

Lecture Notes: Neural Network Architectures

Evelyn Herberg¹

¹ Interdisciplinary Center for Scientific Computing, Ruprecht-Karls-University of Heidelberg,
69120 Heidelberg, Germany

April 2023

Acknowledgement

These lecture notes were written to serve as theoretical background for programming sessions given by Florian Wolf during the SPP1962 Young Researchers' workshop on Deep Learning in March 2023.

I do not claim originality, but merely collected material from various sources and wrote it down in a cohesive way. Any mistakes that I may have introduced are of course at my fault.

Especially, parts of these lecture notes are based on a student seminar at the University of Heidelberg that Roland Herzog and I organized together. I thank the involved students: Deniz Aydin, Laurin Ernst, Yanxin Jia, Xinyu Liang, Hannah Rickmann, Viktor Stein von Kamienski, Xiao Wang and Zixiang Zhou for their contributions.

The main literature for the seminar, and consequently also for these lecture notes, was [15] and [26]. Additional sources are mentioned throughout the document.

Furthermore, I thank Harbir Antil and the CMAI work group at George Mason University for their guidance in understanding Machine Learning from an optimal control point of view. The notation in this document is highly influenced by the CMAI work group, cf. [2, 1, 3].

Please send comments and remarks to evelyn.herberg@iwr.uni-heidelberg.de.

Contents

1	Introduction	3
1.1	Supervised Learning	5
1.2	Unsupervised Learning	5
1.3	Optimization Algorithms	7
1.4	Overfitting and Underfitting	9
1.5	Hyperparameters and Data Set Splitting	12
1.6	Modeling logical functions	12
2	Feedforward Neural Network	14
2.1	Depth and Width	16
2.2	Initialization	16
2.3	Batch Normalization	17
2.4	Classification Tasks	18
2.5	Backpropagation	21
3	Convolutional Neural Network	26
3.1	Convolution	27
3.2	Convolutional Layer	29
3.3	Detector Layer	32
3.4	Pooling Layer	32
3.5	Local Response Normalization	34
4	ResNet	36
4.1	Different ResNet Versions	38
4.2	ResNet18	39
4.3	Transfer Learning	40
5	Recurrent Neural Network	41
5.1	Variants of RNNs	43
5.2	Long term dependencies	44
5.2.1	Gated Recurrent Unit	45
5.2.2	Long Short Term Memory	45
5.3	Language processing	47

1. Introduction

Machine Learning (ML) denotes the field of study in which algorithms infer from given data how to perform a specific task, without being explicitly programmed for the task (Arthur Samuel, 1959). Here, we consider a popular subset of ML algorithms: **Neural Networks**. The inspiration for a Neural Network (NN) originates from the human brain, where biological neurons (nerve cells) respond to the activation of other neurons they are connected to. At a very simple level, neurons in the brain take electrical inputs that are then channeled to outputs. The sensitivity of this relation also depends on the strength of the connection, i.e. a neuron may be more responsive to one neuron, then to another.

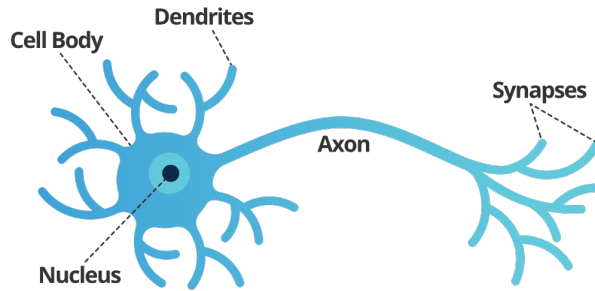


Figure 1. Brain Neuron Structure: electrical inputs are received through dendrites and transmitted via the axon to other cells. There are approximately 86 billion neurons in the human brain. Image modified from: <https://www.smartsheet.com/neural-network-applications>.

For a single neuron/node with input $u \in \mathbb{R}^n$, a mathematical model, named the **perceptron** [27], can be described as

$$y = \sigma \left(\sum_{i=1}^n W_i u_i + b \right) = \sigma(W^\top u + b), \quad (1)$$

where y is the **activation** of the neuron/node, W_i are the **weights** and b is the **bias**.

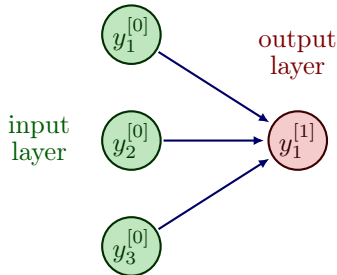


Figure 2. Schematic representation of the perceptron with a three dimensional input. For generality we denote the input by $y^{[0]} = u$ and the output by $y^{[1]} = y$. The weights W_i are applied on the arrows and the bias is added in the node $y_1^{[1]}$.

The function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is called **activation function**. Originally, in [27], it was proposed to choose the Heaviside function as activation function to model whether a neuron fires or not, i.e.

$$\sigma(y) = \begin{cases} 1 & \text{if } y \geq 0, \\ 0 & \text{if } y < 0. \end{cases}$$

However, over time several other activation functions have been suggested and are being used. Typically, they are monotone increasing to remain in the spirit of the original idea, but continuous.

Popular activation functions are, cf. [26, p.90]

$$\begin{aligned} \sigma(y) &= \frac{1}{1 + \exp(-y)} && \text{sigmoid (logistic),} \\ \sigma(y) &= \tanh(y) = \frac{\exp(y) - \exp(-y)}{\exp(y) + \exp(-y)} && \text{hyperbolic tangent,} \\ \sigma(y) &= \max\{y, 0\} && \text{rectified linear unit (ReLU),} \\ \sigma(y) &= \max\{\alpha y, y\} && \text{leaky ReLU.} \end{aligned}$$

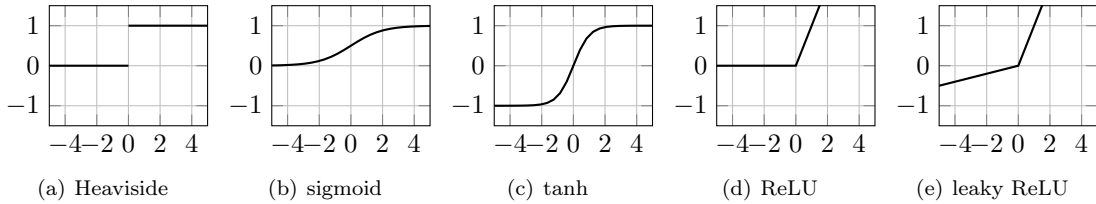


Figure 3. Popular activation functions. Leaky ReLU is displayed for $\alpha = 0.1$.

Remark 1.1 *The nonlinearity of activation functions is an integral part of the Neural Networks success. Since concatenations of linear functions result again in a linear function, see e.g. [26, p.90], the complexity that can be achieved by using linear activation functions is limited.*

While the sigmoid function approximates the Heaviside function continuously, and is differentiable, it contains an exponential operation, which is computationally expensive. Similar problems arise with the hyperbolic tangent function. However, the fact that tanh is closer to the identity function often helps speed up convergence, since it resembles a linear model, as long as the values are close to zero. Another challenge that needs to be overcome is vanishing derivatives, which is visibly present for Heaviside, sigmoid and hyperbolic tangent. In contrast, ReLU is not bounded on positive values, while also being comparatively cheap to compute, because linear computations tend to be very well optimized in modern computing. Altogether, these advantages have resulted in ReLU (and variants thereof) becoming the most widely used activation function currently. As a remedy for the vanishing gradient on negative values, leaky ReLU was introduced. When taking derivatives of ReLU one needs to account for the non-differentiability at 0, but in numerical practice this is easily overcome.

With the help of Neural Networks we want to solve a task, cf. [15, Section 5.1]. Let the performance of the algorithm for the given task be measured by the **loss function** L , which needs to be adequately modeled. By \mathcal{F} we denote the Neural Network. The variables that will be learned are the weights W and biases b of the Neural Network. Hence, we can formulate the following optimization problem, cf. [1, 2, 3]

$$\min_{W, b} \mathcal{L}(y, u, W, b) \quad \text{s.t.} \quad y = \mathcal{F}(u, W, b). \quad (P)$$

One possible choice for \mathcal{F} has already been given in (1), the perceptron. In the subsequent sections we introduce and analyze various other Neural Network architectures. They all have in common that they contain weights and biases, so that the above problem formulation remains sensible.

Before we move on to different network architectures, we discuss the modeling of the loss function. Learning tasks can be divided into two subgroups: Supervised and Unsupervised learning.

1.1. Supervised Learning

In supervised learning we have given data u with known supervision $S(u)$ (also called labels), so that the task is to match the output y of the Neural Network to the supervision. These problems are further categorized depending on the known supervision, e.g. for $S(u) \in \mathbb{N}$ it is called a classification and for $S(u) \in \mathbb{R}$ a regression. Furthermore, the supervision $S(u)$ can also take more complex forms like a black and white picture of 256×256 pixels represented by $[0, 1]^{256}$, a higher dimensional quantity, a sentence, etc. These cases are called structured output learning.

Let us consider one very simple example, cf. [15, Section 5.1.4].

Example 1.2 Linear Regression

We have a given set of inputs $u^{(i)} \in \mathbb{R}^d$ with known supervisions $S(u^{(i)}) \in \mathbb{R}$ for $i = 1, \dots, N$. In this example we only consider weights $W \in \mathbb{R}^d$ and no bias. Additionally, let $\sigma = \text{id}$. The perceptron network simplifies to

$$y^{(i)} = W^\top u^{(i)},$$

and the learning task is to find W , such that $y^{(i)} \approx S(u^{(i)})$. This can be modeled by the **mean squared error (MSE)** function

$$\mathcal{L}(\{y^{(i)}\}_i, \{u^{(i)}\}_i, W) := \frac{1}{2N} \sum_{i=1}^N \|y^{(i)} - S(u^{(i)})\|^2.$$

By convention we will use $\|\cdot\| = \|\cdot\|_2$ throughout the lecture. The chosen loss function is quadratic, convex and non-negative. We define

$$U := \begin{pmatrix} (u^{(1)})^\top \\ \vdots \\ (u^{(N)})^\top \end{pmatrix} \in \mathbb{R}^{N \times d}, \quad S := \begin{pmatrix} S(u^{(1)}) \\ \vdots \\ S(u^{(N)}) \end{pmatrix} \in \mathbb{R}^N,$$

so that we can write $\mathcal{L}(W) = \frac{1}{2} \|UW - S\|_2^2$. Minimizing this function will deliver the same optimal weight W as minimizing the MSE function defined above. We can now derive the gradient

$$\nabla_W \mathcal{L}(W) = U^\top UW - U^\top S$$

and immediately find the stationary point $W = (U^\top U)^{-1} U^\top S$.

1.2. Unsupervised Learning

In unsupervised learning, only the input data u is given and we have no knowledge of supervisions or labels. The algorithm is supposed to learn e.g. a structure or relation in the data. Some examples are k-clustering and principal component analysis (PCA). Modeling the loss function specifies the task and has a direct influence on the learning process. For illustration of this concept, we introduce the k-means algorithm, see eg. [26, Chapter 10], which is used for clustering.

Example 1.3 We have a set of given data points

$$\left\{u^{(i)}\right\}_{i=1}^N \in \mathbb{R}^d,$$

and a desired number of clusters $k \in \mathbb{N}$ with $k \leq N$ and typically $k \ll N$. Every data point is supposed to be assigned to a cluster. Iteratively every data point is assigned to the cluster with the nearest centroid, and we redefine cluster centroids as the mean of the vectors in the cluster. The procedure is specified in Algorithm 1 and illustrated for an example in Figure 4, which can be found e.g. in [26, Chapter 10]. The loss function (also called distortion function in this setup) can be defined as

$$\mathcal{L}(c, \mu) := \sum_{i=1}^N \|u^{(i)} - \mu_{c(i)}\|^2,$$

which is also a model of the quantity that we try to minimize in Algorithm 1. We have a non-convex set of points in \mathbb{R}^d , so the algorithm may converge to a local minimum. To prevent this, we run the algorithm many times, compare the resulting clusterings using the loss function, and choose the one with the minimal value attained in the loss function.

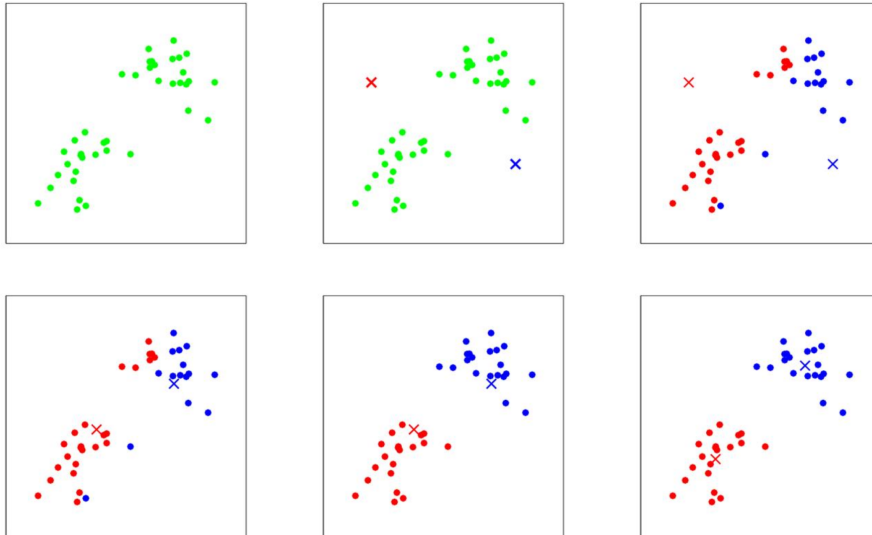


Figure 4. Visualization of k-means algorithm for $k = 2$ clusters. Data points $u^{(i)}$ are indicated as dots, while cluster centroids μ_j are shown as crosses. Top, from left to right: Dataset, random initial cluster centroids μ_1 (red) and μ_2 (blue), every data point is assigned to either the red or blue cluster. Bottom, from left to right: cluster centroids are redefined, every data point is reassigned, cluster centroids are redefined again. Image source: [26, Chapter 10].

We will see various other loss functions \mathcal{L} throughout the remainder of this lecture, all of them specifically tailored to the task at hand.

In the case of Linear Regression, we have a closed form derivative, so we are able to find the solution by direct calculus, while for k-means clustering the optimization was done by a tailored iteration. For general problems we will need a suitable optimization algorithm. We move on to introduce a few options.

Algorithm 1 k-means clustering

Require: Initial cluster centroids μ_1, \dots, μ_k

```

while not converged do
  for  $i = 1 : N$  do
     $c^{(i)} := \arg \min_j \|u^{(i)} - \mu_j\|^2$ 
  end for
  for  $j = 1 : k$  do
     $\mu_j \leftarrow \frac{\sum_{i=1}^N 1_{\{c^{(i)}=j\}} u^{(i)}}{\sum_{i=1}^N 1_{\{c^{(i)}=j\}}}$ 
  end for
end while

```

1.3. Optimization Algorithms

Here, for simplicity we define θ , which collects all variables, i.e. weights W and bias b and write the loss function as

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}^{(i)}(\theta),$$

which we want to minimize. Here, $\mathcal{L}^{(i)}$ indicates the loss function evaluated for data point i , for example with a MSE loss $\mathcal{L}^{(i)}(\theta) = \frac{1}{2} \|y^{(i)} - S(u^{(i)})\|^2$.

First, let us recall the standard **gradient descent** algorithm, see e.g. [6, Section 9.3], which is also known as steepest descent or batch gradient descent.

Algorithm 2 Gradient Descent

Require: Initial point θ^0 , step size $\tau > 0$, counter $k = 0$.

```

while Stopping criterion not fulfilled do
   $\theta^{k+1} = \theta^k - \tau \cdot \nabla \mathcal{L}(\theta^k),$ 
   $k \leftarrow k + 1.$ 
end while

```

Possible stopping criterion are e.g. setting a maximum number of iterations k , reaching a certain exactness $\|\mathcal{L}(\theta)\| < \epsilon$ with a small number $\epsilon > 0$, or a decay in change $\|\theta^{k+1} - \theta^k\| < \epsilon$. Determining a suitable step size is integral to the success of the gradient descent method, especially since this algorithm uses the same step size τ for all components of θ , which can be a large vector in applications. It may happen that in some components the computed descent direction is only providing descent in a small neighborhood, therefore requiring a small step size τ . It is also possible to employ a line search algorithm. However, this is not common in Machine Learning currently. Instead, typically a small step size is chosen, so that it will (hopefully) be not too large for any component of θ , and then it may be adaptively increased. Furthermore, let us remark that the step size is often called **learning rate** in a Machine Learning context.

Additionally, a grand challenge in Machine Learning tasks is that we have huge data sets, and the gradient descent algorithm has to iterate over all data points in every iteration, since $\mathcal{L}(\theta)$

contains all data points, which causes a tremendous computational cost. This motivates the use of the **stochastic gradient descent** algorithm, cf. [26, Algorithm 1], which only takes one data point into account per iteration.

Algorithm 3 Stochastic Gradient Descent (SGD)

Require: Initial point θ^0 , step size $\tau > 0$, counter $k = 0$, maximum number of iterations K .

while $k \leq K$ **do**

Sample $j \in \{1, \dots, N\}$ uniformly.

$$\theta^{k+1} = \theta^k - \tau \cdot \nabla \mathcal{L}^{(j)}(\theta^k),$$

$k \leftarrow k + 1$.

end while

Since the stochastic gradient descent method only calculates the gradient for one data point, it produces an irregular convergence behavior. Indeed, it does not necessarily converge at all, but for a large number of iterations K it often produces a good approximation. In fact, actually converging in training the Neural Network is often not necessary/desired anyhow, since we want to have a solution that generalizes well to unseen data, rather than fit the given data points perfectly. Actually, the latter may lead to overfitting, cf. Section 1.4. Therefore, SGD is a computationally cheap, reasonable alternative to gradient descent. As a compromise, which generates a less irregular convergence behavior, there also exists **mini batch gradient descent**, cf. [26, Algorithm 2], where every iteration takes into account a subset (mini batch) of the data points.

