

# **PROJECT ML**

University of Hamburg

Department of Mathematics

*Jonas Eckhoff - Timo Greve*

*Max Lewerenz - Giulia Satiko Maesaka - John-Robert Wrage*

## **Machine Learning Methods Group 5**

# Contents

<b>List of figures</b>	<b>III</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 What is ML?</b>	<b>2</b>
2.1 Basis ML 1(maybe Data) . . . . .	2
2.2 Basis ML 2 . . . . .	2
<b>3 Classical learning</b>	<b>5</b>
3.1 Supervised . . . . .	5
3.1.1 Regression . . . . .	5
3.1.2 Classification . . . . .	5
3.2 Support Vector Machine . . . . .	7
3.2.1 Introduction . . . . .	7
3.2.2 Mathematics behind SVMs . . . . .	7
3.2.3 In our case . . . . .	8
3.3 Unsupervised . . . . .	8
3.3.1 Clustering . . . . .	9
3.3.2 Pattern search . . . . .	9
3.3.3 Dimension Reduction . . . . .	9
<b>4 Neural Networks and Deep Learning</b>	<b>10</b>
4.1 Convolutional Neural Networks CNN . . . . .	10
4.2 Recurrent Neural Networks RNN . . . . .	10
4.3 Generative adversarial networks GAN . . . . .	11
4.4 Autoencoders . . . . .	11
4.5 Perceptrons MLP . . . . .	12
<b>5 Ensemble Methods</b>	<b>13</b>
5.1 Bagging . . . . .	13
5.2 Boosting . . . . .	13
5.3 Stacking . . . . .	13
<b>6 Reinforcement learning</b>	<b>14</b>
<b>Bibliography</b>	<b>15</b>

## List of Figures

1	The flow of developing an engine based on machine learning. . . . .	3
2	ML-Overview . . . . .	4
3	Supervised vs unsupervised . . . . .	5
4	K-Nearest Neighbors . . . . .	6
5	K-Nearest Neighbors . . . . .	6
6	Dacision Tree . . . . .	7
7	Clustering . . . . .	9
8	Example of a CNN architecture . . . . .	10
9	Example of an autoencoder architecture . . . . .	12
10	Bagging/Boosting . . . . .	13

# 1 Introduction

What is this paper for? What are the main goals?

E.g. this is for ourself. We want to build an overview about the common algorithms, how they work and their advantages/disadvantages.

## 2 What is ML?

Machine learning is the field of study responsible for developing and understanding algorithms that allow computer systems to perform certain tasks independently. The type of tasks considered are those for which any rule-based approach is unfeasible, either because a precise set of rules is unclear, for example in pattern recognition, or because computing according to the rules is too complex, for example in playing Go.

Given a data set, the goal of a machine learning method can be of two types, one is to use the data to find a function, which will then match new input data to a value; another is to find a subdivision of the data set based on similarities. The methods used for the first type of goal are named *supervised learning* methods and the given data set, called *training set*, must be structured in pairs input-output. For the second type of goal, one uses *unsupervised learning* methods.

In a supervised learning method, the output value is either quantitative or qualitative. For the problem of predicting the price of a house given its area and location, the output is a quantitative value; these are called *regression* problems. On the other hand, if the task is to address an emotion to a given face expression, the output is a qualitative value encoding a certain emotion; these are called *classification* problems.

Supervised learning methods are structured in two parts. Firstly it consists of a predetermined set of functions, called *hypothesis space*, from which a final function will be chosen. Secondly the method needs a predetermined way of searching through the hypothesis space, therefore it requires a predetermined *loss function* and *optimization algorithm*. The choice of hypothesis space, loss function and optimization algorithm depends on the task and different methods need to be implemented and evaluated. For the task of recognizing facial emotion, it is known that between the methods so far developed, the best results are achieved by convolutional neural networks. Still, this project tests different methods one more time, for the sake of experience.

### 2.1 Basis ML 1(maybe Data)

We could speak a little bit more about data in this topic?

Finally I would suggest to finish with a Table.

