

---

# Gradient descent in Gaussian random fields as a toy model for high-dimensional optimisation in deep learning

---

**Mariano Chouza**

Tower Research Capital, London

**Stephen Roberts**

University of Oxford

**Stefan Zohren**

University of Oxford

## Abstract

In this paper we model the loss function of high-dimensional optimization problems by a Gaussian random field, or equivalently a Gaussian process. Our aim is to study gradient descent in such loss functions or energy landscapes and compare it to results obtained from real high-dimensional optimization problems such as encountered in deep learning. In particular, we analyze the distribution of the improved loss function after a step of gradient descent, provide analytic expressions for the moments as well as prove asymptotic normality as the dimension of the parameter space becomes large. Moreover, we compare this with the expectation of the global minimum of the landscape obtained by means of the Euler characteristic of excursion sets. Besides complementing our analytical findings with numerical results from simulated Gaussian random fields, we also compare it to loss functions obtained from optimisation problems on synthetic and real data sets by proposing a “black box” random field toy-model for a deep neural network loss function.

## 1 Introduction

For almost a decade there have been significant advances in many areas of machine learning by applying deep learning techniques, such as in the context of image recognition [1, 2], generative adversarial networks [3] and in reinforcement learning, most notably in the development of AlphaGo [4] (see also [5] for an overview). The amount of progress, combined with some issues that were found such as robust adversarial examples [6, 7], have led to interest in getting a better understanding of the underlying process. One contributing factor for the success of deep neural networks has to do with the nature and complexity of its loss function or energy landscape. In particular, it was

found that the loss function of deep neural networks has very similar properties to random fields or Gaussian processes [8, 9, 10]. For example, it was seen that the Hessian of such a loss function is mostly governed by the spectrum of a random matrix [8] which can be used to show that local minima are located in a band close to the global minimum.

Given the above evidence that loss functions of deep neural networks share many properties with those of random energy landscapes, we want to investigate further optimization procedures in such landscapes. In particular, we model the loss function of high-dimensional optimization problem as a Gaussian random field (GRF) [11], which can also be viewed as a Gaussian Process (GP) [12], and study, both theoretically as well as experimentally, the performance of gradient descent in such landscapes as well as properties of the global minimum. Recently, there has been a revived interest in studying distributional properties of GRFs within the research community working on GPs. Examples include the study of the distribution of arc length in GPs [13], as well as expected improvements in batch optimization [14]. We aim to fill a gap in the literature by studying distributional aspects of improvements in gradient descent in such landscapes, including proving asymptotic normality of the improved field value. Interesting scalings can be obtained by studying the optimal learning rate and comparing it with that of random search as well as the location of the global minimum both as a functions of the dimension of the parameter space.

As in [8] but differing from [9] and [10], we explicitly consider the field to be a function of the input and the parameters. Concretely, we choose the loss function to be a Gaussian random field  $\phi(\mathbf{x})$  with squared exponential correlation  $k(r)$  and constant mean  $\mu$ , where  $r = \|\mathbf{x} - \mathbf{x}'\|$  is the distance between two points. As the differences between different parameters will then be just be scaling and translation, we choose  $\mu = 0$  and  $k(r) = \exp(-r^2/2)$  for definiteness. As an example, Figure 1 illustrates a two-dimensional slice of a 500-dimensional realization of the field.

We are concerned with analyzing properties of gradient de-

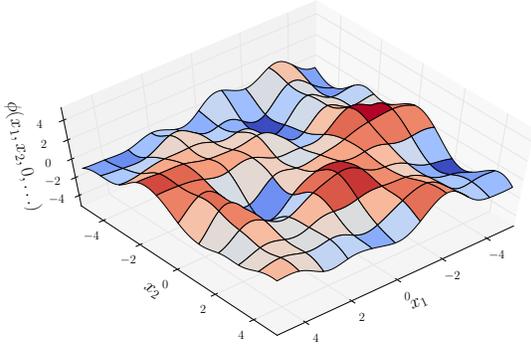


Figure 1: Two-dimensional slice of a realization in  $\mathbb{R}^{500}$  of the previously defined random field.

scent in the above random energy landscape,

$$\mathbf{x}_1 = \mathbf{x}_0 - \eta \nabla \phi(\mathbf{x}_0). \quad (1)$$

Here  $\mathbf{x}_1$  is the updated point,  $\mathbf{x}_0$  is the initial point where we start our gradient descent and  $\eta$  is the learning rate. As we will frequently use the values of the field and its gradient at both points, we introduce the short-hand notation  $\Phi_0 \equiv \phi(\mathbf{x}_0)$ ,  $\Phi_1 \equiv \phi(\mathbf{x}_1)$ ,  $\Xi_0 \equiv \nabla \phi(\mathbf{x}_0)$  and  $\Xi_1 \equiv \nabla \phi(\mathbf{x}_1)$  as illustrated in Figure 2.

After a short introduction to Gaussian Processes (GPs) in the next section, we present our main theoretical results in Section 3. Firstly, in Section 3.1 we obtain a formal expression for the distribution of  $\Phi_1$ , as well as provide analytic expressions for its expected value and variance as a function of the dimension  $N$  of our parameter space. In Section 3.2, we use the expected value of  $\Phi_1$ , to compute the optimal learning rate. When comparing the optimal learning rate with that of random search, it is seen how random search gives superior results for small dimensions while gradient descent outperforms random search in larger dimensions. In Section 3.3 we prove asymptotic normality of the rescaled random variable  $\Phi_1$ . For most practical applications, the Gaussian approximation of the distribution of  $\Phi_1$  is sufficiently close to the true distribution. In the following Section 3.4 we compare the expected value of  $\Phi_1$  with that of the global minimum in a unit ball which we estimate by means of analyzing the Euler characterise of excursion sets. The latter is found to have the same scaling with the dimension  $N$  as the expected value of  $\Phi_1$  but with a slightly larger per-factor. Besides the above theoretical results, we also provide numerical results and simulation experiments in Section 4. More precisely, in Section 4.1 we numerically simulate GRFs of dimensions up to  $N = 500$  and verify the theoretical results obtained in the previous sections. Furthermore, we also investigate gradient descent on a toy model of GFRs which models the loss functions of deep neural networks on synthetic as well as real-life datasets and compare those with the findings on

