

# FYS-STK4155 - Applied data analysis and machine learning

## Project 3

Even M. Nordhagen

December 7, 2018

- Github repository containing programs and results:

<https://github.com/evenmn/FYS-STK4155>

### **Abstract**

The aim of this project is to divide various sounds into ten categories, inspired by the Urban Sound Challenge. For that we have investigated the performance of logistic regression and various neural networks, like Feed-forward Neural Networks (FNNs), Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs).

When it comes to error estimation, accuracy score is a natural choice. Various activation functions were investigated, and regularization was added for the logistic case. Since we mostly used ADAM as optimization tool...

To extract features from our dataset (dimensionality reduction), we used a spectrogram in the CNN case, and Mel Frequency Cepstral Coefficients (MFCCs) for FNN and RNN. We also tried two RNN networks, mainly the Long Short-Term Memory (LSTM) network and the Gated Recurrent Unit (GRU) network.

All the neural networks were more or less able to recognize the training set, but it was a significant difference in test accuracy. FNN provided the highest test accuracy of 94%, followed by LSTM (90%), GRU (88%) and CNN (85%). Our linear model, logistic regression, was only able to classify 59% of the test set correctly.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Sound analysis</b>	<b>4</b>
2.1	The Urban Sound Challenge (USC) . . . . .	4
2.2	Time domain . . . . .	5
2.3	Frequency domain . . . . .	6
2.4	Extract features . . . . .	7
2.4.1	Spectrograms . . . . .	7
2.4.2	Mel Frequency Cepstral Coefficients (MFCCs) . . . . .	7
2.5	Examples from USC . . . . .	7
<b>3</b>	<b>Classification methods</b>	<b>9</b>
3.1	Logistic regression . . . . .	9
3.2	Feed-forward Neural Networks (FNN) . . . . .	10
3.3	Convolutional Neural Networks (CNN) . . . . .	12
3.4	Recurrent Neural Networks (RNN) . . . . .	13
<b>4</b>	<b>Optimization</b>	<b>14</b>
4.1	Gradient Descent (GD) . . . . .	14
4.2	ADAM . . . . .	14
<b>5</b>	<b>Activation</b>	<b>15</b>
5.1	Logistic . . . . .	16
5.2	ReLU . . . . .	16
5.3	ELU . . . . .	16
<b>6</b>	<b>Code</b>	<b>16</b>
6.1	Code structure . . . . .	16
6.2	Implementation . . . . .	16
<b>7</b>	<b>Results</b>	<b>16</b>
7.1	Logistic regression . . . . .	16
7.2	Feed-forward Neural Networks . . . . .	17
7.3	Convolutional Neural Networks . . . . .	17
7.4	Recurrent Neural Networks . . . . .	18
<b>8</b>	<b>Discussion</b>	<b>18</b>
<b>9</b>	<b>Conclusion</b>	<b>18</b>

# 1 Introduction

In the everyday life we are continuously surrounded by sounds, and usually the human brain is able to recognize what kind of sound it is based on experience. Artificial neural networks are again based on studies of the human brain, and if the artificial neurons work as the biological ones, there should be possible training a neural network recognizing sounds as well.

Lately, immense efforts have been put into this subject with the purpose of translating voice into text, leaded by technology companies like Google, Microsoft, Amazon and so on, who develop voice controlled virtual assistants. The technology is promising, the time saving using voice commands vs. keyboard commands is potentially huge and the market is enormous. Some even claim that virtual assistants will run our lives within 20 years. [Fel18] If that is right is still uncertain, but what is certain is that the technology has great potential.

In this final project we will, based on the Urban Sound Challenge, sort sounds into classes using various classification methods. We use the same idea as when the great technology companies recognize voice, but we are going to differentiate between sounds made by **air conditioners**, **car horns**, **children**, **dogs**, **drills**, **engines**, **guns**, **jackhammers**, **sirens** and **street musicians**. In order to do that,

## 2 Sound analysis

In sound analysis, one can either analyze the raw sound files directly or convert them to simpler files and hope we have not lost any essential features. A big advantage with the latter, is that the dimensionality and complexity of the data set is significantly reduced, but we will also lose some information. In this section we will first introduce the Urban Sound Challenge, and then look at the sounds from a time perspective, frequency perspective and see how to extract information. Finally, we take a look at some of the samplings and see if we can classify them with the naked eye.

### 2.1 The Urban Sound Challenge (USC)

The Urban Sound Challenge is a classification contest provided by Analytics Vidhya with the purpose of introducing curious people to a real-world classification problem. After registered, one is provided with a dataset containing sounds from ten classes. For the training data set, the classes (targets) are given, but there is also a test dataset where the targets are unknown. Our task is to classify the test dataset correctly, and by uploading our results to Analytics Vidhya's webpage, they will return the accuracy score. Participants are expected to submit their answers by 31st of December 2018, and there will be a leaderboard. We are allowed to use all possible tools, including open-source libraries like TensorFlow/Keras and Scikit-Learn. For more practical information, see [Vid17].