Algorithm 4 Mini Batch Gradient Descent

Require: Initial point θ^0 , step size $\tau > 0$, counter $k = 0$, maximum number of iterations K , batch size $b \in \mathbb{N}$.

while $k \leq K$ **do**

Sample b examples j_1, \dots, j_b uniformly from $\{1, \dots, N\}$

$$\theta^{k+1} = \theta^k - \tau \cdot \frac{1}{b} \sum_{i=1}^b \nabla \mathcal{L}^{(j_i)}(\theta^k),$$

$k \leftarrow k + 1$.

end while

Finally, we introduce a sophisticated algorithm for stochastic optimization called **Adam**, [23], see Algorithm 5. It is also a gradient-based method, and as an extension of the previous methods it employs adaptive estimates of so-called moments. Good default settings in Adam for the tested machine learning problems are $\tau = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$, cf. [23]. Typically, the stochasticity of $\tilde{\mathcal{L}}(\theta)$ will come from using mini batches of the data set, as in Mini Batch Gradient Descent, Algorithm 4.

Algorithm 5 Adam. All operations on vectors are element-wise. $(g^k)^2$ indicates the element-wise square $g^k \odot g^k$, and $(\beta_1)^k, (\beta_2)^k$ denote the k -th power of β_1 and β_2 , respectively.

Require: Initial point θ^0 , step size $\tau > 0$, counter $k = 0$, exponential decay rates for the moment estimates $\beta_1, \beta_2 \in [0, 1)$, $\epsilon > 0$, stochastic approximation $\widetilde{\mathcal{L}}(\theta)$ of the loss function.

$m_1^0 \leftarrow 0$ (Initialize first moment vector)

$m_2^0 \leftarrow 0$ (Initialize second moment vector)

while θ^k not converged **do**

$g^{k+1} = \nabla_{\theta} \widetilde{\mathcal{L}}(\theta^k)$

$m_1^{k+1} = \beta_1 \cdot m_1^k + (1 - \beta_1) \cdot g^{k+1}$

$m_2^{k+1} = \beta_2 \cdot m_2^k + (1 - \beta_2) \cdot (g^{k+1})^2$

$m_1^{k+1} \leftarrow \frac{m_1^{k+1}}{(1 - (\beta_1)^k)}$

$m_2^{k+1} \leftarrow \frac{m_2^{k+1}}{(1 - (\beta_2)^k)}$

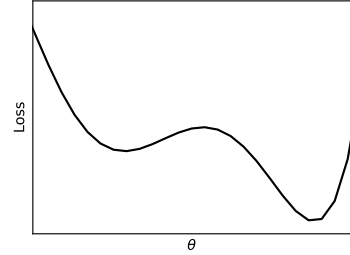
$\theta^{k+1} = \theta^k - \tau \cdot \frac{m_1^{k+1}}{(\sqrt{m_2^{k+1} + \epsilon})}$

$k \leftarrow k + 1$

end while

Remark 1.4 *In any case we need to be cautious when interpreting results, since independent of the chosen algorithm, we are dealing with a non-convex loss function, so that we can only expect convergence to stationary points.*

Figure 5. Simple example of a non-convex loss function with a local and a global minimum.



In the following section we discuss how fitting the given data points and generalizing well to unseen data can be contradictory goals.

1.4. Overfitting and Underfitting

As an example we discuss supervised learning with polynomials of degree r , cf. [12, Section 1.3.3].

Example 1.5 *Define*

$$p(u, W) := \sum_{j=0}^r W_j u^j = W^\top u,$$

with $u = (u^0, \dots, u^r)^\top \in \mathbb{R}^{r+1}$ the potencies of data point u , and $W := (W_0, \dots, W_r)^\top \in \mathbb{R}^{r+1}$. The polynomial p is linear in W , but not in u . As in Linear Regression (Example 1.2), we do not consider bias b here. Our goal is to compute weights W , given data points $u^{(i)}$ with supervisions

$S(u^{(i)})$, so that p makes good predictions on data it hasn't seen before. We again employ the MSE loss function

$$\mathcal{L}(W) = \frac{1}{2N} \sum_{i=1}^N \|p(u^{(i)}, W) - S(u^{(i)})\|^2$$

As before, we write the loss in matrix-vector notation

$$\mathcal{L}(W) = \frac{1}{2N} \|UW - S\|^2$$

where

$$U := \begin{pmatrix} u_0^{(1)} & u_1^{(1)} & \dots & u_r^{(1)} \\ \vdots & \vdots & & \vdots \\ u_0^{(m)} & u_1^{(m)} & \dots & u_r^{(m)} \end{pmatrix}, \quad S := \begin{pmatrix} S(u^{(1)}) \\ \vdots \\ S(u^{(m)}) \end{pmatrix}$$

The minimizer W can be directly calculated, cf. Example 1.2.

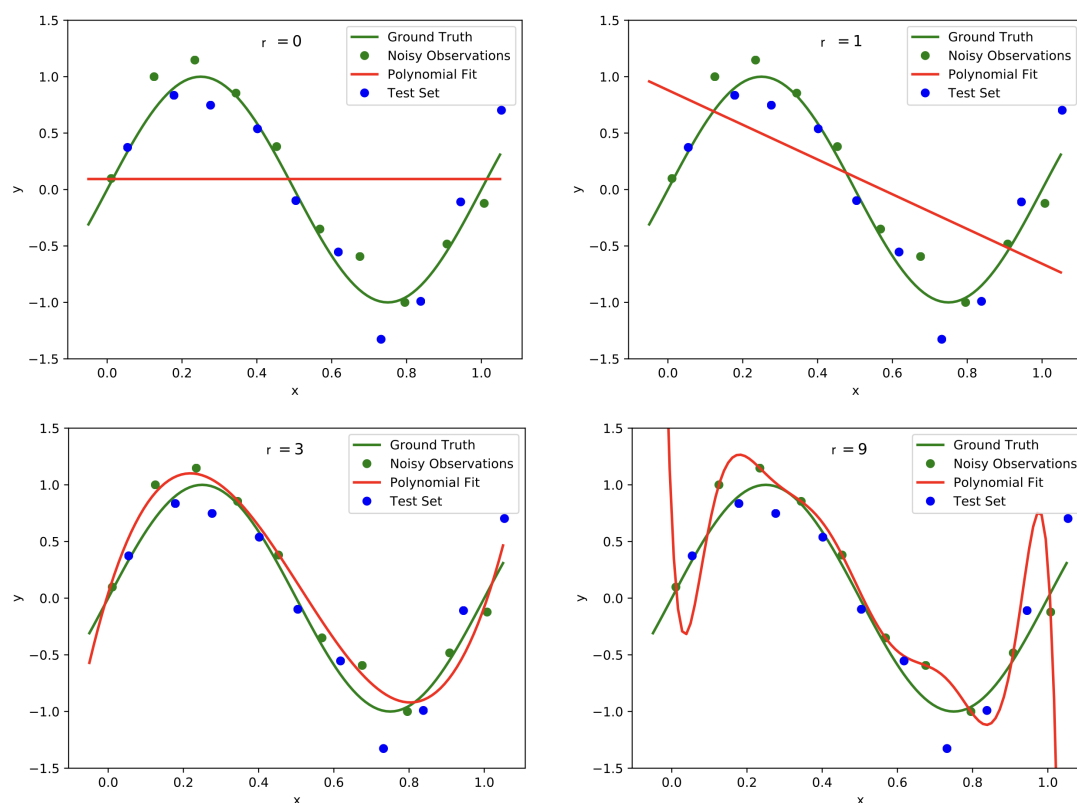


Figure 6. Plots of polynomials of various degrees r (red graph) fitted to the noisy data points (green dots) based on the ground truth (green graph). The model should extend well to the test set data (blue dots). We observe underfitting in the top row for $r = 0$ (left) and $r = 1$ (right). In the bottom left with $r = 3$ reasonable results are achieved, while $r = 9$ in the bottom right leads to overfitting. Image modified from: [12, Fig. 1].

To measure the performance of the polynomial curve fitting we compute the error on data points that were not used to determine the best polynomial fit, because we aim for a model that will generalize well. To this end, finding a suitable degree for the polynomial that we are fitting over the data points is crucial. If the degree is too low, we will encounter **underfitting**, see Figure 6 top row. This means that the complexity of the polynomial is too low and the model does not even fit the data points. A remedy is to increase the degree of the polynomial, see Figure 6 bottom left. However, increasing the degree too much may lead to **overfitting**, see Figure 6 bottom right. The data points are fit perfectly, but the curve will not generalize well.

We can characterize overfitting and underfitting by using some statistics, cf. [26, Section 8.1]. A point estimator $g : \mathcal{U}^N \rightarrow \Theta$ (where \mathcal{U} denotes the data space, and Θ denotes the parameter space) is a function which makes an estimation of the underlying parameters of the model. For example, the estimate for $\theta = W$ from Example 1.2: $\hat{\theta} = (U^\top U)^{-1} U^\top S$ (which we will denote with a hat in this subsection to emphasize that it is an estimation) is an example of a point estimator. We assume that the data from \mathcal{U}^N is i.i.d, so that $\hat{\theta}$ is a random variable. We can define the variance and the bias

$$\text{Var}(\hat{\theta}) := \mathbb{E}(\hat{\theta}^2) - \mathbb{E}(\hat{\theta})^2, \text{Bias}(\hat{\theta}) := \mathbb{E}(\hat{\theta}) - \theta,$$

with \mathbb{E} denoting the expected value. A good estimator has both, low variance and low bias. We can characterize overfitting with low bias and high variance, and underfitting with high bias and low variance. The bias-variance trade-off is illustrated in Figure 7. Hence, we can make a decision based on mean squared error of the estimates

$$\text{MSE}(\hat{\theta}) := \mathbb{E}[(\hat{\theta} - \theta)^2] = \text{Var}(\hat{\theta}) + \text{Bias}(\hat{\theta})^2.$$

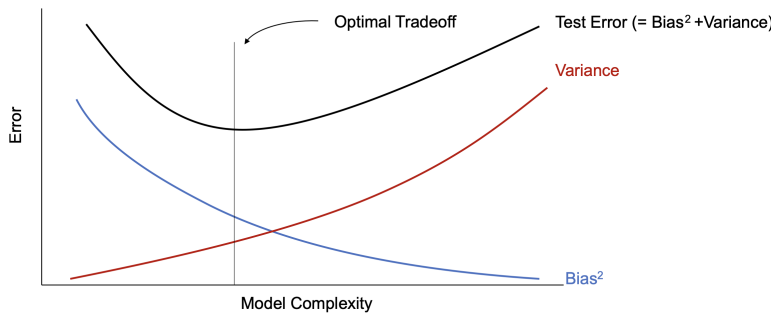


Figure 7. Bias-variance trade-off. Image source: [26, Fig. 8.8].

In general, it can be hard to guess a suitable degree for the polynomial beforehand. We could compute a fitting curve for different choices of r and then compare the error on previously unseen data points of the validation data set, cf. Section 1.5, to determine which one generalizes best. This will require solving the problem multiple times which is unfavorable, especially for large data sets. Also, the polynomial degree can only be set discretely. Another, continuous way is to introduce a penalization term in the loss function

$$\mathcal{L}_\lambda(\theta) := \mathcal{L}(\theta) + \lambda \|\theta\|^2.$$

This technique is also called **weight decay**, cf. [15, Section 5.2.2]. We can also use other norms, e.g. $\|\cdot\|_1$. Here, we can choose a large degree r and for λ big enough, we will still avoid overfitting, because many components of θ will be (close to) zero. Nonetheless, we need to be cautious with the choice of λ . If it is too big, we will face again the problem of underfitting.

We see that choosing values for the degree r and the penalization parameter λ poses challenges, and will discuss this further in the next section.

1.5. Hyperparameters and Data Set Splitting

We call all quantities that need to be chosen before solving the optimization problem **hyperparameters**, cf. [15, Section 5.3]. Let us point out that hyperparameters are not learnt by the optimization algorithm itself, but nevertheless have an impact on the algorithms performance. Examples of hyperparameters include the polynomial degree r , the scalar λ , all parameters in the optimization algorithms (Section 1.3) like the step size τ , and also the architecture of the Neural Network, and many more.

The impact of having a good set of hyperparameters can be tremendous, however finding such a set is not trivial. First of all, we split our given data into three sets. **training** data, **validation** data and **test** data (a 4:1:1 ratio is common). We have seen training and test data before. The data points that we are using as input to solve the optimization problem are called training data, and the unseen data points, which we use to evaluate whether the model generalizes well, are called test data. Since we don't want to mix different causes of error, we also introduce the validation data set. This will be used to compare different choices of hyperparameter configurations, i.e. we train the model on the training data for different hyperparameters, compute the error on the validation data set, choose the hyperparameter setup with the lowest error and finally evaluate the model on the test set. The reasoning behind this is that if we would use the test data set to determine the hyperparameter values, the test error may be not meaningful, because the hyperparameters have been optimized for this specific test set. Since we are using the validation set, we will have the test set with previously unseen data available to determine the generalization error without giving our network an advantage.

Still, imagine you need to choose 5 hyperparameters and have 4 possible values that you want to try for each hyperparameter. This amounts to $4^5 = 1024$ combinations you have to run on the training data and evaluate on the validation set. In real applications the number of hyperparameters and possible values can be much larger, so that it is nearly infeasible to try every combination, but rather common to change one hyperparameter at a time. Luckily, some hyperparameters also have known good default values, like the hyperparameters for Adam Optimizer, Algorithm 5. Apart from that it is a tedious, manual work to try out, monitor and choose suitable hyperparameters.

Finally, we discuss the limitations of shallow Neural Networks, i.e. networks with only one layer.

1.6. Modeling logical functions

Let us consider a shallow Neural Network with input layer $y^{[0]} \in \mathbb{N}^2$ and output layer $y^{[1]} \in \mathbb{N}$.

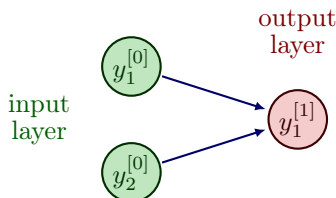


Figure 8. A simple perceptron (shallow Neural Network) with a two dimensional input.

We model true by the value 1 and false by the value 0, which results in the following truth table for the logical "OR" function.

input $y_1^{[0]}$	input $y_2^{[0]}$	$y_1^{[0]}$ OR $y_2^{[0]}$ (output $y_1^{[1]}$)
0	0	0
1	0	1
0	1	1
1	1	1

Table 1. Truth table for the logical "OR" function.

With Heaviside activation function, we have

$$y_1^{[1]} = \begin{cases} 1, & \text{if } W_1 y_1^{[0]} + W_2 y_2^{[0]} + b \geq 0, \\ 0, & \text{else.} \end{cases}$$

The goal is now to choose W_1, W_2, b so that we match the output from the truth table for given input. Obviously, $W_1 = W_2 = 1$ and $b = -1$ is a possible choice that fulfills the task. Similarly, one can find values for W_1, W_2 and b to model the logical "AND" function.

Next, let us consider the logical "XOR" function with the following truth table.

input $y_1^{[0]}$	input $y_2^{[0]}$	$y_1^{[0]}$ XOR $y_2^{[0]}$ (output $y_1^{[1]}$)
0	0	0
1	0	1
0	1	1
1	1	0

Table 2. Truth table for the logical "XOR" function.

In fact, the logical "XOR" function can not be represented by the given shallow Neural Network, since the data is not linearly separable, see e.g. [15, Section 6.1]. This motivates the introduction of additional layers in between the input and output layer, i.e. we choose a more complex function \mathcal{F} in the learning problem (P).

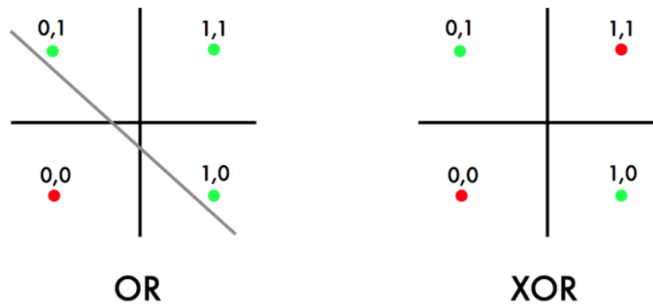


Figure 9. This illustration shows that the logical "OR" function is linearly separable, while the logical "XOR" function is not. Image modified from: <https://dev.to/jbahire/demystifying-the-xor-problem-1blk>.

2. Feedforward Neural Network

Introducing **hidden layers**, i.e. layers between the input and output layer, leads to **Feedforward Neural Networks (FNNs)**, also called **multilayer perceptrons (MLPs)**, cf. [15, Section 6]. Essentially, they are multiple perceptrons organized in layers, $\ell = 0, \dots, L$, where every perceptron takes the output from the previous perceptron as input. The number of layers L is called the **depth** of the network, while the number of neurons per layer n_ℓ is the **width** of the network. The input layer is denoted with $y^{[0]} = u \in \mathbb{R}^{n_0}$ and not counted in the depth of the network. A FNN is called **deep** if it has at least two hidden layers. We now indicate the weights from layer ℓ to $\ell + 1$ by $W^{[\ell]} \in \mathbb{R}^{n_{\ell+1} \times n_\ell}$ and the bias vector by $b^{[\ell]} \in \mathbb{R}^{n_{\ell+1}}$ for $\ell = 0, \dots, L - 1$. To simplify notation, we extend the activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ to vector valued inputs, by applying it component-wise, so that $\sigma : \mathbb{R}^n \rightarrow \mathbb{R}^n$, with $(y_1, \dots, y_n)^\top \mapsto (\sigma(y_1), \dots, \sigma(y_n))^\top$. The FNN layers can be represented in the same way as perceptrons

$$y^{[\ell]} = f_\ell(y^{[\ell-1]}) = \sigma^{[\ell]}(W^{[\ell-1]}y^{[\ell-1]} + b^{[\ell-1]}) \quad \text{for } \ell = 1, \dots, L, \quad (2)$$

where the activation function $\sigma^{[\ell]}$ may differ from layer to layer. We call $y^{[\ell]}$ the **feature vector** of layer ℓ . Compactly, we can write a FNN as a composition of its layer functions, cf. [1, 2]

$$y^{[L]} = \mathcal{F}(u) = f_L \circ f_{L-2} \circ \dots \circ f_1(u).$$

This formulation reinforces the choice of nonlinear activation function σ , cf. Remark 1.1. Otherwise, the output $y^{[L]}$ is linearly dependent on the input u and hidden layers can be eliminated. Hence, with linear activation function, solving rather simple tasks like modeling the logical "XOR" function will not be possible. However, sometimes it can be favorable to have one linear layer.

