DNN-based Topology Optimisation: Spatial Invariance and Neural Tangent Kernel

Benjamin Dupuis

Chair of Statistical Field Theory
Ecole Polytechnique Fédérale de Lausanne
Lausanne, Switzerland
benjamin.dupuis@epfl.ch

Arthur Jacot

Chair of Statistical Field Theory
Ecole Polytechnique Fédérale de Lausanne
Lausanne, Switzerland
arthur.jacot@epfl.ch

Abstract

We study the Solid Isotropic Material Penalisation (SIMP) method with a density field generated by a fully-connected neural network, taking the coordinates as inputs. In the large width limit, we show that the use of DNNs leads to a filtering effect similar to traditional filtering techniques for SIMP, with a filter described by the Neural Tangent Kernel (NTK). This filter is however not invariant under translation, leading to visual artifacts and non-optimal shapes. We propose two embeddings of the input coordinates, which lead to (approximate) spatial invariance of the NTK and of the filter. We empirically confirm our theoretical observations and study how the filter size is affected by the architecture of the network. Our solution can easily be applied to any other coordinates-based generation method.

1 Introduction

Topology optimisation [4], also known as structural optimisation, is a method to find optimal shapes subject to some constraints. It has been widely studied in the field of computational mechanics. Here we are interested in the particular case of the Solid Isotropic Material Penalisation (SIMP) method [19, 1], which is a very common method in this field.

Recently some authors have used Deep Neural Networks (DNNs) to perform topology optimisation. We can differentiate two different approaches in the use of DNNs with SIMP. The first approach consists in generating with the classical algorithms a dataset of optimised shapes and train a DNN on this dataset to produce new optimal shapes [3, 32]. Variations of this approach use Generative Adversarial Networks (GAN) [22, 29] to effectively reproduce classical topology optimisation.

In the second approach, the density is generated pointwise by a DNN, which is trained with gradient descent to optimise the density field with respect to the physical constraints, as proposed in [12] to use the power of deep models without giving up exact physics. We focus on the approach of [7, 6] where the density field is generated by a Fully-Connected Neural Network (FCNN) taking the coordinates of a grid as inputs. Surprisingly, [7] observes that the DNN-generated density fields do not feature checkerboard artifacts, which are common in vanilla SIMP. A traditional method to avoid checkerboard patterns is to add a filter [30, 5], but it is not needed for DNN-generated density fields.

In this paper, we analyse theoretically how the use of a DNN to generate the density field affects the learning. Our main theoretical tool is the Neural Tangent kernel (NTK) introduced in [14] to describe the dynamics of wide neural networks [14, 2, 18, 13].

While this paper focuses on linear elasticity and SIMP, our analysis can extended to other physical problems such as heat transfer [20], or any model where an image is generated by a DNN taking the pixel coordinates as inputs (like in [21]).

1.1 Our contribution

In this paper we study topology optimisation with neural networks. The physical density is represented by a neural network taking an embedding of spatial coordinates as inputs, i.e. the density at a point $x \in \mathbb{R}^d$ is given by $f_{\theta}(\varphi(x))$ for θ the parameters of the network and φ an embedding. We use theoretical tools, in particular the Neural Tangent Kernel (NTK), to understand how the architecture and hyperparameters of the network affect the optimisation of the density field:

- We show that in the infinite width limit (when the number of neurons in the hidden layers grows to infinity), topology optimisation with a DNN is equivalent to topology optimisation with a density filter equal to the "square root" of the NTK. Filtering is a commonly used technique in topology optimisation, aimed to remove checkerboard patterns.
- In topology optimisation as in other physical optimisation problems, it is crucial to guarantee some spatial invariance properties. If the coordinates are taken as inputs of the network directly, the NTK (and the corresponding filter) is not translation invariant, leading to non-optimal shapes and visual artifacts. We present two methods to ensure the spatial invariance of the NTK: embedding the coordinates on the (hyper-)torus or using a random Fourier features embedding (similar to [33]).
- In traditional topology optimisation, the filter size must be tuned carefully. When optimising with a DNN, the filter size depends on the embedding of the coordinates and the architecture of the network. We define a filter radius for the NTK, which plays a similar role as the classical filter size and discuss how it is affected by the choice of embedding, activation function, depth and other hyperparameters like the importance of bias in the network. This tradeoff can also be analysed in terms of the spectrum of the NTK, explaining why neural networks naturally avoid checkerboard patterns.

We confirm and illustrate these theoretical observations with numerical experiments. Our implementation of the algorithm will be made public at https://github.com/benjiDupuis/DeepTopo.

2 Presentation of the method

In this paper, we use a DNN to generate the density field used by the Solid Isotropic Material Penalisation (SIMP) method. Our implementation of SIMP is based on [1] and [19]. In this section we introduce the traditional SIMP method and our neural network setting.

2.1 SIMP method

We consider a regular grid of N elements where the density of element i is denoted $y_i \in [0,1]$, informally the value y_i represents the presence of material at a point i. Our goal is to optimise over the density $y \in \mathbb{R}^N$ to obtain a shape that can withstand forces applied at certain points, represented by a vector F.

The method uses finite element analysis to define a stiffness matrix $K(y) \in S_N^{++}(\mathbb{R})$ from the density y and computes the displacement vector U(y) (which represent the deformation of the shape at all points i as a result of the applied forces F) by solving a linear system K(y)U(y) = F. In our implementation, we performed it either by using sparse Cholesky factorisation [10, 8] or BICGSTAB method [34] (this last one can be used for a high number of pixels).

The loss function is then defined as the compliance $C(y) = U(y)^T K(y) U(y)$, under a volume constraint of the form $\sum_{i=1}^N y_i = V_0$, with $0 \le V_0 \le N$ (see [1, 19]).

2.2 A modified SIMP approach

Several methods exist to optimise the density field $y \in \mathbb{R}^N$, such as gradient descent or the so-called Optimality Criteria (OC) [36]. We propose here an optimisation method inspired from [12] which we will refer as the Modified Filtering method (MF). The advantage of this method is that it can be used with or without DNNs, hence allowing comparison between these two approaches. We first present here the model without DNNs.

In our method, the densities y_i^{MF} are given by:

$$\forall i \in \{1, ..., N\}, \ y_i^{\text{MF}} = \sigma(x_i + \bar{b}(X)), \quad \text{with } \bar{b}(X) \text{ such that } \sum_{i=1}^N y_i^{\text{MF}} = V_0, \tag{1}$$

for $X=(x_1,...,x_N)\in\mathbb{R}^N$ and the sigmoid $\sigma(x)=\frac{1}{1+e^{-x}}$. We will denote this operation as: $Y^{\mathrm{MF}}=\Sigma(X)$. The sigmoid ensures that densities are in [0,1] and the choice of the optimal bias $\overline{b}(X)$ ensures that the volume constraint is satisfied.

Filtering: If the vector X is optimised directly with gradient descent, SIMP often converges toward checkerboard patterns, i.e. some high frequency noise in the image, which is a common issue with SIMP [1]. To overcome this issue a common technique is to use filtering [30]. In this paper, we consider low-pass density filters of the form: $X = T\bar{X}$ where T represents a convolution on the grid, \bar{X} are the design variables and X is the vector in equation 1. The loss function of this method is then naturally defined as: $\bar{X} \longmapsto C(\Sigma(T\bar{X}))$.

The gradient $\nabla_Y C$ is easily obtained by the self-adjointness of the variational problem [36, 16]. We recover $\nabla_X C$ from $\nabla_Y C$ using an implicit differentiation technique [11]. The following proposition is a consequence of implicit function theorem and chain rules:

Proposition 2.1. Let \dot{S} be the vector with entries $\dot{\sigma}(x_i + \bar{b}(X))$. We have $\nabla_X C = D_X \nabla_Y C$ with:

$$D_X := -\frac{1}{|\dot{S}|_1} \dot{S} \dot{S}^T + Diag(\dot{S}). \tag{2}$$

where $|.|_1$ denotes the l^1 norm of a vector. Furthermore D_X is a symmetric positive semi-definite matrix whose null-space is the space of constant vectors and has eigenvalues smaller than $\frac{1}{4}$.

2.3 Proposed algorithm: SIMP with Neural networks

Fully-Connected Neural Networks (FCNN) are characterised by the number of layers L+1, the numbers of neurons in each layer $(n_0, n_1, ..., n_L)$ and an activation function $\mu: \mathbb{R} \longrightarrow \mathbb{R}$, here we will use the particular case $n_L=1$. The activations $a^l \in \mathbb{R}^{n_l}$ and preactivations $\tilde{a}^l \in \mathbb{R}^{n_l}$ are defined recursively for all layers l, using the so-called NTK parameterisation [14]:

$$a^{0}(x) = x, \quad \tilde{a}^{l+1}(x) = \frac{\alpha}{\sqrt{n_{l}}} W^{l} a^{l}(x) + \beta b^{l}, \quad a^{l+1}(x) = \mu(\tilde{a}^{l+1}(x)),$$
 (3)

for some hyperparameters $\alpha, \beta \in [0,1]$ representing the contribution of the weights and bias terms respectively. The parameters $\theta = (\theta_p)_p$, consisting in weight matrices W^l and bias vectors b^l are drawn as i.i.d. standard normal random variables $\mathcal{N}(0,1)$. We denote the output of the network as $f_{\theta}(x) = \tilde{a}^L(x)$.

Remark: To ensure that the variance of the neurons at initialization is the equal to 1 at all layers, we choose α and β such that $\alpha^2 + \beta^2 = 1$ and use a standardised non-linearity, i.e. $\mathbb{E}_{X \sim \mathcal{N}(0,1)}[\mu(X)^2] = 1$ ([15]).

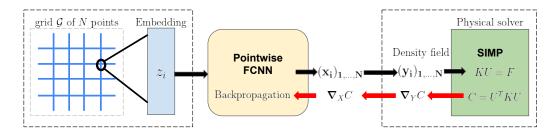


Figure 1: Illustration of our method

In our approach, the pre-densities $X^{\mathrm{NN}}(\theta)=(x_1^{\mathrm{NN}},...,x_N^{\mathrm{NN}})$ are generated by a neural network as $x_i^{\mathrm{NN}}=f_{\theta}(z_i)$ where $z_i\in\mathbb{R}^{n_0}$ is either the coordinates of the grid elements (in this case $n_0=d$) or

an embedding of those coordinates. We then apply the same transformation Σ to obtain the density field $Y^{\rm NN}(\theta) = \Sigma(X^{\rm NN}(\theta))$. Our loss function is then defined as:

$$\theta \longmapsto C(Y^{\mathrm{NN}}(\theta)) = C\big(\Sigma(X(\theta))\big).$$

The design variables are now the parameters θ of the network. The gradient $\nabla_{\theta}C$ w.r.t. to the parameters is computed by first using Proposition A.1 to get $\nabla_{Y^{\rm NN}}C$ followed by traditional backpropagation.

