

# January 22-26: Advanced machine learning and data analysis for the physical sciences

Morten Hjorth-Jensen<sup>1,2</sup>

Department of Physics and Center for Computing in Science Education,  
University of Oslo, Norway<sup>1</sup>

Department of Physics and Astronomy and Facility for Rare Isotope Beams,  
Michigan State University, East Lansing, Michigan, USA<sup>2</sup>

Jan 23, 2024

© 1999-2024, Morten Hjorth-Jensen. Released under CC Attribution-NonCommercial 4.0 license

## Overview of second week, January 22-26

1. Mathematics of neural networks
2. Writing own code (bring back to life your NN code)
3. Discussion of first data set and paths for project 1

# Videos on Neural Networks

- ▶ Video on Neural Networks
- ▶ Video on the back propagation algorithm

# Mathematics of deep learning

## Two recent books online

1. The Modern Mathematics of Deep Learning, by Julius Berner, Philipp Grohs, Gitta Kutyniok, Philipp Petersen, published as Mathematical Aspects of Deep Learning, pp. 1-111. Cambridge University Press, 2022
2. Mathematical Introduction to Deep Learning: Methods, Implementations, and Theory, Arnulf Jentzen, Benno Kuckuck, Philippe von Wurstemberger

## Reminder on books with hands-on material and codes

- ▶ Sebastian Rashcka et al, Machine learning with Sickit-Learn and PyTorch
- ▶ David Foster, Generative Deep Learning with TensorFlow
- ▶ Bali and Gavras, Generative AI with Python and TensorFlow 2

All three books have GitHub addresses from where one can download all codes. We will borrow most of the material from these three texts as well as from Goodfellow, Bengio and Courville's text [Deep Learning](#)

## Reading recommendations

1. Rashkca et al., chapter 11, jupyter-notebook sent separately, from [GitHub](#)
2. Goodfellow et al, chapter 6 and 7 contain most of the neural network background.

# Mathematics of deep learning and neural networks

Neural networks, in its so-called feed-forward form, where each iterations contains a feed-forward stage and a back-propagation stage, consist of series of affine matrix-matrix and matrix-vector multiplications. The unknown parameters (the so-called biases and weights which determine the architecture of a neural network), are updated iteratively using the so-called back-propagation algorithm. This algorithm corresponds to the so-called reverse mode of automatic differentiation.

## Basics of an NN

A neural network consists of a series of hidden layers, in addition to the input and output layers. Each layer  $l$  has a set of parameters  $\Theta^{(l)} = (\mathbf{W}^{(l)}, \mathbf{b}^{(l)})$  which are related to the parameters in other layers through a series of affine transformations, for a standard NN these are matrix-matrix and matrix-vector multiplications. For all layers we will simply use a collective variable  $\Theta$ .

It consist of two basic steps:

1. a feed forward stage which takes a given input and produces a final output which is compared with the target values through our cost/loss function.
2. a back-propagation state where the unknown parameters  $\Theta$  are updated through the optimization of the their gradients. The expressions for the gradients are obtained via the chain rule, starting from the derivative of the cost/function.

These two steps make up one iteration. This iterative process is continued till we reach an eventual stopping criterion.



## Overarching view of a neural network

The architecture of a neural network defines our model. This model aims at describing some function  $f(\mathbf{x})$  which aims at describing some final result (outputs or target values) given a specific input  $\mathbf{x}$ . Note that here  $\mathbf{y}$  and  $\mathbf{x}$  are not limited to be vectors.

The architecture consists of

1. An input and an output layer where the input layer is defined by the inputs  $\mathbf{x}$ . The output layer produces the model output  $\tilde{\mathbf{y}}$  which is compared with the target value  $\mathbf{y}$
2. A given number of hidden layers and neurons/nodes/units for each layer (this may vary)
3. A given activation function  $\sigma(\mathbf{z})$  with arguments  $\mathbf{z}$  to be defined below. The activation functions may differ from layer to layer.
4. The last layer, normally called **output** layer has normally an activation function tailored to the specific problem
5. Finally we define a so-called cost or loss function which is used to gauge the quality of our model.

