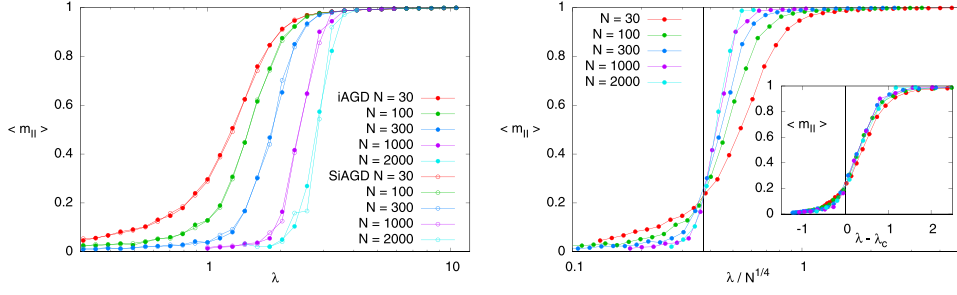


Figure 1. Left: iAGD and SiAGD achieve the same accuracy detecting the signal in tensor PCA with $k = 3$. Right: their algorithmic threshold scales as $\lambda_c \simeq 0.37N^{1/4}$ (only data for iAGD are shown). Inset: the final overlap with the signal mostly depends on $\lambda - \lambda_c$.

identify the algorithmic thresholds for both the full and simplified versions of iAGD, on the other hand we wish to directly test the connection between iAGD and SiAGD performance and the properties of the energy landscape. Numerical results for AGD with finite R and their comparison with iAGD are reported in SM7. As discussed in the previous section, it was shown that there exist no spurious minima [5] such that its overlap with the signal satisfies $\lambda m^{k-2} > C_k$ (for $k = 3$ one finds $C_3 \simeq 0.425\,815$, see SM for further details). In the following we are going to show numerically that such condition is directly related to the algorithmic threshold of iAGD. The results we present are obtained for runs of iAGD and SiAGD on problems of sizes $N = 30, 100, 300, 1000, 2000$. They are then averaged over a number M of different disorder realizations such that $NM = 1.2 \times 10^5$.

- *Algorithmic threshold and threshold phenomenon.* In figure 1 (left panel) we show the mean overlap with the signal, m_{II} , achieved at the end of the algorithm (either for iAGD or SiAGD) as a function of the signal to noise ratio λ . In the right panel we show that a threshold phenomenon (a phase transition) is taking place in the large N limit on the scale $\lambda \sim N^{(k-2)/4} = N^{1/4}$. It is worth noticing, as shown in the left panel, that both versions of the algorithm, iAGD and SiAGD, do achieve the same final mean overlap with the signal. For this reason in the right panel we have re-scaled only the data obtained via iAGD. In the right panel we also mark with a vertical line our best estimation for the critical threshold $\lambda_c \simeq 0.37N^{1/4}$. Finally, the inset shows the same results plotted as a function of $\lambda - 0.37N^{1/4}$. This highlights that the size of the critical window around the algorithmic threshold $\lambda_c \simeq 0.37N^{1/4}$ is almost N independent.
- *Comparison between iAGD and SiAGD.* Although the final overlap achieved by the two versions of the algorithm is the same, the dynamics followed by the algorithms in the first regime is very different (see previous section for the distinction of two regimes in the dynamics). While in the SiAGD algorithm the center of mass takes a straight path to the surface of the sphere of radius \sqrt{N} , in the iAGD algorithm the center of mass moves according to the mean gradient at each time and thus follows a curved trajectory determined by the landscape. *A priori* it is unclear which dynamics is better; we offer an insight by measuring the evolution of the center of mass during and at the end of the first regime.

