

Intrinsic Dimensionality and Neural Networks

Dominique Barbe

The University of Melbourne, Victoria 3010, Australia

Abstract. Neural networks exhibit impressive performances, yet their nature is still poorly understood. We explored two directions in an attempt to increase the knowledge about them. We ran experiments, looking at how the dimensionality of the data changes through the network. And we defined an infinite model, more suited to theoretical analysis.

Keywords: neural networks, intrinsic dimensionality

1 Introduction

This document is an attempt to summarize the work realized at the University of Melbourne, under the supervision of Prof. James Bailey, from the 30/01/2017 to the 30/06/2017. I will go through the experiments I did, the theoretical model for neural networks developed (with the help of Michael E. Houle), and the intuition that guided me here and there.

Should you have any question regarding this document (or anything I worked on, like the code), feel free to contact me at `dominique.barbe@ens-rennes.fr`.

2 Preliminaries

This document assumes some familiarity with neural networks and the Local Intrinsic Dimensionality. We give a quick recap of these two notions here.

2.1 An overview of feedforward neural networks

Consider the following quantities:

- Let $(n, m) \in \mathbb{N}^2$.
- Let $\mathbf{W} \in \mathbb{R}^{m \times n}$ be the *weight matrix* (or *weights*).
- Let $\mathbf{b} \in \mathbb{R}^m$ be the *bias vector*.
- Let $a : \mathbb{R} \rightarrow \mathbb{R}$ be the *activation function*; we also note a the component-wise application of a .

Consider the following function:

$$\begin{aligned} \mathbb{R}^n &\rightarrow \mathbb{R}^m \\ f : \mathbf{x} &\mapsto a(\mathbf{W} \cdot \mathbf{x} + \mathbf{b}) \end{aligned} \tag{1}$$

We call the function f a layer; it is parameterized by its weight matrix \mathbf{W} and its bias vector \mathbf{b} . A *k-layers feedforward neural network* is a function defined by the composition of k layers.

For a more detailed introduction to neural networks, one can refer to [14].

2.2 The Intrinsic Dimensionality

Definition 1 ([8]). Let \mathbf{X} be a random distance variable. For any r such that $F_{\mathbf{X}}(r) > 0$, the intrinsic dimensionality (ID) of \mathbf{X} at r is given by

$$\text{ID}_{F_{\mathbf{X}}}(r) \triangleq \lim_{\epsilon \rightarrow 0^+} \frac{\ln F_{\mathbf{X}}((1+\epsilon)r) - \ln F_{\mathbf{X}}(r)}{\ln((1+\epsilon)r) - \ln r} = \frac{r \cdot F'_{\mathbf{X}}(r)}{F_{\mathbf{X}}(r)}, \quad (2)$$

wherever the limit exists. The second equality follows by applying l'Hôpital's rule to the limits, provided that $F_{\mathbf{X}}$ is positive and differentiable over an open interval containing r .

Under this distributional interpretation, the original data set determines a sample of distances from a given point. The intrinsic dimensionality of this distance distribution $F_{\mathbf{X}}$ is estimated. The definition of $\text{ID}_{F_{\mathbf{X}}}$ can be extended to the case where $r = 0$, by taking the limit of $\text{ID}_{F_{\mathbf{X}}}(r)$ as $r \rightarrow 0^+$, whenever this limit exists:

$$\text{ID}_{F_{\mathbf{X}}}(0) \triangleq \lim_{r \rightarrow 0^+} \text{ID}_{F_{\mathbf{X}}}(r). \quad (3)$$

The intuition and the formal definition we gave involved distances distributions. However, this definition can be extended to any multivariate real-valued function that is non-zero over some open ball around x (the point where the ID is computed).

An important result is the Local ID Representation theorem, which states conditions under which a function F is fully characterized by ID_F .

Local ID Representation Theorem

Theorem 1 ([7]). Let $F : \mathbb{R} \rightarrow \mathbb{R}$ be a real-valued function. Let $v \in \mathbb{R}$ be a value for which $\text{ID}_F(v)$ exists. Let $x, w \in \mathbb{R}$ be values for which the quantities x/w and $F(x)/F(w)$ exist and are positive. Then if F is non-zero and continuously differentiable everywhere over $[\min\{x, w\}, \max\{x, w\}]$, except perhaps at v itself,

$$\frac{F(x)}{F(w)} = \left(\frac{x}{w}\right)^{\text{ID}_F(v)} \cdot G_{F,v,w}(x), \quad \text{where} \quad (4)$$

$$G_{F,v,w}(x) \triangleq \exp\left(\int_x^w \frac{\text{ID}_F(v) - \text{ID}_F(t)}{t} dt\right), \quad (5)$$

whenever the integral exists.

Moreover, let $c > 1$ be a constant. We have:

– if $v \neq 0$,

$$\lim_{\substack{w \rightarrow 0 \\ |x-v| \leq c|w-v|}} G_{F,0,w}(r) = 1 \quad (6)$$

– if $v = 0$,

$$\lim_{\substack{w \rightarrow 0^+ \\ 0 < w/c \leq r \leq cw}} G_{F,0,w}(r) = \lim_{\substack{w \rightarrow 0^- \\ 0 < w/c \leq r \leq cw}} G_{F,0,w}(r) = 1 \quad (7)$$

Informally, given a function F and three values, a , b and c relatively close to each others. We have:

$$\frac{F(a)}{F(b)} \simeq \left(\frac{a}{b}\right)^{\text{ID}_F(c)} \quad (8)$$

In other words, we can switch from a reasoning on the value of F to simply ratios of inputs a/b , the ID doing the link between the two. The notion of intrinsic dimensionality is analogous to the derivative, but allows working with quotients instead of differences:

$$F(a) - F(b) \simeq (a - b) \cdot F'(c) \quad (9)$$

3 Experiments with networks

3.1 Introduction

Most of the internship consisted in a serie of experiments, trying to measure the complexity of the data (using the ID) at different steps through the network and at different moments during training. None of the experiments gave us new, interesting insights. We will go through most of them.

3.2 Experimental setup

An auto-encoder is a neural network that tries to reproduce its input (in other words, it tries to approximate the identity function). It consists of two part: the encoder — a first serie of layers of decreasing size — and the decoder — a serie of layers of increasing size.