The data sets consist of a total number of 8732 sound samplings (5434 training samplings and 3298 test samplings) with a constant sampling rate of 22050Hz. The length of each sampling is maximum 4 seconds, and they are distributed between ten classes, namely

- air conditioner
- car horn
- children playing
- dog bark
- drilling
- engine idling
- gun shot
- jackhammer
- siren
- street music.

The error is evaluated with the accuracy score, which is just how much of the data set that is classified correctly:

$$\text{Accuracy} = \frac{\sum_{i=1}^N I(y_i = t_i)}{N}. \quad (1)$$

## 2.2 Time domain

Sounds are longitudinal waves which actually are different pressures in the air. They are usually represented as a function in the time domain, which means how the pressure vary per time. This function is obviously continuous, but since computers represent functions as arrays, we cannot save all the information.

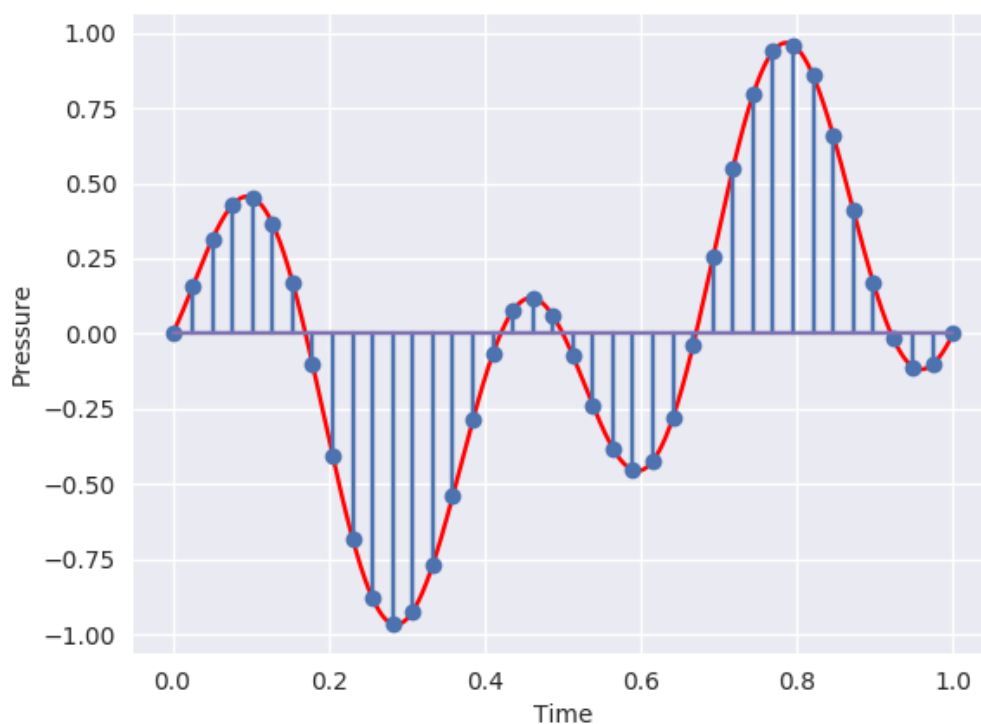


Figure 1: When sampling sound, we need to choose a high enough sampling. Here, the red line is the true sound and the blue dots are sampling points. The sampling rate here is 40 Hz, which is quite low and just for illustration.

How much information we loose depends on the sampling rate, which is the number of sampling points per second (see figure (1)). A rule of thumb is that one should have twice as high sampling rate as the highest sound frequency to keep the most important information. For instance, a human ear can perceive frequencies in the range 20-20000Hz, so around a sampling rate around 40kHz should be sufficient to keep all the information. Ordinary CD's use a sampling rate of 44.1kHz, but for speech recognition, 16kHz is enough to cover the frequency range of human speech. [Gei16]

## 2.3 Frequency domain

Sometimes, a frequency domain gives a better picture than the time domain. On one hand, one loses the time dependency, but on the other one gets information about which frequencies that are found in the wave. To go from the time domain to the frequency domain, one needs to use Fourier transformations, defined by

$$\hat{f}(x) = \int_{-\infty}^{\infty} f(t) \exp(-2\pi i x t) dx \quad (2)$$

where  $f(t)$  is the time function and  $\hat{f}(x)$  is the frequency function. Fortunately we do not need to solve this analytically, Fast Fourier Transformations (FFT) does the job numerically in no time. In figure (2), FFT is used on a time picture to obtain the corresponding frequency picture.