In the left panel of figure 2 we report the mean overlap $\langle m_I \rangle$ achieved at the end of the first phase by the iAGD and SiAGD algorithms. We clearly see that the dynamics followed by the iAGD algorithm reaches a larger overlap. Therefore a natural question arises: how can SiAGD achieve the same accuracy in detection than iAGD although it starts from a lower value of

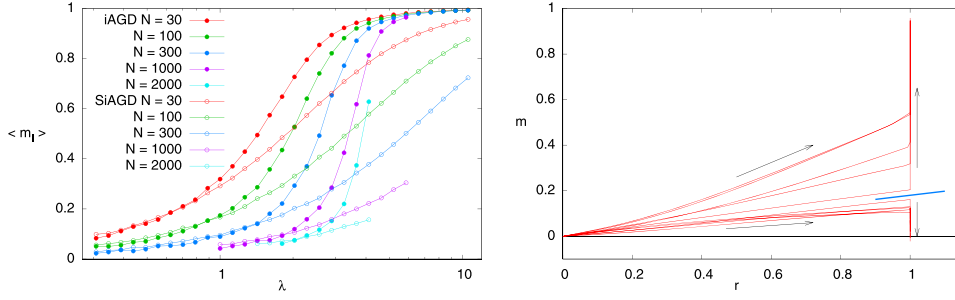


Figure 2. Left: at the end of the first dynamical regime iAGD achieves an overlap with the signal larger or equal to the one achieved by SiAGD. Right: A schematic picture of the trajectories followed by iAGD, represented by $r = \|\mathbf{x}_{CM}\|_2 / \sqrt{N}$ and the overlap m with the signal.

m_I ? While trying to answer this question, we notice an important difference between the two dynamics in the first phase: although both depend on the landscape, they feel the landscape in a quite different way. In the SiAGD algorithm the mean gradient is computed only once at the beginning. Then a straight path is followed until the center of mass hits the sphere. In this sense the algorithm in its first regime should be considered as a strongly out of equilibrium process that feels little of the original landscape and thus ends on a point on the sphere whose energy has not been optimized. SiAGD then secures its own connection to the landscape only in the second regime, where it continues with usual gradient descent that starts from this high energy configuration.

iAGD starts in a similar way computing the mean gradient when the center of mass is close to the origin. At this initial stage, the averaging process reaches its highest efficiency in ironing out the landscape as the replicas are completely uncorrelated. The gradient on the center of mass is much less affected by noise with respect to the one of single replicas. However, as soon as the center of mass starts to approach the sphere of radius \sqrt{N} the cloud of replicas shrinks, thus sampling a progressively smaller region of the landscape, until the mean gradient converges continuously to the standard gradient. Thus we expect iAGD to reach a point on the sphere of lower energy than SiAGD. This is explicitly shown in SM6. In summary iAGD and SiAGD algorithms reach the same accuracy in signal detection, although they land on the sphere on very different points, with iAGD reaching larger overlaps and lower energies.

- *Landscape and dynamics.* In the right panel of figure 2 we show the trajectories followed by the center of mass during the execution of the iAGD algorithm solving 10 problems of size $N = 300$ with $\lambda = 2$: we plot the overlap of the center of mass with the signal $m = (\mathbf{x}_{CM}, \mathbf{v})/N$ versus the normalized norm of the center of mass $r = \|\mathbf{x}_{CM}\|_2 / \sqrt{N}$. Recall that when $r = 1$ iAGD reduces to standard GD. Observing the plot it should be clear that there is a threshold value for the overlap on the sphere (marked by a thick blue line) such that when the algorithm hits the sphere above (below) that threshold value, then GD is able (not able) to recover the signal.

Moreover we notice that the trajectories of the runs that eventually detect the signal tend to bend upwards already in the first dynamical regime.

To better illustrate the threshold phenomenon in m_I we show in the left panel of figure 3 a scatter plot of the final overlap m_{II} versus λm_I . Clouds of points have different sizes for two reasons: for the smaller problems we have studied more samples and finite size effects

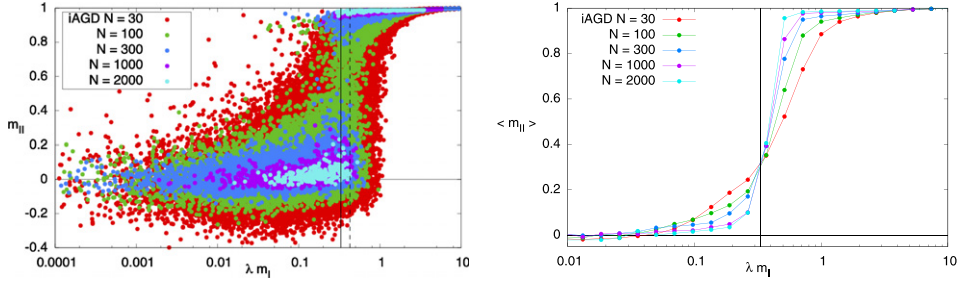


Figure 3. Signal detection is possible if λm_I is larger than the threshold value reported with a full vertical line, estimated from averaged data in the right panel. Complexity is null on the right of the dashed line.