# The optimization problem

The cost function is a function of the unknown parameters  $\Theta$  where the latter is a container for all possible parameters needed to define a neural network

If we are dealing with a regression task a typical cost/loss function is the mean squared error

$$C(\Theta) = \frac{1}{n} \left\{ (\mathbf{y} - \mathbf{X}\theta)^T (\mathbf{y} - \mathbf{X}\theta) \right\}.$$

This function represents one of many possible ways to define the so-called cost function. Note that here we have assumed a linear dependence in terms of the parameters  $\Theta$ . This is in general not the case.

## Parameters of neural networks

For neural networks the parameters  $\Theta$  are given by the so-called weights and biases (to be defined below).

The weights are given by matrix elements  $w_{ij}^{(l)}$  where the superscript indicates the layer number. The biases are typically given by vector elements representing each single node of a given layer, that is  $b_j^{(l)}$ .

## Other ingredients of a neural network

Having defined the architecture of a neural network, the optimization of the cost function with respect to the parameters  $\Theta$ , involves the calculations of gradients and their optimization. The gradients represent the derivatives of a multidimensional object and are often approximated by various gradient methods, including

1. various quasi-Newton methods,
2. plain gradient descent (GD) with a constant learning rate  $\eta$ ,
3. GD with momentum and other approximations to the learning rates such as
  - ▶ Adaptive gradient (ADAGRAD)
  - ▶ Root mean-square propagation (RMSprop)
  - ▶ Adaptive gradient with momentum (ADAM) and many other
4. Stochastic gradient descent and various families of learning rate approximations

## Other parameters

In addition to the above, there are often additional hyperparameters which are included in the setup of a neural network. These will be discussed below.

# Universal approximation theorem

The universal approximation theorem plays a central role in deep learning. [Cybenko \(1989\)](#) showed the following:

Let  $\sigma$  be any continuous sigmoidal function such that

$$\sigma(z) = \begin{cases} 1 & z \rightarrow \infty \\ 0 & z \rightarrow -\infty \end{cases}$$

Given a continuous and deterministic function  $F(\mathbf{x})$  on the unit cube in  $d$ -dimensions  $F \in [0, 1]^d$ ,  $\mathbf{x} \in [0, 1]^d$  and a parameter  $\epsilon > 0$ , there is a one-layer (hidden) neural network  $f(\mathbf{x}; \Theta)$  with  $\Theta = (\mathbf{W}, \mathbf{b})$  and  $\mathbf{W} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^n$ , for which

$$|F(\mathbf{x}) - f(\mathbf{x}; \Theta)| < \epsilon \quad \forall \mathbf{x} \in [0, 1]^d.$$

## Some parallels from real analysis

For those of you familiar with for example the [Stone-Weierstrass theorem](#) for polynomial approximations or the convergence criterion for Fourier series, there are similarities in the derivation of the proof for neural networks. See whiteboard notes for week 2.

## The approximation theorem in words

**Any continuous function  $y = F(\mathbf{x})$  supported on the unit cube in  $d$ -dimensions can be approximated by a one-layer sigmoidal network to arbitrary accuracy.**

[Hornik \(1991\)](#) extended the theorem by letting any non-constant, bounded activation function to be included using that the expectation value

$$\mathbb{E}[|F(\mathbf{x})|^2] = \int_{\mathbf{x} \in D} |F(\mathbf{x})|^2 p(\mathbf{x}) d\mathbf{x} < \infty.$$

Then we have

$$\mathbb{E}[|F(\mathbf{x}) - f(\mathbf{x}; \Theta)|^2] = \int_{\mathbf{x} \in D} |F(\mathbf{x}) - f(\mathbf{x}; \Theta)|^2 p(\mathbf{x}) d\mathbf{x} < \epsilon.$$



## More on the general approximation theorem

None of the proofs give any insight into the relation between the number of hidden layers and nodes and the approximation error  $\epsilon$ , nor the magnitudes of  $\mathbf{W}$  and  $\mathbf{b}$ .

Neural networks (NNs) have what we may call a kind of universality no matter what function we want to compute.