Remark: Note the absence of filter T in the above equations, indeed we will show how neural networks naturally avoid checkerboard patterns, making the use of filtering obsolete.

Initial density field: The SIMP method is usually initialised with a constant density field [1]. Since the neural network is initialized randomly, the initial density field is random and non-constant. To avoid this problem, we subtract the initial density field and add a well-chosen constant:

$$\forall i \in \{1, ..., N\}, \ x_i(\theta) = \bar{f}_{\theta(t)}(z_i) = f_{\theta(t)}(z_i) - f_{\theta(t=0)}(z_i) + \log\left(\frac{V_0}{N - V_0}\right). \tag{4}$$

We used equation 4 to compute $X(\theta)$ in our numerical experiments.

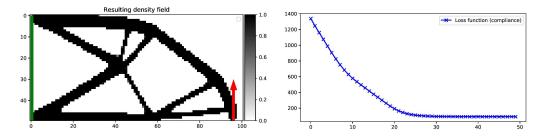


Figure 2: Example of result of our method with applied forces (red arrow) and a fixed boundary (green). Here we used a Gaussian embedding (see section 4 for details).

3 Theoretical Analysis

3.1 Analogy between the Neural Tangent Kernel and filtering techniques

In our paper, we use the Neural Tangent Kernel (NTK [14]) as the main tool to analyse the training behaviour of the FCNN. In our setting (where $n_L=1$) the NTK is defined as:

$$\forall z, z' \in \mathbb{R}^{n_0}, \ \Theta_{\theta}^L(z, z') = \sum_{n} \frac{\partial f_{\theta}}{\partial \theta_p}(z) \frac{\partial f_{\theta}}{\partial \theta_p}(z') = (\nabla_{\theta} f_{\theta}(z) | \nabla_{\theta} f_{\theta}(z')).$$

This is a positive semi-definite kernel. Given some inputs $z_1,...,z_N$ we define the NTK Gram matrix as: $\tilde{\Theta}_{\theta}^L := \left(\Theta^L(z_i,z_j)\right)_{1 < i,j < N} \in \mathbb{R}^{N \times N}$.

Assuming a small enough learning rate, the evolution of the network under gradient descent is well approximated by the gradient flow dynamics $\partial_t \theta(t) = -\nabla_\theta C(t)$. The evolution of the output of the network $X^{\text{NN}}(\theta)$ can then easily be expressed in terms of the NTK Gram matrix [15] for a loss \mathcal{L} :

$$\partial_t X^{\mathrm{NN}}(\theta(t)) = -\tilde{\Theta}^L_{\theta(t)} \nabla_{X^{\mathrm{NN}}} \mathcal{L}.$$

From this equation we can derive the evolution of the physical density field $Y^{\rm NN}$ in our algorithm:

Proposition 3.1. If the network is trained under this gradient flow, then by applying chain rules, we can prove that the density field follows the equation:

$$\partial_t Y^{NN}(\theta(t)) = -D_X(t) \tilde{\Theta}^L_{\theta(t)} D_X(t) \nabla_Y C(Y^{NN}(\theta(t))). \tag{5}$$

The analogy between the NTK and filtering techniques comes from the following observation. With Modified Filtering with a filter T, we show similarly that the density field $Y^{\rm MF}$ evolves as

$$\partial_t Y^{\text{MF}}(t) = -D_X(t)TT^T D_X(t) \nabla_Y C(Y^{\text{MF}}(t)). \tag{6}$$

We see that the NTK Gram matrix and the squared filter TT^T play exactly the same role. An important difference however is that the NTK is random at initialisation and evolves during training.

This difference disappears for large widths (when n_1, \ldots, n_{L-1} are large), since the NTK converges to a deterministic and time independent limit $\tilde{\Theta}_{\infty}^L$ as $n_1, \ldots, n_{L-1} \to \infty$ [14]. Furthermore, in contrast to the finite width NTK (also called empirical NTK), we have access to a closed form formula for the limiting NTK $\tilde{\Theta}_{\infty}^L$ (given in the appendix).

In the infinite width limit, the evolution of the physical densities is then expressed in terms of the limiting NTK Gram matrix $\tilde{\Theta}^L_{\infty}$:

$$\partial_t Y^{\text{NN}}(\theta(t)) = -D_X(t)\tilde{\Theta}_{\infty}^L D_X(t) \nabla_Y C(Y^{\text{NN}}(\theta(t))). \tag{7}$$

From now on we will focus on this infinite-width limit, comparing the NTK Gram matrix $\tilde{\Theta}_{\infty}^{L}$ and the squared filter TT^{T} . Recent results [18, 2, 13] suggest that this limit is a good approximation when the width of the network is sufficiently large. For more details see the appendix, where we compare the empirical NTK with its limiting one and plot its evolution in our setting.

3.2 Spatial invariance

In physical problems such as topology optimisation, it is important to ensure that certain physical properties are respected by the model. We focus in this section on the translation and rotation invariance of topology optimisation: if the force constraints are rotated or translated, the resulting shape should remain the same (up to rotation and translation), as in Figure 4 (b.1 and b.2).

In Modified Filtering method, this property is guaranteed if the filter T is translation and rotation invariant. In contrast the limiting NTK is in general invariant under rotation [14] but not translation. As Figure 4 shows, this leads to some problematic artifacts. The NTK can be made translation and rotation invariant by first applying an embedding $\varphi:\mathbb{R}^d\longrightarrow\mathbb{R}^{n_0}$ with the properties that for any two coordinates $p,p',\varphi(p)^T\varphi(p')$ only depends on the distance $\|p-p'\|_2$. Since the rotation invariance of the NTK implies that $\Theta^L_\infty(z,z')$ depends only on the scalar products z^Tz',zz^T and $z'z'^T$, we have that $\Theta^L_\infty(\varphi(p),\varphi(p'))$ depends only on $\|p-p'\|$ as needed.

The issue is that for finite n_0 there is no non-trivial embedding φ with this property:

Proposition 3.2. Let $\varphi : \mathbb{R}^d \to \mathbb{R}^{n_0}$ for d > 2 and any finite n_0 . If φ satisfies $\varphi(x)^T \varphi(x') = K(||x - x'||)$ for some continuous function K then both φ and K are constant.

To overcome this issue, we present two approaches to approximate spatial invariance with finite embeddings: an embedding on a (hyper)-torus and a random feature [24] embedding based on Bochner theorem [27].

3.2.1 Embedding on a hypertorus

In this subsection we consider the following embedding of a $n_x \times n_y$ regular grid on a torus:

$$\mathbb{R}^2 \ni p = (p_1, p_2) \longmapsto \varphi(p) = r(\cos(\delta p_1), \sin(\delta p_1), \cos(\delta p_2), \sin(\delta p_2)), \tag{8}$$

where $\delta > 0$ is a discretisation angle (our default choice is $\delta = \frac{\pi}{2 \max(n_x, n_y)}$). One can use similar formulas for d > 2 (leading to an hyper-torus embedding), we used d = 2 in equation 8 for simplicity.

This embedding leads to an exact translation invariance and an approximate rotation invariance:

$$\varphi(p)^{T}\varphi(p') = r^{2}(\cos(\delta(p_{1} - p'_{1})) + \cos(\delta(p_{2} - p'_{2}))) = r^{2}\left(2 - \frac{\delta^{2}}{2}\|p - p'\|_{2}^{2}\right) + \mathcal{O}\left(\delta^{4}\|p - p'\|_{4}^{4}\right).$$

As a result, the limiting NTK $\Theta_{\infty}(\varphi(p), \varphi(p'))$ is translation invariant and approximately rotation invariant (for small δ and/or when p, p' are close to each other). Moreover, if we look at the limiting NTK on the whole torus, we obtain that the gram matrix $\tilde{\Theta}_{\infty}$ is a discrete convolution on the input grid, with nice properties summed up in the following proposition:

Proposition 3.3. We can always extend our $n_x \times n_y$ grid and choose δ such that the embedded grid covers the whole torus (typically $\delta = \frac{\pi}{2 \max(n_x, n_y)}$ and take a $n \times n$ grid with $n = 4 \max(n_x, n_y)$).

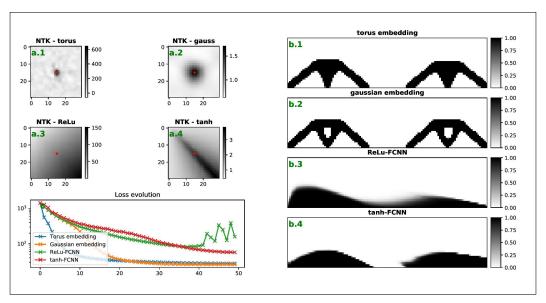


Figure 4: Left: empirical NTK of FCNNs with both embedding (a.1, a.2, see Section 4.1 for details) or without embedding (a.3 with ReLu, a.4 with tanh). Right: Corresponding shape obtained after training. Note that methods without spatial invariance particularly struggles with this symmetric load case (b.3, b.4) while both "embedded methods" respect the symmetry (b.1, b.2). We also observed that training with non-embedded methods is very unstable

Then the Gram matrix $\tilde{\Theta}_{\infty}$ of the limiting NTK is a 2D discrete convolution matrix. Moreover the NTK Gram matrix has a positive definite square root $\sqrt{\tilde{\Theta}_{\infty}}$ which is also a discrete convolution matrix.

As we know, the eigenvectors of such a convolution matrix are the 2D Fourier vectors. The corresponding eigenvalues are the discrete Fourier transforms of the convolution kernel.

The square root of the NTK Gram matrix $\sqrt{\tilde{\Theta}_{\infty}}$ then corresponds to the filtering matrix T in our analogy. Figure 3 shows that on the full torus, the matrix square root $\sqrt{\tilde{\Theta}_{\theta}}$ indeed looks like a typical smoothing filter.

As Figure 4 shows, the torus embedding method gives good numerical results and respect the symmetry of the applied forces F.