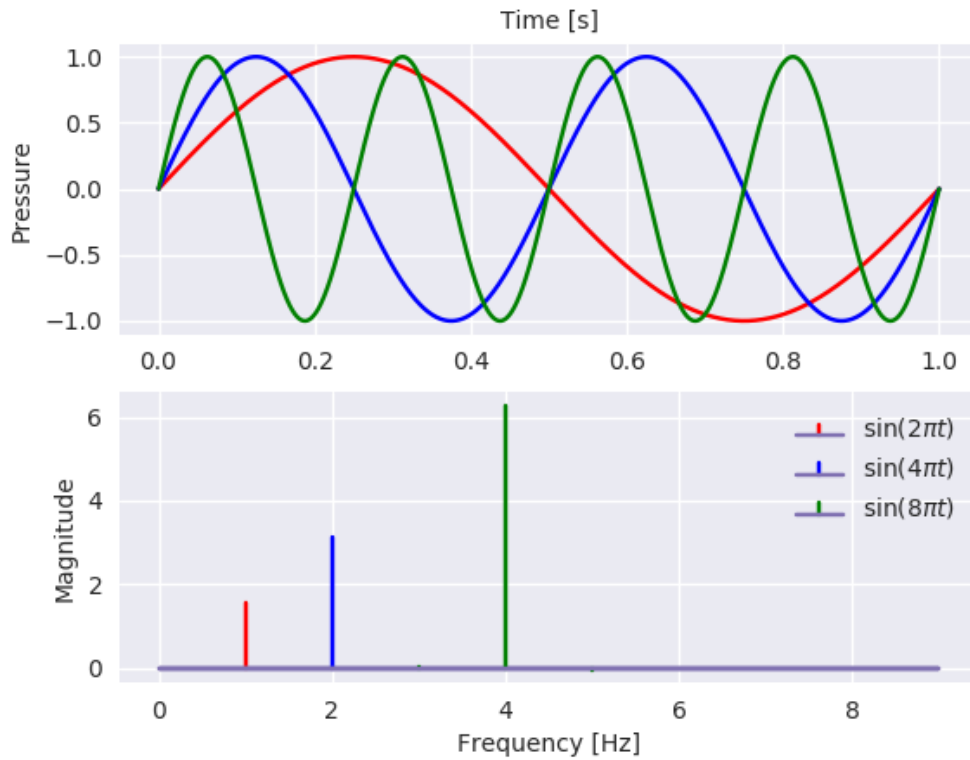


Figure 2: Time picture of three harmonic functions;  $\sin(2\pi t)$ ,  $\sin(4\pi t)$  and  $\sin(8\pi t)$  (upper part) and the corresponding frequency pictures (lower part).

## 2.4 Extract features

There are many ways we can extract features from sound, but the most interesting are those that are based on how humans extract information. First the spectrogram method will be described briefly, and we thereafter turn to the mel frequency ceptral coefficients method.

### 2.4.1 Spectrograms

As we have seen above, we are missing the frequency information in the time domain, and time information in the frequency domain. Is it possible to get the best of both worlds? This is what we try when creating a spectrogram, which is a time-frequency representation.

The procedure goes like this: we divide the spectrum into  $N$  frames and Fourier transform each frame. Thereafter we collect all the frequency domains in a  $N \times M$ -matrix, where  $M$  is the length of each frequency domain. We are then left with an image which hopefully hold the important information, and since Hidden Markov Models are known to implicitly model spectrograms for speech recognition, we have a reason to believe it will work. [Pra15]

### 2.4.2 Mel Frequency Cepstral Coefficients (MFCCs)

MFCCs are features widely used in automatic speech and speaker recognition. [Lyo13] They are inspired by how the human ear extracts features, and has been the main feature extraction tool for sound recognition since it was introduced in the 1980's.

It actually takes advantage of the spectrogram, and use it to reveal the most dominant frequencies. Thereafter, we sum the energy in each frame and take the logarithm, which is the way an ear works. Those energies are then what we call the Mel Frequency Ceptral Coefficients. Unlike the spectrogram, we are now left with a single array of information from the spectrum.

## 2.5 Examples from USC

We will end this section with a few examples of what a sound file may look like. In figure (3), we compare one of the dog bark sounds with one of the jackhammer sounds. On the top, we have the time domain, in the middle the spectrogram is presented and on the bottom we have plotted the MFCCs. As we can see, they are quite different, and distinguish between those two should not be too hard. The spectrograms also really makes sense, but it is hard for a human eye to extract information from the MFCCs.