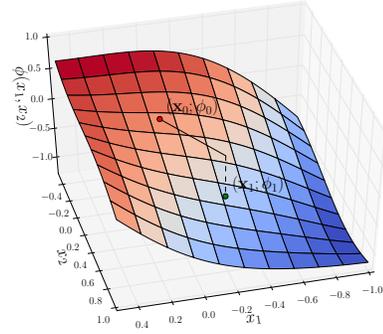


Figure 2: Gradient descent first step in a two-dimensional realization of the random field.

GRFs which are in good agreement.

## 2 Gaussian random fields and Gaussian processes

The Gaussian random field (GRF) we introduced in the previous section can be seen to be equivalent to saying that the field or loss function  $\phi$  is given by a Gaussian Process (GP)

$$\phi \sim \text{GP}(0, K) \quad (2)$$

with zero mean and kernel  $K(\mathbf{x}_0, \mathbf{x}_1) = k(|\mathbf{x}_0 - \mathbf{x}_1|)$  with  $k(r) = \exp(-r^2/2)$ . For a comprehensive review on GPs the reader is referred to [12]. The GP can be understood as an infinite dimensional extension of a multivariate Gaussian distributions such that joint distributions of any finite number of points are again a multivariate Gaussians. For our problem at hand it means that the joint distribution of  $\Phi_0 \equiv \phi(\mathbf{x}_0)$  and  $\Phi_1 \equiv \phi(\mathbf{x}_1)$ , given  $\mathbf{x}_0, \mathbf{x}_1$ , is Gaussian with mean zero and covariance  $K(\mathbf{x}_0, \mathbf{x}_1)$ . Having a kernel which depends on the distance makes closer points more correlated where in fact the correlation goes to one if the distance goes to zero. This essentially ensures that  $\phi$  will be a continuous function. This makes GPs popular choices for prior distributions over continuous functions in Bayesian statistics. One frequently occurring quantity to compute in this context is the posterior distribution of the field  $\phi$  at a new point  $\mathbf{x}_1$  given the value  $\Phi_0$  observed at  $\mathbf{x}_0$  which can be obtained by conditioning the joint distribution. The conditional distribution is indeed Gaussian,  $[\Phi_1]_{\text{cond}} = \Phi_1 | \Phi_0, \mathbf{x}_0, \mathbf{x}_1 \sim \mathcal{N}(\mu, \Sigma)$ , with conditional mean  $\mu$  and conditional covariance  $\Sigma$  given by,

$$\begin{aligned} \mu &= K(\mathbf{x}_0, \mathbf{x}_1)^T K(\mathbf{x}_0, \mathbf{x}_0)^{-1} \Phi_0 \\ \Sigma &= K(\mathbf{x}_1, \mathbf{x}_1) - K(\mathbf{x}_0, \mathbf{x}_1)^T K(\mathbf{x}_0, \mathbf{x}_0)^{-1} K(\mathbf{x}_0, \mathbf{x}_1) \end{aligned} \quad (3)$$

This is a well-known result for GPs and a similar result holds true when conditioning on more than one point. The implications of this relation are important since it essen-

tially means that we can efficiently compute posterior updates of probability distributions at the expense of a few matrix operations.

### 3 Theoretical results

#### 3.1 Distribution of the field after one step of gradient descent

In the previous section we saw that the joint distribution of  $(\Phi_1, \Phi_0)^T$  is a Gaussian and that the conditional distribution can be easily obtained. Moreover, the joint distribution of the  $2N+2$  dimensional vector  $(\Phi_1, \Xi_1, \Phi_0, \Xi_0)^T$  is also a multivariate Gaussian with mean  $\mathbf{0}$  and a covariance matrix  $\Sigma$  which can easily be expressed in terms of the kernel function  $k$  and its derivative,

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \quad (4)$$

$$\Sigma_{11} = \Sigma_{22} = \mathbb{1} \quad (5)$$

$$\begin{aligned} \Sigma_{12} &= \Sigma_{21}^T \\ &= e^{-\Delta \mathbf{x}^2/2} \left( \mathbb{1} - \begin{bmatrix} 0 & \Delta \mathbf{x}^T \\ -\Delta \mathbf{x} & \Delta \mathbf{x} \Delta \mathbf{x}^T \end{bmatrix} \right) \end{aligned} \quad (6)$$

where  $\mathbb{1} \equiv \mathbb{1}_{(N+1) \times (N+1)}$  and  $\Delta \mathbf{x} \equiv \mathbf{x}_1 - \mathbf{x}_0$ .