Remark 2.1 *In practice, it is not uncommon that the last layer of a FNN is indeed linear, i.e. $y^{[L]} = W^{[L-1]}y^{[L-1]}$. As long as the previous layers are nonlinear this does not hinder the expressiveness of the FNN, and is typically used to attain a desired output dimension. In essence, $W^{[L-1]}$ can be seen as a reformatting, in this case.*

However, in the remainder of this section we will consider the FNN architecture as introduced in (2). Let us now try again to represent the "XOR" logical function, this time by a FNN with one hidden layer.

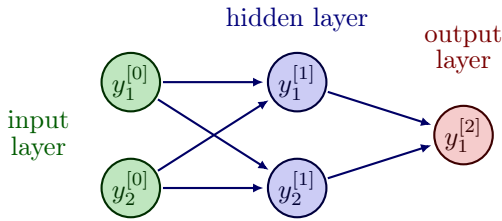


Figure 10. A feedforward network with a two dimensional input and one hidden layer with 2 nodes, i.e. $n_0 = n_1 = 2, n_2 = 1$ and $L = 2$.

The variable choices

$$W^{[0]} = \begin{pmatrix} 1 & 1 \\ -1 & -1 \end{pmatrix}, \quad b^{[0]} = \begin{pmatrix} -1 \\ 1 \end{pmatrix}, \quad W^{[1]} = (1 \quad 1), \quad b^{[1]} = -2,$$

solve the task and lead to the following truth table.

input $y_1^{[0]}$	input $y_2^{[0]}$	$y_1^{[1]}$	$y_2^{[1]}$	$y_1^{[0]}$ XOR $y_2^{[0]}$ (output $y_1^{[2]}$)
0	0	0	1	0
1	0	1	1	1
0	1	1	1	1
1	1	1	0	0

Table 3. Truth table for the logical "XOR" function modeled by the FNN from Figure 10 and given variable choices as above.

Next, we formulate an optimization problem similar to (P) for multiple layers with the above introduced notation, cf. [1, 2, 3]

$$\begin{aligned} \min_{\{W^{[\ell]}\}_\ell, \{b^{[\ell]}\}_\ell} \mathcal{L} \left(\{y^{[L(i)]}\}_i, \{u^{(i)}\}_i, \{W^{[\ell]}\}_\ell, \{b^{[\ell]}\}_\ell \right) & \quad (P_\ell) \\ \text{s.t.} \quad y^{[L(i)]} = \mathcal{F} \left(u^{(i)}, \{W^{[\ell]}\}_\ell, \{b^{[\ell]}\}_\ell \right). & \end{aligned}$$

If we collect again all variables in one vector θ this will have the following length:

$$\underbrace{n_0 \cdot n_1 + \dots + n_{L-1} \cdot n_L}_{\text{weights}} + \underbrace{n_1 + \dots + n_L}_{\text{biases}}$$

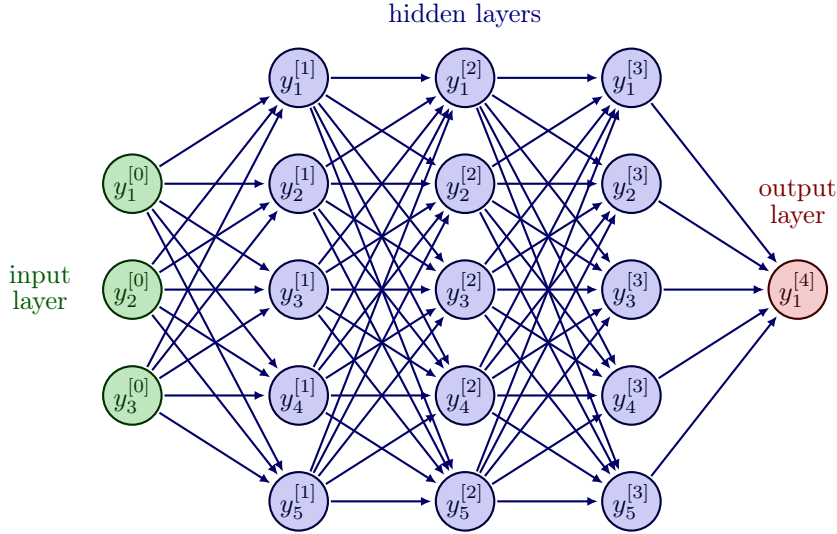


Figure 11. A feedforward network with 3 hidden layers, layer widths $n_0 = 3, n_1 = n_2 = n_3 = 5, n_4 = 1$ and depth $L = 4$. Collecting all variables of this network in a vector will give $\theta \in \mathbb{R}^{86}$.

The question that immediately arises is: How to choose the network architecture, i.e. depth and width of the FNN? The following discussion is based on [15, Section 6.4].

2.1. Depth and Width

The universal approximation theorem, see e.g. [10], states that any vector-valued, multivariate, measurable (in particular, continuous) function $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_L}$ can be approximated with arbitrary small error by a Neural Network with one hidden layer. Hence, a first approach may be to choose a network with depth $L = 2$ and increase the width until the desired accuracy is reached.

However, this poses certain problems. First of all the universal approximation theorem does not imply that a training algorithm will actually reach the desired approximation, but rather that some set of parameters exists, that satisfies the requirement. The training algorithm might for example only find a local minimum or choose the wrong function as a result of overfitting. Another problem may be the sheer size of the layer required to achieve the wanted accuracy. In the worst case the network will need an exponential number of hidden units, with one hidden unit for each possible combination of inputs. Thus in practice, one is discouraged to use only one hidden layer, but instead to use deep networks.

Various families of functions are efficiently approximated by deep networks with a smaller width. If one desires to approximate the same function to the same degree of accuracy, the number of hidden units typically grows exponentially. This stems from the fact, that each new layer allows the network to make exponentially more connections, thus allowing a wider output of target functions. Another reason why one might choose deeper networks is due to the intuition, that our desired function may well be a composition of multiple functions. Each new layer adds a nonlinear layer function to our network, thus making it easier for the FNN to approximate composite functions. Also, heuristically we observe that deep networks typically outperform shallow networks.

Nonetheless, one main problem with deep FNNs is, that the gradient used for training is the product of the partial derivatives of each layer, as we will see in Section 2.5. If these derivatives have small values, then the gradient for earlier layers can become very small. Thus training has a smaller if not even a negligible effect on the first layers when the network is too deep. This is especially a problem for sigmoid activation function, since its derivative is bounded by $\frac{1}{4}$.

As discussed in Section 1.5, choosing hyperparameters like the depth and width is a non-trivial undertaking and currently, the method of choice is experimenting to find a suitable configuration for the given task.

Additionally, in the optimization algorithms, cf. Section 1.3, we need a starting point θ^0 for the variables. Hence, we discuss how to initialize the weights and biases in the FNN.

2.2. Initialization

Recall that due to the non-convex loss function, we can only expect convergence to stationary points, cf. Remark 1.4. Consequently, the choice of the initial point θ^0 can have great impact on the algorithm, since two different initial points can lead to two different results. An unsuitable initial point may even prevent convergence altogether. Similar to choosing hyperparameters, for the choice of initial points there exist several well-tested strategies, cf. [12, Section 4.2.2], but it is still an active field of research.

The naive approach would be to initialize $\theta = 0$ or with some other constant value. Unfortunately, this strategy has major disadvantages. With this initialization, all weights per layer in the Neural Network have the same influence on the loss function and will therefore have the same gradient. This leads to all those neurons evolving symmetrically throughout training, so that different

neurons will not learn different things, which significantly reduces the expressiveness of the FNN. Let us remark that it is fine to initialize the biases $b^{[\ell]}$ with zero, as long as the weights $W^{[\ell]}$ are not initialized constant. Hence, it is sufficient to discuss initialization strategies for the weights.

We know now that the weights should be initialized in a way that they differ from each other to ensure **symmetry breaking**, i.e. preventing the neurons from evolving identically. One way to achieve this is **random initialization**. However, immediately the next question arises: How to generate those random values?

For example, weights can be drawn from a Gaussian distribution with mean zero and some fixed standard deviation. Choosing a small standard deviation, e.g. 0.01, may cause a problem known as vanishing gradients for deep networks, since the small neuron values will be multiplied with each other in the computation of gradients due to the chain rule, leading to exponentially decaying products. As a result learning can be very slow or even diverge. On the other hand, choosing a large standard deviation, e.g. 0.2, can result in exploding gradients, which is essentially the opposite problem, where the products grow exponentially and learning can become unstable, oscillate or even produce "NaN" values for the variables. Furthermore, in combination with saturating activation functions like sigmoid or tanh exploding variable values can lead to saturation of the activation function, which then leads to vanishing gradients and again hinder learning, cf. Figure 3.

To find a good intermediate value for the standard deviation, **Xavier initialization** has been proposed in [14], where the standard deviation value is chosen depending on the input size of the layer, i.e.

$$\frac{1}{\sqrt{n_\ell}}$$

for $W^{[\ell]}, \ell = 0, \dots, L - 1$. Note that since input sizes can vary, the Gaussian distribution that the initial values are drawn from, will also vary. This choice of weights in combination with $b^{[\ell]} = 0$ leads to $\text{Var}(y^{[\ell+1]}) = \text{Var}(y^{[\ell]})$. However, the Xavier initialization assumes zero centered activation functions, which is not fulfilled for sigmoid and all variants of ReLU. As a remedy, the **He initialization** has been proposed in [18], tailored especially to ReLU activation functions. Here, the standard deviation is also chosen depending on the input size of the layer, namely

$$\sqrt{\frac{2}{n_\ell}}.$$

Additionally, there also exists an approach to normalize feature vectors throughout the network.

2.3. Batch Normalization

Assume that we are employing an optimization algorithm, which passes through the data points in batches of size b and that the nodes in hidden layers follow a normal distribution. Then **batch normalization** [22] aims at normalizing the feature vectors in a hidden layer over the given batch, to stabilize training, especially for unbounded activation functions, such as ReLU. It can be seen as insertion of an additional layer in a deep neural network, and this layer type is already pre-implemented in learning frameworks like tensorflow and pytorch:

- `tf.keras.layers.BatchNormalization`,
- `torch.nn.BatchNorm1d` (also 2d and 3d available).

Say, we have given the current feature vectors $y_j^{[\ell]} \in \mathbb{R}^{n_\ell}$ of hidden layer ℓ for all elements in the batch, i.e. $j = 1, \dots, b$. The batch normalization technique first determines the mean and variance

$$\mu^{[\ell]} = \frac{1}{b} \sum_{j=1}^b y_j^{[\ell]}, \quad (\sigma^2)^{[\ell]} = \frac{1}{b} \sum_{j=1}^b \|y_j^{[\ell]} - \mu^{[\ell]}\|^2.$$

Subsequently, the current feature vectors $y_j^{[\ell]}, j = 1, \dots, b$ are normalized via

$$\hat{y}_j^{[\ell]} = \frac{y_j^{[\ell]} - \mu^{[\ell]}}{\sqrt{(\sigma^2)^{[\ell]} + \epsilon}},$$

where $\epsilon \in \mathbb{R}$ is a constant that helps with numerical stability. Finally, the output of the batch normalization layer is computed by

$$y_j^{[\ell+1]} = W^{[\ell]} \hat{y}_j^{[\ell]} + b^{[\ell]} \quad \forall j = 1, \dots, b.$$

As usual, the weight $W^{[\ell]} \in \mathbb{R}^{n_{\ell+1} \times n_\ell}$ and bias $b^{[\ell]} \in \mathbb{R}^{n_{\ell+1}}$ are variables of the neural network and will be learned.

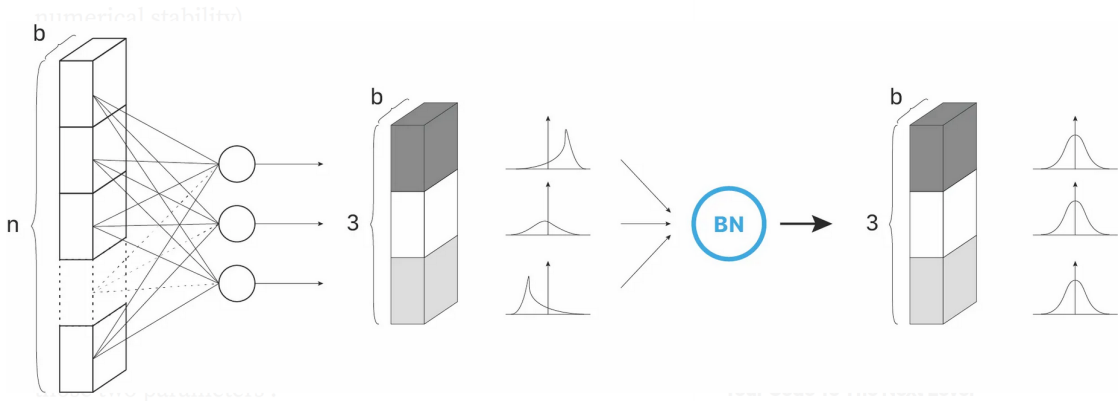


Figure 12. Illustration of Batch Normalization applied to a hidden layer with 3 nodes and batch size b . We assume that every node can be modeled by a normal distribution. Image Source: <https://towardsdatascience.com/batch-normalization-in-3-levels-of-understanding-14c2da90a338>.

Remark 2.2 *The BN layer is a trainable layer, since it contains learnable variables.*

Let us now consider an example task that is commonly solved with FNNs.

2.4. Classification Tasks

As mentioned in 1.1, classification is a supervised learning task with labels $S(u) \in \mathbb{N}$. First, we consider **binary classification**, cf. e.g. [12, Section 2.1], where we classify the data into two

categories, i.e. a spam filter that determines whether an email is spam or not. We could construct our network so that it indicates the category which it concludes is the most likely. However, it can be useful to know the probability assigned to the output, so that we know how certain the decision is. Thus, we aim to have outputs between 0 and 1, so that they sum up to 1. In the case of binary classification, the second output is directly determined by the first. Consequently, it suffices to have a one dimensional output layer $y^{[L]} \in [0, 1]$, which should predict

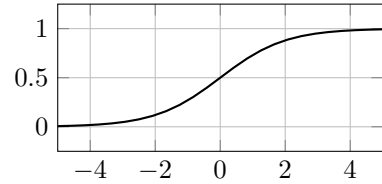
$$P(y^{[L]} = 1 | u, \theta),$$

i.e. the probability of the output being category 1 (e.g. spam) given the input u and variables θ .

Assume that we have already set up a feedforward network up to the second to last layer $y^{[L-1]} \in \mathbb{R}$. It remains to choose the activation function $\sigma^{[L-1]}$ that enters the computation of $y^{[L]}$ and to model the loss function \mathcal{L} .

Since we want $y^{[L]} \in [0, 1]$, a common approach is to use the sigmoid activation function.

$$\sigma(y) = \frac{1}{1 + \exp(-y)} = \frac{\exp(y)}{\exp(y) + 1} \in (0, 1).$$



Let us remark that the cases $y^{[L]} \in \{0, 1\}$ are not possible with this choice of activation function, but we are only computing approximations anyhow.

Next, we construct a loss function. To this end, we assume that the training data is a sample of the actual relationship we are trying to train, thus it obeys a probability function, which we want to recover. The main idea for the loss function is to maximize the likelihood of the input parameters, i.e. if the probability $P(y^{[L]} = S(u) | u, \theta)$ of generating the known supervision $S(u)$ is high for the input data u , the loss should be small, and vice versa. To model the probability we choose the Bernoulli distribution, which models binary classification:

$$P(y^{[L]} = S(u) | u, \theta) = (y^{[L]})^{S(u)}(1 - y^{[L]})^{(1-S(u))}.$$

To achieve small loss for large probabilities, we apply the logarithm and then maximize this function, so that the optimal network variables $\bar{\theta}$ can be determined as follows

$$\begin{aligned} \bar{\theta} &= \operatorname{argmax}_{\theta} \sum_{i=1}^N \log \left(P \left(y^{[L](i)} = S(u^{(i)}) | u^{(i)}, \theta \right) \right) \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^N \log \left((y^{[L](i)})^{S(u^{(i)})} (1 - y^{[L](i)})^{(1-S(u^{(i)}))} \right) \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^N S(u^{(i)}) \cdot \log(y^{[L](i)}) + (1 - S(u^{(i)})) \cdot \log(1 - y^{[L](i)}) \\ &= \operatorname{argmin}_{\theta} \sum_{i=1}^N \underbrace{-S(u^{(i)}) \cdot \log(y^{[L](i)}) - (1 - S(u^{(i)})) \cdot \log(1 - y^{[L](i)})}_{=: \mathcal{L}(y^{[L](i)}, S(u^{(i)})), \text{ Binary Cross Entropy Loss}}, \end{aligned}$$

where $y^{[L](i)}$ is a function of the network variables θ .