I would structure it like: Algorithm, Main Idea(in like 3-5 sentences), possible application(image processing, clustering, etc.) advantages, disadvantages

### 2.2 Basis ML 2

TBA

Best parameters set found on development set:

```
{'C': 10, 'gamma': 0.0001, 'kernel': 'rbf'}
```

Grid scores on development set:

```
0.695 (+/-0.055) for {'C': 1, 'gamma': 0.001, 'kernel': 'rbf'}
0.299 (+/-0.001) for {'C': 1, 'gamma': 0.0001, 'kernel': 'rbf'}
0.708 (+/-0.037) for {'C': 10, 'gamma': 0.001, 'kernel': 'rbf'}
0.711 (+/-0.028) for {'C': 10, 'gamma': 0.0001, 'kernel': 'rbf'}
0.684 (+/-0.036) for {'C': 100, 'gamma': 0.001, 'kernel': 'rbf'}
0.695 (+/-0.028) for {'C': 100, 'gamma': 0.0001, 'kernel': 'rbf'}
0.669 (+/-0.044) for {'C': 1000, 'gamma': 0.001, 'kernel': 'rbf'}
0.669 (+/-0.042) for {'C': 1000, 'gamma': 0.0001, 'kernel': 'rbf'}
0.628 (+/-0.036) for {'C': 1, 'kernel': 'linear'}
0.620 (+/-0.021) for {'C': 10, 'kernel': 'linear'}
0.604 (+/-0.036) for {'C': 100, 'kernel': 'linear'}
0.609 (+/-0.034) for {'C': 1000, 'kernel': 'linear'}
```

Detailed classification report:

The model is trained on first 2000 entries of development set.  
The scores are computed on the first 200 entries of evaluation set.

	precision	recall	f1-score	support
3	0.61	0.93	0.74	115
6	0.67	0.19	0.30	84
accuracy			0.62	199
macro avg	0.64	0.56	0.52	199
weighted avg	0.63	0.62	0.55	199

Figure 1: The flow of developing an engine based on machine learning.