The values of the field and its gradient at  $\mathbf{x}_1$ ,  $\Phi_1$  and  $\Xi_1$ , are going to be Gaussian random variables conditioned on the values at  $\mathbf{x}_0$ ,  $\Phi_0$  and  $\Xi_0$ , similar to the example presented in the previous section,

$$\begin{bmatrix} \Phi_1 \\ \Xi_1 \end{bmatrix}_{cond} \sim \mathcal{N} \left( \Sigma_{12} \Sigma_{22}^{-1} \begin{bmatrix} \Phi_0 \\ \Xi_0 \end{bmatrix}, \mathbb{1} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \right). \quad (7)$$

When multiplying out the terms and replacing  $\Delta \mathbf{x} = -\eta \Xi_0$  from (1), we obtain that  $\Phi_1$  follows a conditional normal distribution with mean and variance,

$$\begin{aligned} m_1(\varphi_0, \xi_0^2) &:= \mathbb{E} [\Phi_1 | \Phi_0 = \varphi_0, \Xi_0^2 = \xi_0^2] \\ &= e^{-\frac{\eta^2}{2} \xi_0^2} (\varphi_0 - \eta \xi_0^2) \end{aligned} \quad (8)$$

$$\begin{aligned} v_1(\xi_0) &:= \text{Var} [\Phi_1 | \Phi_0 = \varphi_0, \Xi_0^2 = \xi_0^2] \\ &= 1 - e^{-\eta^2 \xi_0^2} (1 + \eta^2 \xi_0^2). \end{aligned} \quad (9)$$

As  $\Phi_0$  will have a  $\mathcal{N}(0, 1)$  distribution and  $\Xi_0^2$  will have a  $\chi^2$ -distribution with  $N$  degrees of freedom, being the sum of the squares of  $N$  independent normally distributed components, we can write the overall distribution as

$$\begin{aligned} f_{\Phi_1}(\varphi_1) &= \int_{-\infty}^{+\infty} d\varphi_0 \int_0^{+\infty} d\xi_0^2 \frac{f_{\Phi_0}(\varphi_0) f_{\Xi_0^2}(\xi_0^2)}{(2\pi v_1(\xi_0^2))^{1/2}} \times \\ &\times \exp \left( -\frac{(\varphi_1 - m_1(\varphi_0, \xi_0^2))^2}{2 v_1(\xi_0^2)} \right), \end{aligned} \quad (10)$$

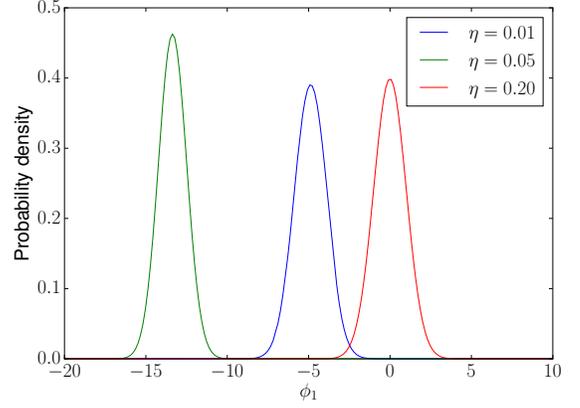


Figure 3: Numerically integrated probability density functions of  $\Phi_1$  plotted for some values of  $\eta$  with  $N = 500$ .

where  $f_{\Phi_0}(\varphi_0)$  is the standard normal probability density function (PDF),  $f_{\Xi_0^2}(\xi_0^2)$  is the PDF for a chi-squared random variable with  $N$  degrees of freedom and  $m_1, v_1$  are the conditional mean and variance of  $\Phi_1$  as given in (8)-(9). For illustration purposes, we plotted the resulting PDF by means of numerical integration as shown in Figure 3.

A closed form expression for the distribution of  $\Phi_1$  seems out of reach, however, we can calculate its moments as well as show asymptotic normality later. For now, we focus on the first moments which can be easily derived,

$$\begin{aligned} \mathbb{E}[\Phi_1] &= \mathbb{E}_{\varphi_0, \xi_0^2} [m_1(\varphi_0, \xi_0^2)] = -\eta \mathbb{E}_{\xi_0^2} \left[ \xi_0^2 e^{-\frac{\eta^2}{2} \xi_0^2} \right] \\ &= -N\eta (\eta^2 + 1)^{-N/2-1}, \end{aligned} \quad (11)$$

where the last expectation value is computed by integrating over the PDF of the  $\chi^2$ -distribution. Similarly, for the variance one obtains

$$\begin{aligned} \text{Var}(\Phi_1) &= \mathbb{E}_{\xi_0^2} [v_1(\xi_0)] \\ &= 1 - \mathbb{E}_{\xi_0^2} \left[ e^{-\eta^2 \xi_0^2} (1 + \eta^2 \xi_0^2) \right] \\ &= 1 + N\eta^2 (1 + 2\eta^2)^{-\frac{N}{2}-2} (N + 1 - 2\eta^2) \\ &\quad - N^2 \eta^2 (\eta^2 + 1)^{-N-2}. \end{aligned} \quad (12)$$

It can be observed that the mean and variance of  $\Phi_1$  converge to those of  $\Phi_0$  in the limit where  $\eta \rightarrow 0$ , as we stay in the same point and also when  $\eta \rightarrow +\infty$ , as the gradient only gives significant information in a neighborhood of  $\mathbf{x}_0$  of size  $\mathcal{O}(1)$ .

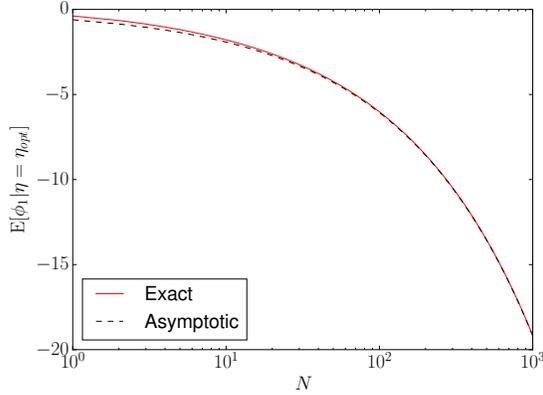


Figure 4: Expected value for  $\Phi_1$  as a function of  $N$  when using the optimal learning rate.