tend to disperse the points more for smaller sizes. We clearly see that for large enough N the data points form two different and well separated clouds: the lower one corresponds to samples where iAGD has been unable to detect the signal, while the upper one corresponds to samples where signal detection was achieved. The choice of using a scaled overlap λm_I for the abscissa is dictated by the observation that the complexity of local minima depends only on the variable $\lambda m_I^{k-2} = \lambda m_I$ (for $k = 3$) in the large N limit and it is null with high probability for $\lambda m_I > C_3 = 0.425815$ (marked by a dashed vertical line in the plot). The full vertical line marks the location of the threshold estimated from the data shown in the right panel of figure 3: in the large N limit if iAGD reaches an overlap satisfying $\lambda m_I \gtrsim 0.33$ then it detects the signal with high probability. We have thus found that the numerically estimated threshold is slightly lower than the one where spurious minima disappear. This can be due to multiple reasons: first the result [5] used to estimate the number of minima only provides an upper bound, a quenched Kac–Rice computation [6] would be needed to obtain the exact value. Second, very recently it has been shown that landscape-based algorithm, such as GD, can succeed even in presence of spurious minima [20]. Moreover it has been also shown that the minima where these landscape-based dynamics end may depend on the starting energy and the most attracting minima are not the most numerous ones [21]. The inspection of this issue in further details is left for future work.

6. Conclusions and discussion

We have proposed a new algorithm which is a generalization of gradient descent and uses the idea that by averaging the gradient of uncorrelated copies of the system one can substantially reduce the roughness of the landscape. One of its main advantages is its generality; in fact, AGD can be straightforwardly and directly applied to many hard inference problems without any prior knowledge.

The spiked tensor problem has provided the perfect testing ground to study its performance, showing that it is slightly better (in the prefactor) than the state-of-the-art algorithms and much better (in the N scaling) than other algorithms based on the landscape.

It is worth discussing the superiority of AGD to AMP. The latter is often the provably best algorithm for models where variables interact in a dense and asymptotically very weak way. However the spiked tensor is one of those problems where AMP is sub-optimal therefore other algorithms can do better. An important message from our work is that the information locally available to AMP is much smaller than the one that can be collected with many uncorrelated

replicas allowing AGD to reach much better performances. Another way of understanding the deep difference between AGD and AMP-like algorithms is to consider the way the elements of the tensor are used by these different algorithms: in AMP all elements are used with a similar weight, while our algorithm gives much more weight to elements having pairs of identical indices. The two algorithms are extracting different information from the same tensor.

We have studied different versions of the algorithm—AGD, iAGD, SiAGD—because on the one hand the versions with infinite R can be solved analytically for the spiked tensor problem, on the other hand the version that presents the best performances is the one with finite R (see SM7) and shows that AGD has the potential to achieve unprecedented results already when working with a limited number of replicas.

Not only AGD outperforms the best available algorithm for signal recovery in the spiked tensor problem (comparison with homotopy is shown in SM8) but we believe it can be straightforwardly extended to more general problems in machine learning. For example we expect that problems where in general a low-rank signal, say a P -dimensional signal, is hidden by additive noise (in this manuscript we considered the $P = 1$ case) can be solved by AGD as the average gradient or the average lowest Hessian eigenvector would point toward the P -dimensional subspace containing the signal. Additional work on the application of AGD to problems of this kind is required to substantiate these claims.

