



# Predicting the Outcome of Soccer Matches

for the

**Master of Science**

from the Course of Studies Informationssysteme

at the Technische Hochschule in Ulm

by

**Lisa Boos, Hemlata Prajapati, Khaled Jallouli, Till Hoffmann**

July 2020

**First Reviewer**

Prof. Dr. Goldstein

**Second Reviewer**

Prof. Dr. Herbort

## **Author's declaration**

Hereby I solemnly declare:

1. that this document, titled *Predicting the Outcome of Soccer Matches* is entirely the product of my own scholarly work, unless otherwise indicated in the text or references, or acknowledged below;
2. I have indicated the thoughts adopted directly or indirectly from other sources at the appropriate places within the document;
3. this document has not been submitted either in whole or part, for a degree at this or any other university or institution;
4. I have not published this document in the past.

I am aware that a dishonest declaration will entail legal consequences.

Ulm, July 2020

---

Lisa Boos, Hemlata Prajapati, Khaled Jallouli, Till Hoffmann

## Abstract

These days Machine Learning, Neural Networks or Data Mining are common buzzwords in the world of computer science. But a lot of people do not know what they are actually talking about, when using these kind of words. Therefore it makes sense to learn about neural networks and machine learning in general during the course of a masters degree in information technology. Fortunately there are many data sets online to use for analysis to get a basic understanding of the topic at hand. The European Soccer Database from Kaggle [6] is such a data set.

During the team project of the information systems master at the Technische Hochschule Ulm it was our task to use the european soccer database for analytics in order to be able to predict the outcome of soccer matches.

From the database we extracted several different features, like amount of goals shot by each team, amount of wins, draws and losses for each team as well as shot-accuracy and shot-efficiency per team and the overall ball possession of each soccer-team per match. Using an algorithm which we adapted for our special data [2] we created a sliding window which aggregates the extracted features over the last 10 soccer-games for each match in the database in chronological order. Because a lot of data samples in the database were incomplete regarding some of the extracted features, we created 5 different versions of the sliding window. Each containing a different amount of features and therefore a different amount of data samples. This gives us the possibility to test classifiers with various different models and later pick the solution that works best. The sliding window option 1 uses the least amount of features (13 features) and has the highest number of data samples (20823 data samples). Option 2 has 21 features and 7033 data samples, option 3 uses 29 features and 7033 samples, option 4 contains 25 features and 6996 samples and option 5 has the most features (33) and 6996 data samples.

We used the 5 sliding window options to train and test various classifiers such as a Decision Tree, a Multi-Layer Perceptron and various different basic sequential neural nets. For the Decision Tree and the Multi-Layer Perceptron we used the scikit-learn API. For the basic sequential neural nets we used the Tensorflow/Keras framework, which is state of the art when doing data analytics tasks.

The decision tree using the first sliding window option and a depth of 4 gives us a testaccuracy of 52,95%. The Multi-Layer Perceptron uses the sliding window option 1 aswell and has a test-accuracy of 53,45%. It has 2 hidden layers with 52 and 32 neurons. The best version of the Keras Sequential Neural Network was trained with sliding window option 3, has 2 hidden layers (21, 21 neurons) and a test-accuracy of 56,25%. It is also the best classifier in general we were able to find until now. At first sight model accuracies

barely over 50% don't seem very good, but considering that the home team has a base chance of 46% to win the match and soccer is still a gambling game up to some point those accuracies aren't too bad at all.

For the goal prediction part we used two different techniques, regression and multi class classification. Regression is showing here better results with a test accuracy of [WE HAVE TO FILL IN] %. Nevertheless multi class classification showed a decent accuracy with around 30%.

For a good visualization we implemented a homepage for predicting soccer games with an API and a backend. This has different functionalities like retraining the models, showing previous predictions and showing upcoming matches.

Nevertheless, the team will continue its work on the project striving for better models by using additional features and varying the amount of features, adapting the split of training- and test-data, trying other kinds of classifiers and tweaking the existing models by changing parameters and input functions.