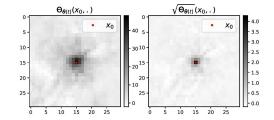


Figure 3: Representation of one line of $\tilde{\Theta}_{\theta}$ on the full torus and of its square root. We used $\beta=0.2$ and $\omega=3$ (see Section 4.1) here to make the filter visible on the whole torus.

3.2.2 Random embeddings for radial kernels

Another approach to approximate a rotation and translation invariant embedding is to use random Fourier features [24], which is a general method to approximate shift invariant kernels of the form k(x,y) = k(x-y). By Bochner theorem [27], any continuous non-zero radial kernel $k(x-y) = K(\|x-y\|)$ can be written as the (scaled) Fourier transform of a probability measure \mathbb{Q} on \mathbb{R}^d :

$$k(r) = k(0) \int_{\mathbb{R}^d} e^{i\omega \cdot r} d\mathbb{Q}(\omega).$$

For radial kernels, we formulate random Fourier features embeddings $\varphi : \mathbb{R}^d \to \mathbb{R}^{n_0}$ as follows:

$$\varphi(p)_i = \sqrt{2k(0)}\sin(w_i^T p + \frac{\pi}{4} + b_i),$$

for i.i.d. samples $w_1,...,w_{n_0}\in\mathbb{R}^d$ from \mathbb{Q} (which is also invariant by rotation) and i.i.d. samples $b_1,...,b_{n_0}\in\mathbb{R}^d$ from any symmetric probability distribution (or uniform laws on $[0,2\pi]$). By the law of large numbers for large n_0 , we have the approximation $\frac{1}{n_0}\varphi(p)^T\varphi(p')\simeq k(p-p')$.

Gaussian embedding: Depending on the kernel k that we want to approximate, it may be difficult to sample from the distribution \mathbb{Q} . The simplest case is for a Gaussian kernel $k(d) = e^{-\frac{1}{2\ell^2}d^2}$, where the distribution \mathbb{Q} of the weights w_i is $\mathcal{N}(0,\frac{1}{\ell^2}I_d)$, i.e. the entries w_{ij} are all i.i.d. $\mathcal{N}(0,\frac{1}{\ell^2})$ Gaussians. For this reason this is the embedding that we will use in our numerical experiments. Note the similarity between this type of embedding and an untrained first layer of a FCNN with sine activation function, weights w_i and bias b_i .

Moreover, the following result shows that we can still define a "square root" of the NTK with those types of embedding and thus complete the analogy with equation 6.

Proposition 3.4. Let φ be an embedding as described above for a positive radial kernel $k \in L^1(\mathbb{R}^d)$ with k(0) = 1, $k \geq 0$. Then there is a filter function $g : \mathbb{R} \to \mathbb{R}$ and a constant C such that for all p, p':

$$\lim_{n_0 \to \infty} \Theta_{\infty}(\varphi(p), \varphi(p')) = C + (g \star g)(p - p'), \tag{9}$$

where Θ_{∞} is the limiting NTK of a network with a Lipschitz, non-constant and standardised activation function μ . (Here \star denotes the convolution product).

As the matrix D_X in equation 7 cancels out the constant frequency (proposition A.1), the constant C doesn't matter, i.e. $D_X \tilde{\Theta}_{\infty}^{(L)} D_X = D_X \left(\tilde{\Theta}_{\infty}^{(L)} - C \right) D_X$.

4 Experimental analysis

4.1 Setup

Most of our experiments were conducted with a torus embedding or a gaussian embedding. For the SIMP algorithm, we adapted the code described in [1, 19]. Here are the hyperparameters used in the experiments.

For the Gaussian embedding, we used $n_0=1000$ and a length scale $\ell=4$. This embedding was followed by one hidden linear layer of size 1000 with standardized ReLu $(x\mapsto\sqrt{2}\max(0,x))$ and a bias parameter $\beta=0.5$.

For the torus embedding we set the torus radius to $r=\sqrt{2}$ (to be on a standard sphere) and the discretisation angle to $\delta=\frac{\pi}{2\max(n_x,n_y)}$ (to cover roughly half the torus, which is a good trade-off between rotation invariance and kernel size), where $n_x\times n_y$ is the size of the grid. It was followed by 2 linear layers of size 1000 with $\beta=0.1$. The ReLu activation is not well-suited in this case because it induces filters that are too wide. The large radius of the NTK kernel can be understood in relation with the order/chaos regimes [28, 23], as observed in [15] the ReLU lies in the ordered regime when $\beta>0$, leading to a "wide" kernel, a narrower kernel can be achieved with non-linearities which lie in the chaotic regime instead. We used a cosine activation of the form $x\mapsto\cos(\omega x)$, which has the advantage that the width of the filter can be adjusted using the ω hyperparameter, see Section 7. When not stated otherwise we used $\omega=5$.

Even though our theoretical analysis is for gradient flow, we obtain similar results with other optimizers such as RPROP [25] (learning rate 10^{-3}) and ADAM [17] (learning rate 10^{-3}). RPROP gave the fastest results, possibly because it is well-suited for batch learning [26]. Vanilla gradient descent can be very slow due to the vanishing of the gradients when the image becomes almost binary (due to the sigmoid), we therefore gradually increased the learning rate during training to compensate.

4.2 Spectral analysis

In SIMP convolution with a low pass filter ensures that low frequencies are optimised faster than high frequencies, to avoid checkerboards.

With the embeddings proposed in the last two subsections, the limiting NTK takes the form of a convolution over the input space \mathbb{R}^d . Figure 5 represents the eigenvalues and eigenimages of the

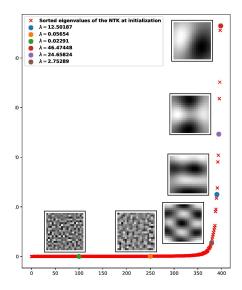


Figure 5: Sorted eigenvalues of the empirical NTK with some eigenvectors (reshaped as images). Obtained with a Gaussian embedding.

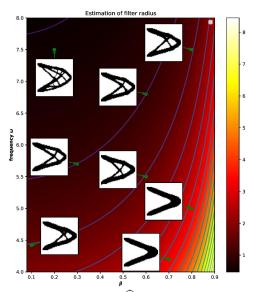


Figure 6: Colormap of $\widehat{R}_{1/2}$ in the (β,ω) plane, torus embedding. Level lines and shapes obtained for different radius are represented.

NTK Gram matrix $\tilde{\Theta}_{\theta(t)}$. Even though this plot is done for a finite width network and a finite random embedding, we see that the eigenimages look like 2D Fourier modes. The fact that the low frequencies have the largest eigenvalues supports the similarity between the NTK and a low pass filter.

This may explain why neural networks naturally avoid checkerboard patterns: the low frequencies of the shape are trained faster than the high frequencies which lead to checkerboard patterns.

4.3 Filter radius

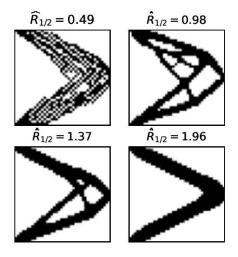


Figure 7: Shape obtained for different values of $\widehat{R}_{1/2}$ with a Gaussian embedding for different values of $\ell \in \{0.5, 1, 1.4, 2\}$.

In the classical SIMP algorithm, the choice of the radius of the filter T is critical. It controls the appearance of checkerboards or intermediate densities.

When using DNNs, there is no explicit choice of filter radius, since the filter depends on the embedding and the architecture of the network. In Section 3.2 we have shown that the NTK is approximately invariant, it can hence be expressed as:

$$\Theta_{\theta(t)}^{L}(\varphi(p), \varphi(p')) \simeq \Phi_{\infty}(\|p - p'\|),$$

where Φ_{∞} can be analytically expressed with the embedding and the limiting NTK (see appendix for a detailed example).

The kernels we consider do not have compact support in general, we therefore focus instead on the radius at half-maximum of Φ_{∞} :

$$\Phi_{\infty}(\widehat{R}_{1/2}) = \frac{1}{2} (\Phi_{\infty}(0) + \inf_{r} \Phi_{\infty}(r)).$$

Note that for simplicity we are computing here the radius of the squared filter, since obtaining a closed form

formula for the square root of the NTK is more difficult. For Gaussian filters the radius of the squared filter is $\sqrt{2}$ times that of the original, suggesting that the filter radius is well estimated by $\frac{1}{\sqrt{2}} \widehat{R}_{1/2}$.

The quantity $\widehat{R}_{1/2}$ is a function of the hyperparameters of the network (α, β, L) , see appendix) and of the embedding (the lengthscale ℓ). Using the formula for $\widehat{R}_{1/2}$, these hyperparameters can be tuned to obtain a specific filter radius.

With the Gaussian embedding, the radius of the filter can easily be adjusted by changing the length-scale ℓ of the embedding. As illustrated in Figure 7.

With the torus embedding, we instead have to change the hyperparameters of the network to adjust the radius of the filter. With the ReLU activation function, the radius is very large which makes it impossible to obtain precise shape. The solution we found is to use a cosine activation $x \mapsto \cos(\omega x)$ with hyper-parameter ω . Figure 6 shows how the radius decreases as ω increases. The β parameter has the opposite effect, as increasing it increases the radius. For different values of ω and β , we obtain a variety of radius and plot the resulting shapes. This plot also illustrates the role of the radius in the determination of the resulting shape. The fact that cosine activation leads to an adjustable NTK radius could explain why periodic activation function help in the representation of high frequency signal as observed in [31].

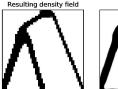
The effect of depth is more complex. For large depths L the NTK either approaches a constant kernel in the so-called order regime (with infinite radius) or a Kronecker delta kernel in the so-called chaos regime (with zero radius) [23, 28, 15]. Depending on whether we are in the order or chaos regime (which is determined by the activation function μ and the parameters α , β), increasing the depth can either increase or decrease the radius.

We conducted an experimental study of the influence of this parameter on the geometry of the final shape. We observed that its complexity (number of holes, high frequencies) is highly controlled by $\widehat{R}_{1/2}$. We see in Figure 7 and 6 some examples of shape obtained for several values of $\widehat{R}_{1/2}$.

4.4 Up-sampling

Since the density field is generated by a DNN, it can be evaluated at any point in \mathbb{R}^d , hence allowing upsampling. As Figure 8 shows, with our method we obtain a smooth and binary shape. Something interesting happens when the network is trained without an embedding: when upsampling we observe some visual artifacts plotted in Figure 9. We believe that it is due to the lack of spatial invariance.