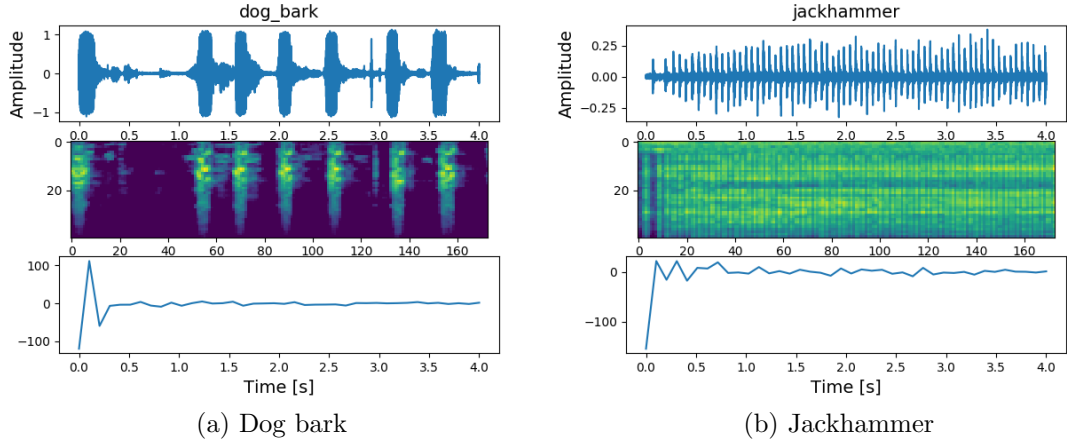


Figure 3: Time domains (top), spectrograms (middle) and MFCCs (bottom) for two random sound files, more specifically dog bark on the left and jackhammer on the right.

However, sometimes two time domains in the same category look totally different, for example if one compares the jackhammer in (4) with the one in (3). Will it be possible for our program to understand that they are in the same category? Also sometimes two time domain from different categories look similar, as we can see in (4). Can we really distinguish between them?

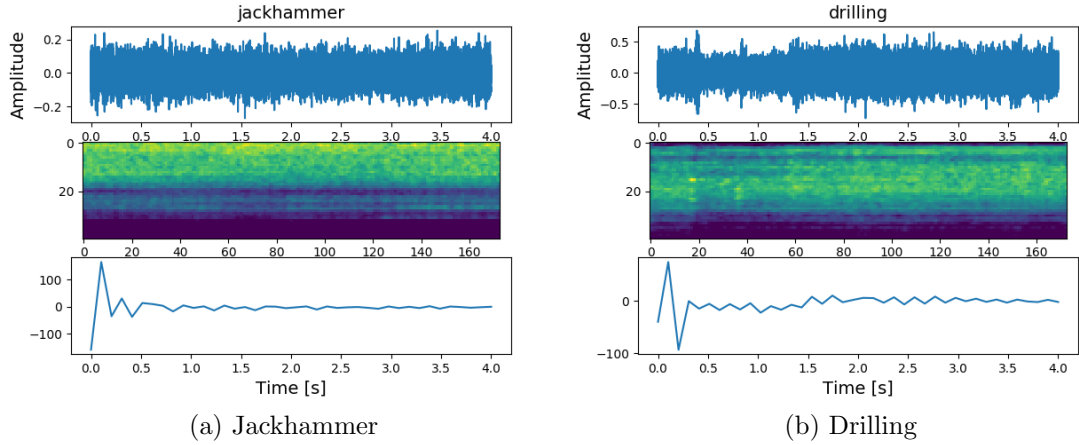


Figure 4: Time domains (top), spectrograms (middle) and MFCCs (bottom) for two random sound files, more specifically jackhammer on the left and drilling on the right.



### 3 Classification methods

Classification is very important in everyday life, and we often use classification even without thinking about it. A simple example is when we distinguish between people, which is usually an easy task for a human. For a computer, on the other hand, classification is difficult, but fortunately some great methods are developed mainly based on neural networks. In this context we will investigate several methods, spanning from **Logistic regression** to various neural networks like **Feed-forward Neural Networks**, **Convolutional Neural Networks** and **Recurrent Neural Networks**.

#### 3.1 Logistic regression

Logistic regression is a linear model, which means that each input node is multiplied with only one weight to get the net output. Since we want the probability of each class as a number between 0 and 1, we typically use the logistic function as activation function, and one hot encoder is often used to represent classes as arrays. In figure (5), logistic regression is illustrated as a single layer perceptron, with inputs on the left and side and outputs on right hand side.