In practice, we minimize the cross-entropy, since it is equivalent to maximizing the likelihood, but stays within our given frame of minimization problems. Let us assume that our data either has the label $S(u) = 1$ (spam) or $S(u) = 0$ (not spam), then the binary cross entropy loss is

$$\mathcal{L}(y^{[L]}, S(u)) = \begin{cases} -\log(y^{[L]}), & \text{if } S(u) = 1, \\ -\log(1 - y^{[L]}), & \text{if } S(u) = 0. \end{cases}$$

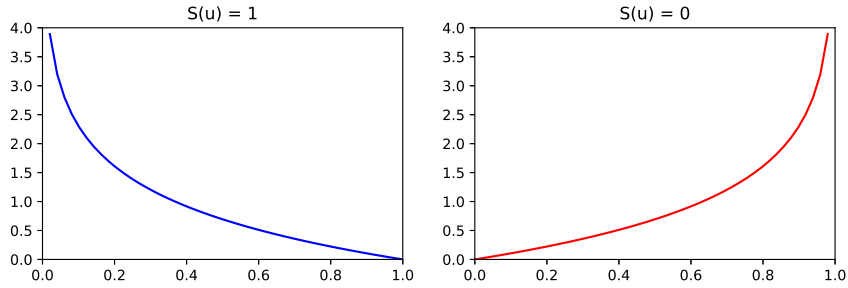


Figure 13. Binary Cross Entropy Loss for label $S(u) = 1$ (left) and label $S(u) = 0$ (right).

We see in Figure 13 that the cross entropy loss for $S(u) = 1$ goes to zero, for $y^{[L]} \nearrow 1$, and grows for $y^{[L]} \searrow 0$, as desired. On the other hand, the cross entropy loss for $S(u) = 0$ goes to zero for $y^{[L]} \searrow 0$ and grows for $y^{[L]} \nearrow 1$.

Altogether, we have the following loss function for the binary classification task

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{[L(i)]}(\theta), S(u^{(i)})).$$

Multiclass classification is a direct extension of binary classification. Here, the goal is to classify the data into multiple (at least 3) categories. A prominent example is the MNIST data set, where black and white pictures of handwritten digits (of size 28×28 pixels) are supposed to be classified as $\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$, cf. Figure 14. A possible network architecture is illustrated in Figure 15.

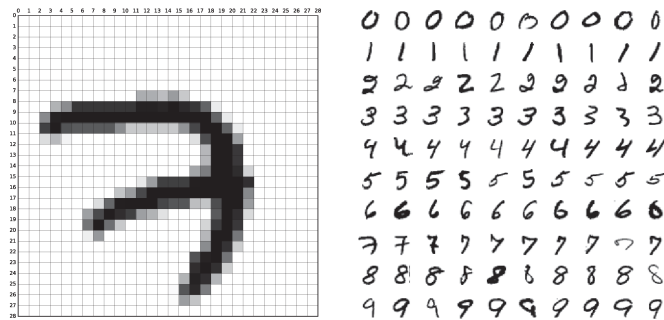


Figure 14. Example of the MNIST database. Sample belonging to the digit 7 (left) and 100 samples from all 10 classes (right). Image Source: [4, Fig. 1].

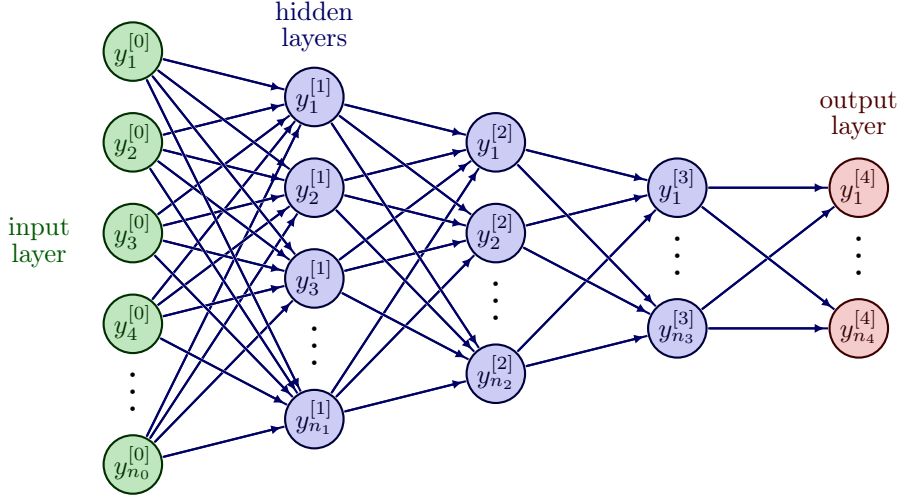


Figure 15. A feedforward network with 3 hidden layers, i.e. depth $L = 4$. For the MNIST data set we have $n_0 = 784$ and $n_4 = 10$.

For this task, we need a generalization of the sigmoid activation function, which will take $y^{[L-1]} \in \mathbb{R}^n$ and map it to $y^{[L]} \in [0, 1]^n$, so that we have $\sum_{i=1}^n y_i^{[L]} = 1$, where $n_{L-1} = n_L = n$ is the number of classes. A suitable option is the softmax function, which is given component-wise by

$$\text{softmax}(y)_i = \frac{\exp(y_i)}{\sum_{j=1}^n \exp(y_j)} \in (0, 1), \quad \text{for } i = 1, \dots, n.$$

Keep in mind that e.g. in the MNIST case we have labels $S(u) \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$, and in general we have multiple labels $S(u) \in \mathbb{N}$. We have seen before that maximizing the log-likelihood is a suitable choice for classification tasks. Since we want to formulate a minimization problem, we choose the negative log-likelihood as loss function

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left(P \left(y^{[L](i)} = S(u^{(i)}) \mid u^{(i)}, \theta \right) \right).$$

In this section we have seen yet another model for the loss function \mathcal{L} , and from Section 1.3 we know that in any case we will need the gradient $\nabla \mathcal{L}(\theta)$ to update our variables θ . Let us discuss how frameworks like pytorch and tensorflow obtain this information.

2.5. Backpropagation

The derivations are based on [15, Section 6.5] and [26, Section 7.3]. When a network, e.g. a FNN, takes an input u , passes it through its layers and finally computes an output $y^{[L]}$, the network **feeds forward** the information, which this is called **forward propagation**. Then a loss is assigned to the output and we aim at employing the gradient of the loss function $\nabla \mathcal{L}(\theta)$ to update the network variables $\theta \in \mathbb{R}^K$. In a FNN we have $K = n_0 \cdot n_1 + \dots + n_{L-1} \cdot n_L + n_1 + \dots + n_L$.

The gradient is then given by

$$\nabla \mathcal{L}(\theta) = \begin{pmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial \theta_K} \end{pmatrix}.$$

To develop an intuition about the process, we discuss the following simple example.

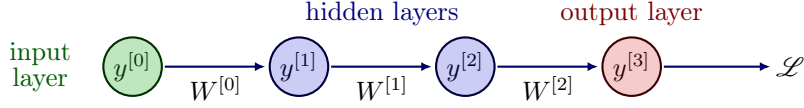


Figure 16. A feedforward network with 2 hidden layers, one node per layer and depth $L = 3$.

Example 2.3 Consider a very simple FNN with one node per layer and assume that we only consider weights $W^{[\ell]} \in \mathbb{R}$ and no biases. For the network in Figure 16 we have $\theta = (W^{[0]}, W^{[1]}, W^{[2]})^\top$,

$$y^{[3]} = \sigma^{[3]}(W^{[2]}y^{[2]}) = \sigma^{[3]}(W^{[2]}\sigma^{[2]}(W^{[1]}y^{[1]})) = \sigma^{[3]}(W^{[2]}\sigma^{[2]}(W^{[1]}\sigma^{[1]}(W^{[0]}y^{[0]}))),$$

and

$$\mathcal{L}(\theta) = \mathcal{L}(y^{[3]}(\theta)) = \mathcal{L}\left(\sigma^{[3]}(W^{[2]}\sigma^{[2]}(W^{[1]}\sigma^{[1]}(W^{[0]}y^{[0]})))\right).$$

Computing the components of the gradient, we employ the chain rule to obtain e.g.

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W^{[0]}} &= \frac{\partial \mathcal{L}}{\partial y^{[3]}} \cdot \frac{\partial y^{[3]}}{\partial W^{[0]}} \\ &= \frac{\partial \mathcal{L}}{\partial y^{[3]}} \cdot \frac{\partial y^{[3]}}{\partial y^{[2]}} \cdot \frac{\partial y^{[2]}}{\partial W^{[0]}} \\ &= \frac{\partial \mathcal{L}}{\partial y^{[3]}} \cdot \frac{\partial y^{[3]}}{\partial y^{[2]}} \cdot \frac{\partial y^{[2]}}{\partial y^{[1]}} \cdot \frac{\partial y^{[1]}}{\partial W^{[0]}}, \end{aligned}$$

and in general for depth L we get

$$\frac{\partial \mathcal{L}}{\partial W^{[\ell]}} = \frac{\partial \mathcal{L}}{\partial y^{[L]}} \cdot \prod_{j=L}^{\ell+2} \frac{\partial y^{[j]}}{\partial y^{[j-1]}} \cdot \frac{\partial y^{[\ell+1]}}{\partial W^{[\ell]}}.$$

Essentially, to calculate the effect of a variable on the loss function we iterate backwards through the network, multiplying the derivatives of each layer. This is called **back propagation**, often abbreviated as **backprop**.

To obtain a computationally efficient version of back propagation, we exploit the fact that parts of the derivatives can be recycled, broadly speaking. E.g. $\frac{\partial \mathcal{L}}{\partial y^{[L]}}$ is a part of all derivatives. So, if we compute the derivative by $W^{[L-1]}$ first, we already have this component available and can reuse it in the computation of the derivative by $W^{[L-2]}$, etc.

In order to formalize the effective computation of derivatives in a backpropagation algorithm, we decompose the forward propagation into two parts, cf. e.g. [26, Section 7.3.2].

$$\begin{aligned} z^{[\ell]} &= W^{[\ell-1]}y^{[\ell-1]} + b^{[\ell-1]} && \in \mathbb{R}^{n_\ell}, \\ y^{[\ell]} &= \sigma^{[\ell]}(z^{[\ell]}) && \in \mathbb{R}^{n_\ell}. \end{aligned}$$

This was not necessary in Example 2.3, because we only consider weights and no biases. Furthermore, we assume that the loss function \mathcal{L} takes the final output $y^{[L]}$ as an input. Especially, no other feature vectors $y^{[\ell]}$ for $\ell \neq L$ enter the loss function directly. This is the case e.g. for mean squared error, cf. Example 1.2, and cross entropy, cf. Section 2.4.

In general, we now have by chain rule for all $\ell = 0, \dots, L-1$

$$\frac{\partial \mathcal{L}}{\partial W^{[\ell]}} = \frac{\partial \mathcal{L}}{\partial y^{[L]}} \cdot \prod_{j=L}^{\ell+2} \left(\frac{\partial y^{[j]}}{\partial z^{[j]}} \cdot \frac{\partial z^{[j]}}{\partial y^{[j-1]}} \right) \cdot \frac{\partial y^{[\ell+1]}}{\partial z^{[\ell+1]}} \cdot \frac{\partial z^{[\ell+1]}}{\partial W^{[\ell]}}, \quad (3)$$

$$\frac{\partial \mathcal{L}}{\partial b^{[\ell]}} = \frac{\partial \mathcal{L}}{\partial y^{[L]}} \cdot \prod_{j=L}^{\ell+2} \left(\frac{\partial y^{[j]}}{\partial z^{[j]}} \cdot \frac{\partial z^{[j]}}{\partial y^{[j-1]}} \right) \cdot \frac{\partial y^{[\ell+1]}}{\partial z^{[\ell+1]}} \cdot \frac{\partial z^{[\ell+1]}}{\partial b^{[\ell]}}. \quad (4)$$

However, we have to understand these derivatives in detail. First of all, let us introduce the following definition from [21].

Definition 2.4 Let $A, B \in \mathbb{R}^{m \times n}$ be given matrices, then $A \odot B \in \mathbb{R}^{m \times n}$, with entries

$$(A \odot B)_{i,j} := (A)_{i,j} \cdot (B)_{i,j}, \quad \text{for } i = 1, \dots, m, j = 1, \dots, n,$$

is called the **Hadamard product** of A and B .

Furthermore, we define the derivative of the component-wise activation function $\sigma : \mathbb{R}^m \rightarrow \mathbb{R}^m$ as follows

$$\sigma' : \mathbb{R}^m \rightarrow \mathbb{R}^m, \quad z = \begin{pmatrix} z_1 \\ \vdots \\ z_m \end{pmatrix} \mapsto \begin{pmatrix} \sigma'(z_1) \\ \vdots \\ \sigma'(z_m) \end{pmatrix} = \sigma'(z).$$

Let us introduce two special cases of multi-dimensional chain rule, cf. [26, p.98], which will prove helpful to calculate the derivatives.

1. Consider $a = \sigma(z) \in \mathbb{R}^m$, where σ is a component-wise function, e.g. an activation function, $z \in \mathbb{R}^m$, and $f = f(a) \in \mathbb{R}$. Then, it holds

$$\frac{\partial f}{\partial z} = \frac{\partial f}{\partial a} \odot \sigma'(z) \in \mathbb{R}^m. \quad (5)$$

2. Consider $z = Wy + b \in \mathbb{R}^m$ and $f = f(z) \in \mathbb{R}$, with $W \in \mathbb{R}^{m \times n}$ and $y \in \mathbb{R}^n$. Then, it holds

$$\frac{\partial f}{\partial y} = W^\top \cdot \frac{\partial f}{\partial z} \in \mathbb{R}^n, \quad (6)$$

$$\frac{\partial f}{\partial W} = \frac{\partial f}{\partial z} \cdot y^\top \in \mathbb{R}^{m \times n}, \quad (7)$$

$$\frac{\partial f}{\partial b} = \frac{\partial f}{\partial z} \in \mathbb{R}^m. \quad (8)$$

We can now start working our way backwards through the network to get all derivatives. Assume that we know $\frac{\partial \mathcal{L}}{\partial y^{[L]}} \in \mathbb{R}^{n_L}$, which will depend in detail on the choice of loss function. We can employ (5) with $f = \mathcal{L}, z = z^{[L]}, a = y^{[L]}$ to compute

$$\bar{z}^{[L]} := \frac{\partial \mathcal{L}}{\partial z^{[L]}} = \frac{\partial \mathcal{L}}{\partial y^{[L]}} \odot (\sigma^{[L]})'(z^{[L]}) \in \mathbb{R}^{n_L}.$$

Here, we employ a typical notation from automatic differentiation (AD), i.e. the gradient of the loss with respect to a certain variable is denoted by the name of that variable with an overbar. Now, we know

$$\bar{y}^{[L-1]} := \frac{\partial \mathcal{L}}{\partial y^{[L-1]}} = \bar{z}^{[L]} \cdot \frac{\partial z^{[L]}}{\partial y^{[L-1]}} \in \mathbb{R}^{n_{L-1}},$$

i.e. we can reuse the previously derived gradient. Furthermore, from (6) with $f = \mathcal{L}$, $y = y^{[L-1]}$, $W = W^{[L-1]}$, $z = z^{[L]}$ we deduce

$$\bar{y}^{[L-1]} = (W^{[L-1]})^\top \bar{z}^{[L]}.$$

Subsequently, we use $\bar{y}^{[L-1]}$ to compute $\bar{z}^{[L-1]}$, and so forth. In this way we can keep iterating to build up the products in (3) and (4).

In every layer, $\ell = 0, \dots, L-1$, we also want to determine $\bar{W}^{[\ell]} := \frac{\partial \mathcal{L}}{\partial W^{[\ell]}}$ and $\bar{b}^{[\ell]} := \frac{\partial \mathcal{L}}{\partial b^{[\ell]}}$. We show this exemplary for $\ell = L-1$. It holds

$$\begin{aligned} \bar{W}^{[L-1]} &= \bar{z}^{[L]} \cdot \frac{\partial z^{[L]}}{\partial W^{[L-1]}} \in \mathbb{R}^{n_L \times n_{L-1}}, \\ \bar{b}^{[L-1]} &= \bar{z}^{[L]} \cdot \frac{\partial z^{[L]}}{\partial b^{[L-1]}} \in \mathbb{R}^{n_L}. \end{aligned}$$

Making use of (7) and (8) with the same choices as in the computation of $\bar{y}^{[L-1]}$ and $b = b^{[L-1]}$, we get

$$\begin{aligned} \bar{W}^{[L-1]} &= \bar{z}^{[L]} (y^{[L-1]})^\top, \\ \bar{b}^{[L-1]} &= \bar{z}^{[L]}. \end{aligned}$$

With this technique, we have an iterative way to efficiently calculate all gradients needed for the variable update in the optimization method, cf. Section 1.3.

Remark 2.5

- (i) *It is even more elegant and efficient to update $W^{[\ell]}$ and $b^{[\ell]}$ during backpropagation, i.e. on the fly. This way we do not need to store the gradient and can overwrite the weight and bias variables. It is only necessary to save the current gradients for the next loop, so we could rewrite the backpropagation algorithm with temporary gradient values. This only works if the stepsize / learning rate τ is previously known and fixed for all variables, since for a line search we would need to know the full gradient and could only update the variables afterwards.*
- (ii) *Considering we have N training data points, the backpropagation algorithm has to take all of them into account. When the loss is a sum of loss functions for each data point, this can be easily incorporated into the algorithm by looping over $i = 1, \dots, N$ and introducing a sum where necessary.*

Altogether, we formulate Algorithm 6, which collects the gradients with respect to the weights and biases in one final gradient vector $\nabla \mathcal{L}(\theta)$.

In frameworks like pytorch and tensorflow, backpropagation is already implemented, e.g. in pytorch the function "autograd" handles the backward pass. Broadly speaking, autograd collects

Algorithm 6 Backpropagation.

Require: Training data set $\{u^{(i)}, S(u^{(i)})\}_{i=1}^N$.

Require: Current weights $W^{[\ell]}$ and biases $b^{[\ell]}$ for $\ell = 0, \dots, L - 1$.

Require: Activation functions $\sigma^{[\ell]}$ for $\ell = 1, \dots, L$.