Acknowledgments

We thank Marco Baity-Jesi, Gerard Ben Arous, Aukosh Jagannath, Florent Krzakala, Marc Mézard, Andrea Montanari and Lenka Zdeborová for interesting and very useful discussions. This work has been conceived and mainly developed at the Kavli Institute for Theoretical Physics within the program entitled ‘The Rough High-Dimensional Landscape Problem’, as such this research was supported in part by the National Science Foundation under Grant No. NSF PHY-1748958. We also acknowledge support by the Simons Foundation collaboration Cracking the Glass Problem (No. 454935 to G Biroli and No. 454949 to G Parisi).

ORCID iDs

Federico Ricci-Tersenghi  <https://orcid.org/0000-0003-4970-7376>

References

- [1] Montanari A and Richard E 2014 A statistical model for tensor PCA *Advances in Neural Information Processing Systems* pp 2897–905
- [2] Lesieur T, Miolane L, Lelarge M, Krzakala F and Zdeborová L 2017 Statistical and computational phase transitions in spiked tensor estimation *IEEE Int. Symp. on Information Theory (ISIT)* IEEE pp 511–5
- [3] Hopkins S B, Shi J and Steurer D 2015 Tensor principal component analysis via sum-of-square proofs *Conf. on Learning Theory* pp 956–1006
- [4] Ben Arous G, Gheissari R and Jagannath A 2018 Algorithmic thresholds for tensor PCA (arXiv:1808.00921)
- [5] Ben Arous G, Song M, Montanari A and Nica M 2017 The landscape of the spiked tensor model (arXiv:1711.05424)

- [6] Ros V, Ben Arous G, Biroli G and Cammarota C 2019 Complex energy landscapes in spiked-tensor and simple glassy models: ruggedness, arrangements of local minima, and phase transitions *Phys. Rev. X* **9** 011003
- [7] Wu Z 1996 The effective energy transformation scheme as a special continuation approach to global optimization with application to molecular conformation *SIAM J. Optim.* **6** 748–68
- [8] Baldassi C, Borgs C, Chayes J, Ingrosso A, Lucibello C, Saggiotti L and Zecchina R 2016 Unreasonable effectiveness of learning neural nets: accessible states and robust ensembles *Proc. Natl Acad. Sci. USA* **113** 7655–62
- [9] Anandkumar A, Deng Y, Ge R and Mobahi H 2017 Homotopy analysis for tensor PCA *Proc. Mach. Learn. Res.* vol 65 pp 1–26
- [10] Alexander S W, Ahmed E A and Moore C 2019 The Kikuchi hierarchy and tensor PCA (arXiv:1904.03858)
- [11] Hillar C J and Lim L-H 2013 Most tensor problems are np-hard *J. ACM* **60** 45
- [12] Donoho D L et al 2006 Compressed sensing *IEEE Trans. Inf. Theory* **52** 1289–306
- [13] Decelle A, Krzakala F, Moore C and Zdeborová L 2011 Inference and phase transitions in the detection of modules in sparse networks *Phys. Rev. Lett.* **107** 065701
- [14] Decelle A, Krzakala F, Moore C and Zdeborová L 2011 Asymptotic analysis of the stochastic block model for modular networks and its algorithmic applications *Phys. Rev. E* **84** 066106
- [15] Chaganty A T and Liang P 2014 Estimating latent-variable graphical models using moments and likelihoods *Int. Conf. on Machine Learning* pp 1872–80
- [16] Yue M L and Li G 2017 Phase transitions of spectral initialization for high-dimensional nonconvex estimation (arXiv:1702.06435)
- [17] Tao T 2012 *Topics in Random Matrix Theory* vol 132 (Providence, RI: American Mathematical Society)
- [18] Edwards S F and Jones R C 1976 The eigenvalue spectrum of a large symmetric random matrix *J. Phys. A: Math. Gen.* **9** 1595–603
- [19] Baik J, Ben Arous G and Pécché S 2005 Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices *Ann. Probab.* **33** 1643–97
- [20] Mannelli S S, Krzakala F, Urbani P and Zdeborova L 2019 Passed and spurious: analysing descent algorithms and local minima in spiked matrix-tensor model (arXiv:1902.00139)
- [21] Folea G, Franz S and Ricci-Tersenghi F 2019 Rethinking mean-field glassy dynamics and its relation with the energy landscape: the awkward case of the spherical mixed p-spin model (arXiv:1903.01421)