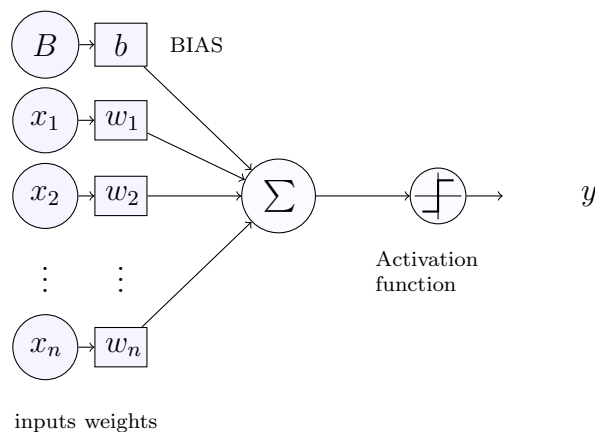


Figure 5: Logistic regression model with  $n$  inputs.

The first thing one might observe, is the bias on top, which is added to avoid the weights to explode or vanish. For instance, if we want the output to be 1, but the inputs are really small, the weights need to be large to give 1 without the bias. When adding the bias, we do not need that large weights, and the perceptron turns out to be more stable. The bias node value needs to be given, but the value does not really matter since the bias weight will be adjusted with respect to it.

The net output is then found to be

$$z = \sum_{i=1}^n x_i \cdot w_i + b \equiv \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{w}_i \quad (3)$$

where  $\mathbf{x} \equiv [1, x]$ ,  $\mathbf{w} \equiv [b, w]$  and the bias node value is set to 1. We get the real output by sending the net output through the activation function  $f$ ,  $y = f(z)$ .

What we really want to find is the optimal weight values, which is not necessary unambiguous. Initially the weight are set to random, preferably small values, and they are updated with the formula

$$\mathbf{w}_i^+ = \mathbf{w}_i - \eta \cdot [f(\mathbf{x}_i \cdot \mathbf{w}_i) - t_i] \mathbf{x}_i \quad (4)$$

with  $\eta$  as the learning rate,  $t_i$  as a target and  $\mathbf{w}_i^+$  as an updated weight. This formula was derived from the cross entropy cost function and makes use of gradient descent, see [Nor18].

### 3.2 Feed-forward Neural Networks (FNN)

Feed-forward neural networks work mostly in the same way as the single layer perceptron, but the difference is that the perceptron model is not single anymore. With two sets of weights, we call it a double layer perceptron, and with more layers it is called a multi layer perceptron model. In figure (6), we have illustrated a double layer perceptron.

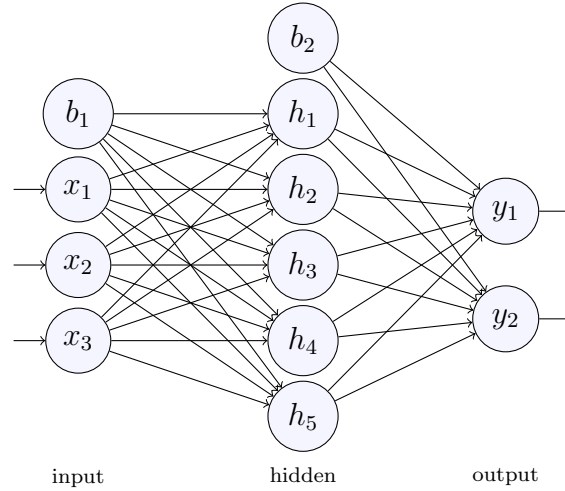


Figure 6: Neural network with 3 input nodes, 5 hidden nodes and 2 output nodes, in addition to the bias nodes.

As one can see, it differs from the single perceptron model in the way that it has one more layer. This makes the model nonlinear, and we are able to solve more complex problems than with the linear one. Actually, the double layer perceptron is able to approximate any continuous function according to the **the universal approximation theorem**.

Generally, each layer, included the output layer, consists of multiple nodes which gives us a matrix of weights between all the layers. The outputs are found in the same way as for the single perceptron, where the output at layer  $l$  is the dependent on the output at layer  $l - 1$ :

$$z_j^{(l)} = \sum_{i=1}^{h_{l-1}} z_i^{(l-1)} \cdot w_{ij}. \quad (5)$$

where  $h_{l-1}$  is the number of nodes in layer  $l - 1$  and the outputs  $z^{l-1}$  again are assumed to take the bias node values. We activate each layer with an activation function, which can be layer unique, similarly as for the single perceptron:  $y^l = f(z^l)$ . There is a bunch of activation function to choose from, and we consider it so important that we branched it into its own section, see section 5.