Require: Loss function $\mathcal{L}(y^{[L]})$ and its gradient $\nabla \mathcal{L}(y^{[L]}) = \frac{\partial \mathcal{L}}{\partial y^{[L]}} \in \mathbb{R}^{n_L}$.

$$y^{[0](i)} = u^{(i)} \in \mathbb{R}^{n_0} \quad \text{for } i = 1, \dots, N.$$

for $\ell = 1, \dots, L$ **do**

$$z^{[\ell](i)} = W^{[\ell-1]} y^{[\ell-1]} + b^{[\ell-1]} \in \mathbb{R}^{n_\ell} \quad \text{for } i = 1, \dots, N,$$

$$y^{[\ell](i)} = \sigma^{[\ell]}(z^{[\ell](i)}) \in \mathbb{R}^{n_\ell} \quad \text{for } i = 1, \dots, N.$$

end for

$$\text{Compute loss } \mathcal{L} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{[L](i)}) \in \mathbb{R}.$$

$$\bar{y}^{[L](i)} = \frac{1}{N} \cdot \nabla \mathcal{L}(y^{[L](i)}) \in \mathbb{R}^{n_L} \quad \text{for } i = 1, \dots, N.$$

for $\ell = L, L - 1, \dots, 1$ **do**

$$\bar{z}^{[\ell](i)} = \bar{y}^{[\ell](i)} \odot (\sigma^{[\ell]})'(z^{[\ell](i)}) \in \mathbb{R}^{n_\ell} \quad \text{for } i = 1, \dots, N,$$

$$\bar{y}^{[\ell-1](i)} = (W^{[\ell-1]})^\top \bar{z}^{[\ell](i)} \in \mathbb{R}^{n_{\ell-1}} \quad \text{for } i = 1, \dots, N,$$

$$\bar{W}^{[\ell-1]} = \sum_{i=1}^N \bar{z}^{[\ell](i)} (y^{[\ell-1](i)})^\top \in \mathbb{R}^{n_\ell \times n_{\ell-1}},$$

$$\bar{b}^{[\ell-1]} = \sum_{i=1}^N \bar{z}^{[\ell](i)} \in \mathbb{R}^{n_\ell}.$$

end for

the data and all executed operations in a directed acyclic graph. In this graph the inputs are the leaves, while the outputs are the roots. Now to automatically compute the gradients, the graph can be traced from roots to leaves, employing the chain rule. This coincides with the computations that we just derived by hand. For more details we refer to https://pytorch.org/tutorials/beginner/blitz/autograd_tutorial.html.

3. Convolutional Neural Network

In this section, based on [9],[15, Section 9], we consider Neural Networks with a different architecture: **convolutional neural networks** (CNNs or ConvNets). They were first introduced by Kunihiko Fukushima in 1980 under the name "neocognitron" [11]. Famous examples of convolutional neural networks today are "LeNet" [25], see Figure 17 and "AlexNet" [24].

As a motivation, consider a classification task where the input is an image of size $n_{0,1} \times n_{0,2}$ pixels. We want to train a Neural Network so that it can decide, e.g. which digit is written in the image (MNIST data set). We have seen in Figure 15 that the image with $n_{0,1} = n_{0,2} = 28$ has been reshaped (vectorized, flattened) into a vector in $\mathbb{R}^{n_{0,1} \cdot n_{0,2}} = \mathbb{R}^{784}$, so that we can use it as an input for a regular FNN. However, this approach has several disadvantages:

1. Vectorization causes the input image to lose all of its spatial structure, which could have been helpful during training.
2. Let e.g. $n_{0,1} = n_{0,2} = 1000$, then $n_0 = 10^6$ and the weight matrix $W^{[0]} \in \mathbb{R}^{n_1 \times 10^6}$ contains an enormous number of optimization variables. This can make training very slow or even infeasible.

On the contrary, convolutional neural networks are designed to exploit the relationships between neighboring pixels. In fact, the input of a CNN is typically a matrix or even a three-dimensional tensor, which is then passed through the layers while maintaining this structure. CNNs take small patches, e.g. squares or cubes, from the input images and learn features from them. Consequently, they can subsequently recognize these features in other images, even when they appear in other parts of the image.

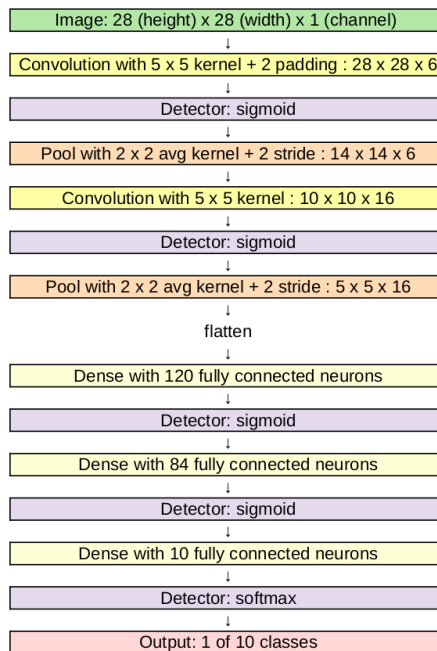


Figure 17. Architecture of LeNet-5.

In Figure 17 we see the architecture of "LeNet-5". The inputs are images, where we have 1 channel, because we consider grayscale images. At first we have two sequences of convolution layer (yellow), Section 3.2, detector layer (violet), Section 3.3, and pooling layer (orange), Section 3.4. These layers retain the multidimensional structure of the input. Since this network is built for a classification tasks, the output should be a vector of 10. Consequently, the multi-dimensional output of a hidden layer is flattened, i.e. vectorized, and the remaining layers are fully connected layers (bright yellow) as we have seen in FNNs.

In other, larger architectures, like AlexNet, cf. Figure 24, to avoid overfitting with large fully connected layers, a technique called **dropout** is applied. The key idea is to randomly drop units with a given probability and their connections from the neural network during training, for more details we refer to [28].

Remark 3.1

- (i) We will view convolution, detector and pooling layers as separate layers. However, it is also possible to define a convolutional layer to consist of a convolution, detector and pooling stage, cf. Figure 18. This can be a source of confusion when referring to convolutional layers, which we should be aware of.
- (ii) Throughout the remainder of this section we omit layer indices ℓ to simplify notation, and we indicate the data with capital letter Y to clarify that they are matrices or tensors.

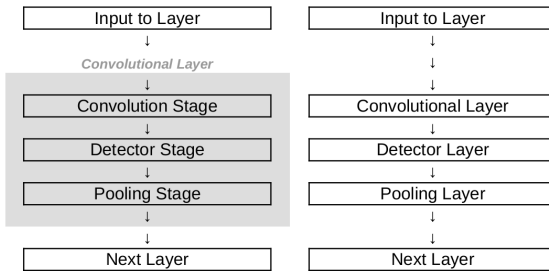


Figure 18. Convolutional layer (gray) consisting of stages (left) compared to viewing the operations as separate layers (right). We use the terminology as depicted on the right hand side, and refer to convolutional, detector and pooling layers as separate layers.

Before we move on to a detailed introduction of the different layer types in CNNs, let us recall the mathematical concept of a convolution.

3.1. Convolution

As explained in [15, Section 9.1], in general, convolution describes how one function influences the shape of another function. But it can also be used to apply a weight function to another function, which is how convolution is used in convolutional neural networks.

Definition 3.2 Let $f, g : \mathbb{R}^n \rightarrow \mathbb{R}$ be two functions. If both f and g are integrable with respect to Lebesgue measure, we can define the **convolution** as:

$$c(t) = (f * g)(t) = \int f(x)g(t - x) dx,$$

for some $t \in \mathbb{R}^n$. Here, f is called the input and g is called the kernel. The new function $c : \mathbb{R}^n \rightarrow \mathbb{R}$ is called the feature map.

However, for convolutional neural networks we need the discrete version.

Definition 3.3 Let $f, g : \mathbb{Z}^n \rightarrow \mathbb{R}$ be two discrete functions. The **discrete convolution** is then defined as:

$$c(t) = (f * g)(t) = \sum_{x \in \mathbb{Z}^n} f(x)g(t - x),$$

for some $t \in \mathbb{Z}^n$.

A special case of the discrete convolution is setting f and g to n -dimensional vectors and using the indices as arguments. We illustrate this approach in the following example.

Example 3.4 Let X and Y be two random variable each describing the outcome of rolling a dice. The probability mass functions are defined as:

$$f_X(t) = f_Y(t) = \begin{cases} \frac{1}{6}, & \text{if } t \in \{1, 2, 3, 4, 5, 6\}, \\ 0, & \text{if } t \in \mathbb{Z} \setminus \{1, 2, 3, 4, 5, 6\}. \end{cases}$$

We aim at calculating the probability that the sum of both dice rolls equals nine. To this end, we take the vectors of all possible outcomes and arrange them into two rows. Here, we flip the second vector and slide it to the right, such that the numbers which add to nine align.

1	2	3	4	5	6		
		6	5	4	3	2	1

Now, we replace the outcomes with their respective probabilities, multiply the adjacent components and add up the results.

$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$		
		$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

This gives

$$f_{X+Y}(9) = \frac{1}{36} + \frac{1}{36} + \frac{1}{36} + \frac{1}{36} = \frac{1}{9},$$

i.e. the probability that the sum of the dice equals nine is $\frac{1}{9}$.

In fact, all the steps we have just done are equivalent to calculating a discrete convolution:

$$f_{X+Y}(9) = \sum_{x=1}^6 f_X(x)f_Y(9-x) = (f_X * f_Y)(9)$$

3.2. Convolutional Layer

For the convolutional layers in CNNs we define convolutions for matrices, cf. e.g. [15, (9.4)]. This can be extended to tensors straight forward.

Definition 3.5 Let $Y \in \mathbb{R}^{n_1 \times n_2}$ and $K \in \mathbb{R}^{m_1 \times m_2}$ be given matrices, such that $m_1 \leq n_1$ and $m_2 \leq n_2$. The **convolution** of Y and K is denoted by $Y * K$ with entries

$$[Y * K]_{i,j} := \sum_{k=1}^{m_1} \sum_{l=1}^{m_2} K_{k,l} Y_{i+m_1-k, j+m_2-l},$$

for $1 \leq i \leq n_1 - m_1 + 1$ and $1 \leq j \leq n_2 - m_2 + 1$. Here, Y is called the input and K is called the kernel.

In Machine Learning often the closely related concept of **(cross) correlation**, cf. e.g. [15, (9.6)], is used, and incorrectly referred to as convolution, where

$$[Y \circledast K]_{i,j} := \sum_{k=1}^{m_1} \sum_{l=1}^{m_2} K_{k,l} Y_{i-1+k, j-1+l},$$

for $1 \leq i \leq n_1 - m_1 + 1$ and $1 \leq j \leq n_2 - m_2 + 1$. The (cross) correlation has the same effect as convolution, if you flip both, rows and columns of the kernel K , see the changed indices indicated in red. Since we learn the kernel anyway, it is irrelevant whether the kernel is flipped, thus either concept can be used.

We illustrate the matrix computations with an example.

Example 3.6 For this example we have the data matrix

$$Y = \begin{pmatrix} 1 & 5 & -2 & 0 & 2 \\ 3 & 8 & 7 & 1 & 0 \\ -1 & 0 & 1 & 2 & 3 \\ 4 & 2 & 1 & -1 & 2 \end{pmatrix},$$

and the kernel

$$K = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix}.$$

The computation of $[Y * K]_{1,1}$ can be illustrated as follows

$$[Y * K]_{1,1} = \begin{pmatrix} +1 \cdot 9 & +5 \cdot 8 & -2 \cdot 7 & 0 & 2 \\ +3 \cdot 6 & +8 \cdot 5 & +7 \cdot 4 & 1 & 0 \\ -1 \cdot 3 & +0 \cdot 2 & +1 \cdot 1 & 2 & 3 \\ 4 & 2 & 1 & -1 & 2 \end{pmatrix} = 9 + 40 - 14 + 18 + 40 + 28 - 3 + 0 + 1 = 119.$$

The gray values of Y are not used in the computation. Here, we see that K is flipped when used in the convolution. This also clarifies, how the (cross) correlation can be more intuitive, where

$$[Y \circledast K]_{1,1} = \begin{pmatrix} +1 \cdot 1 & +5 \cdot 2 & -2 \cdot 3 & 0 & 2 \\ +3 \cdot 4 & +8 \cdot 5 & +7 \cdot 6 & 1 & 0 \\ -1 \cdot 7 & +0 \cdot 8 & +1 \cdot 9 & 2 & 3 \\ 4 & 2 & 1 & -1 & 2 \end{pmatrix} = 1 + 10 - 6 + 12 + 40 + 42 - 7 + 0 + 9 = 101.$$

In a similar way we can proceed to calculate the remaining values by shifting the kernel over the matrix

$$[Y \circledast K]_{1,2} = \begin{pmatrix} 1 & +5 \cdot 1 & -2 \cdot 2 & +0 \cdot 3 & 2 \\ 3 & +8 \cdot 4 & +7 \cdot 5 & +1 \cdot 6 & 0 \\ -1 & +0 \cdot 7 & +1 \cdot 8 & +2 \cdot 9 & 3 \\ 4 & 2 & 1 & -1 & 2 \end{pmatrix} = 5 - 4 + 0 + 32 + 35 + 6 + 0 + 8 + 18 = 100.$$

Altogether, we get

$$Y * K = \begin{pmatrix} 119 & 120 & 53 \\ 155 & 155 & 102 \end{pmatrix} \quad \text{and} \quad Y \circledast K = \begin{pmatrix} 101 & 100 & 87 \\ 95 & 55 & 58 \end{pmatrix}.$$

Especially, $Y * K \neq Y \circledast K$.

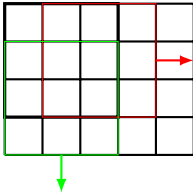


Figure 19. An image of size 4×5 is divided in blocks of size 3×3 by moving one pixel at a time either horizontally or vertically, as shown exemplary in red and green. Here, the black square is denoted by the index $(1, 1)$, the red one by $(1, 2)$ and the green one by $(2, 1)$.

The kernel size, which is typically square, e.g. $m \times m$, is a hyperparameter of the CNN. Furthermore, the convolutional layer has additional hyperparameters that need to be chosen. We have seen that a convolution with a $m \times m$ kernel reduces the dimension from $n_1 \times n_2$ to $n_1 - m + 1 \times n_2 - m + 1$. To retain the image dimension we can use **(zero) padding**, cf. [9], applied to the input Y with $p \in \mathbb{N}_0$. Choosing $p = 0$ yields Y again, whereas $p = 1$ results in

$$\hat{Y} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 5 & -2 & 0 & 2 & 0 \\ 0 & 3 & 8 & 7 & 1 & 0 & 0 \\ 0 & -1 & 0 & 1 & 2 & 3 & 0 \\ 0 & 4 & 2 & 1 & -1 & 2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix},$$

for Y from Example 3.6. Consequently, the padded matrix \hat{Y} is of dimension $(n_1 + 2p) \times (n_2 + 2p)$. To retain the image dimension, we need to choose p so that

$$\begin{aligned} (n_1 + 2p) - m + 1 &= n_1, \\ (n_2 + 2p) - m + 1 &= n_2, \end{aligned}$$

i.e. $p = \frac{m-1}{2}$, which is possible for any odd m .

Furthermore, we can choose the **stride** $s \in \mathbb{N}$, which indicates how far to move the kernel. For example, in Figure 19 the stride is chosen as $s = 1$, while the stride in Figure 20 is $s = 2$.

Let us remark that a stride $s > 1$ reduces the output dimension of the convolution to

$$\left(\frac{n_1 - m}{s} + 1 \right) \times \left(\frac{n_2 - m}{s} + 1 \right).$$

Altogether, we can describe the convolutional layer. It consists of $M \in \mathbb{N}$ filters with identical hyperparameters: kernel size $m \times m$, padding p and stride s , but each of them has its own learnable

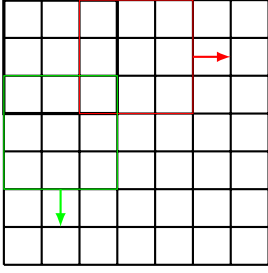


Figure 20. A visualization of the convolution of a 7×7 images with a 3×3 kernel and stride $s = 2$.

kernel K . Consequently, the filters in this layer have $M \cdot m^2$ variables in total. Applying all M filters to an input matrix $Y \in \mathbb{R}^{n_1 \times n_2}$ leads to an output of size

$$\left(\frac{n_1 + 2p - m}{s} + 1 \right) \times \left(\frac{n_2 + 2p - m}{s} + 1 \right) \times M,$$

where the results for all M filters are stacked, cf. [9]. Typically, the **depth** M is chosen as a power of 2, and growing for deeper layers, while height and width are shrinking, cf. Figure 24.

Obviously, the output is a tensor with three dimensions, hence the subsequent layers need to process 3-tensor-valued data. In fact, for colored images already the original input of the network is a tensor. The (cross) correlation operation (and also the convolution operation) can be generalized to this case in the following way.

Assume we have an input tensor of size $n_1 \times n_2 \times n_3$, then we choose a three dimensional kernel of size

$$m \times m \times n_3,$$

i.e. the depth coincides. No striding or padding is applied in the third dimension. Hence, the output is of dimension

$$\left(\frac{n_1 + 2p - m}{s} + 1 \right) \times \left(\frac{n_2 + 2p - m}{s} + 1 \right) \times 1,$$

which can be understood as a matrix by discarding the redundant third dimension, cf. Figure 21. Doing this for M filters, again leads to the output being a 3-tensor.

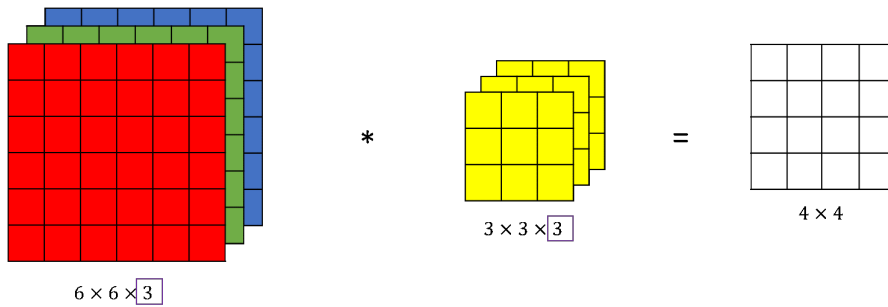


Figure 21. Illustration of a convolution on a tensor, specifically a colored image (with red, green, blue color channels) with a three-dimensional kernel and its result, which is a matrix. Here, no padding $p = 0$ and a single stride $s = 1$ is employed. Image Source: <https://datahacker.rs/convolution-rgb-image/>.