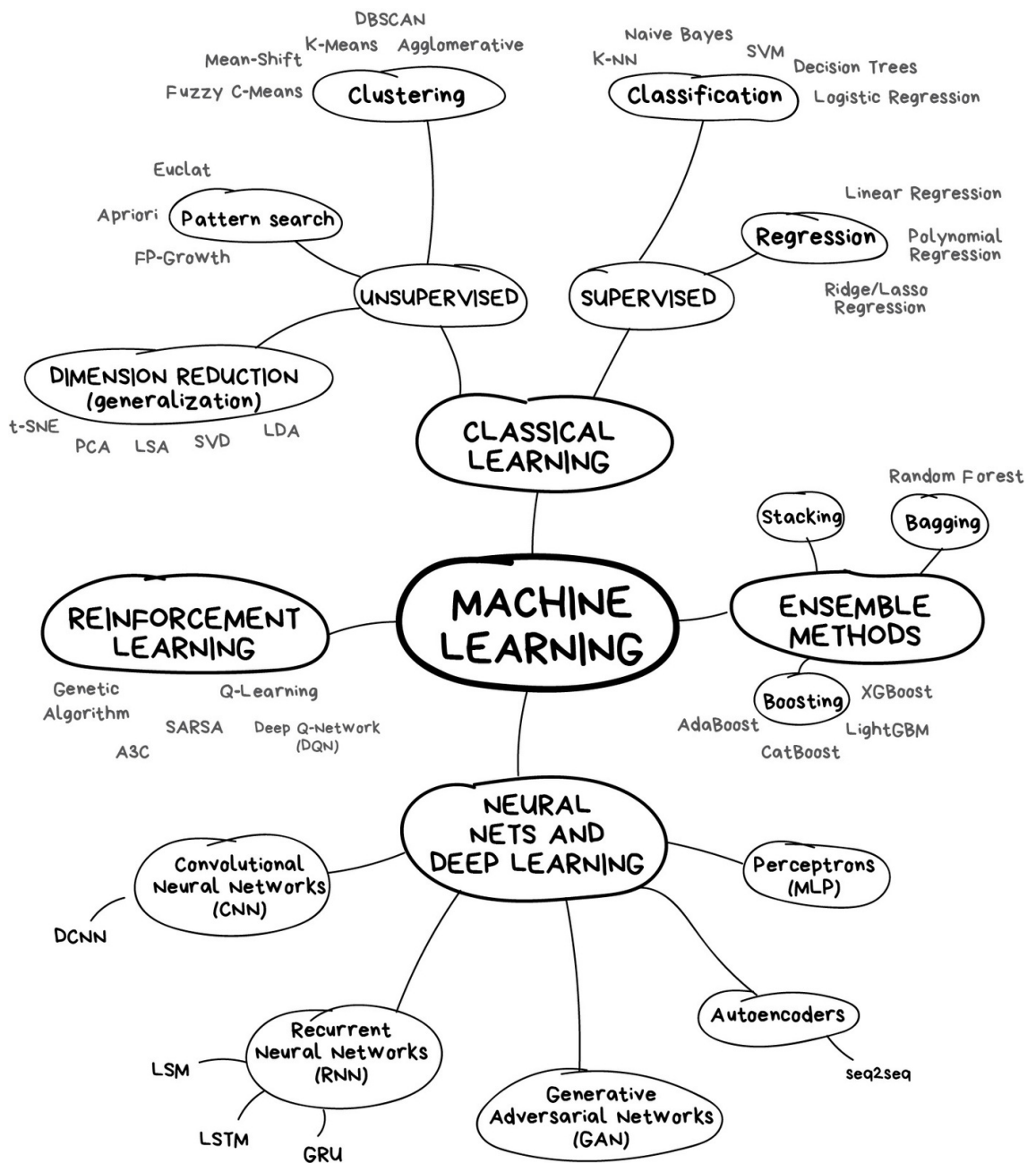


Figure 2: ML-Overview

### 3 Classical learning

Difference between Supervised unsupervised. We need to check what we wrote in chapter 2 before writing this part.

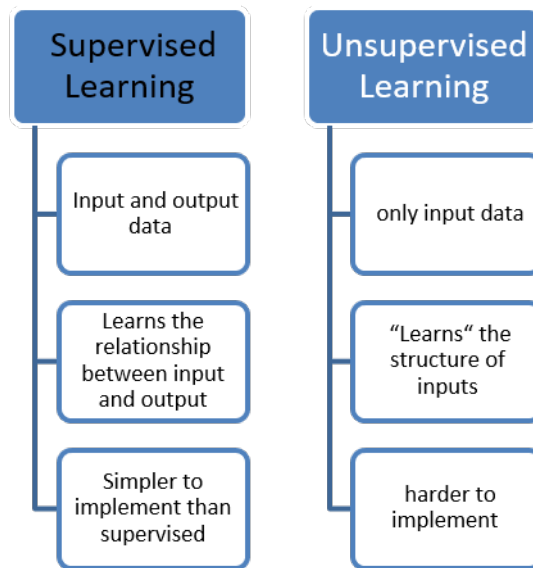


Figure 3: Supervised vs unsupervised

#### 3.1 Supervised

In supervised learning, we have a collection of input-output data pairs. In this case, the machine has a "supervisor" or a "teacher" who gives the machine all the answers, like whether it's a cat in the picture or a dog. The teacher has already divided (labeled) the data into cats and dogs, and the machine is using these examples to learn. One by one. Dog by cat.

##### 3.1.1 Regression

Regression is basically classification where we forecast a number instead of category. Examples are car price by its mileage, traffic by time of the day, demand volume by growth of the company etc.

Like for all other subsections i would suggest that we structure the Algorithms like in the table at Ch.2.(**Algorithm**, **possible application**, etc.), but also add some links etc. Trainingsdaten und berechnet den Abstand zu den jeweiligen Merkmalen. In

##### 3.1.2 Classification

Splits objects based at one of the attributes known beforehand. Separate socks by based on color, documents based on language, music by genre. It is not limited to one attribute. Common algorithms are:  
Naive Bayes

K-Nearest-Neighbours(KNN)

The KNN algorithm assumes that similar things exist in close proximity. In other words, similar things are near to each other.