# Contents

<b>Acronyms</b>	<b>VI</b>
<b>List of Figures</b>	<b>VII</b>
<b>List of Tables</b>	<b>VIII</b>
<b>Listings</b>	<b>IX</b>
<b>1. Introduction</b>	<b>1</b>
1.1. Motivation . . . . .	1
1.2. Goals . . . . .	1
1.3. Procedure . . . . .	3
<b>2. Data Understanding</b>	<b>6</b>
<b>3. Feature Selection and Extraction</b>	<b>8</b>
<b>4. Data Preprocessing</b>	<b>10</b>
4.1. Sliding Windows or Data Options . . . . .	10
4.2. Data Profiling Part 2 . . . . .	13
4.3. Prediction of the final result . . . . .	14
4.4. Prediction of the scored goals . . . . .	16
4.5. Functions and Input Data For Neural Networks . . . . .	19
<b>5. Modeling</b>	<b>21</b>
5.1. Predicting the winner . . . . .	21
5.2. Predicting the final goals . . . . .	28
<b>6. Evaluation</b>	<b>37</b>
6.1. Prediction of the final result . . . . .	37
6.2. Predicting the final goals . . . . .	38
<b>7. Deployment</b>	<b>42</b>
7.1. Backend . . . . .	42
7.2. API . . . . .	45
7.3. Frontend . . . . .	47
<b>8. Conclusion</b>	<b>48</b>
<b>Bibliography</b>	<b>49</b>

<b>A. Appendix</b>	<b>50</b>
A.1. Architectures for various neural nets . . . . .	50
A.2. Daily Scrum Logs first therm . . . . .	54
A.3. Daily Scrum Logs second therm . . . . .	66

# **Acronyms**

# List of Figures

1.1.	The six stages of CRISP-DM [7] . . . . .	3
2.1.	Entity Relationship Diagram European Soccer Database . . . . .	7
4.1.	Data for Option 1 . . . . .	11
4.2.	Sliding Window Option 1 . . . . .	11
4.3.	Sliding Window Option 2 . . . . .	12
4.4.	Sliding Window Option 3 . . . . .	12
4.5.	Sliding Window Option 4 . . . . .	13
4.6.	Pandas Profiling Sliding Window 3 . . . . .	14
4.7.	Results Sliding Window Option 2 . . . . .	15
4.8.	Scored Goals Sliding Window Option 2 . . . . .	16
4.9.	The initial number of rows per each scored goal(s) . . . . .	18
4.10.	Encode and Decode functions . . . . .	18
4.11.	Encoding results . . . . .	19
4.12.	Excerpt of Input Data for Neural Networks as Numpy Array . . . . .	20
5.1.	Decision Tree (max_depth=4) . . . . .	23
5.2.	MLP Cost function . . . . .	23
5.3.	Evolution of the loss function of the models over time . . . . .	25
5.4.	Baseline model for regression using deep learning . . . . .	31
5.5.	Display of Validation and Training Loss . . . . .	32
5.6.	Predictions on Test Dataset . . . . .	33
5.7.	Multi Class Classification loss graph . . . . .	35
5.8.	Multi Class Classification accuracy graph . . . . .	36
7.1.	Architecture of the backend . . . . .	42

# List of Tables

5.1.	Test-accuracy of various models with different parameters . . . . .	26
5.2.	Test-accuracies for variation of hidden layers and neurons for sliding window option 1 . . . . .	26
5.3.	Test-accuracies for variation of hidden layers and neurons for sliding window option 2 . . . . .	26
5.4.	Test-accuracies for variation of hidden layers and neurons for sliding window option 3 . . . . .	27
5.5.	Test-accuracies for variation of hidden layers and neurons for sliding window option 4 . . . . .	27
5.6.	Test-accuracies for variation of hidden layers and neurons for sliding window option 5 . . . . .	27
5.7.	Quality Model for Away Team and Home Team . . . . .	33
6.1.	Comparison of classifiers . . . . .	37
6.2.	Round and Decode predicted values . . . . .	38
6.3.	Calculate the degree difference . . . . .	39
6.4.	Quality Model for models with different hidden units for dataset sliding02 .	40
6.5.	Quality Model for models with different hidden units for dataset sliding03 .	40
A.1.	Test Variation of Hidden-Layers and Neurons for Neural Nets . . . . .	50

# Listings

3.1. XML structure of shot-statistics . . . . .	8
3.2. XML structure of ball-possession-statistics . . . . .	9
4.1. SQL code for Sliding Window . . . . .	10
4.2. Python code for matches_all.csv . . . . .	10
4.3. Python code for normalization . . . . .	15
4.4. Python code for encoding classes . . . . .	16
4.5. Scored goals Python code for normalization . . . . .	17
4.6. Scored goals Python code for encoding classes . . . . .	18
5.1. Python code for simple Keras Sequential Model Instantiation . . . . .	24
5.2. Python code for multi class classification . . . . .	33
5.3. Python code for multi class classification model evaluation . . . . .	34
7.1. JSON structure of API call on /soccerGamesClassification . . . . .	46

# 1. Introduction

The project is about predicting the outcome of soccer games and training a neural network. For the training the database from Kaggle in source [6] has been used.