Most of the experiments use auto-encoders, as these networks do not require labels for the data. They try to find a generic model for their training set (generic because their is other goal than reconstructing the dataset). How they model the data is what we want to analyze. We wanted to see how much a network would simplify the data it receives: does the data reconstructed by the network has a low or a high dimensionality?

There are two kind of plots: LID distributions, and mean LID evolution. Most of the experiments involve the following process: reconstructing the dataset (i.e taking its image by the network), and estimating the LID for all/sampled points. The LIDs collected are either plotted as a distribution, or averaged and plotted (against the training epoch most of the time).

We use the MNIST [1] and CIFAR-10 [2] datasets for all experiments.

3.3 Output complexity

We compare different autoencoders. For each of them, we take the image of a dataset by a network, and compute the LID distribution of the transformed dataset. The LID is estimated using the MLE ($k = 100$ neighbours are used). It is computed for every point in the training set, and the probability distribution

function of the distribution of the LID is plotted. The results are shown in Figure 1 and 2.

The following networks are used:

- The networks named "D x " are networks using one dense hidden layer, with x hidden neurons. Rectified linear units are used for the activation in the hidden layer, and sigmoid units for the output.
- The networks named "C x " are convolutional neural network, using x features maps of 5 by 5 pixels. Rectified linear units are used for the activation, followed by a max pooling layer of size 2 by 2. The output layer also uses sigmoid functions for the activation.

Remark: I didn't spend that much time tweaking the training of the networks. Their performances may be good, but still not exploit to the fullest their individual architectures. This can explain why, for example, the networks D8 and D16 perform similarly on the MNIST dataset, despite one of them being more powerful.

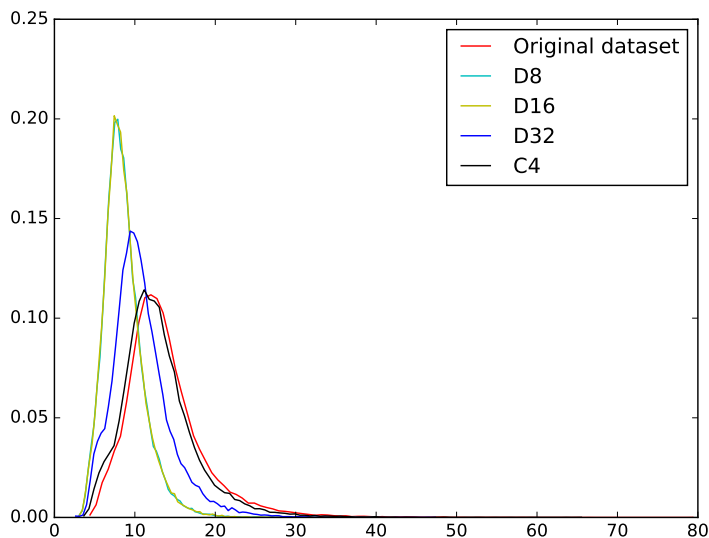


Fig. 1: LID distribution for the MNIST dataset

3.4 Output complexity during training

We examined the evolution of the output complexity during training. For this series of experiments we plotted the training accuracy and the mean LID against the training epoch.

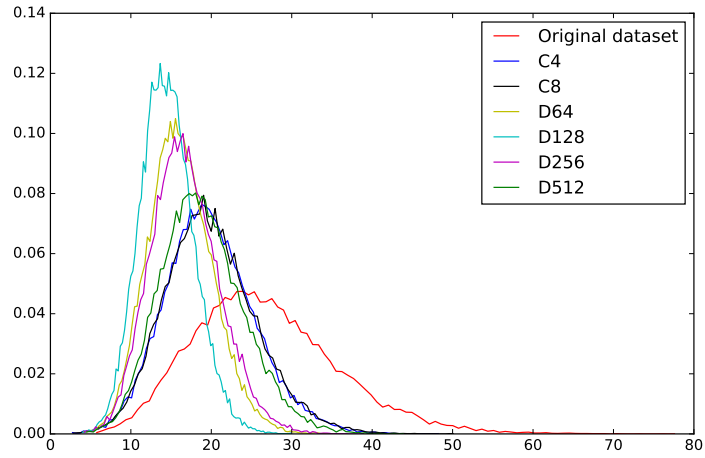


Fig. 2: LID distribution for the CIFAR-10 dataset

A first experiment study the short-term dynamic. Here, a 2-layers auto-encoder with 32 hidden units is used. The activation functions are ReLU for the hidden layer and sigmoid for the output layer (for re-scaling purpose). The ID is computed on the test set, estimated using 1000 samples with the Maximum Likelihood Estimation [9]. The learning rate is set low (.002), as convergence is not more important here than the dynamic. Results are shown on Figure 3.

The second experiment studies the long-term dynamic. Again, a 2-layer auto-encoder is used, with 32 hidden units for MNIST and 64 for CIFAR. The learning rates are respectively $lr = 0.005$ and $lr = 0.001$. Results are shown on Figure 4.