Note that this second experiment was done with batch norm, as described in [7], since for this problem it was difficult to obtain a good shape with a vanilla ReLU-FCNN. With our embeddings, we can achieve complex shapes without batch-norm.





A



Figure 8: Density field obtained with a Torus embedding (left) and up sampling of factor 6 of the same network (right).

Figure 9: Exemple of up-sampling of a FCNN (ReLu FCNN with batchnorms) without embedding, exhibing typical visual artifacts.

5 Conclusion

Using the NTK, we were able to give a simple theoretical description of topology optimisation with DNNs, showing a similarity to traditional filtering techniques. This theory allowed us to identify a problem: since the NTK is not translation invariant, the spatial invariance of topology optimisation is not respected, leading to visual artifacts and non-optimal shapes. We propose a simple solution to this problem: adding a spatial invariant embedding to the coordinates before the DNN.

Using this method, our models are able to learn efficient shapes while avoiding checkerboard patterns. We give tools to adjust the implicit filter size induced by the hyperparameters, to give control over the complexity of the final shape. Using the learned network, we can easily perform good quality

up-sampling. The techniques described in this paper can easily be translated to any other problem where spatial invariance is needed.

The NTK is a simple yet powerful tool to analyse a practical method such as SIMP when combined with a DNN. Morover it can be used to make informed choices of the DNN's architecture and hyperparameters.

Acknowledgments and Disclosure of Funding

There is no funding or competing interests associated to this work.

References

- [1] Erik Andreassen, Anders Clausen, Mattias Schevenels, Boyan Lazarov, and Ole Sigmund. Efficient topology optimization in matlab using 88 lines of code. *Structural and Multidisciplinary Optimization*, 43:1–16, 11 2011.
- [2] Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *CoRR*, abs/1904.11955, 2019.
- [3] Saurabh Banga, Harsh Gehani, Sanket Bhilare, Sagar Patel, and Levent Kara. 3d topology optimization using convolutional neural networks. *CoRR*, abs/1808.07440, 2018.
- [4] Bendsoe and Sigmund. Topology optimization: Theory, methods and applications. *Springer Science and Business*, April 2013.
- [5] Martin Bendsøe. Bendsoe, m.p.: Optimal shape design as a material distribution problem. structural optimization 1, 193-202. *Structural Optimization*, 1:193–202, 01 1989.
- [6] Aaditya Chandrasekhar and K. Suresh. Length scale control in topology optimization using fourier enhanced neural networks. 2020.
- [7] Aaditya Chandrasekhar and Krishnan Suresh. Tounn: Topology optimization using neural networks. *Structural and Multidisciplinary Optimization*, 2020.
- [8] Yanqing Chen, Timothy A. Davis, and William W. Hager. Algorithm 887: Cholmod, supernodal sparse cholesky factorization and update/downdate. *ACM Transactions on Mathematical Software*, pages 1–14, 2008.
- [9] Amit Daniely, Roy Frostig, and Yoram Singer. Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity. *CoRR*, abs/1602.05897, 2016.
- [10] Timothy A. Davis. User guide for cholmod: a sparse cholesky factorization and modification package. 2009.
- [11] Andreas Griewank and Christèle Faure. Reduced functions, gradients and hessians from fixed-point iterations for state equations. *Numerical Algorithms*, 30:113–139, 06 2002.
- [12] Stephan Hoyer, Jascha Sohl-Dickstein, and Sam Greydanus. Neural reparameterization improves structural optimization. *CoRR*, abs/1909.04240, 2019.
- [13] Jiaoyang Huang and Horng-Tzer Yau. Dynamics of deep neural networks and neural tangent hierarchy. *CoRR*, abs/1909.08156, 2019.
- [14] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *CoRR*, abs/1806.07572, 2018.
- [15] Arthur Jacot, Franck Gabriel, and Clément Hongler. Order and chaos: NTK views on DNN normalization, checkerboard and boundary artifacts. CoRR, abs/1907.05715, 2019.
- [16] Jiaqi Jiang and Jonathan A. Fan. Global optimization of dielectric metasurfaces using a physics-driven neural network. *Nano Letters*, 19(8):5366–5372, Jul 2019.
- [17] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. 2017.
- [18] Jaehoon Lee, Lechao Xiao, Samuel S Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. *Journal of Statistical Mechanics: Theory and Experiment*, 2020(12):124002, Dec 2020.

- [19] Kai Liu and Andres Tovar. An efficient 3d topology optimization code written in matlab. *Structural and Multidisciplinary Optimization*, 50, 12 2014.
- [20] Gilles Marck, Maroun Nemer, Jean-Luc Harion, Serge Russeil, and Daniel Bougeard. Topology optimization using the simp method for multiobjective conductive problems. *Numerical Heat Transfer Part B-fundamentals - NUMER HEAT TRANSFER PT B-FUND*, 61:439–470, 06 2012.
- [21] Alexander Mordvintsev, Nicola Pezzotti, Ludwig Schubert, and Chris Olah. Differentiable image parameterizations. *Distill*, 2018. https://distill.pub/2018/differentiable-parameterizations.
- [22] Zhenguo Nie, Tong Lin, Haoliang Jiang, and Levent Burak Kara. Topologygan: Topology optimization using generative adversarial networks based on physical fields over the initial domain. *CoRR*, abs/2003.04685, 2020.
- [23] Ben Poole, Subhaneil Lahiri, Maithra Raghu, Jascha Sohl-Dickstein, and Surya Ganguli. Exponential expressivity in deep neural networks through transient chaos. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, Advances in Neural Information Processing Systems 29, pages 3360–3368. Curran Associates, Inc., 2016.
- [24] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2008.
- [25] M. Riedmiller and H. Braun. A direct adaptive method for faster backpropagation learning: the rprop algorithm. In *IEEE International Conference on Neural Networks*, pages 586–591 vol.1, 1993.
- [26] Martin Riedmiller and Heinrich Braun. A direct adaptive method for faster backpropagation learning: The rprop algorithm. pages 586–591, 1993.
- [27] W. Rudin. Fourier Analysis on Groups. Wiley Classics Library. Wiley, 1990.
- [28] Samuel S. Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep information propagation. 2017.
- [29] M.-H. Herman Shen and Liang Chen. A new CGAN technique for constrained topology design optimization. CoRR, abs/1901.07675, 2019.
- [30] Ole Sigmund. Morphology-based black and white filters for topology optimization. *Structural and Multidisciplinary Optimization*, 33:401–424, 04 2007.
- [31] Vincent Sitzmann, Julien N. P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *CoRR*, abs/2006.09661, 2020.
- [32] Ivan Sosnovik and Ivan V. Oseledets. Neural networks for topology optimization. CoRR, abs/1709.09578, 2017.
- [33] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *CoRR*, abs/2006.10739, 2020
- [34] H. Vorst. Bi-cgstab: A fast and smoothly converging variant of bi-cg for the solution of nonsymmetric linear systems. *SIAM J. Sci. Comput.*, 13:631–644, 1992.
- [35] Lechao Xiao, Jeffrey Pennington, and Samuel S. Schoenholz. Disentangling trainability and generalization in deep learning. *CoRR*, abs/1912.13053, 2019.
- [36] Luzhong Yin and Wei Yang. Optimality criteria method for topology optimization under multiple constraints. *Computers and Structures*, 79(20):1839–1850, 2001.

A Derivation of the algorithm

In this section we show how to derive the equations used in our algorithm, especially the ones corresponding to implicit differentiation [11]. Let us recall that we consider a vector $X \in \mathbb{R}^N$ and compute a vector $Y = \Sigma(X) \in [0,1]^N$ (either Y^{MF} or Y^{NN}) by:

$$\forall i \in \{1,...,N\}, \ y_i = \sigma(x_i + \bar{b}(X)), \quad \text{such that: } \sum_{i=1}^N y_i = V_0, \quad \sigma(x) = \frac{1}{1 + e^{-x}},$$

Where X denotes $(x_1, ..., x_N)$.

We want to show that this operation is well defined and find a formula to recover $\nabla_X C$ from a given $\nabla_Y C$. More precisely we have the following result.

Proposition A.1 (Proposition A.1 in the paper). Let $X \in \mathbb{R}^N$, the operation $Y = \Sigma(X)$ is well defined. Moreover, let \dot{S} be the vector of the $\dot{\sigma}(x_i + \bar{b}(X))$. Then we have $\nabla_X C = D_X \nabla_Y C$ with:

$$D_X := -\frac{1}{|\dot{S}|_1} \dot{S} \dot{S}^T + Diag(\dot{S}). \tag{10}$$

 D_X is a symmetric positive semi-definite matrix whose kernel corresponds to constant vectors and has eigenvalues smaller than $\frac{1}{2}$.

Proof: Let us consider the function $F: \mathbb{R}^N \times \mathbb{R} \longrightarrow \mathbb{R}$ defined by: $F(z,b) = \sum_{i=1}^N \sigma(z_i+b)$. It is clear that F(X,.) is strictly increasing on \mathbb{R} from 0 to N. Then $\exists ! \bar{b} \in \mathbb{R}$ such that $F(X,\bar{b}) = V_0$.