Remark 3.7

- (i) Convolutional layers have the advantage that they have less variables than fully connected layers applied to the flattened image. For example consider a grayscale image of size 28×28 as input in LeNet, cf. Figure 17. The first convolution has 6 kernels with 5×5 entries. Due to padding with $p = 2$ the output is $28 \times 28 \times 6$, so the image size is retained. Additionally, before applying a detector layer, we add a bias per channel, so 6 bias variables in this case. In total, we have to learn **156** variables. Now, imagine this image is flattened to a 784×1 vector and fed into a FNN with fully connected layer, where we also want to retain the size, i.e. the first hidden layer has 784 nodes. This results in a much larger number of variables:

$$\underbrace{784 \cdot 784}_{\text{weight}} + \underbrace{784}_{\text{bias}} = \mathbf{615440}.$$

- (ii) Directly related, a disadvantage of convolutional layers is that every output only sees a subset of all input neurons, cf. e.g. [15, Section 9.2]. We denote this set of seen inputs by **effective receptive field** of the neuron. In an FNN with fully connected layers the effective receptive field of a neuron is the entire input. However, the receptive field of a neuron increases with depth of the network, as illustrated in Figure 22.

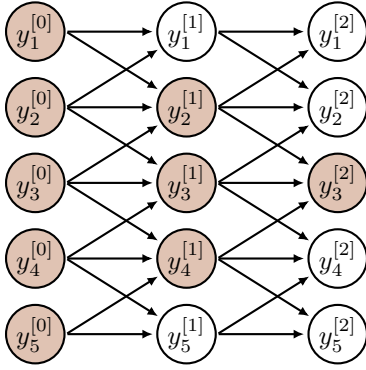


Figure 22. Simplified CNN architecture with an input layer $y^{[0]}$ and two subsequent convolutional layers $y^{[1]}$ and $y^{[2]}$, each with a one dimensional kernel of size $m = 3$, stride $s = 1$ and zero padding $p = 1$. The colored nodes are the receptive field of the neuron $y_3^{[2]}$.

3.3. Detector Layer

In standard CNN architecture after a convolutional layer, a detector layer is applied. This simply means performing an activation function. To this end, we extend the activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ to matrix and tensor valued inputs by applying it component-wise, as we did for vectors before in Section 2, e.g. for a 3-tensor $Y = \{Y_{i,j,k}\}_{i,j,k}$ with $i = 1, \dots, n_1, j = 1, \dots, n_2, k = 1, \dots, n_3$ we get

$$(\sigma(Y))_{i,j,k} = \sigma(Y_{i,j,k}).$$

3.4. Pooling Layer

After the detector layer, typically a pooling layer (also called downsampling layer) is applied, cf. e.g. [15, Section 9.3]. This layer type is responsible for reducing the first two dimensions (height and width) and usually does not interfere with the third dimension (depth) of the data Y , but rather is applied for all channels independently. Consequently, the depth of the output coincides with the depth of the input and we omit the depth in our discussion.

As in convolutional layers, pooling layers have a filter size $m \times m$, stride s and padding p . However, almost always $p = 0$ is chosen. The most popular values for for the filter size and stride are $m = s = 2$. Again, with an input of size $n_1 \times n_2$ the output dimension is

$$\left(\frac{n_1 + 2p - m}{s} + 1 \right) \times \left(\frac{n_2 + 2p - m}{s} + 1 \right) \quad m=s=2, p=0 \quad \frac{n_1}{2} \times \frac{n_2}{2}.$$

One common choice is **Max Pooling** (or Maximum Pooling), where the largest value is selected, cf. e.g. [9]. Below we see an example of max pooling with a 2×2 kernel, stride $s = 2$ and no padding.

$$\left(\begin{array}{cc|cc} 1 & 3 & 0 & -7 \\ -2 & 4 & 1 & -1 \\ \hline 0 & 1 & 8 & -3 \\ 2 & 0 & 4 & 5 \end{array} \right) \xrightarrow{\max} \left(\begin{array}{cc|cc} 4 & 1 & & \\ \hline 2 & 8 & & \end{array} \right)$$

Another common choice is **Average Pooling**, where we take the mean of all values. Below we see an example of average pooling (abbreviated: "avg") with a 2×2 kernel,

$$K = \begin{pmatrix} 0.25 & 0.25 \\ 0.25 & 0.25 \end{pmatrix},$$

stride $s = 2$ and no padding.

$$\left(\begin{array}{cc|cc} 1 & 3 & 0 & -7 \\ -2 & 4 & 1 & -1 \\ \hline 0 & 1 & 8 & -3 \\ 2 & 0 & 4 & 5 \end{array} \right) \xrightarrow{\text{avg}} \left(\begin{array}{cc|cc} 1.50 & -1.75 & & \\ \hline 0.75 & 3.50 & & \end{array} \right)$$

The effect of average pooling applied to an image is easily visible: It blurs the image. In the new image every pixel is an average of a pixel and its neighboring $m^2 - 1$ pixels, see Figure 23. Depending on the choice of stride s and padding p , the blurred image may also have less pixels.



Figure 23. Original image (left) and blurred image produced by average pooling (right) with a 5×5 kernel, stride $s = 1$ and zero padding with $p = 2$. Image Source: Laurin Ernst.

Remark 3.8

- (i) Pooling layers do not contain variables to learn.
- (ii) We have seen that when using CNNs, we make the following assumptions:
- (a) Pixels far away from each other do not need to interact with each other.
 - (b) Small translations are not relevant.

If these assumptions do not hold, employing a CNN can result in underfitting.

3.5. Local Response Normalization

Similar to batch normalization, Local Response Normalization (LRN) [24, Section 3.3] stabilizes training with unbounded activation functions like ReLU. This strategy was first introduced within the "AlexNet" architecture [24], cf. Figure 24, because contrary to previous CNNs like "LeNet-5", which used sigmoid activation, "AlexNet" employs ReLU activation.

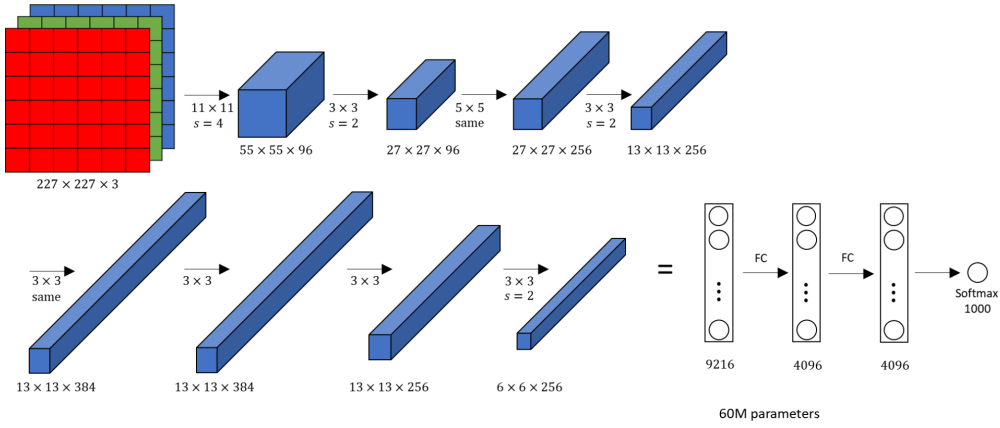


Figure 24. Architecture of AlexNet. ReLU activation is employed in the hidden layers. Image Source: <https://datahacker.rs/deep-learning-alexnet-architecture/>.

The inter-channel LRN, as introduced in [24, Section 3.3], see also Figure 25 a), is given by

$$\hat{Y}_{i,j,k} = \frac{Y_{i,j,k}}{\left(\kappa + \gamma \sum_{m=\max(1, k-\frac{n}{2})}^{\min(M, k+\frac{n}{2})} (Y_{i,j,m})^2 \right)^\beta}.$$

Here, $Y_{i,j,k}$ and $\hat{Y}_{i,j,k}$ denote the activity of the neuron before and after normalization, respectively. The indices i, j, k indicate the height, width and depth of Y . We have $i = 1, \dots, n_1$, $j = 1, \dots, n_2$ and $k = 1, \dots, M$, where M is the number of filters in the previous convolutional layer. The values $\kappa, \gamma, \beta, n \in \mathbb{R}$ are hyperparameters, where κ is used to avoid singularities, and γ and β are called normalization and contrasting constants, respectively. Furthermore, n dictates how many surrounding neurons are taken into consideration, see also Figure 25. In [24] $\kappa = 2, \gamma = 10^{-4}, \beta = 0.75$ were chosen.

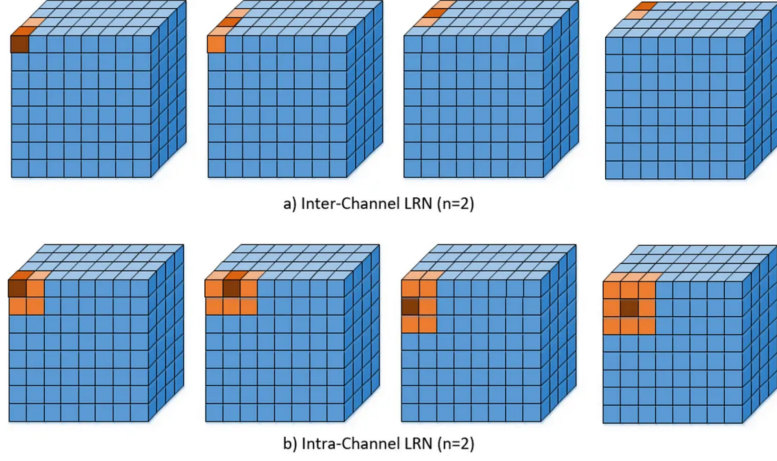


Figure 25. Illustration of local response normalization. Inter-channel version a) as introduced in [24, Section 3.3] and intra-channel version b). Both for $n = 2$. For clarification: The red pixel in the top left cube is $Y_{1,1,1}$, while the red pixel in the top row, second from the left cube is $Y_{1,1,2}$ and the red pixel in the bottom row, second from the left cube is $Y_{1,2,1}$. Image source: <https://towardsdatascience.com/difference-between-local-response-normalization-and-batch-normalization-272308c034ac>.

In the case of intra-channel LRN, the neighborhood is extended within the same channel. This leads to the following formula

$$\hat{Y}_{i,j,k} = \frac{Y_{i,j,k}}{\left(\kappa + \gamma \sum_{p=\max(1,i-\frac{n}{2})}^{\min(n_1,i+\frac{n}{2})} \sum_{q=\max(1,j-\frac{n}{2})}^{\min(n_2,j+\frac{n}{2})} (Y_{p,q,k})^2 \right)^{\beta}}.$$

Remark 3.9 *The LRN layer is non-trainable, since it only contains hyperparameters and no variables.*

4. ResNet

We have seen in Example 2.3 that for a FNN with depth L we have the derivative

$$\frac{\partial \mathcal{L}}{\partial W^{[\ell]}} = \frac{\partial \mathcal{L}}{\partial y^{[L]}} \cdot \prod_{j=L}^{\ell+2} \frac{\partial y^{[j]}}{\partial y^{[j-1]}} \cdot \frac{\partial y^{[\ell+1]}}{\partial W^{[\ell]}}.$$

In the case that we consider a very deep network, i.e. large L , the product in the derivative can be problematic, [5, 14], especially if we take derivatives with respect to variables from early layers. Two cases may occur:

1. If $\frac{\partial y^{[j]}}{\partial y^{[j-1]}} < 1$ for all j , the product, and hence the whole derivative, tends to zero for growing L . This problem is referred to as **vanishing gradient**.
2. On the other hand, if $\frac{\partial y^{[j]}}{\partial y^{[j-1]}} > 1$ for all j , the product, and hence the whole derivative, tends to infinity for growing L . This problem is referred to as **exploding gradient**.

Residual Networks (ResNets) have been developed in [17, 19] with the intention to solve the vanishing gradient problem. Employing the same notation as in FNNs, simplified ResNet layers can be represented in the following way

$$y^{[\ell]} = y^{[\ell-1]} + \sigma^{[\ell]}(W^{[\ell-1]}y^{[\ell-1]} + b^{[\ell-1]}) \quad \text{for } \ell = 1, \dots, L, \quad (9)$$

with $y^{[0]} = u$ the input data. Essentially, a ResNet is a FNN with an added **skip connection**, i.e. $+y^{[\ell-1]}$, cf. Figure 26

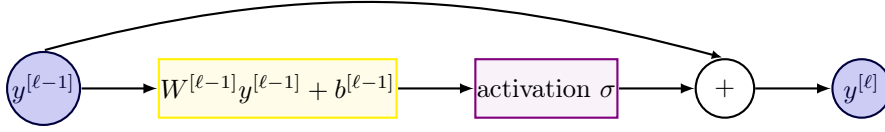


Figure 26. Illustration of a simplified ResNet layer.

Remark 4.1 *The ResNet layers in the current form (9) only work, if all feature vectors $y^{[\ell]}$ have the same dimension n_ℓ , so that we can add them up. To allow for different layer sizes, we need to insert projection operators $P_{\ell-1}^\ell \in \mathbb{R}^{n_\ell \times n_{\ell-1}}$, cf. [2, Section 4], i.e.*

$$y^{[\ell]} = P_{\ell-1}^\ell y^{[\ell-1]} + \sigma^{[\ell]}(W^{[\ell-1]}y^{[\ell-1]} + b^{[\ell-1]}) \quad \text{for } \ell = 1, \dots, L,$$

We now revisit the simple FNN with two hidden layers from Example 2.3, and add skip connections to make it a ResNet, see Figure 27.

Example 4.2 *Consider a simple ResNet with one node per layer and assume that we only consider weights $W^{[\ell]} \in \mathbb{R}$ and no biases. For the network in Figure 27 we have $\theta = (W^{[0]}, W^{[1]}, W^{[2]})^\top$, and*

$$\mathcal{L}(\theta) = \mathcal{L}(y^{[3]}(\theta)).$$

We define for $\ell = 1, \dots, L$

$$a^{[\ell]} := \sigma^{[\ell]}(W^{[\ell-1]}y^{[\ell-1]}),$$

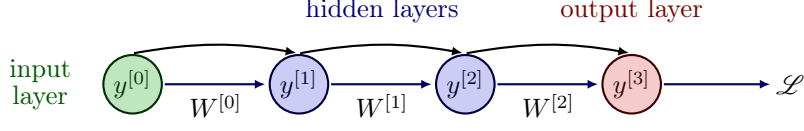


Figure 27. A ResNet with 2 hidden layers, one node per layer and depth $L = 3$.

so that in the ResNet setup

$$y^{[\ell]} = y^{[\ell-1]} + a^{[\ell]}.$$

Computing the components of the gradient, we employ the chain rule to obtain e.g.

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W^{[0]}} &= \frac{\partial \mathcal{L}}{\partial y^{[3]}} \cdot \frac{\partial y^{[3]}}{\partial W^{[0]}} \\ &= \frac{\partial \mathcal{L}}{\partial y^{[3]}} \cdot \frac{\partial}{\partial W^{[0]}} (y^{[2]} + a^{[3]}) \\ &= \frac{\partial \mathcal{L}}{\partial y^{[3]}} \cdot \left(\frac{\partial y^{[2]}}{\partial W^{[0]}} + \frac{\partial a^{[3]}}{\partial y^{[2]}} \cdot \frac{\partial y^{[2]}}{\partial W^{[0]}} \right) \\ &= \frac{\partial \mathcal{L}}{\partial y^{[3]}} \cdot \left(\mathbb{I} + \frac{\partial a^{[3]}}{\partial y^{[2]}} \right) \cdot \frac{\partial y^{[2]}}{\partial W^{[0]}} \\ &= \frac{\partial \mathcal{L}}{\partial y^{[3]}} \cdot \left(\mathbb{I} + \frac{\partial a^{[3]}}{\partial y^{[2]}} \right) \cdot \left(\mathbb{I} + \frac{\partial a^{[2]}}{\partial y^{[1]}} \right) \cdot \frac{\partial y^{[1]}}{\partial W^{[0]}}, \end{aligned}$$

where \mathbb{I} denotes the identity. In general for depth L , we get

$$\frac{\partial \mathcal{L}}{\partial W^{[\ell]}} = \frac{\partial \mathcal{L}}{\partial y^{[L]}} \cdot \prod_{j=L}^{\ell+2} \left(\mathbb{I} + \frac{\partial a^{[j]}}{\partial y^{[j-1]}} \right) \cdot \frac{\partial y^{[\ell+1]}}{\partial W^{[\ell]}}. \quad (10)$$

If we generalize the derivative (10) to ResNet architectures, where we do not only consider weights $W^{[\ell]}$, see e.g. [2, Theorem 6.1], the structure of the product in the derivative remains the same, i.e. it also contains an identity term.

Remember that for FNNs it holds $y^{[j]} = a^{[j]}$, i.e. the fraction in the product coincides in both cases. However, due to the added identity, even if

$$\frac{\partial a^{[j]}}{\partial y^{[j-1]}} < 1$$

holds for all j , we will not encounter vanishing gradients in the ResNet architecture. The exploding gradients problem can still occur.