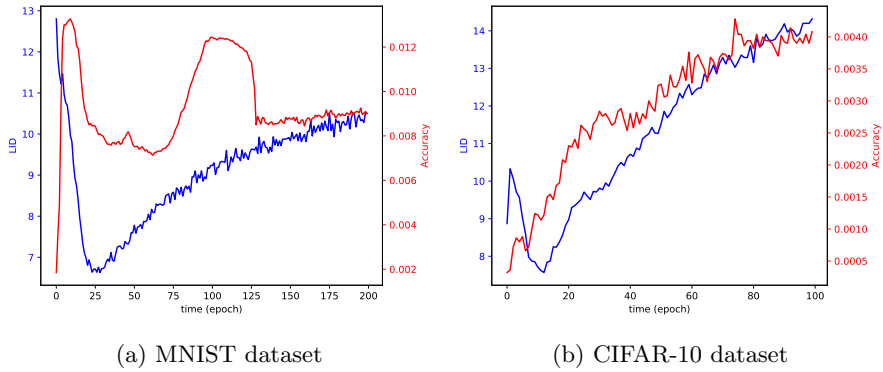


Fig. 3: Accuracy and output ID plotted against the training epoch

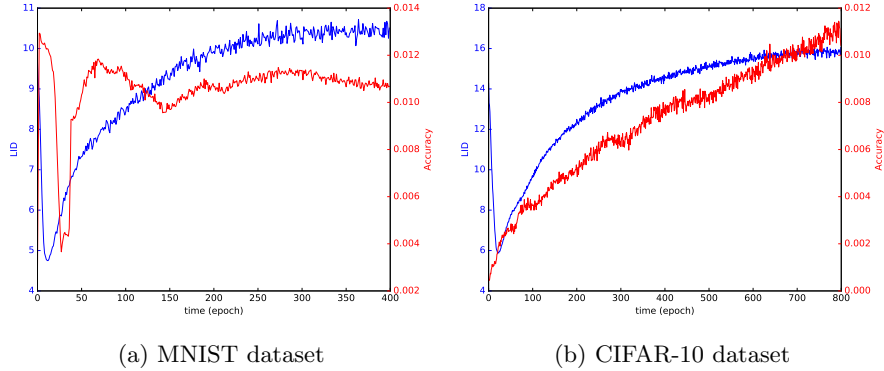


Fig. 4: Accuracy and output ID plotted against the training epoch

These plots show a decrease followed by an increase in output ID. Interpreting the ID as a complexity measure, these plots can be interpreted as such: the parameters being random at first, the data undergo random transformations, and thus the output complexity starts high. As patterns are learned, the network fits a simple model to the data, decreasing the output complexity. Finally, finer details are learned; the approximation getting better, the complexity increases again, getting closer to the one of the original dataset.

4 Experiments with neurons

All the networks studied here have the same architecture unless told otherwise:

- feedforward autoencoders with one hidden layer.
- the hidden layer uses rectified linear units.
- the output layer uses sigmoid units.

Three datasets are used: MNIST [1], Cifar-10 [2] and Hasy [3]. All three are images sets. For each dataset, all but 10000 images are used for the training set, and the remaining 10000 for the validation set.

4.1 Introduction

Rectified linear units allow for an easy partition of a dataset in two: either an input has a null activation (**low-activation set**) or a strictly positive activation (**high activation set**). We want to investigate the properties of these two sets, for a given neuron in a given network.

First, we ran the following experiment: we took an autoencoder trained on the MNIST dataset. For each neuron, we partitioned the validation set in two (using the previously described method, that is computing the high and low-activation sets). We then computed the following metrics:

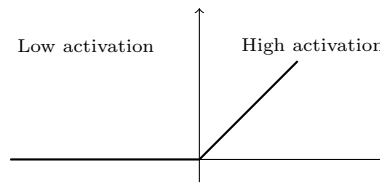


Fig. 5: Rectified Linear Unit

- size of the two sets.
- dominant class of the two sets.
- proportion of the dominant class.
- mean LID of the two sets (using ALIDv4, $k=12$, sampled on 500 points).

Note: the LID of the original images is computed. The LID is computed for each set independently (that is: neighbours are drawn from the same set, not the whole dataset).

Below are some examples of such metrics for a few neurons (taken from the same network):

- Neuron (a):
 - low-activation set: 159 images. Class 6 is dominant (49% of the images). LID = 11.87.
 - high-activation set: 9841 images Class 1 is dominant (11% of the images). LID = 14.91.
- Neuron (b):
 - low-activation set: 4796 images. Class 3 is dominant (17% of the images). LID = 15.55.
 - high-activation set: 5204 images Class 6 is dominant (18% of the images). LID = 14.14.
- Neuron (c):
 - low-activation set: 9612 images. Class 1 is dominant (11% of the images). LID = 15.37.
 - high-activation set: 388 images Class 7 is dominant (36% of the images). LID = 12.44.

We observed the following correlations:

- The two sets can be unbalanced. One is huge and the other is small. The huge set has high LID (close to the one of the original dataset); the small set has low ID. The huge set doesn't have a clear dominant class (often 11%-12% of its images for 10 classes); the small set has a clear dominant class.
- The two sets can be balanced. They then tend to have the same LID, and a dominant class of the same proportion.

From here, we concluded that:

- Neurons can be classified into three categories:

- *normal neurons* are unbalanced, having a small low-activation set and a huge high-activation one (They are the most common, hence the name).
 - *balanced neurons* are balanced. Their two sets are of the same size.
 - *inverted neurons* are unbalanced, having a huge low-activation set and a small high-activation one.
- Some neurons clearly belong to one of those three categories; some are harder to classify (neither really unbalanced nor balanced).
 - From now on, we will only compute the size of the high-activation set of the different neurons we analyze. The empiric correlations described earlier show that it is enough to determine the type of a given neuron.

We then tried to answer the following question: what has an influence on the distribution of the type of neurons? We explored three parameters: the initialization of the weights, the size of the hidden layer and the training database.

4.2 Influence of the initialization

We trained 7 different networks using the same architecture (256 hidden neurons), the same training set (MNIST here) and the same training process (i.e. number of epochs and learning rate). The only difference between them are the random weight initialization.

For each network, we compute the size of the high-activation set for each of its neuron. Then, the distribution of these sizes is plotted. The results are displayed on Figure 6.

The initialization doesn't seem to have a significant effect on the distribution of the type of neurons. Further experiments with the other 2 datasets confirmed this trend (the plots are not shown). Here, for every network trained, most neurons are *normal*, a small number are *inverted*, and almost none are *balanced*.

4.3 Influence of the size of the hidden layer

We trained 5 different networks, using between 64 and 1024 hidden neurons (still with the MNIST dataset). The results are displayed on Figure 7.

Note: the input dimension is 784 here; 1024 is already beyond that.

Again, the number of hidden units doesn't seem to have a significant effect on the distribution of the type of neurons.

Still, networks using a high number of neurons may have a smoother distribution. Maybe this is an artifact of the discretization required to trace the plot. Maybe there is a limit distribution when the number of hidden neurons tends toward infinity. We may investigate this question later, but we postponed it, as auto-encoder with too many hidden units do not make sense.

4.4 Influence of the dataset

We trained networks using 256 hidden neurons on the three datasets (MNIST, Cifar-10 and Hasy). Again, the distribution of the sizes of the high-activation sets are plotted together. The results are displayed on Figure 8.