As $\partial_b F(X, \bar{b}) > 0$, by the implicit functions theorem, there exists a neighbourhood V of X in \mathbb{R}^N , a neighbourhood U of \bar{b} in \mathbb{R} and a function $\bar{b}: V \longrightarrow \mathbb{R}$ of class C^1 such that:

$$\forall (z,b) \in V \times U, \ F(z,b) = V_0 \iff b = \bar{b}(z).$$

Moreover we also get from the implicit function theorem that:

$$\frac{\partial \bar{b}}{\partial x_i}(X) = -\left(\frac{\partial F}{\partial b}(X,\bar{b})\right)^{-1} \frac{\partial F}{\partial x_i}(X,\bar{b}) = -\left(\sum_{i=1}^N \dot{\sigma}(x_i + \bar{b})\right)^{-1} \dot{\sigma}(x_i + \bar{b}),$$

and we can apply chain rules:

$$\begin{split} \frac{\partial C}{\partial x_i} &= \sum_{j=1}^N \frac{\partial C}{\partial y_j} \frac{\partial y_j}{\partial x_i} \\ &= \sum_{i=1}^N \frac{\partial C}{\partial y_j} \dot{\sigma}(x_j + \bar{b}(x)) \left(\frac{\partial \bar{b}}{\partial x_i} + \delta_{ij} \right), \end{split}$$

So that equation 10 makes sense. Now, if we denote $\dot{S}=(a_1,...,a_N)$, let us recall that we defined $a_i=\dot{\sigma}x_i+\bar{b}(X)$ where σ is the sigmoid function. By taking any $u\in\mathbb{R}^N$, we remark that:

$$(D_X u)_i = \frac{a_i}{|\dot{S}|_1} \sum_{i=1}^N a_j (u_i - u_j).$$
(11)

We easily deduce from equation 11 that $\ker(D_X) = \operatorname{span}(1_N)$ and that $D_X \in S_N^+(\mathbb{R})$. Indeed:

$$\forall u \in \mathbb{R}^{N}, \quad u^{T}(D_{X})u = -\frac{1}{|\dot{S}|_{1}}u^{T}\dot{S}\dot{S}^{T}u^{T} + \sum_{i=1}^{N}a_{i}u_{i}^{2}$$

$$= \frac{1}{|\dot{S}|_{1}} \left\{ -\left(\sum_{i=1}^{N}a_{i}u_{i}\right)^{2} + \left(\sum_{i=1}^{N}a_{i}\right)\left(\sum_{i=1}^{N}a_{i}u_{i}^{2}\right)\right\}$$

$$= \frac{1}{|\dot{S}|_{1}} \sum_{1 \leq i, j \leq N} a_{i}a_{j}u_{i}(u_{i} - u_{j})$$

$$= \frac{1}{|\dot{S}|_{1}} \sum_{1 \leq i < j \leq N} a_{i}a_{j}(u_{i} - u_{j})^{2} \geq 0.$$

Eigenvalues: We already know that 0 is an eigenvalue with multiplicity 1. So let $u \neq 0$ in \mathbb{R}^N and $\lambda > 0$ such that: $D_X u = \lambda u$. Then we easily show:

$$\forall i \in [1, N], \quad \frac{a_i - \lambda}{a_i} u_i = \frac{1}{|\dot{S}|_1} \sum_{j=1}^N a_j u_j =: \langle u \rangle_a.$$

If $\langle u \rangle_a = 0$, then necessarily $\lambda \in \{a_1,...,a_N\}$ If $\langle u \rangle_a \neq 0$, then we can assume (by normalising u) that $\langle u \rangle_a = 1$ and we have $u_i = \frac{a_i}{a_i - \lambda}$. Then we can replace $u_i = \frac{a_i}{a_i - \lambda}$ in the equation $\langle u \rangle_a = 1$:

$$\sum_{j=1}^N a_j = \sum_{j=1}^N \frac{a_j^2}{a_j - \lambda}, \quad \text{which by substraction leads to} \quad F(\lambda) := \sum_{j=1}^N \frac{a_j}{a_j - \lambda} = 0,$$

By studying the function F, we see that $\forall \lambda > \max_i(a_i), F(\lambda) < 0$. Therefore an eigenvalue always satisfies the inequality:

$$\lambda \le \max\{a_1, ..., a_N\} \le \|\dot{\sigma}\|_{\infty} = \frac{1}{4},$$

The last inequality coming from the fact that $a_i = \dot{\sigma}(x_i + \bar{b}(X))$, as mentionned earlier.

Remark: As shown above an important property of the matrix D_X is that it cancels out constants, which allows us to consider the limiting NTK up to some constant. The fact that the eigenvalues of D_X are in $[0,\frac{1}{4}]$ can help to avoid exploding gradients.

Equations of evolution

We quickly show how equations 5, 6 and 7 of the paper are derived. The proofs are mainly based on chain rules.

Let us first remark that the matrix D_X introduced above actually corresponds to the jacobian matrix $\nabla_X \Sigma$ of the application $\Sigma: \mathbb{R}^N \longrightarrow [0,1]^N$. So we can immediately applied chain rules to $Y^{\mathrm{NN}} = \Sigma(X(\theta))$ and get:

$$\begin{split} \frac{\partial Y^{\text{NN}}}{\partial t} &= D_{X(\theta(t))} \frac{\partial X(\theta(t))}{\partial t} \\ &= -D_{X(\theta(t))} \tilde{\Theta}^L_{\theta(t)} \nabla_{X_{\theta(t)}} C \quad \text{(Gradient Descent)} \\ &= -D_{X(\theta(t))} \tilde{\Theta}^L_{\theta(t)} D_{X(\theta(t))} \nabla_{Y^{\text{NN}}} C(\theta(t)) \quad \text{(By proposition A.1)}. \end{split}$$

Similarly, for the MF method, we set $X=T\bar{X}$ and obtain:

$$\begin{split} \frac{\partial Y^{\rm MF}}{\partial t} &= D_{X(t)} \frac{\partial X(t)}{\partial t} \\ &= D_{X(t)} T \frac{\partial \bar{X}(t)}{\partial t} \quad \text{(Linearity)} \\ &= -D_{X(t)} T \nabla_{\bar{X}} C \quad \text{(Gradient descent)} \\ &= -D_{X(t)} T T^T \nabla_X C \quad \text{(Chain rule)} \\ &= -D_{X(t)} T T^T D_{X(t)} \nabla_{Y^{\rm MF}} C. \end{split}$$

Details about embeddings

C.1 Torus embedding

The aim of this section is to give details about properties of the limiting NTK in case of Torus embedding. As a reminder we consider the following embedding:

$$\mathbb{R}^2 \ni p = (p_1, p_2) \longmapsto \varphi(p) = r(\cos(\delta p_1), \sin(\delta p_1), \cos(\delta p_2), \sin(\delta p_2));$$

In particular we show the following proposition which basically says that $\tilde{\Theta}_{\infty}$ is in that case a discrete convolution and derive from there its spectral properties and construct its positive semi-definite square root

Proposition C.1 (Proposition 3.3 in the paper). We can always extend our $n_x \times n_y$ grid and choose δ such that the embedded grid covers the whole torus (typically $\delta = \frac{\pi}{2\max(n_x,n_y)}$ and take a $n \times n$ grid with $n = 4\max(n_x,n_y)$). Then the Gram matrix $\tilde{\Theta}_{\infty}$ of the limiting NTK is a 2D discrete convolution matrix. Moreover the NTK Gram matrix has a positive definite square root $\sqrt{\tilde{\Theta}_{\infty}}$ which is also a discrete convolution matrix.

proof: We assume that we extend the grid in a $n \times n$ grid with $n \ge n_x, n_y$. Now we take $\delta = \frac{2\pi}{n}$ and we consider the limiting NTK Gram matrix on $\varphi(\llbracket n, n \rrbracket \times \llbracket n, n \rrbracket)$.

As $\Theta_{\infty}(\varphi(p), \varphi(p'))$ depends only on p-p', we can see the limiting NTK Gram Matrix as a discrete convolution kernel \mathcal{K} acting on $\mathbb{Z}/n\mathbb{Z} \times \mathbb{Z}/n\mathbb{Z}$:

$$\Theta_{\infty}((k,k'),(j,j')) = \mathcal{K}(k-k',j-j'),$$

For (k, k'), $(j, j') \in \mathbb{Z}/n\mathbb{Z} \times \mathbb{Z}/n\mathbb{Z}$.

We see $\tilde{\Theta}_{\infty}$ as a n^2 square matrix with each index in $\mathbb{Z}/n\mathbb{Z} \times \mathbb{Z}/n\mathbb{Z}$.

We introduce the Fourier vectors $\Omega_m = (e^{-i2\pi \frac{mk}{n}})_{0 \le k \le n_x - 1}$. As $\tilde{\Theta}_{\infty}$ is a 2D convolution matrix, we classically have the following results:

The eigenvectors of $\tilde{\Theta}_{\infty}$ are exactly given by:

$$\Omega_m \otimes \Omega_M$$

for $0 \le m \le n_x - 1$ and $0 \le M \le n_y - 1$, \otimes denotes the Kronecker product. The corresponding eigenvalue is given by the discrete Fourier transform $\widehat{\mathcal{K}}(m,M)$ with:

$$\widehat{\mathcal{K}}(m,M) = \sum_{j=0}^{n-1} \sum_{j'=0}^{n-1} e^{-i2\pi \frac{mj}{n}} e^{-i2\pi \frac{Mj'}{n}} \mathcal{K}(j,j').$$

Moreover, as the matrix $\tilde{\Theta}_{\infty}$ is positive definite (from the positive definiteness of the NTK, [14]) those eigenvalues verify $\hat{\mathcal{K}}(m,M) \geq 0$ and it makes sense to write the square root of the NTK Gram Matrix as the inverse Fourier transform of the $\sqrt{\hat{\mathcal{K}}(m,M)}$:

$$\sqrt{\widetilde{\Theta}_{\infty}}((k,k'),(j,j')) = \frac{1}{n^2} \sum_{m=0}^{n-1} \sum_{M=0}^{n-1} e^{i2\pi \frac{m(j-k)}{n}} e^{i2\pi \frac{M(j'-k')}{n}} \sqrt{\widehat{\mathcal{K}}(m,M)},$$
(12)

It is easy to see that the matrix defined by equation 12 is symmetric and positive semi-definite. Indeed we can write $\sqrt{\tilde{\Theta}_{\infty}}((k,k'),(j,j'))=g(k-j,k'-j')$ with g the Fourier transform of a positive vector.

Moreover it follows from the (discrete) convolution theorem that $\sqrt{\tilde{\Theta}_{\infty}}^2 = \tilde{\Theta}_{\infty}$. Therefore $\sqrt{\tilde{\Theta}_{\infty}}((k,k'),(j,j'))$ is indeed the positive semi-definite matrix square root of $\tilde{\Theta}_{\infty}$.

Thus the square root of the NTK Gram matrix can be seen as a convolution filter as well (it is invariant by translation as a function of (k-j,k'-j')).

C.2 Dimension of radial embeddings

In this section we prove that feature maps associated to continuous radial kernels are either trivial or of infinite dimension. this result is what motivates discussion in section 3.2 of the paper.