It does not mean that an NN can be used to exactly compute any function. Rather, we get an approximation that is as good as we want.

## Class of functions we can approximate

The class of functions that can be approximated are the continuous ones. If the function  $F(\mathbf{x})$  is discontinuous, it won't in general be possible to approximate it. However, an NN may still give an approximation even if we fail in some points.

## Setting up the equations for a neural network

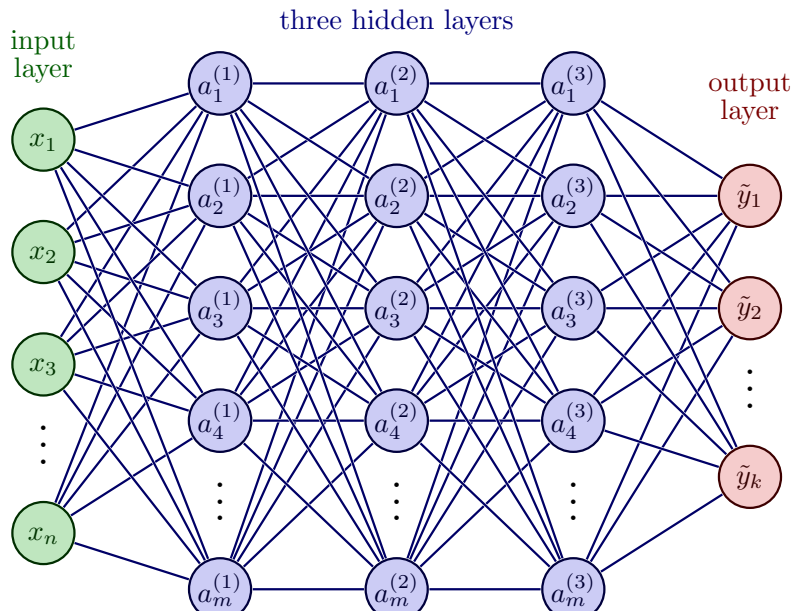
The questions we want to ask are how do changes in the biases and the weights in our network change the cost function and how can we use the final output to modify the weights and biases?

To derive these equations let us start with a plain regression problem and define our cost function as

$$\mathcal{C}(\Theta) = \frac{1}{2} \sum_{i=1}^n (y_i - \tilde{y}_i)^2,$$

where the  $y_i$ s are our  $n$  targets (the values we want to reproduce), while the outputs of the network after having propagated all inputs  $\mathbf{x}$  are given by  $\tilde{\mathbf{y}}_i$ .

## Layout of a neural network with three hidden layers



## Definitions

With our definition of the targets  $\mathbf{y}$ , the outputs of the network  $\tilde{\mathbf{y}}$  and the inputs  $\mathbf{x}$  we define now the activation  $z_j^l$  of node/neuron/unit  $j$  of the  $l$ -th layer as a function of the bias, the weights which add up from the previous layer  $l - 1$  and the forward passes/outputs  $\hat{a}^{l-1}$  from the previous layer as

$$z_j^l = \sum_{i=1}^{M_{l-1}} w_{ij}^l a_i^{l-1} + b_j^l,$$

where  $b_k^l$  are the biases from layer  $l$ . Here  $M_{l-1}$  represents the total number of nodes/neurons/units of layer  $l - 1$ . The figure in the whiteboard notes illustrates this equation. We can rewrite this in a more compact form as the matrix-vector products we discussed earlier,

$$\hat{\mathbf{z}}^l = \left( \hat{\mathbf{W}}^l \right)^T \hat{\mathbf{a}}^{l-1} + \hat{\mathbf{b}}^l.$$