We will see in Section 4.1 that there exist several versions of ResNets. However, from a mathematical point of view the simplified version (9) is especially interesting, because it can be related to ordinary differential equations (ODEs), as first done in [16]. Inserting a parameter $\tau^{[\ell]} \in \mathbb{R}$ in front of the activation function σ and rearranging the terms of the forward propagation delivers

$$\begin{aligned} y^{[\ell]} &= y^{[\ell-1]} + \tau^{[\ell]} \sigma(W^{[\ell-1]} y^{[\ell-1]} + b^{[\ell-1]}) \\ \Rightarrow \frac{y^{[\ell]} - y^{[\ell-1]}}{\tau^{[\ell]}} &= \sigma(W^{[\ell-1]} y^{[\ell-1]} + b^{[\ell-1]}). \end{aligned}$$

Here, we consider the same activation function σ for all layers. Now, the left hand side of the equation can be interpreted as a finite difference representation of a time derivative, where $\tau^{[\ell]}$ is the time step size and $y^{[\ell]}, y^{[\ell-1]}$ are the values attained at two neighboring points in time. This relation between ResNets and ODEs is also studied under the name of Neural ODEs, [7]. It is also possible to learn the time step size $\tau^{[\ell]}$ as an additional variable, [2].

Let us now introduce the different ResNet versions from the original papers, [17, 19].

4.1. Different ResNet Versions

In contrast to the simplified ResNet layer version (9) that we introduced, original ResNet architectures [17] consist of **residual blocks**, cf. Figure 28. Here, different layers are grouped together into one residual block and then residual blocks are stacked to form a ResNet.

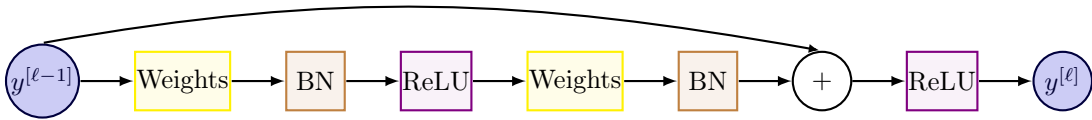


Figure 28. Illustration of a residual block as introduced in [17].

In the residual block, cf. Figure 28, we have the following layer types:

- Weights: fully connected or convolutional layer,
- BN: Batch Normalization layer, cf. Section 2.3,
- ReLU: activation function $\sigma = \text{ReLU}$.

Clearly, the residual block and the simplified ResNet layer both contain a skip connection, which is the integral part of ResNets success, since it helps avoid the vanishing gradient problem. However, the residual block is less easy to interpret from a mathematical point of view and can not directly be related to ODEs.

In frameworks like Tensorflow and Pytorch, if you encounter a network architecture called "ResNet", it will usually be built by stacking residual blocks of this original form, Figure 28.

Subsequently, in a follow up paper, [19], several other options to sort the occurring layers in a residual block have been introduced. The option which performed best in numerical tests (see also Figure 30) is illustrated in Figure 29.

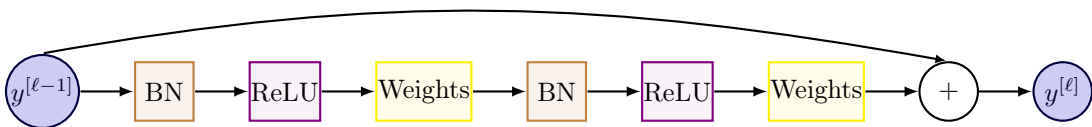


Figure 29. Illustration of a full pre-activation residual block as proposed in [19, Fig.1(b)].

Remark 4.3 The authors of [19] call the residual block in Figure 29 the **full pre-activation residual block**, since both activation functions are exercised before (pre) the skip connection. Meanwhile, in the original residual block there is also a post-activation, i.e. an activation function after (post) the skip connection. In this sense, the simplified ResNet layer (9) can be termed a *pre-activation ResNet layer*.

A ResNet built with full pre-activation residual blocks can be found e.g. in Tensorflow under the name "ResNetV2". In the literature, there also exist other variants of the simplified ResNet layer, e.g. with a weight matrix applied outside the activation function.

In Figure 30 we see a comparison of a 1001-layer ResNet built with original residual blocks and a 1001-layer ResNet built with full pre-activation (proposed) residual blocks. This result clearly demonstrates the advantage of full pre-activation for very deep networks, since both training loss and test error can be improved with the proposed residual block.

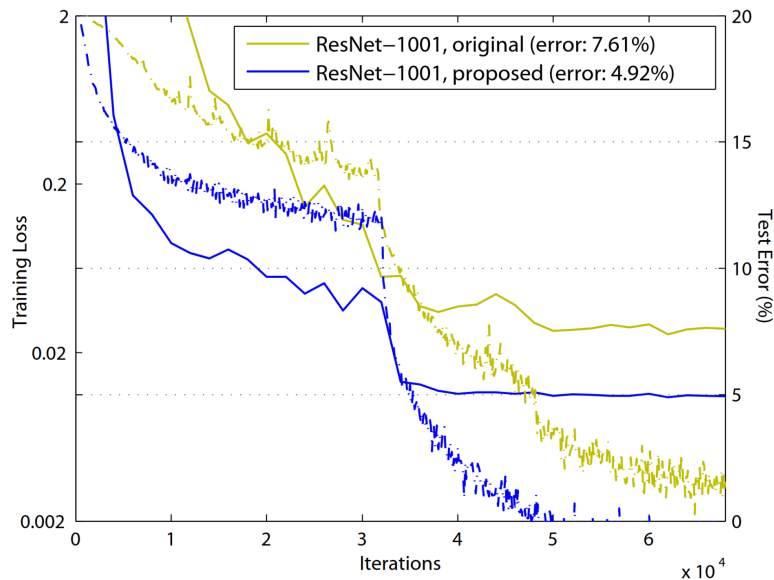


Figure 30. Training loss (dashed line, left y-axis) and test error (solid line, right y-axis) plotted against the iteration counter for a 1001-layer ResNet on the CIFAR-10 dataset. Here, the original residual block (blue) is compared to the proposed full pre-activation residual block (green). Image Source: [19, Fig.1]

4.2. ResNet18

As an example for a ResNet architecture, we look at "ResNet18", cf. Figure 31. Here, 18 indicates the number of layers with learnable weights, i.e. convolutional and fully connected layers. Even though the batch normalization layers also contain learnable weights, they are typically not counted here. This is a ResNet architecture intended for use on image data sets, hence the weights layers are convolutional layers, like in a CNN.

In Block A the input data is pre-processed, while Block B to Block E are built of two residual blocks each. To be certain which type of residual block the network is built of, it is recommended to look into the details of the implementation. However, in most cases the original residual

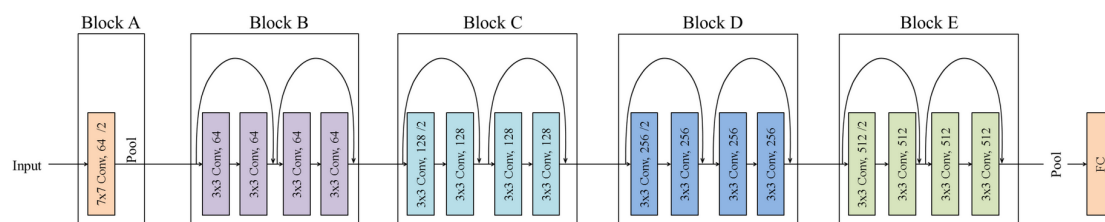


Figure 31. Illustration of "ResNet18" architecture, built of residual blocks. Image Source: [13].

block, cf. Figure 28, is employed. Finally, as usual for a classification task, the data is flattened and passed through a fully connected (FC) layer before the output is generated. Altogether, "ResNet18" has over 10 million trainable parameters, i.e. variables.

4.3. Transfer Learning

An advantage of having so many different ResNet architectures available and also pre-trained is, that they can be employed for **Transfer Learning**, see e.g. [29]. The main idea of transfer learning is to take a model that has been trained on a (potentially large) data set for the same type of task (e.g. image classification), and then adjust the first/last layer to fit your data set. Depending on your task you may need to change one or both layers, for example:

- (i) Your input data has a different structure: adapt the first layer.
- (ii) Your data set has a different set of labels (supervisions): adapt the last layer.

If your input data and labels both coincide with the original task then you don't need to employ transfer learning. You can just use the pre-trained model for your task. When adapting a layer, this layer needs to be initialized. All remaining layers can be initialized with the pre-trained weights, which will most likely give a good starting point. Then you train the adapted network on your data, which will typically take a lot less time than training with a random initialization. Hence, transfer learning can save a significant amount of computing time. Clearly, transfer learning is also possible with other network architectures, as long as the network has been pre-trained.

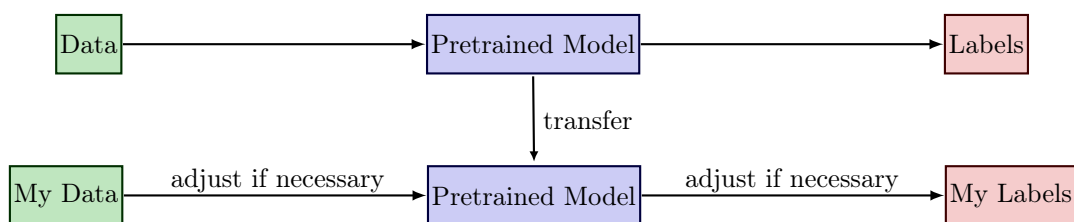


Figure 32. Illustration of transfer learning.

5. Recurrent Neural Network

The Neural Networks we introduced so far rely on the **assumption of independence** among the training and test examples. They process one data point at a time, which is no problem for data sets, in which every data point is generated independently. However, for sequential data that occurs in machine translation, speech recognition, sentiment classification, etc., the dependence is highly relevant to the task.

Recurrent Neural Networks (RNNs), cf. e.g. [15, Section 10] and [12, Section 8.1], are connectionist models that capture the dynamics of sequences via cycles in the network of nodes. Unlike standard FNNs, recurrent neural networks retain a state that can represent information from an arbitrarily long context window.

Example 5.1 (machine translation)

Translate a given english input sentence u , consisting of T_{in} words $u^{<t>}$, $t = 1, \dots, T_{\text{in}}$, e.g.

The	sun	is	shining	today
$u^{<1>}$	$u^{<2>}$	$u^{<3>}$	$u^{<4>}$	$u^{<5>}$

to a german output sentence y , consisting of T_{out} words $y^{<t>}$, $t = 1, \dots, T_{\text{out}}$. Hopefully, the output will be something like

Heute	scheint	die	Sonne
$y^{<1>}$	$y^{<2>}$	$y^{<3>}$	$y^{<4>}$

A comparison of FNN and RNN architecture can be seen in Figure 33. For simplicity of notation we condense all hidden layers of the FNN into a representative computation node h .

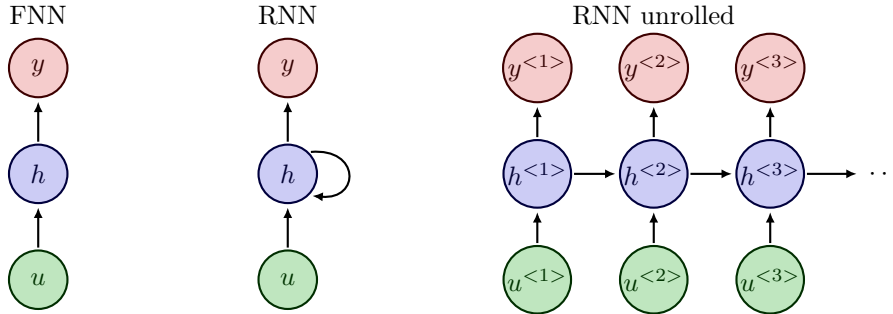


Figure 33. Feedforward Neural Network compared to Recurrent Neural Network with input u , output y and hidden computation nodes h . The index is understood as time instance.

In RNNs the computation nodes h are often called **RNN cells**, cf. [15, Section 10.2]. A RNN cell for a time instance t takes as an input $u^{<t>}$ and $h^{<t-1>}$, and computes the outputs $h^{<t>}$ and $y^{<t>}$, cf. Figure 34. More specifically for all $t = 1, \dots, T_{\text{out}}$

$$h^{<t>} = \sigma(W_{\text{in}} \cdot [h^{<t-1>}; u^{<t>}] + b), \quad (11)$$

$$y^{<t>} = W_{\text{out}} \cdot h^{<t>}. \quad (12)$$

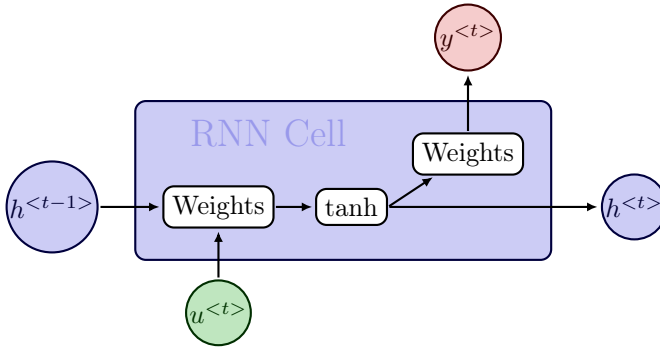


Figure 34. Architecture of a RNN cell.

The equations (11) and (12) describe the forward propagation in RNNs. Here, $[h^{<t-1>}; u^{<t>}]$ denotes the concatenation of the vectors, and $h^{<0>}$ is set to a vector of zeros, so that we do not need to formulate a special case for $t = 1$. Depending on the application, a softmax function may be applied to $W_{\text{out}}h^{<t>}$ to get the output $y^{<t>}$.

It may happen that input and output have different lengths $T_{\text{in}} \neq T_{\text{out}}$, see e.g. Example 5.1. Depending on the task and the structure of the data, there exist various types of RNN architectures, cf. [12, Section 8.1] and Figure 35:

- one to many, e.g. image description (image to sentence),
- many to one, e.g. sentiment analysis (video to word),
- many to many, e.g. machine translation (sentence to sentence), like Example 5.1,
- many to many, e.g. object tracking (video to object location per frame).

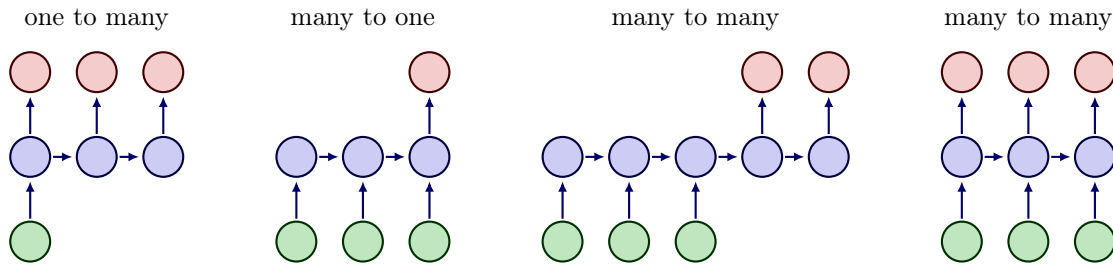


Figure 35. Illustration of different types of RNN architectures.

We note that the weights $W_{\text{in}}, W_{\text{out}}$ and bias b in (11) and (12) do not change over time, but coincide for all temporal layers of the RNN. Sharing the variables allows the RNN to model variable length sequences, whereas if we had specific parameters for each value of the order parameter, we could not generalize to sequence lengths not seen during training. Typically, $\sigma = \tanh$ is chosen in RNNs, and this does also not vary between the layers. To obtain a complete optimization problem (P), we still need a loss function \mathcal{L} , since the RNN represents only the network \mathcal{F} . To this end, each output $y^{<t>}$ is evaluated with a loss function $\mathcal{L}^{<t>}$ and

the final loss is computed by taking the sum over all time instances

$$\mathcal{L}(\theta) = \sum_{t=1}^{T_{\text{out}}} \mathcal{L}^{\langle t \rangle}(y^{\langle t \rangle}(\theta)).$$

Here, as usual, θ contains the weights $W_{\text{in}}, W_{\text{out}}$, and bias b .

5.1. Variants of RNNs

We briefly introduce two popular variants of RNNs.

In many applications the output at time t should be a prediction depending on the whole input sequence, not only the "earlier" inputs $u^{\langle i \rangle}$ with $i \leq t$. E.g., in speech recognition, the correct interpretation of the current sound as a phoneme may depend on the next few phonemes because of co-articulation and potentially may even depend on the next few words because of the linguistic dependencies between nearby words. As a remedy, we can combine a forward-going RNN and a backward-going RNN, which is then called a **Bidirectional RNN**, [15, Section 10.3]. This architecture allows to compute an output $y^{\langle t \rangle}$ that depends on both the past and the future inputs, but is most sensitive to the input values around time t . Figure 36 (left) illustrates the typical bidirectional RNN, with $h^{\langle t \rangle}$ and $g^{\langle t \rangle}$ representing the states of the sub-RNNs that move forward and backward through time, respectively.

Another variant of RNNs is the **Deep RNN**, [15, Section 10.5]. As seen in FNNs, Section 2, multiple hidden layers allow the network to have a higher expressiveness. Similarly, a RNN can be made deep by stacking RNN cells, see Figure 36 (right).