### 3.2 Optimal learning rate and comparison to random search

The optimal learning rate  $\eta_{opt}$  can be easily derived by computing  $\frac{d}{d\eta} \mathbb{E}[\Phi_1] |_{\eta=\eta_{opt}} = 0$ , yielding

$$\eta_{opt} = (N + 1)^{-\frac{1}{2}} \quad (13)$$

Computing the second derivative shows that indeed it is a minimum. We can now obtain the expected value of the field for the optimal learning rate which is given by

$$\begin{aligned} \mathbb{E}[\Phi_1 | \eta = \eta_{opt}] &= -N(N + 1)^{-\frac{1}{2}} \left( \frac{N + 2}{N + 1} \right)^{-\frac{N}{2} - 1} \\ &= -\sqrt{\frac{N}{e}} + \mathcal{O}\left(N^{-\frac{1}{2}}\right). \end{aligned} \quad (14)$$

The expected value of  $\Phi_1$  will improve as the square root of the number of dimensions  $N$ , as shown in Figure 4, and, as we show in Section 3.4, this is within a constant factor of the minimum field value within a unit radius ball. We would not expect significantly better results, as the information provided by the gradient decays very fast for distances greater than the correlation length (1 in our case).

The expected step length will also tend to 1 for large values of  $N$  when using the optimal learning rate. That can be seen by using the previously discussed fact that the squared gradient has a  $\chi^2$ -distribution with  $N$  degrees of freedom and computing the expectation:

$$\begin{aligned} \mathbb{E}[\eta_{opt} \Xi_0] &= \frac{1}{\sqrt{N + 1}} \sqrt{2} \frac{\Gamma(N/2 + 1/2)}{\Gamma(N/2)} \\ &= 1 + \mathcal{O}(N^{-1}). \end{aligned} \quad (15)$$

Taking  $N = 500$ , a single step of gradient descent with the

optimal learning rate gives us an expected value of

$$\mathbb{E}[\Phi_1 | \eta = \eta_{opt}, N = 500] \approx -\sqrt{\frac{500}{e}} \approx -13.56. \quad (16)$$

To put this value into context we compare it to random search. Since any evaluation of the random field would give a value smaller than the above with probability  $F_{\mathcal{N}}(-13.56) \approx 3.46 \cdot 10^{-42}$ , where  $F_{\mathcal{N}}$  is the cumulative distribution function of a standard normal, more than  $10^{41}$  tries would be needed on average to get to a value smaller than that from random search. On the other hand, for  $N = 1$  the expected value after a gradient descent step would only be

$$\mathbb{E}[\Phi_1 | \eta = \eta_{opt}, N = 1] = -\frac{2}{3\sqrt{3}} \approx -0.385. \quad (17)$$

This value or better would be obtained by random search in an average of  $F_{\mathcal{N}}(-0.385)^{-1} \approx 2.85$  tries. The difference exemplifies how gradient descent becomes increasingly powerful when moving to higher dimensional optimization. Below, in Section 3.4, we further investigate how this compares to the value of the global minimum.

### 3.3 Asymptotic normality of the distribution of the field

We now analyze convergence of the distribution of  $\Phi_1$  for  $N \rightarrow \infty$  when we scale the learning rate around its optimal value, namely, under the scaling

$$\eta = \frac{X}{\sqrt{N}}, \quad (18)$$

where  $X$  is the rescaled learning rate. Under this scaling we see that the expected value of  $\Phi_1$

$$\mathbb{E}[\Phi_1] = \mu_N(X) + \dots, \quad \mu_N(X) := -\sqrt{N} X e^{-\frac{X^2}{2}} \quad (19)$$

is of  $\mathcal{O}(N^{1/2})$  while the variance

$$\text{Var}(\Phi_1) = \sigma^2(X) + \dots, \quad \sigma^2(X) := 1 + X^2 e^{-X^2} \quad (20)$$

remains finite. We will now show that

**Theorem 1.** *As the dimension  $N \rightarrow \infty$ , the rescaled field value after a single step of gradient descent converges asymptotically to a normal distribution:*

$$\Phi_1 - \mu_N(X) \xrightarrow{d} \mathcal{N}(0, \sigma^2(X)) \quad (21)$$

where  $X$ ,  $\mu_N(X)$ ,  $\sigma^2(X)$ , are defined in (18) - (20).

*Proof.* To prove the theorem we first compute the moment generating function

$$\begin{aligned} \mathbb{E}[e^{t\Phi_1}] &= \mathbb{E}_{\varphi_0, \xi_0^2} \left[ \exp \left\{ t m_1(\varphi_0, \xi_0^2) + \frac{t^2}{2} v_1(\xi_0) \right\} \right] \\ &= e^{\frac{t^2}{2}} \mathbb{E}_{\xi_0^2} \left[ \exp \left\{ -\frac{t^2}{2} \eta^2 \xi_0^2 e^{-\eta^2 \xi_0^2} + \right. \right. \\ &\quad \left. \left. - t \eta \xi_0^2 e^{-\eta^2 \xi_0^2 / 2} \right\} \right] \end{aligned} \quad (22)$$

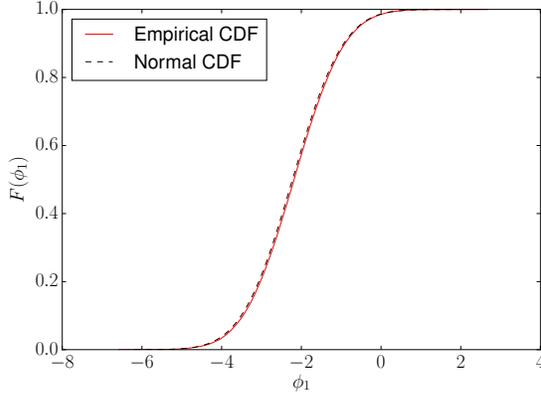


Figure 5: Comparison of the empirical cumulative distribution function of simulated values of  $\phi_1$  with the normal cumulative distribution function with corresponding mean and standard deviation ( $\eta = 0.1 \eta_{opt}$ ).