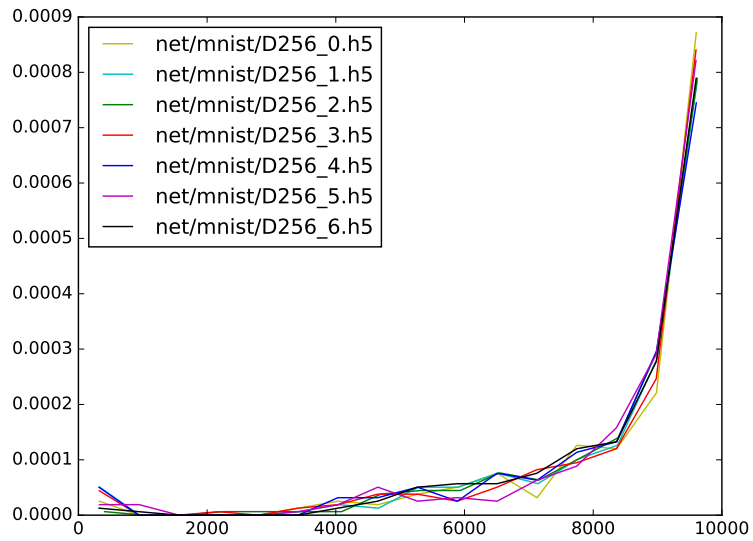


Fig. 6: Comparison with 7 different initialization on the repartition of the neuron types

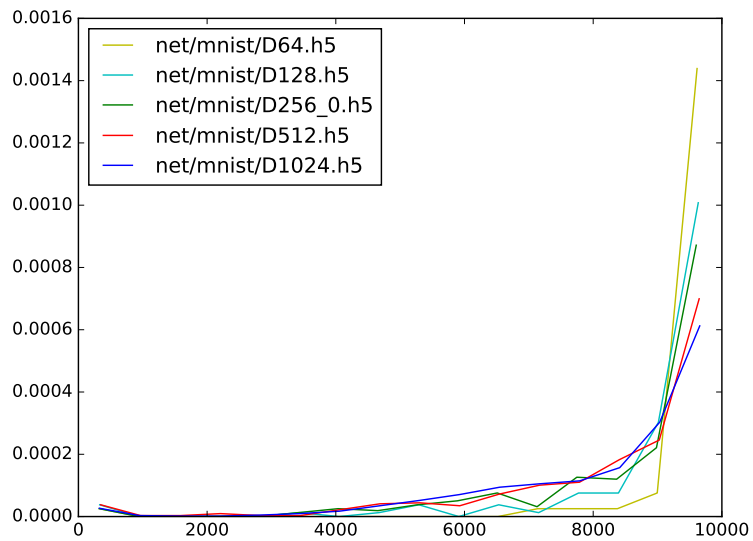


Fig. 7: Comparison with 5 different number of hidden neurons on the repartition of the neuron types

Remark: The distributions are averaged over 7 different initialization. We have seen that the initialization doesn't have a significant impact on the type of the neurons; we can take advantage of this fact to obtain more precise plots by using many networks to analyse more neurons for a given architecture/database.

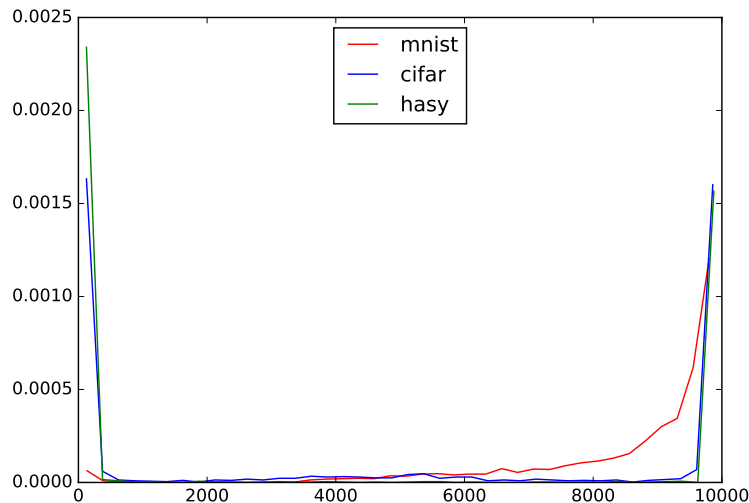


Fig. 8: Comparison with 3 different datasets on the repartition of the neuron types

The dataset has a large impact on the distribution of the type of neurons. MNIST produces mostly *normal* neurons, where Hasy creates mostly *inverted* neurons. Cifar-10 lying in between the two.

5 Some results on ID

Before moving on with the infinite model for neural network, we give some mildly interesting results on ID

5.1 Composition of functions

We give a set of theorems on the behaviour of ID when composing (multivariate) functions. Let $g : \mathbb{R}^n \mapsto \mathbb{R}^m$ and $\mathbf{x} \in \mathbb{R}^n$. We note:

- x_i the i th component of \mathbf{x} .
- $g_j : \mathbb{R}^n \mapsto \mathbb{R}$ the function obtained from g by taking only its j th output.

- ${}_i g : \mathbb{R} \mapsto \mathbb{R}^m$ the function obtained from g obtained by treating all its inputs as constants (equal to the coordinates of \mathbf{x}), except the i th one. We note that $g(\mathbf{x}) = {}_i g(x_i)$.
- ${}_i g_j : \mathbb{R} \mapsto \mathbb{R}$ the function obtained from ${}_i g$ by taking only its j th output. We note that $g_j(\mathbf{x}) = {}_i g_j(x_i)$.

When composing functions, we have the following relations on ID:

Theorem 2. Let $\mathbb{R} \xrightarrow{g} \mathbb{R} \xrightarrow{f} \mathbb{R}$. We have:

$$\text{ID}_{f \circ g}(x) = \text{ID}_g(x) \cdot \text{ID}_f(g(x)) \quad (10)$$

Theorem 3. Let $\mathbb{R}^n \xrightarrow{g} \mathbb{R} \xrightarrow{f} \mathbb{R}$. We have:

$$\text{ID}_{f \circ g}(\mathbf{x}) = \text{ID}_g(\mathbf{x}) \cdot \text{ID}_f(g(\mathbf{x})) \quad (11)$$

Theorem 4. Let $\mathbb{R} \xrightarrow{g} \mathbb{R}^m \xrightarrow{f} \mathbb{R}$. We have:

$$\text{ID}_{f \circ g}(x) = \sum_{j=1}^m \left[\text{ID}_{g_j}(x) \cdot \text{ID}_{j f}(g_j(x)) \right] \quad (12)$$