Let us first recall Bochner theorem ([27]):

Theorem C.1 (Bochner). Let $(x,y) \mapsto k(x-y)$ be a continuous shift invariant positive definite kernel on \mathbb{R}^d . Then it is the Fourier transform of a finite positive measure Λ on \mathbb{R}^d :

$$k(r) = \int_{\mathbb{R}^d} e^{i\omega \cdot r} d\Lambda(\omega).$$

The function k appearing in the above theorem will be called a positive definite function, according to the following definition:

Definition C.1. Let $k : \mathbb{R}^d \longrightarrow \mathbb{R}$, then k is a positive definite function when for all n, all $p_1, \ldots, p_n \in \mathbb{R}^d$ and all $c_1, \ldots, c_n \in \mathbb{R}$ we have:

$$\sum_{1 \le i, j \le n} c_i c_j k(x_i - x_j) \ge 0.$$

Moreover we will denote SO(d) the set of rotations matrices of dimension d and the Fourier transform (for an integrable function ψ):

$$\mathcal{F}\psi(\omega) = \int_{\mathbb{R}^p} \psi(p) e^{-i\omega \cdot p} dp.$$

Let us now recall the result that we want to prove:

Proposition C.2 (Proposition 3.2 in the paper). Let $\varphi : \mathbb{R}^d \to \mathbb{R}^m$ for d > 2 and any finite m. If φ satisfies

$$\varphi(x)^T \varphi(x') = K(\|x - x'\|) \tag{13}$$

for some continuous function K then both φ and K are constant. We will denote $k(x-x'):=K(\|x-x'\|)$.

Proof: We procede in the following way: We consider an embedding φ as described above and we are going to show that, when K is not constant, one can construct arbitrarily big linearly independent families $\varphi(p_1), \ldots, \varphi(p_n)$.

For now let us take pairwise distinct $p_1, \ldots, p_n \in \mathbb{R}^d$ and $c_1, \ldots, c_n \in \mathbb{R}$ such that:

$$\sum_{k=1}^{n} c_k \varphi(p_k) = 0.$$

A clever choice for p_1, \ldots, p_n will be done later.

For any $p \in \mathbb{R}^d$ and any rotation $R \in SO(d)$ we can write:

$$0 = \varphi(p)^{T} \sum_{k=1}^{n} c_{k} \varphi(p_{k}) = \sum_{k=1}^{n} c_{k} K(\|p - p_{k}\|) = \sum_{k=1}^{n} c_{k} K(\|Rp - Rp_{k}\|)$$
$$= \varphi(Rp)^{T} \sum_{k=1}^{n} c_{k} \varphi(Rp_{k}).$$

Since this is true for all p' = Rp we can deduce that for all $p \in \mathbb{R}^d$ and all $R \in SO(d)$ we have:

$$\sum_{k=1}^{n} c_k k(p - Rp_k) = 0.$$

We denote by Λ the finite measure on \mathbb{R}^d given by Bochner's theorem applied on k.

Let us take a test function $\psi \in \mathcal{S}(\mathbb{R}^p)$ in the Schwartz space, we can write successively that for all rotation $R \in SO(d)$:

$$0 = \int_{\mathbb{R}^d} \mathcal{F}\psi(p) \sum_{k=1}^n c_k k(p - Rp_k) dp$$

$$= \int_{\mathbb{R}^d} \mathcal{F}\psi(p) \sum_{k=1}^n c_k \int_{\mathbb{R}^d} e^{i\omega \cdot (p - Rp_k)} d\Lambda(\omega) dp, \quad \text{(Bochner's theorem)}$$

$$= \int_{\mathbb{R}^d} \left(\sum_{k=1}^n c_k e^{-i\omega \cdot (Rp_k)} \right) \int_{\mathbb{R}^d} \mathcal{F}\psi(p) e^{i\omega \cdot p} dp \, d\Lambda(\omega), \quad \text{(Fubini's theorem)}$$

$$= (2\pi)^d \int_{\mathbb{R}^d} \psi(\omega) \sum_{k=1}^n c_k e^{-i\omega \cdot (Rp_k)} d\Lambda(\omega), \quad \text{(Fourier inversion)}$$

As K is not constant, we can find $\omega_0 \in \mathbb{R}^d \setminus \{0\}$ such that for all $\epsilon > 0$ small enough we have $\Lambda(B(\omega_0, \epsilon)) > 0$ (otherwise the finite positive measure Λ would be concentrated on 0 and k would be constant).

Let $R \in SO(d)$, if we assume that $S := \sum_{k=1}^n c_k e^{-i\omega_0 \cdot (Rp_k)} \neq 0$ then we can find a small enough open ball $B(\omega_0, \epsilon)$ on which Re(S) and Im(S) have constant sign and such that: $|Re(S)| \geq c_1 > 0$ or $|Im(S)| \geq c_1 > 0$.

We choose ψ such that $\psi \geq 0$, ψ has compact support in $B(\omega_0, \epsilon)$ and $\psi \geq c_2 > 0$ on $B(\omega_0, \frac{\epsilon}{2})$. Then we obtain a contradiction by writing $0 \geq (2\pi)^d c_1 c_2 \Lambda(B(\omega_0, \frac{\epsilon}{2}))$. (We separate real and imaginary parts).

This implies that:

$$\forall R \in SO(d), \ \sum_{k=1}^{n} c_k e^{-i(R\omega_0) \cdot p_k} = 0, \tag{14}$$

Now we take a particular choice of (p_i) , let $p_k = (k, 0, \dots, 0) \in \mathbb{R}^d$.

Up to rotations, we can assume without loss of generality that $\omega_0=(w,0,\ldots,0)$ with $w\neq 0$. Moreover, we consider the particular case of rotations in the 2D plane generated by $(1,0,\ldots,0)$ and $(0,1,0,\ldots,0)$.

Therefore, equation 14 implies that:

$$\forall \theta \in \mathbb{R}, \ \sum_{k=1}^{n} c_k \left(e^{-iw \cos(\theta)} \right)^k = 0,$$

So that the polynomial $\sum_k c_k z^k$ has an infinite number of roots. Thus $c_1 = \cdots = c_n = 0$.

C.3 Random features embedding

In this section we give some details about the way we define random embeddings, which is very similar but slightly different than in [24].

If the kernel is properly scaled (i.e. k(0) = 1) then Λ defines a probability measure. That's why we introduce a probability measure $\mathbb Q$ and write:

$$k(r) = k(0) \int_{\mathbb{D}^d} e^{i\omega \cdot r} d\mathbb{Q}(\omega) = k(0) \mathbb{E}_{\omega \sim \mathbb{Q}}[e^{i\omega \cdot r}].$$

Now, following the reasoning in [24] we consider:

$$\varphi(p)_i = \sqrt{2k(0)}\sin(\omega \cdot p + \frac{\pi}{4} + b)$$

With $\omega \sim \mathbb{Q}$ and b a random variable with a symmetric law (note that \mathbb{Q} is also symmetric). Then we have:

$$\begin{split} \mathbb{E}[\varphi(p)_{i}\varphi(p')_{i}] &= 2k(0)\mathbb{E}\left[\left(\frac{e^{i\omega\cdot p + \frac{\pi}{4} + b} - e^{-i\omega\cdot p - \frac{\pi}{4}} - b}{2i}\right)\left(\frac{e^{i\omega\cdot p' + \frac{\pi}{4} + b} - e^{-i\omega\cdot p' - \frac{\pi}{4}} - b}{2i}\right)\right] \\ &= -\frac{k(0)}{2}\left(e^{i\frac{\pi}{2}}\mathbb{E}[e^{i\omega\cdot (p+p') + 2b}] + e^{-i\frac{\pi}{2}}\mathbb{E}[e^{-i\omega\cdot (p+p') - 2b}] \\ &\qquad - \mathbb{E}[e^{i\omega\cdot (p-p')}] - \mathbb{E}[e^{-i\omega\cdot (p-p')}]\right) \\ &= k(0)\mathbb{E}[e^{i\omega\cdot (p-p')}] \\ &= k(p-p'). \end{split}$$

Therefore we reduce the variance by drawing i.i.d. samples $\omega_1, \ldots, \omega_{n_0}$ and b_1, \ldots, b_{n_0} as described in section 3 and computing the mean $\frac{1}{n_0}\varphi(p)^T\varphi(p')$. By the strong law of large numbers we have the almost sure convergence:

$$\frac{1}{n_0}\varphi(p)^T\varphi(p')\underset{n_0\to\infty}{\longrightarrow} k(p-p'),$$

Now we can obtain Gaussian embedding by drawing the bias from δ_0 and weights from $\mathcal{N}(0, \frac{1}{\ell^2}I_d)$. from the above formulas we immediately get:

$$k(p-p') = e^{-\frac{\|p-p'\|_2^2}{2\ell^2}}.$$

D **Precise computations of the Neural Tangent Kernel**

We now give more details about the computation of the limiting NTK and detail how we obtain the limiting kernels used in Figures 6 and 7 of the paper.

D.1 Limiting NTK

For this purpose, following several authors ([14], [35], [18]), we need to introduce some gaussian processes and their associated kernels. For a symmetric positive kernel Σ let us define:

$$\begin{cases} \mathcal{T}(\Sigma)(z,z') = \mathbb{E}_{(X,Y) \sim \mathcal{N}(0,\Sigma_{z,z'})} \big[\mu(X) \mu(Y) \big] \\ \dot{\mathcal{T}}(\Sigma)(z,z') = \mathbb{E}_{(X,Y) \sim \mathcal{N}(0,\Sigma_{z,z'})} \big[\dot{\mu}(X) \dot{\mu}(Y) \big] \end{cases} \quad \text{With}: \quad \Sigma_{z,z'} = \begin{pmatrix} \Sigma(z,z) & \Sigma(z,z') \\ \Sigma(z,z') & \Sigma(z',z') \end{pmatrix}.$$

Then we set $\Sigma^1(z,z')=\Theta^1_\infty(z,z')=\beta^2+\frac{\alpha^2}{n_0}z^Tz'$ and we define recursively: $\sigma^{l+1}=\beta^2+\alpha^2\mathcal{T}(\Sigma^l),\quad \dot{\Sigma}^{l+1}=\alpha^2\dot{\mathcal{T}}(\Sigma^l),\quad \Theta^{l+1}_\infty=\dot{\Sigma}^{l+1}\Theta^l_\infty+\Sigma^{l+1}.$

$$\sigma^{l+1} = \beta^2 + \alpha^2 \mathcal{T}(\Sigma^l), \quad \dot{\Sigma}^{l+1} = \alpha^2 \dot{\mathcal{T}}(\Sigma^l), \quad \Theta_{\infty}^{l+1} = \dot{\Sigma}^{l+1} \Theta_{\infty}^l + \Sigma^{l+1}. \tag{15}$$

Using those formulas it is clear that the limiting NTK is invariant under rotation.