Inserting the scaling relation for the learning rate, (18), and using a saddle point expansion of the integral over  $\xi_0^2$  when writing out the expectation one can see that the above expression is given to leading order by

$$\begin{aligned} \mathbb{E} [e^{t\Phi_1}] &= e^{\frac{t^2}{2}} \exp \left\{ -\frac{t^2}{2} X^2 e^{-X^2} + \right. \\ &\quad \left. - tX\sqrt{N}e^{-X^2/2} \right\} + \dots \\ &= \exp \left\{ \frac{\sigma^2(X)t^2}{2} + t\mu_N(X) \right\} + \dots \end{aligned} \quad (23)$$

where the expectation over  $\xi_0^2$  collapsed to its saddle point which to leading order is given by  $\mathbb{E}\xi_0^2 = N$ . Thus

$$\lim_{N \rightarrow \infty} \mathbb{E} \left[ e^{t(\Phi_1 - \mu_N(X))} \right] = e^{\frac{\sigma^2(X)t^2}{2}} \quad (24)$$

which proves the above convergence.  $\square$

Figure 5 shows an example of the normal approximation of the distribution of  $\Phi_1$  for finite  $N$ . Details of the numerical analysis will be presented in Section 4.

### 3.4 Comparison with optimal values

In the previous sections we have analyzed the expected value of the field after a step of gradient descent. A natural follow-up is to ask how does it compare with the global extremum of the field. As we are working with a Gaussian random field with a covariance that decays to zero, we can expect to find values with arbitrarily large magnitude at enough distance, but a more useful comparison can be done by restricting ourselves to a unit ball  $\mathcal{B}_N(\mathbf{x}_0)$  around the random starting point  $\mathbf{x}_0$ .

There is no known analytical expression for the expected value of a Gaussian random field maximum or minimum in

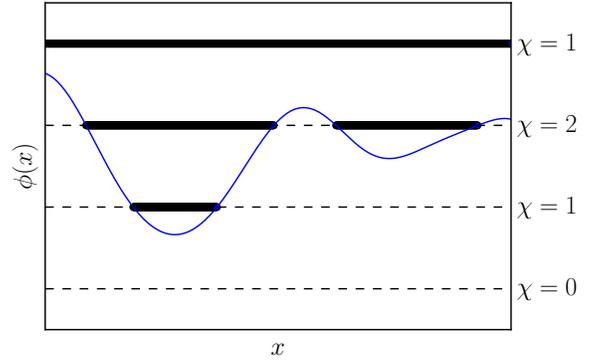


Figure 6: Excursion sets  $A_u$  and their Euler characteristics for different values of  $u$ .

any multidimensional domain, but a number of powerful estimation techniques have been developed [11, 15, 16]. Most of them are based on finding quantities that are related to the extrema and computing their expectations. In our case we will use the exact computation of the expected Euler characteristic, as described in [11], to get an estimate for the number of connected components of an excursion set and use that estimate to get the expected value of the minimum.

In general, the Euler characteristic can be seen as the unique functional  $\chi$  from a family of subsets  $\mathcal{A}$  of a manifold  $\mathcal{M}$  to the integers that has the following properties:

- $\chi(\emptyset) = 0$ .
- $\chi(X) = 1$  if  $X$  is contractible.
- $\chi(X \cup Y) = \chi(X) + \chi(Y)$  if  $X \cap Y = \emptyset$ .

In the following discussion we will only use these properties of the Euler characteristic and assume that all the sets under consideration are included in  $\mathcal{A}$ . A detailed proof of the uniqueness of the Euler characteristic and a description of the family  $\mathcal{A}$  in the context of random fields can be found in [11].

In our case the manifold is the unit ball,  $\mathcal{M} = \mathcal{B}_N(\mathbf{x}_0)$ . Now we can define an excursion set  $A_u$  in  $\mathcal{B}_N(\mathbf{x}_0)$  as the subset of  $\mathcal{B}_N(\mathbf{x}_0)$  composed by the points where the field  $\varphi$  reaches a value of  $u$  or smaller. It is clear then that  $A_u$  will be non-empty if and only if  $u > \min_{\mathbf{x} \in \mathcal{B}_N(\mathbf{x}_0)} \phi(\mathbf{x})$ , allowing us to connect geometrical properties of the excursion set  $A_u$  with the value of the minimum.

As our random field is continuous, if we start  $u$  from a large positive value and gradually decrease it, we expect the excursion set  $A_u$  to start being all of  $\mathcal{B}_N(\mathbf{x}_0)$ , then getting some holes, being disconnected, turning into a few contractible components and finally ending as the empty set, as shown

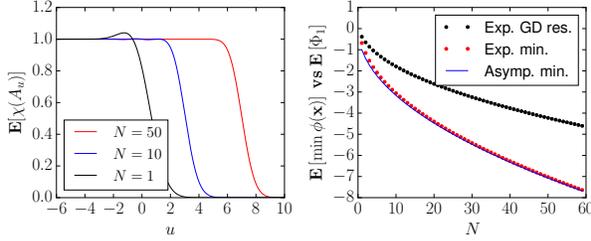


Figure 7: The left figure shows the shape of the expected Euler characteristic as a function of the threshold level  $u$  for multiple values of  $N$ . Those Euler characteristic values are used to estimate the expected values of the minimum (Exp. min.) and its asymptotics (Asymp. min.), shown in the right figure compared against the results of a gradient descent step (Exp. GD res.), both as a function of  $N$ .

in Figure 6. If our excursion set has the form of disjoint contractible components, its Euler characteristic gives us the number of components and we will be able to use its expected value for different values of  $u$  to estimate the value of the minimum.