Theorem 5. Let $\mathbb{R} \xrightarrow{g} \mathbb{R}^m \xrightarrow{f} \mathbb{R}$. We have:

$$\text{ID}_{f \circ g}(\mathbf{x}) = \sum_{j=1}^m \text{ID}_{j f \circ g_j}(\mathbf{x}) \quad (13)$$

Proof. (Theorem 2) For $\mathbb{R} \xrightarrow{g} \mathbb{R} \xrightarrow{f} \mathbb{R}$:

$$\begin{aligned} \text{ID}_{f \circ g}(x) &= \frac{x \cdot (f \circ g)'(x)}{(f \circ g)(x)} \\ &= \frac{x \cdot g'(x) \cdot g(x) \cdot f'(g(x))}{g(x) \cdot f(g(x))} \\ &= \text{ID}_g(x) \text{ID}_f(g(x)) \end{aligned}$$

(Theorem 3) For $\mathbb{R}^n \xrightarrow{g} \mathbb{R} \xrightarrow{f} \mathbb{R}$:

$$\begin{aligned} \text{ID}_{f \circ g}(\mathbf{x}) &= \frac{\mathbf{x} \cdot \nabla(f \circ g)(\mathbf{x})}{(f \circ g)(\mathbf{x})} \\ &= \sum_{i=1}^n \left[\frac{x_i \cdot {}_i g'(x_i) \cdot f'({}_i g(x_i))}{f(g(\mathbf{x}))} \right] \\ &= \sum_{i=1}^n \left[\frac{x_i \cdot {}_i g'(x_i)}{{}_i g(x_i)} \cdot \frac{{}_i g(x_i) \cdot f'({}_i g(x_i))}{f(g(\mathbf{x}))} \right] \\ &= \sum_{i=1}^n \left[\text{ID}_{{}_i g}(x_i) \cdot \frac{g(\mathbf{x}) \cdot f'(g(\mathbf{x}))}{f(g(\mathbf{x}))} \right] \\ \text{ID}_{f \circ g}(\mathbf{x}) &= \text{ID}_g(\mathbf{x}) \text{ID}_f(g(\mathbf{x})) \end{aligned}$$

(Theorem 4) For $\mathbb{R} \xrightarrow{g} \mathbb{R}^m \xrightarrow{f} \mathbb{R}$:

$$\begin{aligned}
\text{ID}_{f \circ g}(x) &= \frac{x \cdot (f \circ g)'(x)}{(f \circ g)(x)} \\
&= \sum_{j=1}^n \frac{x \cdot g'_j(x) \cdot {}_j f(x)}{f(g(x))} \\
&= \sum_{j=1}^n \left[\frac{x \cdot g'_j(x)}{g_j(x)} \cdot \frac{{}_j f'(g_j(x))}{f(g(x))} \right] \\
&= \sum_{j=1}^n \left[\text{ID}_{g_j}(x) \cdot \frac{{}_j f'(g_j(x))}{f(g(x))} \right] \\
\text{ID}_{f \circ g}(x) &= \sum_{j=1}^n \left[\text{ID}_{g_j}(x) \cdot \text{ID}_{{}_j f}(g_j(x)) \right]
\end{aligned}$$

(Theorem 5) For $\mathbb{R}^n \xrightarrow{g} \mathbb{R}^m \xrightarrow{f} \mathbb{R}$:

$$\begin{aligned}
\text{ID}_{f \circ g}(\mathbf{x}) &= \sum_{i=1}^n \text{ID}_{f \circ_i g}(x_i) \\
&= \sum_{i=1}^n \sum_{j=1}^n \left[\text{ID}_{i g_j}(x_i) \cdot \text{ID}_{{}_j f}(g_j(x_i)) \right] \\
&= \sum_{j=1}^n \sum_{i=1}^n \left[\text{ID}_{i g_j}(x_i) \cdot \text{ID}_{{}_j f}(g_j(x_i)) \right] \\
&= \sum_{j=1}^n \sum_{i=1}^n \left[\text{ID}_{i g_j}(x_i) \cdot \text{ID}_{{}_j f}(g_j(\mathbf{x})) \right] \\
&= \sum_{j=1}^n \left[\text{ID}_{g_j}(\mathbf{x}) \cdot \text{ID}_{{}_j f}(g_j(\mathbf{x})) \right] \\
\text{ID}_{f \circ g}(\mathbf{x}) &= \sum_{j=1}^n \text{ID}_{{}_j f \circ g_j}(\mathbf{x})
\end{aligned}$$

Notes on Theorem 5: Theorem 5 give a non-trivial (yet simple to prove) decomposition of the ID of the composition of two functions. The proof involves breaking the overall path taken by the inputs into every individual possible path, and recomposing them.

We can plot the interplay between the different coordinates of $\mathbb{R} \xrightarrow{g} \mathbb{R}^m \xrightarrow{f} \mathbb{R}$ (see Figure 9). The first decomposition is done along each input coordinate. Then the sum is decomposed again using Theorem 4. Finally, it is recomposed using Theorem 3.

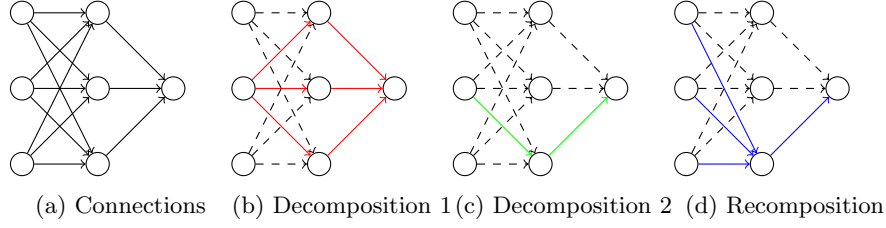


Fig. 9: Paths involved in the proof of Theorem 5.

5.2 The weighted ID of common functions

We try to find explicit formulas for the wID of common transformations involved in neural networks (linear combinations and activation functions).