When neurons have constant variance, the following notion of dual activation function is often very useful:

Definition D.1. Let $\mu: \mathbb{R} \longrightarrow \mathbb{R}$ be a function such that $\mathbb{E}_{X \sim \mathcal{N}(0,1)}[\mu(X)^2] < +\infty$, then its dual function $\hat{\mu}: [-1,1] \longrightarrow \mathbb{R}$ is defined by:

$$\hat{\mu}(\rho) = \mathbb{E}_{(X,Y) \sim \mathcal{N}(0,\Sigma_{\rho})}[\mu(X)\mu(Y)], \quad \textit{With} : \Sigma_{\rho} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}.$$

We will use some properties of the dual function, which are described in [9].

D.2 Another way of seeing Gaussian embedding

As explained above (Section 3.2), the Gaussian embedding can be seen as the first hidden layer of a neural network, with the first layer untrained. Thus it actually corresponds to Σ^2 with the above

Let us consider the activation function $\mu: a \longmapsto \lambda \sin(\omega a + \frac{\pi}{4})$ and denote:

$$\forall x,y \in \mathbb{R}^{n_0}, \ \Sigma_{x,y}^1 = \binom{\beta^2 + \frac{1-\beta^2}{n_0} \|x\|_2^2 \quad \beta^2 + \frac{1-\beta^2}{n_0} x^T y}{\beta^2 + \frac{1-\beta^2}{n_0} x^T y \quad \beta^2 + \frac{1-\beta^2}{n_0} \|y\|_2^2},$$

We are looking at:

$$\Sigma^{2}(x,y) = \beta^{2} + (1 - \beta^{2}) \mathbb{E}_{(X,Y) \sim \mathcal{N}(0,\Sigma_{x,y}^{(1)})} [\mu(X)\mu(Y)].$$

Let $(X,Y) \sim \mathcal{N}(0,\Sigma^1_{x,y})$, then X-Y and X+Y are normal random variables and $\mathbb{V}(X-Y)=$ $\frac{1-\beta^2}{n_0}\|x-y\|_2^2.$ Thus, using properties of characteristic functions we get:

$$\begin{split} \mathbb{E}[\mu(X)\mu(Y)] &= \lambda^2 \mathbb{E}\bigg[\bigg(\frac{e^{i\omega X + \frac{\pi}{4}} - e^{-i\omega X - \frac{\pi}{4}}}{2i}\bigg) \bigg(\frac{e^{i\omega Y + \frac{\pi}{4}} - e^{-i\omega Y - \frac{\pi}{4}}}{2i}\bigg)\bigg] \\ &= -\frac{\lambda^2}{4} \bigg(e^{i\frac{\pi}{2}} \mathbb{E}[e^{i\omega(X+Y)}] + e^{-i\frac{\pi}{2}} \mathbb{E}[e^{-i\omega(X+Y)}] - \mathbb{E}[e^{i\omega(X-Y)}] - \mathbb{E}[e^{-i\omega(X-Y)}]\bigg) \\ &= \frac{\lambda^2}{2} \mathbb{E}[e^{i\omega(X-Y)}] \\ &= \frac{\lambda^2}{2} \exp\bigg\{ -\frac{1}{2}\omega^2 \frac{1-\beta^2}{n_0} \|x-y\|_2^2 \bigg\}. \end{split}$$

D.3 Computation of the NTK used for Figure 7 in the paper

In this section we show how one can derived analytically the function Φ_{∞} described in Section 7. This kind of computation can be used to derive numerically the filter radius $\hat{R}_{1/2}$ and tune the hyperparameters.

We use here a Gaussian embedding φ of size n_0 with lenghtscale ℓ followed by one hidden linear layer (activation function $x \to \sqrt{2} \max(0, x)$) of size n_1 and the output layer $n_2 = 1$. We also take $\alpha^2 + \beta^2 = 1$ in those experiments, to ensure constant variance of the neurons.

By the strong law of large numbers we have for the limiting NTK of the first layer:

$$\Theta^{1}_{\infty}(\varphi(p), \varphi(p')) = \beta^{2} + \frac{1 - \beta^{2}}{n_{0}} \varphi(p)^{T} \varphi(p') \underset{n_{0} \to \infty}{\longrightarrow} \beta^{2} + (1 - \beta^{2}) e^{-\frac{\|p - p'\|_{2}^{2}}{2l^{2}}} =: G(\|p - p'\|).$$

For the second layer, we use the notion of dual function defined above. In the case of the standardized ReLu it is computed in [9]:

$$\hat{r}(\rho) = \rho - \frac{\rho \arccos(\rho) - \sqrt{1 - \rho^2}}{\pi}, \quad \rho \in [-1, 1],$$

and:

$$\dot{\hat{r}}(\rho) = \dot{\hat{r}}(\rho) = 1 - \frac{\arccos(\rho)}{\pi}.$$

So that we can write, with d = ||p - p'||:

$$\Phi_{\infty}(d) = \hat{r}(G(d)) + G(d)\dot{\hat{r}}(G(d)).$$

Therefore Φ_{∞} only depends on ℓ and β . From this expression we can use standard Python libraries to approximate $\hat{R}_{1/2}$ for given values of the hyperparameters.

D.4 Computation of the NTK used for Figure 6 in the paper

Now we derive an approximate of the quantity $\hat{R}_{1/2}$ used in Figure 6 of the paper. This is a little bit more difficult than with Gaussian embedding because the rotation invariance is now only an approximation, even in the infinite-width limit.

With Torus embedding, we have $n_0=4$. The embedding is followed by two hidden linear layers with standardised cosine activation function, and then the last linear layer. We used here $r=\sqrt{2}$ $\delta=\frac{\pi}{80}$ (which is the formula suggested in the paper with $n_x=n_y=40$). As in the case of Gaussian embedding, we set $\alpha^2=1-\beta^2$. This ensures that neurons have constant variance and allows easy analytical computations.

Thanks to the Torus embedding described above, we get for the first layer:

$$\Theta_{\infty}^{1}(\varphi(p), \varphi(p')) = \beta^{2} + \frac{1 - \beta^{2}}{n_{0}} \varphi(p)^{T} \varphi(p')$$

$$= \beta^{2} + \frac{1 - \beta^{2}}{2} \left(\cos(\delta(p_{1} - p'_{1})) + \cos(\delta(p_{2} - p'_{2})) \right)$$

As rotation invariance is not analytically correct here, we look at the limiting NTK in the direction $p_1 = p_2$, which gives:

$$\Sigma^1(\varphi(p),\varphi(p')) = \Theta^1_\infty(\varphi(p),\varphi(p')) = \beta^2 + (1-\beta^2)\cos(\delta r)$$

with $r = |p_1 - p_1'| = |p_2 - p_2'|$.

For the next layers, we use the dual function of the standardised cosine (see [9]) given by:

$$\hat{\mu}(\rho) = \frac{\cosh(\omega^2 \rho)}{\cosh(\omega^2)},$$

and its derivative:

$$\hat{\mu}(\rho) = \omega^2 \frac{\sinh(\omega^2 \rho)}{\cosh(\omega^2)},$$

Then the limiting NTK is simply given by the following formulas:

$$\begin{split} &\Sigma^{l+1}(\varphi(p),\varphi(p')) = \beta^2 + (1-\beta^2)\hat{\mu}(\Sigma^l(\varphi(p),\varphi(p'))), \\ &\dot{\Sigma}^{l+1}(\varphi(p),\varphi(p')) = (1-\beta^2)\hat{\mu}(\dot{\Sigma}^l(\varphi(p),\varphi(p'))), \\ &\Theta^{l+1}_{\infty}(\varphi(p),\varphi(p')) = \Sigma^{l+1}(\varphi(p),\varphi(p')) + \dot{\Sigma}^{l+1}(\varphi(p),\varphi(p'))\Theta^l_{\infty}(\varphi(p),\varphi(p')). \end{split}$$

This way we construct a function $\Phi_{\infty}(r)$ with r an approximation of the radius and we can use it to compute numerically an approximation of $\hat{R}_{1/2}$ as before.

E Square root of the NTK in the case of random embedding

We now prove that we can define a notion of a square root of the NTK. First we need a technical lemma:

Lemma E.1. Let μ be a continuous function such that $\mathbb{E}_{X \sim \mathcal{N}(0,1)}[\mu(X)^2] = 1$, $C \in [0,1]$ a constant and $f \geq 0$ a positive definite function (in the sense of definition C.1) such that $C + f(p) \leq 1$. Then the function

$$F: p \longmapsto \hat{\mu}(C + f(p)) - \hat{\mu}(C),$$

is positive definite, where $\hat{\mu}$ denotes the dual function of μ (see definition D.1).

Proof:

Let us take $p_1, ..., p_m \in \mathbb{R}^d$ and $c_1, ..., c_m \in \mathbb{R}$. We introduce the Hermite expansion $\sum_k a_k h_k$ of μ and write its dual function as (see [9]):

$$\hat{\mu}(\rho) = \sum_{k=0}^{+\infty} a_k^2 \rho^k, \quad \rho \in [-1, 1],$$

Then by Bernoulli's formula:

$$\hat{\mu}(C + f(p_i - p_j)) - \hat{\mu}(C) = \sum_{k=1}^{+\infty} a_k^2 f(p_i - p_j) \sum_{s=0}^{k-1} C^{k-1-s} (C + f(p_i - p_j))^s.$$

Thus by polynomial combination with positive coefficients of positive semi-definite kernels:

$$\sum_{i,j=1}^{m} c_i c_j F(p_i - p_j) = \sum_{k=1}^{+\infty} \sum_{s=0}^{k-1} a_k^2 C^{k-1-s} \sum_{i,j=1}^{m} c_i c_j f(p_i - p_j) (C + f(p_i - p_j))^s \ge 0,$$

Which achieves the proof.

Let us recall the statement that we want to prove:

Proposition E.1 (Proposition 3.4 in the paper). Let φ be an embedding as described in section 3.2.2 of the paper, for a positive radial kernel $k \in L^1(\mathbb{R}^d)$ with k(0) = 1. Then there is a filter function $g : \mathbb{R} \to \mathbb{R}$ and a constant C such that for all p, p':

$$\lim_{n_0 \to \infty} \Theta_{\infty}(\varphi(p), \varphi(p')) = C + (g \star g)(p - p'), \tag{16}$$

where Θ_{∞} is the limiting NTK of a network with Lipschitz, non constant, and standardized activation function μ .