Following [11], we note that it has been proved that the expected Euler characteristic  $\chi$  of the excursion set  $A_u$  under the conditions we described is given by

$$\mathbb{E}[\chi(A_u)] = \sum_{j=0}^N \mathcal{L}_j(\mathcal{B}_N(\mathbf{x}_0)) \rho_j(u), \quad (25)$$

where  $\mathcal{L}_j(\mathcal{B}_N(\mathbf{x}_0))$  are the Lipschitz-Killing curvatures of  $\mathcal{B}_N(\mathbf{x}_0)$  which are given by

$$\mathcal{L}_j(\mathcal{B}_N(\mathbf{x}_0)) = \binom{N}{j} \frac{\omega_N}{\omega_{N-j}}, \quad \omega_j = \frac{\pi^{j/2}}{\Gamma(j/2 + 1)} \quad (26)$$

and  $\rho_j(u)$  is given by

$$\rho_j(u) = (2\pi)^{-(j+1)/2} H_{j-1}(u) e^{-\frac{u^2}{2}}, \quad (27)$$

with  $H_j$  being the Hermite polynomials.

As shown in Figure 6, we anticipate the expected Euler characteristic starting at 1 for large positive values of  $u$  (the whole set has characteristic 1), to be 1 for values of  $u$  that leave a single non-empty excursion set with high probability (a small ‘‘droplet’’ has characteristic 1) and to decrease to 0 when it starts being less probable to find a non-empty excursion set (in other words, when  $u$  is below the minimum).

The expected behavior can be seen in Figure 7, with the threshold increasing in absolute value as we move to higher dimensional spaces. If we estimate the expected minimum as the value where the expected Euler characteristic is 0.5, we find its values are well approximated by  $-\sqrt{N}$ , i.e.

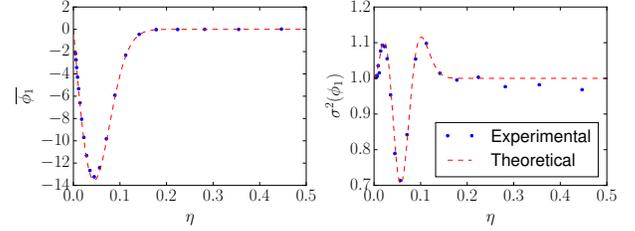


Figure 8: Comparison of sample mean and variance with the theoretically expected values for multiple values of the learning rate  $\eta$  and  $N = 500$ .

$$\mathbb{E} \left[ \min_{\mathbf{x} \in \mathcal{B}(\mathbf{x}_0)} \phi(\mathbf{x}) \right] \approx -\sqrt{N}. \quad (28)$$

Comparing those values with the values of the field after a gradient descent step, as found in Section 3.2, we see they differ by a constant factor of  $\sqrt{e}$ .

Applying this bound to get the expected value of the minima inside the unit ball for  $N = 500$ , we find it will be  $\mathbb{E}[\min \phi] \approx -\sqrt{500} \approx -22.36$ , which should be contrasted with the expected field value after one step of gradient descent,  $\mathbb{E}[\Phi_1] \approx -13.56$ , as obtained in Section 3.2.

## 4 Experimental results

### 4.1 Random field simulation

Our experiments will require generating approximate instances of Gaussian random fields in spaces of high dimensionality, with values of  $N$  reaching 500. Most of the conventional methods for simulating random fields [17, 18] don’t scale well to a large number of dimensions, as they generate explicit grids representing the field values.

We can take the spectral representation [11] of the random field,

$$\phi(\mathbf{x}) = \int e^{i\mathbf{z}^T \mathbf{x}} W(d\mathbf{z}), \quad (29)$$

that can be seen as the Fourier transform of the field expressed as a stochastic integral, and use Monte Carlo sampling to approximate the integration over  $\mathbf{z}$ . In that way, we obtain an approximate instance of the random field expressed as the real part of the sum of  $M$  complex exponentials

$$\phi_{sim}(\mathbf{x}) = \Re \mathbf{w}^T \overline{\exp}(i Z \mathbf{x}), \quad (30)$$

where  $\overline{\exp}$  is component-wise exponentiation,  $\mathbf{w} \sim \mathcal{CN}(\mathbf{0}_M, M^{-1} \mathbf{1}_M)$  is a complex Gaussian random vector and  $Z \in \mathbb{R}^{M \times N}$  is a real Gaussian random matrix with

independently distributed elements  $Z_{mn} \sim \mathcal{N}(0, 1)$ . This can be considered a multidimensional variant of the randomization method described in [19], although in a high dimensional context it is important to ensure that the number of samples  $M$  is significantly higher than the number of dimensions  $N$  to avoid confining the gradient to a low dimensional subspace.

The gradient can then be computed by differentiating the previous expression:

$$\nabla \phi_{sim}(\mathbf{x}) = -\mathfrak{S} Z^T \text{diag}(\mathbf{w}) \overline{\exp}(i Z \mathbf{x}). \quad (31)$$

The value of  $M$ , being the number of samples, will determine how accurate our representation of the random field will be. Higher values will increase the amount of computational resources required and lower values will produce a lower quality realization of the random field. A value of  $2 \cdot 10^4$  was found to give high quality results for  $N \leq 500$  at acceptable computational cost.

We first start by comparing the expected values for the sample mean and variance computed in section 3.1 with experimental results. With the previously discussed representations and for 20 different values of the learning rate  $\eta$ , we can do one step of gradient descent for  $10^4$  different starting points distributed uniformly in  $[-10^6, 10^6]^{500}$ . The resulting sample means and variances are shown in Figure 8 compared with the theoretical expectations and we can see they match them quite accurately.