The straight-line distance (Euclidean distance) is a popular and familiar choice. Note that there are other ways of calculating distance, and one way might be preferable depending on the problem we are solving. Other distances could be Manhattan or Hamming.

**Algorithm:**

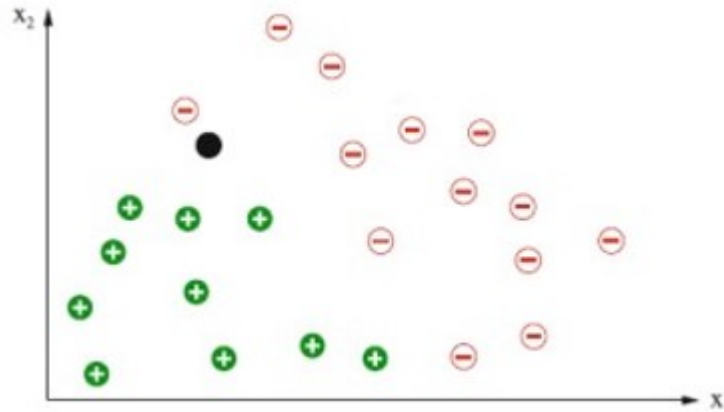


Figure 4: K-Nearest Neighbors

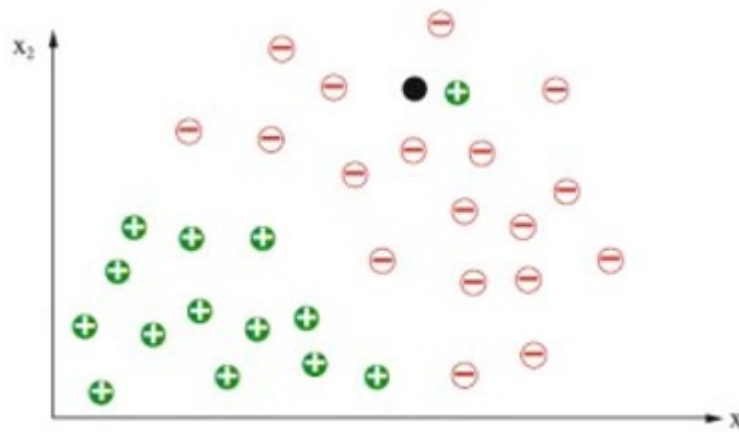


Figure 5: K-Nearest Neighbors

**Possible application**

(classification-) Decision tree

All the data is automatically divided to yes/no questions. This is done by selecting the question such that the best information gain is created (splitting the data via a feature which can separate the different labels best).

**Algorithm:**

K-Nearest-Neighbours(KNN)

**Possible application**



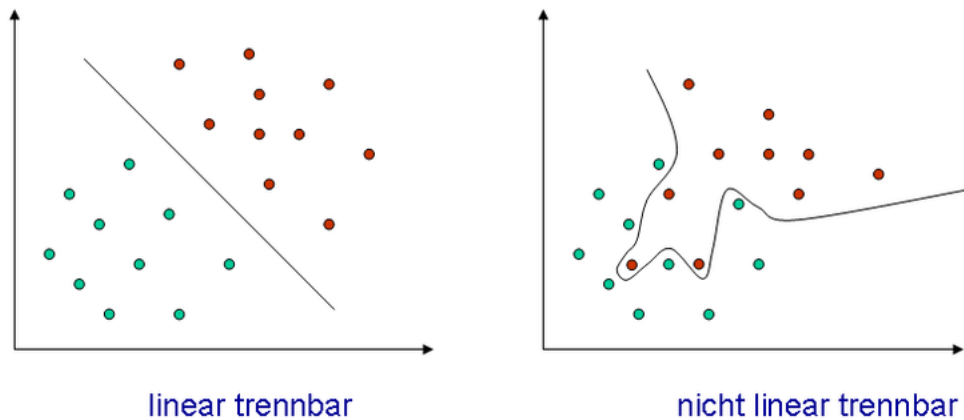
Figure 6: Dacision Tree

Trees are used in ensemble methods (Random forest, adaboost,...). See chapter 5.