## Inputs to the activation function

With the activation values  $\mathbf{z}^l$  we can in turn define the output of layer  $l$  as  $\mathbf{a}^l = f(\mathbf{z}^l)$  where  $f$  is our activation function. In the examples here we will use the sigmoid function discussed in our logistic regression lectures. We will also use the same activation function  $f$  for all layers and their nodes. It means we have

$$a_j^l = f(z_j^l) = \frac{1}{1 + \exp -(z_j^l)}.$$

## Derivatives and the chain rule

From the definition of the activation  $z_j^l$  we have

$$\frac{\partial z_j^l}{\partial w_{ij}^l} = a_i^{l-1},$$

and

$$\frac{\partial z_j^l}{\partial a_i^{l-1}} = w_{ji}^l.$$

With our definition of the activation function we have that (note that this function depends only on  $z_j^l$ )

$$\frac{\partial a_j^l}{\partial z_j^l} = a_j^l(1 - a_j^l) = f(z_j^l)(1 - f(z_j^l)).$$

## Derivative of the cost function

With these definitions we can now compute the derivative of the cost function in terms of the weights.

Let us specialize to the output layer  $l = L$ . Our cost function is

$$\mathcal{C}(\Theta^L) = \frac{1}{2} \sum_{i=1}^n (y_i - \tilde{y}_i)^2 = \frac{1}{2} \sum_{i=1}^n (a_i^L - y_i)^2,$$

The derivative of this function with respect to the weights is

$$\frac{\partial \mathcal{C}(\Theta^L)}{\partial w_{jk}^L} = (a_j^L - y_j) \frac{\partial a_j^L}{\partial w_{jk}^L},$$

The last partial derivative can easily be computed and reads (by applying the chain rule)

$$\frac{\partial a_j^L}{\partial w_{jk}^L} = \frac{\partial a_j^L}{\partial z_j^L} \frac{\partial z_j^L}{\partial w_{jk}^L} = a_j^L (1 - a_j^L) a_k^{L-1}.$$



## Bringing it together, first back propagation equation

We have thus

$$\frac{\partial \mathcal{C}(\Theta^L)}{\partial w_{jk}^L} = (a_j^L - y_j) a_j^L (1 - a_j^L) a_k^{L-1},$$

Defining

$$\delta_j^L = a_j^L (1 - a_j^L) (a_j^L - y_j) = f'(z_j^L) \frac{\partial \mathcal{C}}{\partial (a_j^L)},$$

and using the Hadamard product of two vectors we can write this as

$$\boldsymbol{\delta}^L = f'(\hat{\mathbf{z}}^L) \circ \frac{\partial \mathcal{C}}{\partial (\mathbf{a}^L)}.$$

## Analyzing the last results

This is an important expression. The second term on the right handside measures how fast the cost function is changing as a function of the  $j$ th output activation. If, for example, the cost function doesn't depend much on a particular output node  $j$ , then  $\delta_j^L$  will be small, which is what we would expect. The first term on the right, measures how fast the activation function  $f$  is changing at a given activation value  $z_j^L$ .

## More considerations

Notice that everything in the above equations is easily computed. In particular, we compute  $z_j^L$  while computing the behaviour of the network, and it is only a small additional overhead to compute  $f'(z_j^L)$ . The exact form of the derivative with respect to the output depends on the form of the cost function. However, provided the cost function is known there should be little trouble in calculating

$$\frac{\partial \mathcal{C}}{\partial (a_j^L)}$$

With the definition of  $\delta_j^L$  we have a more compact definition of the derivative of the cost function in terms of the weights, namely

$$\frac{\partial \mathcal{C}}{\partial w_{jk}^L} = \delta_j^L a_k^{L-1}.$$

## Derivatives in terms of $z_j^L$

It is also easy to see that our previous equation can be written as

$$\delta_j^L = \frac{\partial \mathcal{C}}{\partial z_j^L} = \frac{\partial \mathcal{C}}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L},$$

which can also be interpreted as the partial derivative of the cost function with respect to the biases  $b_j^L$ , namely

$$\delta_j^L = \frac{\partial \mathcal{C}}{\partial b_j^L} \frac{\partial b_j^L}{\partial z_j^L} = \frac{\partial \mathcal{C}}{\partial b_j^L},$$

That is, the error  $\delta_j^L$  is exactly equal to the rate of change of the cost function as a function of the bias.