To compare the expected distribution with the empirical one, we repeated the gradient descent step simulation using  $10^5$  points and  $\eta = 0.1 \eta_{opt}$ . The simulated results can be seen in Figure 5 and they also fit very closely with the expected distribution.

## 4.2 Experiments on synthetic and real datasets

In this section we show how this random field model can be used to classify real data. To do that, we introduce a toy model in which we take a standard multilayer network based binary classifier and we replace the entire network by a ‘‘black box’’ loss function, given by a static random field, with no adjustable internal parameters. The  $N_P$ -dimensional parameter vector  $\beta$ , replacing the weights of a normal network, and the  $i$ -th  $N_I$ -dimensional input to be classified  $\mathbf{x}_{input}^i$  are concatenated to get a  $N$ -dimensional vector, with  $N = N_P + N_I$ ,

$$\mathbf{x}^i = \begin{bmatrix} \beta \\ \mathbf{x}_{input}^i \end{bmatrix}, \quad (32)$$

that is the random field input.

It is a normal practice [5] in classifier networks to use softmax as the activation function in the last layer and cross-entropy as the loss function. As we are replacing the rest of

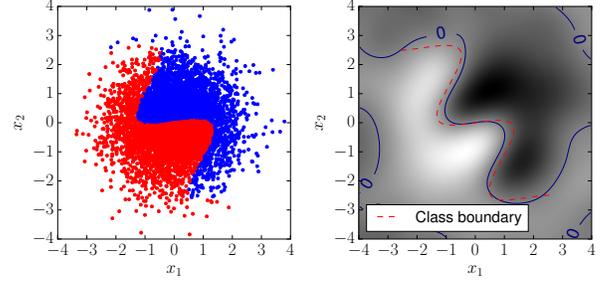


Figure 9: Original dataset and associated values of the random field  $\phi$  for the trained parameters, showing how the classification works.

the network with a random field and using only two classes, the output of the classifier can be written as

$$y^i = \text{sigmoid}(\phi(\mathbf{x}^i)), \quad (33)$$

where  $\mathbf{x}^i$  is the input vector associated with the  $i$ -th input instance and  $\text{sigmoid}(z) = 1/(1 + \exp(-z))$  is the sigmoid function.

As usual in supervised binary classification problems, we associate a true class label  $y_{true}^i \in \{0, 1\}$  to each of our input instances  $\mathbf{x}_{input}^i$  and we try to minimize the cross entropy loss between the true labels  $y_{true}^i$  and the classifier outputs  $y^i$ , i.e.  $L_i = y_{true}^i \log y^i - (1 - y_{true}^i) \log(1 - y^i)$ .

The training process is standard minibatch gradient descent. By analogy with neural networks, where only the weights are updated, only the parameter vector  $\beta$  is updated after each minibatch. Following usual practice, we divide the input data into training and test sets, using the training set to select a value for the parameter vector  $\beta$  and the test set to evaluate the accuracy of the classifier. Furthermore, both sets of input instances are normalized to mean 0 and mean norm 1, matching the scale of the data set distribution to the correlation scale of the random field. This is empirically observed to make a significant difference in classification accuracy.

One way of visualizing the training process is to think of the parameter vector  $\beta$  as selecting a random field slice. Then the gradient descent over the loss function will try to select a slice where the naive Bayes decision surface divides both classes, putting the instances where  $y_{true}^i = 1$  in the positive side and the instances with  $y_{true}^i = 0$  in the negative side of the surface. This process can be seen clearly in Figure 9, where the parameter vector after training can be seen as selecting a 2D slice of the random field  $\phi$ , where the intersection of the naive Bayes decision boundary with the slice is close to the class boundary.

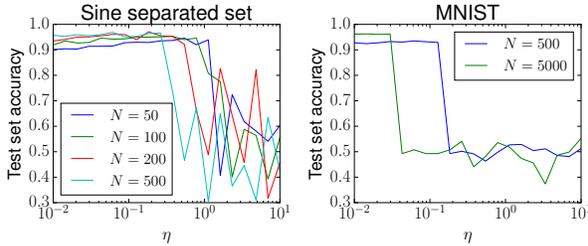


Figure 10: Sine separated and MNIST test set accuracy as a function of the learning rate  $\eta$  and the dimensionality  $N$ .

Note that the here proposed toy model of a “black box” random field that mimics the energy landscape of a deep neural network is different to GP classification [20, 12] or other kernel methods such as support vector machines [21, 22] in the sense that we combine the input and parameter vector into a joint input vector to a GRF which is kept fixed during the learning process.

To test the classification power of these random field instances, we run the training and test process over two simple data sets:

- Normally distributed points in the  $\mathbb{R}^2$  plane separated by a sine function (see Figure 9), with 6000 elements in the training set and 1000 in the test set ( $N_T = 2$ , as we are talking about points in the plane).
- MNIST [23, 24], modified to classify the digits as even or odd and using 60000 elements in the training set and 10000 in the test set ( $N_T = 784$  in this case).

The training is done using a fixed batch size of 128 and 10 epochs, combined with different learning rates and values of  $N$  to evaluate their impact over test set accuracy.

The test set accuracy for MNIST with  $N = 5000$  is over 96% showing that, even though it is far from matching the state of the art, the model has significant classification power. When compared with models with a similar number of parameters, the model is competitive [24].

We can observe in Figure 10 that accuracy doesn’t depend on the learning rate until reaching a critical value and then it drops to random performance. The MNIST drop for  $N = 500$  is at  $\eta \approx 0.13$  and for  $N = 5000$  at  $\eta \approx 0.03$ ,

$$0.03 \cdot \sqrt{5000} \approx 2.12 \approx 2.68 \approx 0.13 \cdot \sqrt{500}, \quad (34)$$

roughly matching the scaling found before in the single step regime.