## 3.2 Support Vector Machine

### 3.2.1 Introduction

Support Vector Machines (SVM) are supervised learning algorithms for classification. They are one of the most widely used supervised learning algorithms. SVMs offer a high accuracy when it comes to classification. The idea behind a SVM is to find a hyperplane that seperates classes of datapoints with a large margin. Where the margin is the smallest distance between the closest datapoint  $x$  of one class and the hyperplane (therefore the name 'support vector'). Because of this feature the SVM is sometimes also called **large margin classifier**. The so called kernel trick can be applied when non linear data has to be classified.



### 3.2.2 Mathematics behind SVMs

For given input data  $X = \{x_1, \dots, x_m\}$ , with  $x_i \in \mathbb{R}^n$  and matching result vector  $Y = \{y_1, \dots, y_m\}$ , with  $y_i \in \{-1, 1\}$ . We want a classifier

$$h : \mathbb{R}^n \rightarrow \{-1, 1\} \text{ such that } h(x) = \begin{cases} 1 & \text{,if } w^T x + b > 0, \\ -1 & \text{,if } w^T x + b < 0. \end{cases} \text{ with } w \in \mathbb{R}^n \text{ and } b \in \mathbb{R}. \text{ We want h to be}$$

correct for most samples. This formulation leads to the following primal problem:

$$\min_{w,b,s} \frac{1}{2} w w^T + C \sum_{i=1}^m s_i,$$

$$\text{s.t.: } y_i(w^T x_i + b) \geq 1 - s_i, s_i \geq 0.$$

Where  $s_i$  are the so called slack variables that should denote the distance from the correct margin if a point

is misclassified.

Intuitively we want to maximize the margin what results in minimizing  $\|w\|^2$  including a penalty when something is misclassified.  $C$  controls the strength of the misclassification penalty. One could view it as an inverse regularization parameter. A large  $C$  results in overfitting and a smaller  $C$  results in a "smoother" fit.

The dual problem to the above Primal problem is:

$$\min_{\lambda} \frac{1}{2} \lambda^T Q \lambda - e^T \lambda$$

s.t.:  $y^T \lambda = 0$  and  $0 \leq \lambda_i \leq C$  for all  $i = 1, \dots, m$ . The entries of the matrix  $Q \in \mathbb{R}^{m \times m}$  are given by  $Q_{i,j} = y_i y_j \langle x_i, x_j \rangle$ .

This shows that the minimization only depends on the scalar product of  $x_i$  and  $x_j$  and we can apply the kernel trick. Instead of the standard scalar product we can use a kernel function  $K(x_i, x_j)$ .

### 3.2.3 In our case

We use the python library sklearn, which can easily be installed via pip. First we get the training data and split it into training and test data.

```
In [5]: file_name = '../Resources/dataset/fer2013.csv'
df = pandas.read_csv(file_name)
print(df)
```

	emotion	pixels	Usage
0	70 80 82 72 58 58 60 63 54 58 60 48 89 115 121...	Training	
1	0 151 150 147 155 148 133 111 140 170 174 182 15...	Training	
2	231 212 156 164 174 138 161 173 182 200 106 38...	Training	
3	4 24 32 36 30 32 23 19 20 30 41 21 22 32 34 21...	Training	
4	6 4 0 0 0 0 0 0 0 0 0 3 15 23 28 48 50 58 84...	Training	
...	...	...	
35882	6 50 36 17 22 23 29 33 39 34 37 37 39 43 48 5...	PrivateTest	
35883	3 178 174 172 173 181 188 191 194 196 199 200 20...	PrivateTest	
35884	0 17 17 16 23 28 22 19 17 25 26 20 24 31 19 27 9...	PrivateTest	
35885	3 30 28 28 29 31 30 42 68 79 81 77 67 67 71 63 6...	PrivateTest	
35886	2 19 13 14 12 13 16 21 33 50 57 71 84 97 108 122...	PrivateTest	