The remaining part now is the weight updates, which are based on backward propagation. We then need to update the last set of weights first, and work back to the input,

$$\begin{aligned} w_{ij}^{(3)} &= w_{ij}^{(3)} - \eta \cdot \delta_j \cdot z_i^{(3)} \\ w_{ij}^{(2)} &= w_{ij}^{(2)} - \eta \sum_{k=1}^{h_3} \delta_k \cdot w_{jk}^{(3)} f'(z_j^{(3)}) \cdot z_i^{(2)} \\ w_{ij}^{(1)} &= w_{ij}^{(1)} - \eta \sum_{k=1}^{h_3} \sum_{l=1}^{h_2} \delta_k \cdot w_{lk}^{(3)} f'(z_l^{(3)}) \cdot w_{jl}^{(2)} f'(z_j^{(2)}) \cdot z_i^{(1)} \end{aligned}$$

where we have used the short hand

$$\delta_j = (t_j - z_j^{(3)}) \cdot f'(z_j^{(3)})$$

and the  $z_i^{(1)}$  are just the input nodes. If we take a close look at the weight updating equations, we observe that there is a pattern. By defining a unique function  $\delta_{ij}^{(l)} = w_{ij}^{(l)} f'(z_i^{(l)})$  for each layer, we can generalize the equations above, and find an expression for an arbitrary number of layers. It might be easier to vectorize first.

### 3.3 Convolutional Neural Networks (CNN)

Convolutional neural networks are known to be good at image classification, but how can we use it on sound classification? The idea is to turn the samples into images, and for that one can for example use a mel spectrogram, as described in the theory section.

CNNs are based on FNNs, but the image is initially feed through a few layers that extract the most important information.

#### Convolutional layer

Initially, the image is sent into a convolutional layer, which is meant to reveal structures and shapes in the image. The way we do it, is to introduce a filter that we slide over the entire image and multiply with all pixels (with overlap). Every time the filter is multiplied with a set of pixels, we sum all the multiplications and add the value to an activation map. The activation map is completed after we have multiplied the filter with the entire image. A typical filter has dimensions 16x16, but depends on the image shape. It is important to choose a filter that is big enough to cover structures.

#### Pooling layer

It is common to insert a pooling layer in-between convolutional layers, but why do we do that? A pooling layer is just a way to reduce the dimensionality of the representation such that we do not need to optimize that many parameters, but also helps to control overfitting. It works the way that we divide the representation (usually an activation map) into regions of equal size and represent each region with one single number. Max pooling is apparently the most popular technique, which just represents each region with the largest number in that region, but it is also possible to use average pooling, min pooling etc.. [Kar13]

As an example, we can divide the image into 2x2 regions, which will reduce the size of the representation with 75%.

#### Dropout

Dropout is widely used in neural network layers, and is another method to prevent overfitting. The way it works is just to drop out units in the current layer.

To understand the idea, we need to take a close look at why overfitting occurs. Overfitting occurs when neighbor nodes collaborate, and become a powerful cluster which fits the model too specifically for the training set. To brake up those evil clusters, we can simply set random nodes to zero. [Sri+14]

### **Fully connected layer**

The last part of a convolutional network is often referred to as a fully connected layer, which is just a FNN. The output from the other layers needs to be flattened out before it is sent into the FNN.

## **3.4 Recurrent Neural Networks (RNN)**

Last, but not least, we will examine recurrent neural networks. In the time domain, the sound at one time is dependent on sounds at other times, but the methods above are not able to utilize that because they lack memory.

In RNNs, information is shared between the different nodes, in that sense they have a short memory. This makes the RNNs applicable when dealing with sequential data, such as handwriting recognition, speech recognition and other data sets where the input location matters.

In its simplest form, a so-called vanilla recurrent neural network, values on hidden nodes are fed back to the same nodes at the next time step, which creates a short-term memory. This is known to give better results on sequential data than ordinary FNNs, but it also has some cons that makes it useless in our case, especially the vanishing gradient problem. Instead, we will turn to more advanced methods which have a long-term memory as well, namely long short-term memory and gated recurrent units.

### **Long Short-Term Memory (LSTM)**

Long short-term memory is an extension of the vanilla network, and is able to remember multiple previous time steps at the same time. It is therefore suited to extract important features which repeat occasionally. The way it works, is that it has three gates: an input gate, forget gate and output gate. Based on the targets, the network decides whether the input is important (pass to output gate) or noise (send to forget gate). Because of the forget gate, the gradients will remain steep and the vanishing gradient problem will not occur. Even though an input is sent to the forget gate, the information will still be stored in the long time memory. [Don18]

### **Gated Recurrent Unit (GRU)**

Gated recurrent unit is a relatively fresh method, which has shown similar performances as the more traditional LSTM. GRUs come with an input gate, update gate and a reset gate, which work in a similar way as the LSTM gates. The update gate determines how much of the past information is relevant for the future and the reset gate decides how much of the past information to forget. [Cho+14]

## 4 Optimization

In data science, we are surrounded by large data sets and functions to extract information from them. After a function describing the data set is found, one typically wants to find its minimum or maximum to optimize various parameters. Since the extremes seldomly are available analytically, one has to rely on numerical optimization methods. There exist loads of methods, but we will only look at two gradient methods: Gradient Descent and ADAM, where the latter is preferred.