## Bringing it together

We have now three equations that are essential for the computations of the derivatives of the cost function at the output layer. These equations are needed to start the algorithm and they are

$$\frac{\partial \mathcal{C}(\hat{W}^L)}{\partial w_{jk}^L} = \delta_j^L a_k^{L-1}, \quad (1)$$

and

$$\delta_j^L = f'(z_j^L) \frac{\partial \mathcal{C}}{\partial (a_j^L)}, \quad (2)$$

and

$$\delta_j^L = \frac{\partial \mathcal{C}}{\partial b_j^L}, \quad (3)$$

## Final back propagating equation

We have that (replacing  $L$  with a general layer  $l$ )

$$\delta_j^l = \frac{\partial \mathcal{C}}{\partial z_j^l}.$$

We want to express this in terms of the equations for layer  $l + 1$ .

## Using the chain rule and summing over all $k$ entries

We obtain

$$\delta_j^l = \sum_k \frac{\partial \mathcal{C}}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} = \sum_k \delta_k^{l+1} \frac{\partial z_k^{l+1}}{\partial z_j^l},$$

and recalling that

$$z_j^{l+1} = \sum_{i=1}^{M_l} w_{ij}^{l+1} a_i^l + b_j^{l+1},$$

with  $M_l$  being the number of nodes in layer  $l$ , we obtain

$$\delta_j^l = \sum_k \delta_k^{l+1} w_{kj}^{l+1} f'(z_j^l),$$

This is our final equation.

We are now ready to set up the algorithm for back propagation and learning the weights and biases.

## Setting up the back propagation algorithm

The four equations provide us with a way of computing the gradient of the cost function. Let us write this out in the form of an algorithm.

**First**, we set up the input data  $\hat{x}$  and the activations  $\hat{z}_1$  of the input layer and compute the activation function and the pertinent outputs  $\hat{a}^1$ .

**Secondly**, we perform then the feed forward till we reach the output layer and compute all  $\hat{z}_l$  of the input layer and compute the activation function and the pertinent outputs  $\hat{a}^l$  for

$l = 1, 2, 3, \dots, L$ .

**Notation:** The first hidden layer has  $l = 1$  as label and the final output layer has  $l = L$ .



## Setting up the back propagation algorithm, part 2

Thereafter we compute the output error  $\hat{\delta}^L$  by computing all

$$\delta_j^L = f'(z_j^L) \frac{\partial \mathcal{C}}{\partial (a_j^L)}.$$

Then we compute the back propagate error for each  $l = L - 1, L - 2, \dots, 1$  as

$$\delta_j^l = \sum_k \delta_k^{l+1} w_{kj}^{l+1} f'(z_j^l).$$

## Setting up the Back propagation algorithm, part 3

Finally, we update the weights and the biases using gradient descent for each  $l = L - 1, L - 2, \dots, 1$  and update the weights and biases according to the rules

$$w_{jk}^l \leftarrow w_{jk}^l - \eta \delta_j^l a_k^{l-1},$$

$$b_j^l \leftarrow b_j^l - \eta \frac{\partial \mathcal{C}}{\partial b_j^l} = b_j^l - \eta \delta_j^l,$$

with  $\eta$  being the learning rate.

## Updating the gradients

With the back propagate error for each  $l = L - 1, L - 2, \dots, 1$  as

$$\delta_j^l = \sum_k \delta_k^{l+1} w_{kj}^{l+1} f'(z_j^l),$$

we update the weights and the biases using gradient descent for each  $l = L - 1, L - 2, \dots, 1$  and update the weights and biases according to the rules

$$w_{jk}^l \leftarrow w_{jk}^l - \eta \delta_j^l a_k^{l-1},$$

$$b_j^l \leftarrow b_j^l - \eta \frac{\partial \mathcal{C}}{\partial b_j^l} = b_j^l - \eta \delta_j^l,$$

## Fine-tuning neural network hyperparameters

The flexibility of neural networks is also one of their main drawbacks: there are many hyperparameters to tweak. Not only can you use any imaginable network topology (how neurons/nodes are interconnected), but even in a simple FFNN you can change the number of layers, the number of neurons per layer, the type of activation function to use in each layer, the weight initialization logic, the stochastic gradient optimizer and much more. How do you know what combination of hyperparameters is the best for your task?