Linear combinations

Theorem 6. Let $(a, b) \in \mathbb{R}^2$. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ $f : x \mapsto a \cdot x + b$. Then:

$$\text{wID}_f(x) = a \cdot x \quad (14)$$

Theorem 7. Let $\mathbf{a} \in \mathbb{R}^n$ and $b \in \mathbb{R}$. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ $f : \mathbf{x} \mapsto \mathbf{a} \cdot \mathbf{x} + b$. Then:

$$\text{wID}_f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{a} \quad (15)$$

$$\text{wID}_j f(x) = x \cdot a_j \quad (16)$$

Sigmoid

Theorem 8. Let σ be the sigmoid function: $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ $\sigma : x \mapsto \frac{e^x}{1+e^x}$. The sigmoid function has the two following properties:

$$\sigma(x) = 1 - \sigma(-x) \quad (17)$$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)) \quad (18)$$

Its weighted ID is:

$$\text{wID}_\sigma(x) = x \cdot \sigma(x) \cdot \sigma(-x) \quad (19)$$

Proof. (First property)

$$1 - \sigma(-x) = \frac{1 + e^{-x}}{1 + e^{-x}} - \frac{e^{-x}}{1 + e^{-x}} = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x} = \sigma(x)$$

(Second property)

$$\begin{aligned}\sigma'(x) &= \frac{e^x \cdot (1 + e^x) - e^x \cdot e^x}{(1 + e^x)^2} = \frac{e^x}{1 + e^x} - \left(\frac{e^x}{1 + e^x} \right)^2 \\ &= \sigma(x) - \sigma(x)^2 = \sigma(x) \cdot (1 - \sigma(x))\end{aligned}$$

(Weighted ID)

$$\text{wID}_\sigma(x) = x \cdot \sigma'(x) = x \cdot \sigma(x) \cdot (1 - \sigma(x)) = x \cdot \sigma(x) \cdot \sigma(-x)$$

Remark: The weighted ID of the sigmoid function has the following shape:

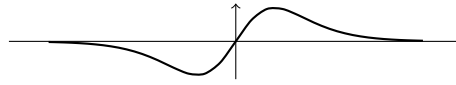


Fig. 10: The weighted ID of the sigmoid function

ReLU

Theorem 9. Let $\text{ReLU} : x \mapsto \max(0, x)$ be the Rectified Linear Unit function. Its weighted ID is:

$$\text{wID}_{\text{ReLU}}(x) = \text{ReLU}(x) \tag{20}$$

Proof. Let $x \in \mathbb{R}^*$.

$$\text{wID}_{\text{ReLU}}(x) = x \cdot \text{ReLU}'(x) = \begin{cases} x \cdot 1 & \text{if } x > 0 \\ x \cdot 0 & \text{if } x < 0 \end{cases} = \text{ReLU}(x)$$

6 A model for infinite neural networks

6.1 Introduction

The universal approximation theorem states that a feedforward neural network with one hidden layer can approximate any continuous function on compacts of \mathbb{R}^n , as the number of hidden neurons tends to infinity [15]. This theorem is our main reason to consider infinite networks, but other results push us in the same directions.

Chiyuan Zhang et al. [13] showed through experiments that deep neural networks easily fit random noise. Yet they are also able to generalize properly from data that make sense. The networks they manipulate are over-parameterized. Despite this fact, they avoid overfitting. This suggests that the capacity (i.e the number of parameters) is not the proper metric to measure the expressive power of a network. This quantity could then be abstracted by considering infinite networks. Moreover, even when neural networks achieve a perfect accuracy on their

training sets, their accuracy at test time still increases as hidden neurons are added. In other words, even when the training examples are perfectly recognized there is still something to learn.

Our intuition is that real (i.e finite) networks are an approximation of infinite objects: neural networks with an infinite number of neurons at every hidden layer. Real networks are the result of a sampling process from the space of the possible weights. Increasing the size of a network then allows for a better approximation. The quality of the approximation (and thus the performances of the network) is determined by the quality of the sampling process. Overfitting would then not be linked to the number of parameters, but rather the capacity of the sampling process to select the right neurons for difficult regions of the input space.

Generally speaking, a continuous point of view seem to be a better support for intuition, design and theoretical analysis in machine learning. This is the point of view adopted with the ID, where one manipulates distributions instead of fixed set of points. When practical work has to be done the ID is estimated, the data being interpreted as samples from an underlying distribution. We wish to achieve a similar result with neural networks: to get rid of the discrete complexity (the connections, the number of neurons, etc), to abstract it into series of continuous transformations.

6.2 Proposed model

The model we propose is inspired by Hazam et al. [16]. They model infinite neural networks with an arbitrary depth to create new kernels mimicking deep learning. In order to compute analytically their kernels, they restrict themselves to Gaussian distributions of the weights. We lift this restriction to get a more generic model.

In this section, we use the following notations.

- let $\langle \cdot \rangle$ be the scalar product in \mathbb{R}^n
- let $a : \mathbb{R} \rightarrow \mathbb{R}$ be an activation function.

We consider a neural network with one hidden layer and one output neuron. We isolate a hidden neuron, and analyze its contribution to the output. This hidden neuron is associated with a weighting of the input $\mathbf{w} \in \mathbb{R}^n$, and a weighting of its output $u \in \mathbb{R}$. Its contribution to the global output is:

$$a(\langle \mathbf{x}, \mathbf{w} \rangle) \cdot u \tag{21}$$

In practice with neural networks, the hidden layers contains a set number of hidden neurons. The i th neuron is associated with a weight $\mathbf{w}_i \in \mathbb{R}^n$. The output of each neuron is weighted by a number $u(\mathbf{w}_i) \in \mathbb{R}$. The final output is a weighted sum of the activation of each neuron:

$$\sum_i a(\langle \mathbf{x}, \mathbf{w}_i \rangle) \cdot u(\mathbf{w}_i) \tag{22}$$

We give a theoretical model for infinite neural networks. Instead of having a fixed, finite number of hidden neurons (and their respective weights), we consider all possible weightings of the input. The final output is again a weighted sum of the activation of each neuron. Consider the following function parameterized by $u : \mathbb{R}^n \rightarrow \mathbb{R}$:

$$F_u : \mathbf{x} \mapsto \int_{\mathbb{R}^n} a(\langle \mathbf{x} \cdot \mathbf{w} \rangle) \cdot u(\mathbf{w}) d\mathbf{w} \quad (23)$$

This infinite network is parameterized by $u : \mathbb{R}^n \rightarrow \mathbb{R}$, a continuous and integrable function weighting each activation, and μ , a probability distribution of the different weights.