### 4.1 Gradient Descent (GD)

The gradient descent method is the simplest optimization method available. It simply moves the steepest slope towards the extreme with constant step length (learning rate), calculating the gradient based on the entire data set. Mathematically, we can write it as

$$\theta_{ij}^+ = \theta_{ij} - \eta \cdot \frac{\partial c(\boldsymbol{\theta})}{\partial \theta_{ij}} \quad (6)$$

where  $\theta_{ij}^+$  is the updated  $\theta_{ij}$  and  $\eta$  is the learning rate, in the same way as in section 3. GD has some major cons: it is likely to be stuck in a local extreme and because of its lack of adaptiveness, it is either really slow or likely to walk past the extreme. An example implementation looks like

```
for iter in range(T):
    for i in range(N):
        y[i] = feedforward(X[i])
        gradient = (y[i] - t[i]) * df(y[i])
        theta -= self.eta * gradient
```

### 4.2 ADAM

ADAM is often the preferred optimization method in machine learning, due to its computational efficiency, little memory occupation and implementation ease. Since it was introduced in a paper by D. P. Kingma and J. Ba in 2015 [KB14], the paper has collected more than 15000 citations! So what is so special about ADAM?

ADAM is a momentum based method that only rely on the first derivative of the cost function. Required inputs are exponential decay rates  $\beta_1$  and  $\beta_2$ , cost function  $c(\boldsymbol{\theta})$  and initial weights  $\boldsymbol{\theta}$ , in addition to the learning rate  $\eta$  and eventually a regularization  $\lambda$ . Iteratively, we first calculate the gradient of the cost function,

and use this to calculate the first and second moment estimates. Further, we bias-correct the moments before we use this to update the weights. An implementation may look like this

```
def ADAM(eta, b1, b2, c, theta):
    m = 0
    v = 0
    for i in range(100):
        grad = d(c)/d(theta)
        m = b1*m + (1 - b1)*grad
        v = b2*v + (1 - b2)*grad*grad
        m_ = m/(1 - b1**i)
        v_ = v/(1 - b2**i)
        theta -= eta * m_/(sqrt(v_) + lambda)
```

## 5 Activation

The activation function is a term that we repeatedly have used, but not really detailed. It is used to activate the outputs, which often need to take some certain properties. For instance, when doing classification, the outputs are probabilities and therefore between 0 and 1. A nonlinear activation function is often preferred to reinforce the non-linearity of the neural network.

Traditionally, the logistic function and the tanh function have been used as activation functions, since they were believed to work in the same way as the human brain. However, in 2012 a new network was introduced known as AlexNet, taking image recognition to a new level. They used a function named Rectified Linear Units (ReLU), which is zero for negative values. Based on this, some successful activation functions have been innovated in the recent years, such as Leaky ReLU and Exponential Linear Unit (ELU). They will be examined one after another below.

Before 2012, the logistic, the tanh and the pure linear functions were the standard activation functions, but then Alex Krizhevsky published an article where he introduced a new activation function called *Rectified Linear Units (ReLU)* which outperformed the classical activation functions. [4] The network he used is now known as AlexNet, and helped to revolutionize the field of computer vision. [5] After that, the ReLU activation function has been modified several times (avoiding zero derivative among others), and examples of innovations are *leaky ReLU* and *Exponential Linear Unit (ELU)*. All those networks are linear for positive numbers, and small for negative numbers. Often, especially in the output layer, a straight linear function is used as well.

In figure (??), *standard ReLU*, *leaky ReLU* and *ELU* are plotted along with the logistic function.

## 5.1 Logistic

## 5.2 ReLU

## 5.3 ELU

# 6 Code

Using the language of machine learning: Python. Useful packages sklearn, keras, TensorFlow. Parallel processing using Dask. Pandas used in reading. Librosa

## 6.1 Code structure

## 6.2 Implementation

# 7 Results

## 7.1 Logistic regression

Table 1: Accuracy, softmax output

Epochs	Regularization							
	0.1		0.001		0.00001		0	
	Train	Test	Train	Test	Train	Test	Train	Test
10	0.3925	0.4112	0.2000	0.2042	0.2543	0.2447	0.2308	0.2429
20	0.4572	0.4802	0.2927	0.2843	0.3320	0.3238	0.4135	0.3919
30	0.5741	0.5842	0.2948	0.2815	0.3387	0.3395	0.4393	0.4186
40	0.5780	0.5603	0.2948	0.2861	0.3433	0.3385	0.4402	0.4103
50	0.5916	0.5998	0.3007	0.2833	0.3389	0.3395	0.4590	0.4453



## 7.2 Feed-forward Neural Networks

Table 2: Accuracy, sigmoid function was used in hidden layers, softmax on output. Dropout was used.