Before writing the proof, let us make some remarks on the assumptions of this proposition and their immediate implications:

- We recall that the fact that μ is "standardised" means here: $\mathbb{E}_{X \sim \mathcal{N}(0,1)}[\mu(X)^2] = 1$.
- As mentioned before (Section 2.3 of the paper) we assume for simplicity that $\alpha^2 = 1 \beta^2$ to ensure constant variance of the neurons (we consider $\beta \in [0, 1)$).
- We denote by A the Lipschitz constant of μ . By Rademacher theorem, we know that μ is almost everywhere differentiable and $\|\dot{\mu}\|_{\infty} \leq A$. The fact that μ is not constant ensures that $\hat{\mu}$ is (strictly) increasing on [0,1).

- Moreover, the Lipschitz assumption also implies that $|\hat{\mu}(1)| \le A^2 < +\infty$ and therefore $\hat{\mu}$ is continuous on [-1,1] by Abel's theorem on entire series.
- The procedure to approximate the kernel k in Section 3.2.2 of the paper assumes that k is
 continuous (to be able to apply Bochner's theorem). It is therefore also the case in this proof.

Proof of the proposition:

Step 1: We want to show by recursion that for all $l \ge 1$ there exists some constant $C_l \in [0, 1)$ such that for all $p, p' \in \mathbb{R}^d$ we have in probability:

$$\Sigma^{l}(\varphi(p), \varphi(p')) \underset{n_0 \to \infty}{\longrightarrow} C_l + f_l(p - p'),$$
 (17)

With f_l a radial positive definite function such that $f_l \ge 0$ and $f_l \in L^1(\mathbb{R}^d)$.

For l = 1, we know that this is true by the law of large numbers:

$$\Sigma^{1}(\varphi(p), \varphi(p')) = \Theta^{1}_{\infty}(\varphi(p), \varphi(p')) = \beta^{2} + \frac{1 - \beta^{2}}{n_{0}} \varphi(p)^{T} \varphi(p')$$

$$\underset{n_{0} \to \infty}{\longrightarrow} \beta^{2} + (1 - \beta^{2})k(p - p'),$$
(18)

We just set $f_1 = (1 - \beta^2)k$. Now we assume $l \ge 2$:

We have by our normalisation assumptions $\Sigma^l(\varphi(p), \varphi(p)) = C_l + f_l(0) = 1$. Using the continuity of $\hat{\mu}$ (see [9] for the properties of $\hat{\mu}$), we have:

$$\Sigma^{l+1}(\varphi(p), \varphi(p')) = \beta^2 + (1 - \beta^2)\hat{\mu}(\Sigma^l(\varphi(p), \varphi(p')))$$

$$\underset{n_0 \to \infty}{\longrightarrow} \beta^2 + (1 - \beta^2)\hat{\mu}(C_l + f_l(p - p')). \tag{19}$$

Using properties of the dual function given in [9], we know that $\hat{\mu}$ is positive, increasing and convex in [0, 1]. Moreover as f_l is radial positive definite we have $f_l \leq f_l(0) = 1 - C_l$. Then by convexity:

$$\hat{\mu}(C_l + f_l(p - p')) = \hat{\mu}\left(\frac{f_l(p - p')}{1 - C_l} + \left(1 - \frac{f_l(p - p')}{1 - C_l}\right)C_l\right)$$

$$\leq \frac{f_l(p - p')}{1 - C_l}\hat{\mu}(1) + \left(1 - \frac{f_l(p - p')}{1 - C_l}\right)\hat{\mu}(C_l).$$

Using that $\hat{\mu}$ is increasing:

$$|\hat{\mu}(C_l + f_l(p - p')) - \hat{\mu}(C_l)| \le \frac{\hat{\mu}(1) - \hat{\mu}(C_l)}{1 - C_l} f_l(p - p'),$$

So that we can rewrite equation 19 in the following form:

$$\Sigma^{l+1}(\varphi(p),\varphi(p')) \underset{n_0 \to \infty}{\longrightarrow} \beta^2 + (1-\beta^2)\hat{\mu}(C_l) + f_{l+1}(p-p'),$$

With
$$f_{l+1}(p-p') = (1-\beta^2)(\hat{\mu}(C_l + f_l(p-p')) - \hat{\mu}(C_l))$$
 and $C_{l+1} = \beta^2 + (1-\beta^2)\hat{\mu}(C_l)$.

The previous inequality, lemma E.1 and the fact that $\hat{\mu}$ is increasing in [0, 1) ensure the properties of f_{l+1} and C_{l+1} .

Step 2: As $\hat{\mu}$ is also positive, continuous, increasing and convex in [0,1], we can obtain a convergence in probability similar to equation 17 but for $\dot{\Sigma}^l$:

$$\dot{\Sigma}^l(\varphi(p),\varphi(p')) \xrightarrow{n \to \infty} B_l + h_l(p-p'),$$

With $B_l \geq 0$, and h_l a positive definite function such that $h_l \in L^1(\mathbb{R}^d)$ and $h_l \geq 0$.

Now we want to show by recursion that for a fixed l:

$$\forall p, p' \in \mathbb{R}^d, \ \Theta_{\infty}^l(\varphi(p), \varphi(p')) \xrightarrow[n_0 \to \infty]{} C_{\mu,\beta,l} + \theta_l(p - p'). \tag{20}$$

With θ_l a positive definite function such that $\theta_l \in L^1(\mathbb{R}^d)$ and $C_{\mu,\beta,l} \geq 0$. Again we know that this is true for l=1 by equation 18.

We have:

$$\Theta_{\infty}^{l+1}(\varphi(p),\varphi(p')) \underset{n_0 \to \infty}{\longrightarrow} (C_{\mu,\beta,l} + \theta_l(p-p')) \dot{\Sigma}^{(l+1)}(p,p') + C_{l+1} + f_{l+1}(p-p').$$

So that we can set:

$$\theta_{l+1}(\varphi(p), \varphi(p')) = C_{\mu,\beta,l} h_{l+1}(p-p') + \theta_l(p-p') \dot{\Sigma}^{l+1}(\varphi(p), \varphi(p')) + f_{l+1}(p-p'),$$

and:

$$C_{\beta,\mu,l+1} = C_{l+1} + C_{\beta,\mu,l}B_l.$$

Using that $|\theta_l(p-p')\dot{\Sigma}^{l+1}(\varphi(p),\varphi(p'))| \leq A^2|\theta_l(p-p')|$ and all the previous results, the recursion works automatically and we have equation 20 for all $l \geq 2$.

Moreover $(p, p') \longmapsto \theta_l(p - p')\dot{\Sigma}^{1+l}(p, p')$ is positive semi-definite as a product of two positive semi-definite kernels. By sum we deduce that θ_{l+1} is positive semi-definite and by recursion we have the result for all θ_l .

Step 3: Now, using integrability of θ_l , we know that its Fourier transform defines a function $q \in L^{\infty}(\mathbb{R}^d)$.

From dominated convergence theorem we deduce that q is continuous.

Therefore in the sense of distributions, the Fourier transform of θ_L is given by a finite positive measure (Bochner's theorem) and also by $q \in L^{\infty}(\mathbb{R}^d)$. We deduce that q is the density of this finite positive measure (the Radon-Nikodym derivative with respect to the Lebesgue measure).

From those arguments we get $q \ge 0$ and $q \in L^1(\mathbb{R}^d)$. We then have the Fourier inversion formula for θ_L :

$$\theta_L(p-p') = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} q(\omega) e^{i\omega \cdot (p-p')} d\omega, \quad \text{with: } q \ge 0$$

Hence it makes sense to define:

$$g = \mathcal{F}^{-1}(\sqrt{q}),$$

In the sense of the Fourier transform of a L^2 function. Then the convolution theorem ensures:

$$\theta_L = q \star q$$
.

Remark: Here we used lemma E.1 and the dual activation function to show that both f_l and θ_l are positive definite. If we only show that $\theta_l \in L^1(\mathbb{R}^d)$ it is still possible to show the same properties of the function q by using positive definites of $C + \theta_L$ and take the Fourier transform in the sense of distributions, which leads to $(2\pi)^d C \delta_0 + q = (2\pi)^d M$ with M a finite positive measure. Then arguments based on test functions and the continuity of q give the result. The advantage of lemma E.1 is that it is a bit more general.

F Additional experimental results

F.1 Plots of the Neural Tangent Kernel

Here are some additional experimental results regarding the comparison between the theoretical (limiting) NTK $\tilde{\Theta}_{\infty}$ and the empirical NTK $\tilde{\Theta}_{\theta(t)}$. Here again the "lines" of the Gram matrices are reshaped as images.

Figure 10 represents the comparison between the limiting NTK and the emprirical NTK with a Gaussian embedding. We can observe that the infinite-width limit seems to be well-respected.

Figure 11 shows the evolution of the NTK during the optimisation process. While the NTK begins to change at the end of training (it is due to the alignment of descent directions, because of the sigmoid we use to control the volume, pre-densities $(x_i)_{1 \leq i \leq N}$ tend to infinity) the NTK stays close to Θ_{∞} during the part of training where the final shape is created. This justifies even more that it is pertinent to study the effect of the NTK on the final geometry.

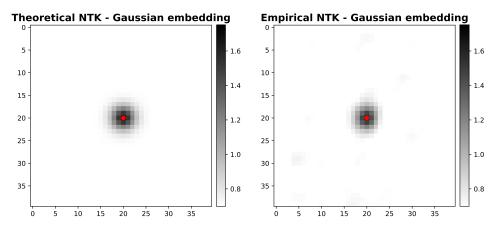


Figure 10: Comparison between one line of the Gram matrix of the empirical NTK $\tilde{\Theta}_{\theta(t)}$ and and of the corresponding limiting NTK $\tilde{\Theta}_{\infty}$. Here we use a Gaussian embedding as described in the paper

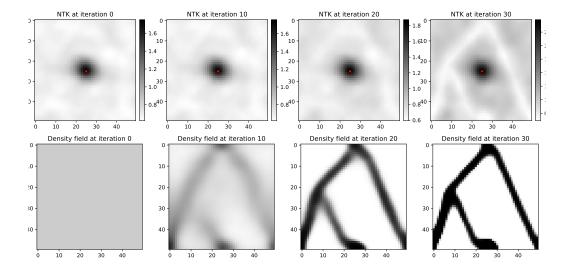


Figure 11: Evolution of the NTK of a network with a Gaussian embedding with hyperparameters as described in Section 4.1. We can see a relative stability of the NTK