[35887 rows x 3 columns]

Next we fit the model to the training data for the chosen parameters and save it with the pickle model.

```
In [ ]: clf = svm.SVC(decision_function_shape = 'ovo', kernel= input_kernel, cache_size=2000, C = in_c, gamma= in_gamma)
clf.fit(X_train, Y_train)
with open(file_name, 'wb') as fid:
    pickle.dump(clf, fid)
```

Now the model can give a prediction on new images.

```
In [55]: clf_loaded.predict(X_test[50:55])
Out[55]: array([3, 3, 5, 3, 4])
```

#### Sources:

- Martin Lotz: *Mathematics of Machine Learning*. Lecture Notes, Warwick (UK), 2020.
- scikit-learn: <https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>  
<https://scikit-learn.org/stable/modules/svm.html>
- Coursera: Machine Learning, by Stanford University  
<https://www.coursera.org/learn/machine-learning/home/welcome>
- Wikipedia: Support vector machine. [https://en.wikipedia.org/wiki/Support\\_vector\\_machine](https://en.wikipedia.org/wiki/Support_vector_machine)

## 3.3 Unsupervised

Because our topic was a classification problem, we did not add content to this chapter. Like mentioned, in unsupervised learning one tries to find structure in unlabeled data.



Figure 7: Clustering

### 3.3.1 Clustering

TBA

### 3.3.2 Pattern search

TBA

### 3.3.3 Dimension Reduction

TBA

## 4 Neural Networks and Deep Learning

### 4.1 Convolutional Neural Networks CNN

A Convolutional Neural Network is a Deep Learning algorithm which is primarily used to classify images. It can more easily capture spatial (and temporal) dependencies since the convolution filters use a linear combination of neighboring pixels to calculate an output value. After a convolution filter is applied, a Max-Pooling filter can be applied, which lowers the amount of nodes in the network. A 2x2 Max-Pooling filter, for example, replaces every 2x2 square with the maximal value in that square, resulting in lowering the amount of nodes after that layer to a quarter. At the end the multidimensional layers are flattened to a one dimensional vector that is connected to the output via a fully connected layer.

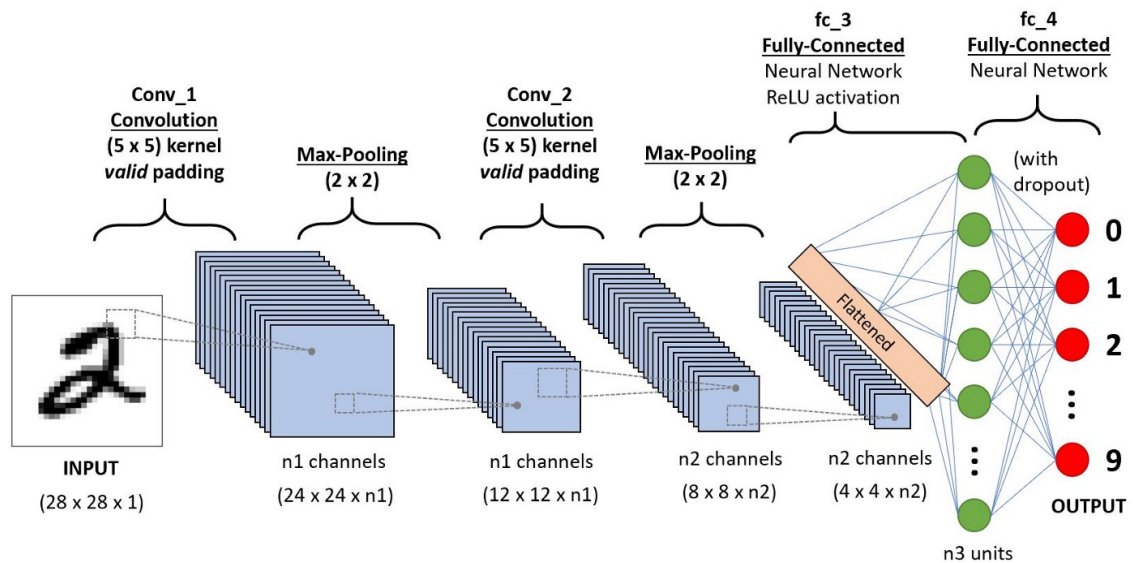


Figure 8: Example of a CNN architecture

Helpful links:

A Comprehensive Guide to Convolutional Neural Networks — the ELI5 way

<https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-net>

Since we want to detect emotions in images, a convolutional neural network would be a good choice. It is the go-to type of network for these kind of computer vision problems.

### 4.2 Recurrent Neural Networks RNN

A recurrent neural network (RNN) is a type of artificial neural network commonly used in speech recognition and natural language processing. It uses feedback loops, that allow information to persist (the network has memory), to process sequences of data.

RNNs are not useful to analyse facial expression in an image.

### 4.3 Generative adversarial networks GAN

Generative adversarial networks (GANs) are algorithmic architectures that use two neural networks, pitting one against the other in order to generate new, synthetic instances of data that can pass for real data. They are used widely in image generation, video generation and voice generation.

One neural network, called the generator, generates new data instances, while the other, the discriminator, evaluates them for authenticity; i.e. the discriminator decides whether each instance of data that it reviews belongs to the actual training dataset or not.

Helpful links:

A Beginner's Guide to Generative Adversarial Networks (GANs)

<https://pathmind.com/wiki/generative-adversarial-network-gan>

GANs can be used to generate images and videos. We could use this type of network to expand training data. Having a well trained GAN would allow us to create new images for specific emotions that look convincingly real in order to prevent overfitting by creating an endless stream of new images. It may also be used to create video training data that can be used to train a neural network with memory to better capture changes in emotion in a video stream. The training of a GAN however takes a very long time compared to a simple CNN which makes it probably unfeasible for us even though it would create interesting possibilities to explore.

### 4.4 Autoencoders

Autoencoders are a specific type of feedforward neural networks where the input is the same as the output. They compress the input into a lower-dimensional code and then reconstruct the output from this representation. The code is a compact “summary” or “compression” of the input. An autoencoder consists of 3 components: encoder, code and decoder. The encoder compresses the input and produces the code, the decoder then reconstructs the input only using this code.

Helpful links:

Applied Deep Learning - Part 3: Autoencoders

<https://towardsdatascience.com/applied-deep-learning-part-3-autoencoders-1c083af>

An autoencoder is used to reduce the dimensionality. Since there is a lot of redundancy in images, especially if they all contain more or less centered faces, it can be expected that an autoencoder could shrink the input size considerably before using a different neural network. This could potentially make the training process a lot quicker. We would however not be able to use a convolutional neural network anymore because the 2 dimensional image structure is lost by removing the redundancy with the autoencoder. Moreover some information may be lost in the process which could affect the overall accuracy of the neural network.

The effectiveness of an autoencoder could possibly be explored if there are given time constraints for the training process. It can also be compared to principal component analysis (PCA) which removes redundancy in a more controlled way.