Such a function computes every possible weighted input $\langle \mathbf{x} \cdot \mathbf{w} \rangle$, then computes the corresponding activations, and finally weights these activations according to u . Neural networks with one hidden layer can then be seen as an approximation of such functions. Instead of considering all possible initial weights \mathbf{w} , only a finite number are chosen; the integral then becomes a finite sum. In other words, the hidden neurons of a neural networks (and their associated weights) can be thought as the result of a sampling process of an underlying distribution.

This model has other differences with the finite one (presented earlier) than adopting a continuous point of view. The output layer is not equipped with an activation function, for two reasons. First, the universal approximation theorem (for finite networks) do not require activation functions on the output layer. Second, the last activation function is commonly used for convenience or interpretation reasons: a sigmoid function would fit the output in $(0, 1)$, a softmax layer fit the output to a probability distribution (the output values are in $(0, 1)$ and sum to 1), etc. To keep the analysis generic, we remove the last activation function.

6.3 An interpretation for binary classification

Equation 23 can be rewritten to have a simple interpretation in the case of binary classification. u being a weighting of all the possible activations, it resembles a probability density function. However, it is not always positive, and its integral is not 1. We thus define the following quantities:

$$u_+ = \max(0, u) \quad U_+ = \int u_+ \quad \mu_+ = \frac{u_+}{U_+} \quad (24)$$

$$u_- = \max(0, -u) \quad U_- = \int u_- \quad \mu_- = \frac{u_-}{U_-} \quad (25)$$

u_+ and u_- are the positive and negative parts of u ; they are the two positive functions verifying $u = u_+ - u_-$. Then, μ_+ and μ_- are two probability density functions, obtained from u_+ and u_- by re-scaling them. We note that these distributions describe disjoint regions of \mathbb{R}^n ; in other words, when one is (strictly)

positive, the other has to be null. We now rewrite Equation 23:

$$F_u(\mathbf{x}) = \int_{\mathbb{R}^n} a(\langle \mathbf{x} \cdot \mathbf{w} \rangle) \cdot u(\mathbf{w}) d\mathbf{w} \quad (26)$$

$$= \int_{\mathbb{R}^n} a(\langle \mathbf{x} \cdot \mathbf{w} \rangle) \cdot u_+(\mathbf{w}) d\mathbf{w} - \int_{\mathbb{R}^n} a(\langle \mathbf{x} \cdot \mathbf{w} \rangle) \cdot u_-(\mathbf{w}) d\mathbf{w} \quad (27)$$

$$= U_+ \int_{\mathbb{R}^n} a(\langle \mathbf{x} \cdot \mathbf{w} \rangle) d\mu_+(\mathbf{w}) - U_- \int_{\mathbb{R}^n} a(\langle \mathbf{x} \cdot \mathbf{w} \rangle) d\mu_-(\mathbf{w}) \quad (28)$$

$$= U_+ \mathbb{E}_{\mathbf{w} \sim \mu_+} [a(\langle \mathbf{x} \cdot \mathbf{w} \rangle)] - U_- \mathbb{E}_{\mathbf{w} \sim \mu_-} [a(\langle \mathbf{x} \cdot \mathbf{w} \rangle)] \quad (29)$$

Thus, the output of the network is the (weighted) averaged activation over one region of the weights space, minus the (weighted) averaged activation over a disjoint region of the weight space.

We assume a to be both positive and increasing (like ReLU or sigmoid). Intuitively, when the input \mathbf{x} is in a dense regions of μ_+ , many scalar products $\langle \mathbf{x} \cdot \mathbf{w} \rangle$ will be high. As a result, many activations will be high (a being positive and increasing). Thus, the expectation $U_+ \mathbb{E}_{\mathbf{w} \sim \mu_+} [a(\langle \mathbf{x} \cdot \mathbf{w} \rangle)]$ will tend to be high. Conversely \mathbf{x} will be in a void region of μ_- ; by a similar argument $U_- \mathbb{E}_{\mathbf{w} \sim \mu_-} [a(\langle \mathbf{x} \cdot \mathbf{w} \rangle)]$ will tend to be low (close to 0). Thus, when the input \mathbf{x} is in a dense regions of μ_+ , the output of the network $F_u(\mathbf{x})$ will tend to be high.

A symmetric argument can be given for inputs in dense regions of μ_- giving a low (negative) output. The function F_u can be seen as binary classifier, where μ_+ describes the density of the first class, associated with positive outputs, and μ_- describes the density of the second class, associated with negative outputs.

6.4 Possible extensions

Adding more layers We gave a model suited for neural networks with one hidden layer. It is possible to extend this construction for any number of layer.

We give an example for the modeling of a 3-layers feedforward neural network: we start from the function F_u , a weighting by a given u of all the possible activations. We consider the activation of this quantity: $a(F_u)$. This is one possible activation for the second (hidden) layer. A weighting of all possible activations gives an infinite model for a 3-layers feedforward neural network. Let $v : \mathbb{R}^{\mathbb{R}} \rightarrow \mathbb{R}$. We define:

$$G_v = \int a(F_u(\mathbf{x})) \cdot v(u) du \quad (30)$$

(We assume that v is chosen such that the above integral is well defined.)

Infinite networks with one layer were parameterized with a real-valued function over \mathbb{R} , similar to a probability density. Networks with two layers are now parameterized with real-valued function of functions, similar to probability measures. This construction can be repeated an arbitrary number of times. The interpretation of such models may be harder though. This direction is left open for future studies.

The effect of perturbations Let $\mathbf{x} \in \mathbb{R}^n$ be a reference location in space, the unit vector $\hat{\mathbf{y}} \in \mathbb{R}^n$ be the direction of perturbation, and $u : \mathbb{R} \rightarrow \mathbb{R}$ a suitable parameter for an infinite network. Consider the following function:

$$\begin{aligned} \mathbb{R}_+ &\rightarrow \mathbb{R} \\ \bar{F}_u : \epsilon &\mapsto F_u(\mathbf{x} + \epsilon \hat{\mathbf{y}}) \end{aligned} \quad (31)$$