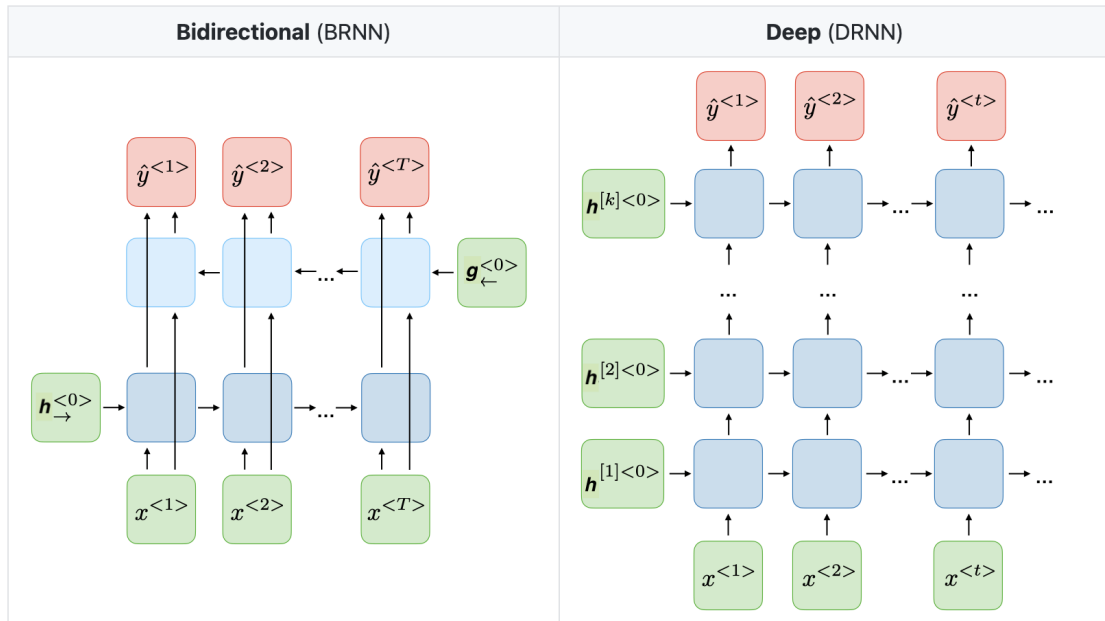


Figure 36. Examples of Bidirectional RNNs and Deep RNNs. Here, the inputs are denoted by x instead of u . Image source: <https://stanford.edu/~shervine/teaching/cs-230/>.

5.2. Long term dependencies

In this section we investigate one of the main challenges, that a RNN can encounter, cf. [15, Section 10.7]. Consider the following illustrative example.

Example 5.2

Predict the next word in the sequence:

1. The cat, which ..., *was* ...
2. The cats, which ..., *were* ...

Here, depending on whether we are talking about one cat or multiple cats the verb has to be adjusted. The "... part in the sentence can be very extensive, so that the dependence becomes long.

The gradient from the output $y^{<t>}$ with large t has to propagate back through many layers to affect weights in early layers. Here, the vanishing gradient and exploding gradient problems (cf. Section 4) may occur and hinder training. The exploding gradient problem can be solved relatively robustly by **gradient clipping**, see e.g. [15, Section 10.11.1]. The idea is quite simple. If a gradient $\partial_{\theta_i} \mathcal{L}$, with respect to some variable θ_i gets too large, we rescale it. I.e. if $\|\partial_{\theta_i} \mathcal{L}\| \geq C \in \mathbb{R}$ for a hyperparameter C , we set

$$\partial_{\theta_i} \mathcal{L} \leftarrow C \cdot \frac{\partial_{\theta_i} \mathcal{L}}{\|\partial_{\theta_i} \mathcal{L}\|}.$$

Let us remark that this is a heuristic approach. In contrast, the vanishing gradient problem is more difficult to solve.

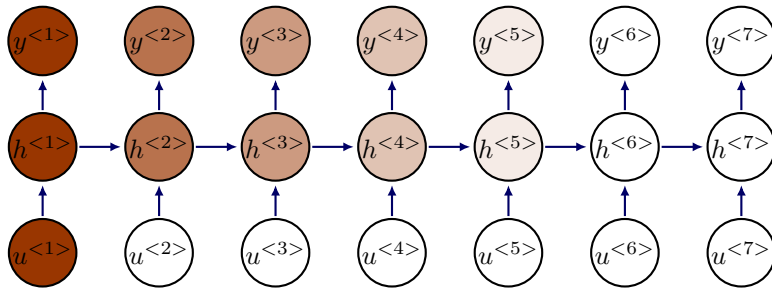


Figure 37. Illustration of vanishing gradient problem for RNNs. The shading of the nodes indicates the sensitivity over time of the network nodes to the input $u^{<1>}$ (the darker the shade, the greater the sensitivity). The sensitivity decays over time.

A common remedy is to modify the RNN cell so that it can capture long term dependencies better, and avoids vanishing gradients. Two popular options are **Gated Recurrent Unit (GRU)** from 2014 [8], and the cell architecture as suggested already in 1997 in **Long Short Term Memory (LSTM)** networks [20]. The core idea in both cell architectures is to add gating mechanisms. These gates have a significant influence on whether, and how severely, the input and previous hidden state influence the output and new hidden state. Additionally, the gating mechanism helps to solve the vanishing gradient problem.

5.2.1. Gated Recurrent Unit

The gated recurrent unit has a reset (or relevance) gate Γ_r and an update gate Γ_u . The computations for one unit are as follows

$$\begin{aligned} \Gamma_r &= \sigma(W_r \cdot [h^{<t-1>; u^{<t>}] + b_r), && \text{reset gate} \\ \Gamma_u &= \sigma(W_u \cdot [h^{<t-1>; u^{<t>}] + b_u), && \text{update gate} \\ \tilde{h}^{<t>} &= \tanh(W_{in} \cdot [\Gamma_r \odot h^{<t-1>; u^{<t>}] + b), && \text{hidden state candidate} \\ h^{<t>} &= \Gamma_u \odot \tilde{h}^{<t>} + (1 - \Gamma_u) \odot h^{<t-1>}, && \text{hidden state} \\ y^{<t>} &= W_{out} \cdot h^{<t>}. && \text{output} \end{aligned}$$

The computations of the hidden state candidate $\tilde{h}^{<t>}$ and the output $y^{<t>}$ resemble the computations in the RNN cell (11) and (12), respectively. However, e.g. if $\Gamma_u = 0$, then the new hidden state will coincide with the previous hidden state and the candidate will not be taken into account. Also, the GRU has significantly more variables per cell, in comparison with the standard RNN cell.

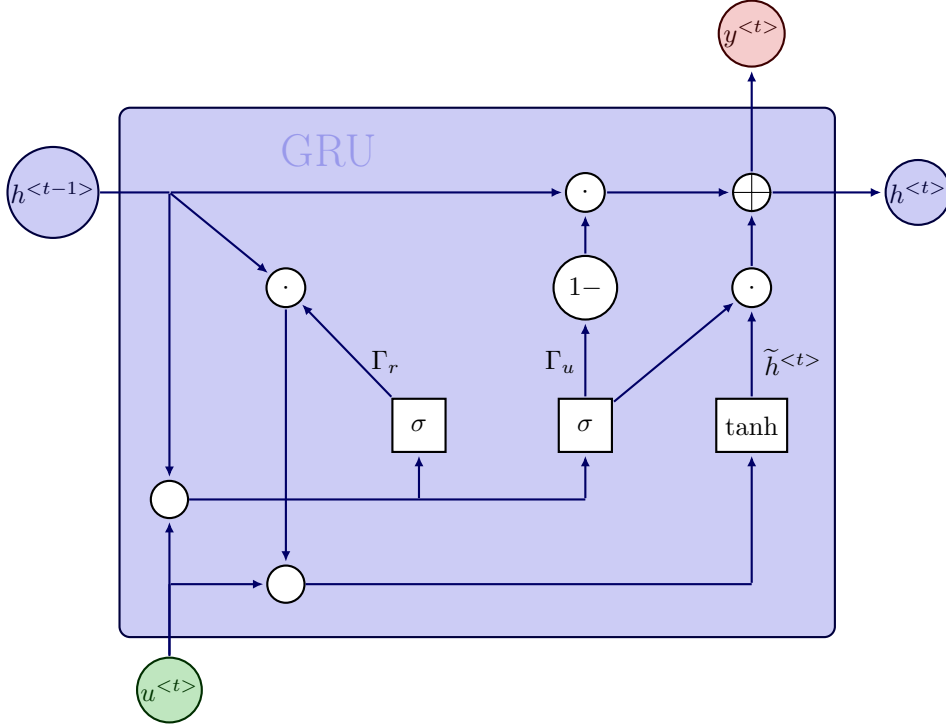


Figure 38. Architecture of a gated recurrent unit. Weights are omitted in this illustration. A white circle illustrates concatenation, while a circle with a dot represents the Hadamard product and a circle with a plus indicates an addition.

5.2.2. Long Short Term Memory

The key to LSTM networks is that in addition to the hidden state, there also exists a cell state $c^{<t>}$, which is propagated through the network. It can be understood like a conveyor belt,

which only has minor interactions and runs down the entire chain of LSTM cells, see Figure 39. This allows information to flow through the network easily. In contrast to GRU, the LSTM cell contains three gates: the forget gate Γ_f , input gate Γ_i and output gate Γ_o .

$$\begin{aligned}
 \Gamma_f &= \sigma(W_f \cdot [h^{<t-1>}; u^{<t>}] + b_f), & \text{forget gate} \\
 \Gamma_i &= \sigma(W_i \cdot [h^{<t-1>}; u^{<t>}] + b_i), & \text{input gate} \\
 \Gamma_o &= \sigma(W_o \cdot [h^{<t-1>}; u^{<t>}] + b_o), & \text{output gate} \\
 \tilde{c}^{<t>} &= \tanh(W_c \cdot [h^{<t-1>}; u^{<t>}] + b_c), & \text{cell state candidate} \\
 c^{<t>} &= \Gamma_f \odot \tilde{c}^{<t-1>} + \Gamma_i \odot \tilde{c}^{<t>}, & \text{cell state} \\
 h^{<t>} &= \Gamma_o \odot \tanh(c^{<t>}), & \text{hidden state} \\
 y^{<t>} &= W_{\text{out}} \cdot h^{<t>}. & \text{output}
 \end{aligned}$$

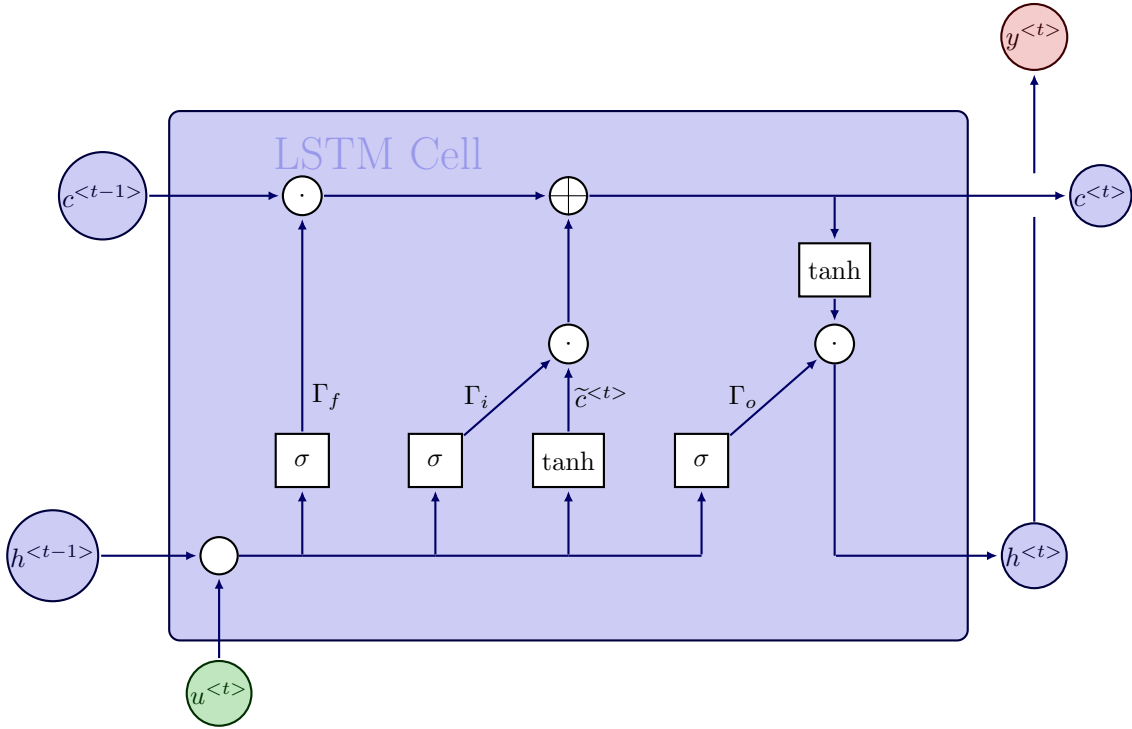


Figure 39. Architecture of a LSTM cell. Weights are omitted in this illustration. A white circle illustrates concatenation, while a circle with a dot represents the Hadamard product and a circle with a plus indicates an addition.

Remark 5.3 Without gates, i.e. $\Gamma_f = \Gamma_i = \Gamma_o = 1$, the LSTM network has a certain similarity with the ResNet structure, which was developed later than the LSTM, in 2016 in [17]. This is not so surprising, since both networks aim at solving the vanishing gradient problem. In fact, propagating the cell state has similar effects on the gradients as introducing skip connections.

5.3. Language processing

An important application of RNNs is language processing, e.g. machine translation, see Example 5.1. In such tasks the words need to be represented, so that the RNN can work with them. Furthermore, we need a way to deal with punctuation marks, and an indicator for the end of a sentence.

To represent the words, we form a dictionary. For the english language we will end up with a vector containing more than 10000 words. Intuitively, we sort the words alphabetically and to simplify computations we use a **one-hot representation**. E.g., if "the" is the 8367th word in the english dictionary vector, we represent the first input

$$u^{<1>} = (0 \quad \dots \quad 0 \quad 1 \quad 0 \quad \dots \quad 0)^\top = e_{8367},$$

with the 8367th unit vector. This allows for an easy way to measure correctness in supervised learning and later on we can use the dictionary to recover the words. Additionally, it is common to create a token for unknown words, which are not in the dictionary. Punctuation marks can either be ignored, or we also create tokens for them. However, the dictionary should at least contain an "end of sentence" token to separate sentences from each other.

References

- [1] H. Antil, T. S. Brown, R. Löhner, F. Togashi, and D. Verma. *Deep Neural Nets with Fixed Bias Configuration*. 2022. arXiv: 2107.01308 [math.OC].
- [2] H. Antil, H. Díaz, and E. Herberg. *An Optimal Time Variable Learning Framework for Deep Neural Networks*. 2022. arXiv: 2204.08528 [math.OC].
- [3] H. Antil, R. Khatri, R. Löhner, and D. Verma. "Fractional deep neural network via constrained optimization". In: *Machine Learning: Science and Technology* 2.1 (2020), p. 015003.
- [4] A. Baldominos, Y. Saez, and P. Isasi. "A Survey of Handwritten Character Recognition with MNIST and EMNIST". In: *Applied Sciences* 9.15 (2019). ISSN: 2076-3417. DOI: 10.3390/app9153169. URL: <https://www.mdpi.com/2076-3417/9/15/3169>.
- [5] Y. Bengio, P. Simard, and P. Frasconi. "Learning long-term dependencies with gradient descent is difficult". In: *IEEE transactions on neural networks* 5.2 (1994), pp. 157–166.
- [6] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [7] R. T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud. "Neural ordinary differential equations". In: *Advances in neural information processing systems* 31 (2018).
- [8] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. *Learning phrase representations using RNN encoder-decoder for statistical machine translation*. 2014. arXiv: 1406.1078.
- [9] *CS231n Convolutional Neural Networks for Visual Recognition*. <https://cs231n.github.io/convolutional-networks/>. 2023.
- [10] G. Cybenko. "Approximation by superpositions of a sigmoidal function". In: *Mathematics of control, signals and systems* 2.4 (1989).
- [11] K. Fukushima. "Neocognitron: A Self-organizing Neural Network Model for a Mechanism of Pattern Recognition Unaffected by Shift in Position". In: *Biological Cybernetics* 36 (1980).
- [12] A. Geiger. *Deep Learning Lecture Notes*. <https://drive.google.com/file/d/16TaFr6d3eZXNkShgJJxaf6CN7x/view>. 2021.

- [13] S. Ghassemi and E. Magli. “Convolutional Neural Networks for On-Board Cloud Screening”. In: *Remote Sensing* 11 (2019). DOI: 10.3390/rs11121417.
- [14] X. Glorot and Y. Bengio. “Understanding the difficulty of training deep feedforward neural networks”. In: *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings. 2010.
- [15] I. Goodfellow, Y. Bengio, and A. Courville. *Deep learning*. MIT press, 2016.
- [16] E. Haber and L. Ruthotto. “Stable architectures for deep neural networks”. In: *Inverse problems* 34.1 (2017). DOI: 10.1088/1361-6420/aa9a90.
- [17] K. He, X. Zhang, S. Ren, and J. Sun. “Deep residual learning for image recognition”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.
- [18] K. He, X. Zhang, S. Ren, and J. Sun. “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification”. In: *Proceedings of the IEEE international conference on computer vision*. 2015.
- [19] K. He, X. Zhang, S. Ren, and J. Sun. “Identity mappings in deep residual networks”. In: *European conference on computer vision*. Springer. 2016.
- [20] S. Hochreiter and J.f Schmidhuber. “Long short-term memory”. In: *Neural computation* 9.8 (1997).
- [21] R. A. Horn. “The hadamard product”. In: *Proc. Symp. Appl. Math.* Vol. 40. 1990, pp. 87–169.
- [22] S. Ioffe and C. Szegedy. “Batch normalization: Accelerating deep network training by reducing internal covariate shift”. In: *International conference on machine learning*. pmlr. 2015, pp. 448–456.
- [23] D. P. Kingma and J. Ba. “Adam: A method for stochastic optimization”. In: *arXiv preprint arXiv:1412.6980* (2014). arXiv: 1412.6980.
- [24] A. Krizhevsky, I. Sutskever, and G. E. Hinton. “ImageNet Classification with Deep Convolutional Neural Networks”. In: *Communications of the ACM* 60.6 (2017).
- [25] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86 (1998).
- [26] A. Ng. *CS229 Lecture Notes*. https://cs229.stanford.edu/notes2022fall/main_notes.pdf. 2022.
- [27] F. Rosenblatt. “The perceptron: a probabilistic model for information storage and organization in the brain.” In: *Psychological review* 65.6 (1958), p. 386.
- [28] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. “Dropout: a simple way to prevent neural networks from overfitting”. In: *The journal of machine learning research* 15.1 (2014).
- [29] L. Torrey and J. Shavlik. “Transfer learning”. In: *Handbook of research on machine learning applications and trends: algorithms, methods, and techniques*. IGI global, 2010, pp. 242–264.