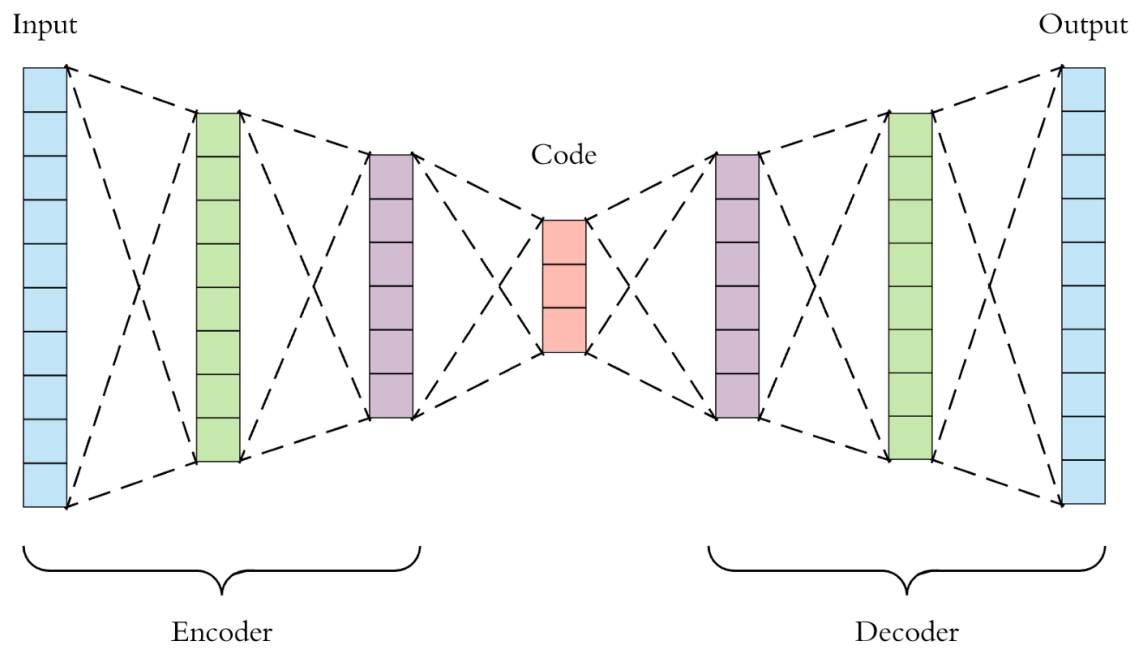


Figure 9: Example of an autoencoder architecture

## 4.5 Perceptrons MLP

## 5 Ensemble Methods

The main idea of ensemble is to combine algorithms to achieve a better accuracy than with a single model. Ensemble methods can be split into heterogeneous and homogeneous ensembles. In homogeneous ensembles while in heterogeneous ensembles different algorithms can be combined. We implemented a homogeneous method called adaboost with decision tree stumps (depth=1).

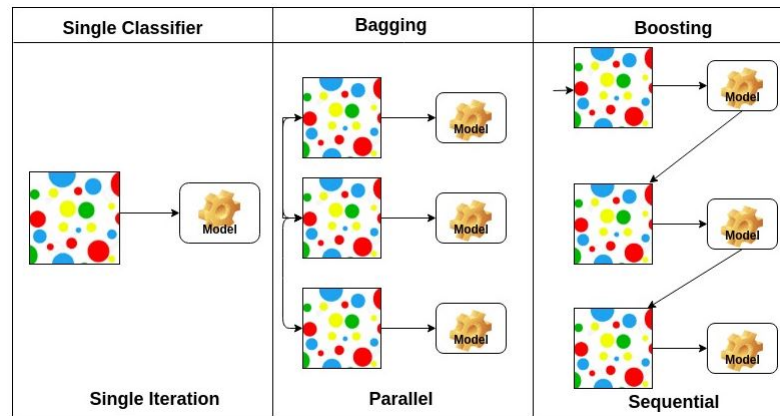


Figure 10: Bagging/Boosting

### 5.1 Bagging

### 5.2 Boosting

### 5.3 Stacking



## **6 Reinforcement learning**

## Bibliography

- [1] DEUTSCHE AKTUARVEREINIGUNG E.V., 2012: Unisex-Tarifierung, abgerufen am 10. November 2019: [https://aktuar.de/fachartikelaktuaraktuell/Lebensversicherung\\_Unisex\\_Aktuaraktuell\\_19.pdf#search=unisex](https://aktuar.de/fachartikelaktuaraktuell/Lebensversicherung_Unisex_Aktuaraktuell_19.pdf#search=unisex)
- [2] ERTEL, W. (2016): *Grundkurs Künstliche Intelligenz*. Springer Vieweg, Wiesbaden. 4. Auflage.