Epochs	Hidden nodes							
	1024		2x1024		3x1024		4x1024	
	Train	Test	Train	Test	Train	Test	Train	Test
20	0.8684	0.8464	0.8603	0.8602	0.8541	0.8602	0.8341	0.8298
40	0.9245	0.8924	0.9204	0.8942	0.9195	0.8868	0.8990	0.8924
60	0.9416	0.9062	0.9508	0.9154	0.9317	0.9016	0.9254	0.8960
80	0.9572	0.8942	0.9556	0.9163	0.9508	0.9172	0.9418	0.9016
100	0.9607	0.9016	0.9579	0.9126	0.9627	0.9181	0.9510	0.9200

We could go further and increase the number of epochs. With 3x1024 hidden nodes and 1000 epochs, we were able to get an accuracy of 0.99 for the training set and 0.94 for the test set.

## 7.3 Convolutional Neural Networks

Table 3

Epochs	Hidden nodes			
	1024	2x1024	3x1024	4x1024
100	0.4371	0.4381	0.3333	
200	0.4379	0.4367	0.3333	
300	0.4371	0.4382	0.3333	
400	0.9930	0.9929	0.9655	
500	0.9935	0.9934	0.9679	

## 7.4 Recurrent Neural Networks

Table 4

Epochs	Hidden nodes			
	1024	2x1024	3x1024	4x1024
100	0.4371	0.4381	0.3333	
200	0.4379	0.4367	0.3333	
300	0.4371	0.4382	0.3333	
400	0.9930	0.9929	0.9655	
500	0.9935	0.9934	0.9679	

## 8 Discussion

## 9 Conclusion

## References

- [Kar13] A. Karpathy. *Convolutional Neural Networks for Visual Recognition*. <http://cs231n.github.io/convolutional-networks/>. [Online; accessed 06-December-2018]. 2013.
- [Lyo13] J. Lyons. *Mel Frequency Cepstral Coefficient (MFCC) tutorial*. <http://practicalcryptography.com/miscellaneous/machine-learning/guide-mel-frequency-cepstral-coefficients-mfccs/>. [Online; accessed 06-December-2018]. 2013.
- [Cho+14] K. Cho et al. “Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation”. In: *CoRR* abs/1406.1078 (2014). arXiv: 1406.1078. URL: <http://arxiv.org/abs/1406.1078>.
- [KB14] D. P. Kingma and J. Ba. “Adam: A Method for Stochastic Optimization”. In: *CoRR* abs/1412.6980 (2014). arXiv: 1412.6980. URL: <http://arxiv.org/abs/1412.6980>.
- [Sri+14] N. Srivastava et al. “Dropout: A Simple Way to Prevent Neural Networks from Overfitting”. In: *Journal of Machine Learning Research* 15 (2014), pp. 1929–1958. URL: <http://jmlr.org/papers/v15/srivastava14a.html>.
- [Pra15] K. Prahallad. “Speech Technology: A Practical Introduction.” In: (2015).
- [Gei16] A. Geitgey. *How to do Speech Recognition with Deep Learning*. <https://medium.com/@ageitgey/machine-learning-is-fun-part-6-how-to-do-speech-recognition-with-deep-learning-28293c162f7a>. [Online; accessed 06-December-2018]. 2016.
- [Vid17] Analytics Vidhya. *Practice Problem: Urban Sound Challenge*. <https://datahack.analyticsvidhya.com/contest/practice-problem-urban-sound-classification/>. [Online; accessed 06-December-2018]. 2017.
- [Don18] N. Donges. *Recurrent Neural Networks and LSTM*. <https://towardsdatascience.com/recurrent-neural-networks-and-lstm-4b601dd822a5>. [Online; accessed 07-December-2018]. 2018.
- [Fel18] R. Feloni. *20 years from now, regular people will have the personal assistants currently reserved for the rich and famous - through tech*. <https://nordic.businessinsider.com/ai-assistants-will-run-our-lives-20-years-from-now-2018-1?r=US&IR=T>. [Online; accessed 06-December-2018]. 2018.
- [Nor18] E. M. Nordhagen. *FYS-STK4155 - Applied data analysis and machine learning*. 1st Edition. 2018.