$\bar{F}_u(\epsilon)$ is the output of the network for an input \mathbf{x} modified by an amount ϵ in a direction $\hat{\mathbf{y}}$. We want to study how much the output of the network can change for small perturbations. We thus compute the ID of \bar{F}_u around 0. Let $A_{\mathbf{w}} : \epsilon \mapsto a(\langle \mathbf{x} + \epsilon \hat{\mathbf{y}} \cdot \mathbf{w} \rangle)$. We have:

$$\text{wID}_{\bar{F}_u}(\epsilon) = \epsilon \cdot \left(\int_{\mathbb{R}^n} A_{\mathbf{w}}(\epsilon) \cdot u(\mathbf{w}) d\mathbf{w} \right)'(\epsilon) \quad (\text{See Equation 2}) \quad (32)$$

$$= \int_{\mathbb{R}^n} \epsilon \cdot A'_{\mathbf{w}}(\epsilon) \cdot u(\mathbf{w}) d\mathbf{w} \quad (33)$$

$$= \int_{\mathbb{R}^n} \text{wID}_{A_{\mathbf{w}}}(\epsilon) \cdot u(\mathbf{w}) d\mathbf{w} \quad (34)$$

The quantities $\bar{F}_u(\epsilon)$ (output of the network), $A_{\mathbf{w}}(\epsilon)$ (activation of a neuron) and $\text{wID}_{A_{\mathbf{w}}}(\epsilon)$ (weighted ID of a function with an analytic form) can all be computed. Thus, $\text{wID}_{\bar{F}_u}$ could be computed in a theoretical case (if u is known) or approximated in a practical case (by approximating u from a finite network). The effect of perturbations could then be studied with the Local ID Representation Theorem:

$$\frac{\bar{F}_u(\epsilon_1)}{\bar{F}_u(\epsilon_2)} \simeq \left(\frac{\epsilon_1}{\epsilon_2} \right)^{\text{ID}_{\bar{F}_u}(0)} \quad (35)$$

Further research in this direction could help understanding the efficiency of adversarial perturbations on neural networks [17]. (Now, this is completely speculative, but:) the quotient ϵ_1/ϵ_2 can be rewritten using the LID of the dataset at x , using again the Local ID Representation Theorem. Let ϵ_1 be chosen such that it corresponds to the first nearest neighbour by expectation of x ($F_{\mathbf{X}}(\epsilon_1)$)

Expressive power We already know the following results:

- (finite) neural networks can approximate any continuous function [15]
- (finite) neural networks have a perfect finite sample expressivity [13]: they can learn any finite set of input/output, provided that the number of hidden neurons is large enough.

A natural question to ask is: what is the expressivity of this new model. We adopt the point of view of the ID for this question: the data are not a finite set of input/output pairs; rather, the training set is the result of a sampling process from an underlying distribution. There are a few possibilities from there:

- is the suggested model (both with one layer or more) capable of modeling any function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$? Can it do so with just a restricted class of functions (continuous functions, functions with a compact support, etc)?

- if this is not possible, can it approximate any function, or some class of function? In which sense would it approximate (\mathbb{L}^1 ? \mathbb{L}^2 ? \mathbb{L}^∞ ?)?

7 Conclusion

Although most of the internship consisted in experimenting with neural networks, the most promising results are on the theoretical side, with the infinite model. Abstracting the discrete details of both the architecture and the database yield simple objects to work with. We believe that further research in this direction could lead to a better understanding of the learning process, to a better understanding of the adversarial setting and to a better intuition regarding design principles.

References

1. Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278-2324, November 1998. <http://yann.lecun.com/exdb/mnist/>
2. Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. 2009. <https://www.cs.toronto.edu/~kriz/cifar.html>
3. Martin Thoma. The HASYv2 dataset. *CoRR*, abs/1701.08380, 2017. <http://arxiv.org/abs/1701.08380>
4. Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. <https://arxiv.org/abs/1409.0575>
5. David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel and Demis Hassabis. Mastering the game of Go with deep neural networks and tree search. doi:10.1038/nature16961 <https://storage.googleapis.com/deepmind-media/alphago/AlphaGoNaturePaper.pdf>
6. Explaining How a Deep Neural Network Trained with End-to-End Learning Steers a Car. Explaining How a Deep Neural Network Trained with End-to-End Learning Steers a Car. *CoRR*, arXiv:1704.07911, 2017. <https://arxiv.org/abs/1704.07911>
7. M. E. Houle. Inlierness, outlierness, hubness and discriminability: an extreme-value-theoretic foundation. Technical Report NII-2015-002E, NII, Mar 2015. http://www.nii.ac.jp/TechReports/public_html/15-002E.pdf
8. Michael E. Houle. Dimensionality, Discriminability, Density & Distance Distributions. In *ICDMW*, pages 468–473, 2013.
9. L. Amsaleg, O. Chelly, T. Furon, S. Girard, M. E. Houle, K. Kawarabayashi, and M. Nett. Estimating local intrinsic dimensionality. In *KDD*, pages 29–38, 2015. <http://mistis.inrialpes.fr/~girard/Fichiers/p29-amsaleg.pdf>
10. D. R. Karger and M. Ruhl. Finding nearest neighbors in growth-restricted metrics. In *STOC*, pages 741–750, 2002.

11. M. E. Houle, H. Kashima, and M. Nett. Generalized expansion dimension. In *ICDMW*, pages 587–594, 2012.
12. Linnainmaa, Seppo. Taylor expansion of the accumulated rounding error. *BIT Numerical Mathematics*. 16 (2): 146160. doi:10.1007/bf01931367
13. Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht and Oriol Vinyals. Understanding deep learning requires rethinking generalization ICLR, 2017. <http://arxiv.org/abs/1611.03530>
14. A., Nielsen, Michael. Neural Networks and Deep Learning (2015-01-01). <http://neuralnetworksanddeeplearning.com/chap2.html>
15. Gybenko, G. Approximation by superposition of sigmoidal functions. *Mathematics of Control, Signals and Systems*, 2(4), 303-314, 1989. <http://www.cs.cmu.edu/~bhiksha/courses/deeplearning/Fall.2016/pdfs/Cybenko.pdf>
16. Tamir Hazan and Tommi S. Jaakkola. Steps Toward Deep Kernel Methods from Infinite Neural Networks. *CoRR*, abs/1508.05133, 2015. <http://arxiv.org/abs/1508.05133>
17. Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow and Rob Fergus. Intriguing properties of neural networks ICLR, abs/1312.6199, 2014. <http://arxiv.org/abs/1312.6199>
18. Alexey Kurakin, Ian J. Goodfellow and Samy Bengio Adversarial examples in the physical world *CoRR*, abs/1607.02533 2016. <http://arxiv.org/abs/1607.02533>