## 5 Conclusion

The successes of deep learning as well as some unexpected weaknesses, such as the difficulty of combining good generalization and resistance to adversarial examples, have led to a significant research effort aiming to understand why their training process performs so well in high dimensional problems. The complex structure of deep neural networks error landscapes makes it difficult to understand how the optimization process is working, but observing its performance in a simple random field model can help to clarify some of the reasons behind its successes and limitations.

In this work we aim to get a better understanding of gradient descent as a tool for high dimensional optimization by obtaining theoretical and empirical results about its performance over Gaussian random fields. Following a brief introduction, we establish some asymptotic results about the distribution of field values reached after a single step of gradient descent. Those results are then compared with a theoretical estimate of the extreme field values at a similar distance. Finally, we compare the previously obtained theoretical results with experimental simulations, while also showing that the “black box” Gaussian random field model is capable of solving realistic classification tasks.

We show our theoretical results about the distribution of values after a gradient descent step in Section 3. Starting in Section 3.1, we obtain the first and second moments as a function of the learning rate  $\eta$  and the number of dimensions  $N$  of the parameter space. In the following Section 3.2 we use the previously derived expressions to get the optimal value for the learning rate as a function of the number of dimensions, finding that  $\eta_{\text{opt}}(N) \approx N^{-1/2}$  for large values of  $N$ . Using that optimal learning rate we show that the expected value of the field after a gradient descent step is approximately  $\mathbb{E}[\Phi_1] \approx -(N/e)^{1/2}$ , comparing very favorably with the values that can be obtained through a random search when  $N \gg 1$ , as those are independent of the dimensionality of the space. Closing our analysis of the distribution of values, we prove in Section 3.3 that in the high dimensional limit the distribution of the values after the gradient descent step is approximately normal, with a variance that is independent of the dimensionality and a mean that is proportional to  $-N^{1/2}$ . Finally, in Section 3.4 we show using the expected Euler characteristic of excursion sets that the expected minimum inside the unit ball will only differ from the expected value we obtain through one step of gradient descent with the optimal learning rate by a factor of  $\sqrt{e}$  in the  $N \gg 1$  limit.

In Section 4 we start by showing how we simulate a high-dimensional random field and comparing the experimental gradient descent results with the previous theoretical results in Section 4.1, finding them to be in good agreement. Finally, we show that the model we obtain by replacing a neural network by a Gaussian random field can be trained by

gradient descent, obtaining competitive results in a simple synthetic dataset and in MNIST, once we take into account that the model is only using 5000 parameters.

The introduced “black box” GRF model is successful at combining nontrivial classification performance in realistic datasets with being simple enough to be susceptible to exact theoretical analysis. A possible line of future investigation would be to look at other aspect of deep neural networks through the lens of our toy model such as the interpretability of hidden layer neurons in image classification tasks or transfer learning. That could be combined with extending these results to the normal multistep minibatch training process.

## References

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [2] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- [3] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [4] Satinder Singh, Andy Okun, and Andrew Jackson. Artificial intelligence: Learning to play go from scratch. *Nature*, 550(7676):550336a, 2017.
- [5] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [6] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. <http://arxiv.org/abs/1412.6572>, 2014.
- [7] Justin Gilmer, Luke Metz, Fartash Faghri, Samuel S Schoenholz, Maithra Raghu, Martin Wattenberg, and Ian Goodfellow. Adversarial spheres. <http://arxiv.org/abs/1801.02774>, 2018.
- [8] Anna Choromanska, Mikael Henaff, Michael Mathieu, Gérard Ben Arous, and Yann LeCun. The loss surfaces of multilayer networks. In *Artificial Intelligence and Statistics*, pages 192–204, 2015.
- [9] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep neural networks as Gaussian processes. <http://arxiv.org/abs/1711.00165>, 2017.
- [10] Samuel S Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep information propagation. <http://arxiv.org/abs/1611.01232>, 2016.
- [11] Robert J Adler and Jonathan E Taylor. *Random fields and geometry*. Springer Science & Business Media, 2009.
- [12] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press, 2005.
- [13] Justin D. Bewsher, Alessandra Tosi, Michael A. Osborne, and Stephen J. Roberts. Distribution of Gaussian process arc lengths. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017.
- [14] Nikitas Rontsis, Michael A. Osborne, and Paul J. Goulart. Distributionally ambiguous optimization techniques for batch Bayesian optimization. <http://arxiv.org/abs/1707.04191>, 2017.
- [15] David Aldous. *Probability approximations via the Poisson clumping heuristic*, volume 77. Springer Science & Business Media, 2013.
- [16] Jean-Marc Azaïs and Mario Wschebor. *Level sets and extrema of random processes and fields*. John Wiley & Sons, 2009.
- [17] Edmund Bertschinger. Multiscale Gaussian random fields and their application to cosmological simulations. *The Astrophysical Journal Supplement Series*, 137(1):1, 2001.
- [18] Annika Lang and Jürgen Potthoff. Fast simulation of Gaussian random fields. *Monte Carlo Methods and Applications*, 17(3):195–214, 2011.
- [19] Peter R Kramer, Orazgeldi Kurbanmuradov, and Karl Sabelfeld. Comparative analysis of multiscale Gaussian random field simulation algorithms. *Journal of Computational Physics*, 226(1):897–924, 2007.
- [20] C. K. I. Williams and D. Barber. Bayesian classification with Gaussian processes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(12):1342, 1998.
- [21] Bernhard E. Boser, Isabelle M. Guyon, and Vladimir N. Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory – COLT ’92*, page 144, 1992.

- [22] Bernhard Schölkopf and Alexander J. Smola. *Learning with Kernels*. MIT Press, 2002.
- [23] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278, 1998.
- [24] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. The MNIST database. <http://yann.lecun.com/exdb/mnist/>.