- ▶ You can use grid search with cross-validation to find the right hyperparameters.

However, since there are many hyperparameters to tune, and since training a neural network on a large dataset takes a lot of time, you will only be able to explore a tiny part of the hyperparameter space.

- ▶ You can use randomized search.
- ▶ Or use tools like [Oscar](#), which implements more complex algorithms to help you find a good set of hyperparameters quickly.

## Hidden layers

For many problems you can start with just one or two hidden layers and it will work just fine. For the MNIST data set you can easily get a high accuracy using just one hidden layer with a few hundred neurons. You can reach for this data set above 98% accuracy using two hidden layers with the same total amount of neurons, in roughly the same amount of training time.

For more complex problems, you can gradually ramp up the number of hidden layers, until you start overfitting the training set. Very complex tasks, such as large image classification or speech recognition, typically require networks with dozens of layers and they need a huge amount of training data. However, you will rarely have to train such networks from scratch: it is much more common to reuse parts of a pretrained state-of-the-art network that performs a similar task.

## Which activation function should I use?

The Back propagation algorithm we derived above works by going from the output layer to the input layer, propagating the error gradient on the way. Once the algorithm has computed the gradient of the cost function with regards to each parameter in the network, it uses these gradients to update each parameter with a Gradient Descent (GD) step.

Unfortunately for us, the gradients often get smaller and smaller as the algorithm progresses down to the first hidden layers. As a result, the GD update leaves the lower layer connection weights virtually unchanged, and training never converges to a good solution. This is known in the literature as **the vanishing gradients problem**.

In other cases, the opposite can happen, namely the the gradients can grow bigger and bigger. The result is that many of the layers get large updates of the weights the algorithm diverges. This is the **exploding gradients problem**, which is mostly encountered in recurrent neural networks. More generally, deep neural networks suffer from unstable gradients, different layers may learn at widely different speeds

## Is the Logistic activation function (Sigmoid) our choice?

Although this unfortunate behavior has been empirically observed for quite a while (it was one of the reasons why deep neural networks were mostly abandoned for a long time), it is only around 2010 that significant progress was made in understanding it.

A paper titled [Understanding the Difficulty of Training Deep Feedforward Neural Networks](#) by Xavier Glorot and Yoshua Bengio found that the problems with the popular logistic sigmoid activation function and the weight initialization technique that was most popular at the time, namely random initialization using a normal distribution with a mean of 0 and a standard deviation of 1.

They showed that with this activation function and this initialization scheme, the variance of the outputs of each layer is much greater than the variance of its inputs. Going forward in the network, the variance keeps increasing after each layer until the activation function saturates at the top layers. This is actually made worse by the fact that the logistic function has a mean of 0.5, not 0 (the hyperbolic tangent function has a mean of 0 and behaves slightly better than the logistic function in deep networks).

## The derivative of the Logistic function

Looking at the logistic activation function, when inputs become large (negative or positive), the function saturates at 0 or 1, with a derivative extremely close to 0. Thus when backpropagation kicks in, it has virtually no gradient to propagate back through the network, and what little gradient exists keeps getting diluted as backpropagation progresses down through the top layers, so there is really nothing left for the lower layers.

In their paper, Glorot and Bengio propose a way to significantly alleviate this problem. We need the signal to flow properly in both directions: in the forward direction when making predictions, and in the reverse direction when backpropagating gradients. We don't want the signal to die out, nor do we want it to explode and saturate. For the signal to flow properly, the authors argue that we need the variance of the outputs of each layer to be equal to the variance of its inputs, and we also need the gradients to have equal variance before and after flowing through a layer in the reverse direction.



## Insights from the paper by Glorot and Bengio

One of the insights in the 2010 paper by Glorot and Bengio was that the vanishing/exploding gradients problems were in part due to a poor choice of activation function. Until then most people had assumed that if Nature had chosen to use roughly sigmoid activation functions in biological neurons, they must be an excellent choice. But it turns out that other activation functions behave much better in deep neural networks, in particular the ReLU activation function, mostly because it does not saturate for positive values (and also because it is quite fast to compute).