## 1.1. Motivation

Machine Learning and Neural Networks is a main topic in the software development field. Many companies start to try analyze different kind of data with this approach. Because of this reason as master students in information systems it is very interesting to gather knowledge about machine learning and neural networks. For this matter a prediction of soccer matches is a very nice procedure to do. For this topic there is not many additional knowledge to gather, what you need for analyzing the data and prepare them for prediction. Additionally there is a free database with data of matches for a couple of years. In the beginning the knowledge about machine learning is very low and the goal is to improve the handling of python in combination with machine learning techniques. This includes pre-processing data with Pandas, normalizing data sets and extracting the right features, as well as developing adequate models for the machine-learning algorithm. Additionally to the gaining of knowledge we like to create an algorithm, which predicts the outcome of soccer games with a decent accuracy.

## 1.2. Goals

The main goals of this project can be divided into 2 parts:

- **Gaining knowledge**

- Improve skills in python

For the the project and for future work with neural networks it is necessary to become familiar with python. There are many things, which are different in python than in other programming languages. All these things have to be discovered and it is necessary to learn how to deal with the python syntax. Additionally its important to gain knowledge about some library, which you can use with python and make it easier to solve some problems.

- Gaining knowledge about data pre-processing

It is very important to edit the data in a way, that it is ready for a neural network. Algorithms can for example not deal with strings, only with numbers. For this matter, before there has to be some knowledge gained, that the data can be processed in a proper way.

- Gaining knowledge about neural networks

Neural networks is one of the most famous topics in the technical world these times. So as master students in computer science it is urgent to gain some knowledge about neural networks and machine learning. For the project it is necessary too to understand the overall technology of the whole topic.

- **Outcome of the project**

- Finding the right features

As a first step it should be done a feature selection. For this the database has to be analyzed and there should be choose some first features by gut feeling. This features has to be checked, if they are independent of each other. During the project it should be always reconsidered, if it make sense to add some more features and checked whether it improves the accuracy.

- Normalizing the features in a proper way

The features have to be normalized before using them for a prediction model. For this procedure it is necessary to find algorithms or write some. The normalization of the data is a major step in the project development.

- Finding a good model for the prediction

There are different libraries available for machine learning. The recommended library is Tensorflow in combination with Keras. Additional there has to be research for other libraries and approaches. The different models are compared and in the end the best model will be choose. It is important to look to the prediction of the actual winner than the prediction of the goals in a different way. So this step has to be done for both ideas.

- Getting a decent accuracy with the prediction

The accuracy for the winner should be higher than 50% in the end. With more than 50% you would theoretically always earning money, if you would bet everything, which the model is predicting. The main goal is to improve the accuracy during the whole project. For the predicting of the goals the accuracy can be lower than 50% because it is harder to predict the exact number of goals which were shot by each team.

- Creating a Homepage

We want to provide a easy to use Homepage which makes it easy for people to look through our former predictions and predict upcoming soccer games. For this matter it is necessary to create a frontend for showing the data in a easy to understand and nice way, an API and a backend.

The first project team had the members Khaled Jallouli, Martin Schmidt, Sergej Dechant and Lisa Boos. The second Team Khaled Jallouli and Lisa Boos from the first team and the two new team members Hemlata Prajapati and Till Hoffmann. The knowledge gaining part was for both of the teams from major interest, except of the ‘Gaining knowledge in data pre-processing’, this part was only done by the first team. In the outcome of the project, ‘Finding the right features’, ‘Normalizing the features’ was only done by the old team. The ‘Finding a good model for prediction’ was done by the first team for the winner prediction and for the goal prediction by the new team. The quality criteria was a team work of only the new team as same as the ‘Creating a Homepage’ part. The team will change for the next step again, Hemlata Prajapati and Till Hoffmann will continue the project with maybe new students and Khaled Jallouli and Lisa Boos will leave the team.

### 1.3. Procedure

For the procedure process we decided to proceed in order to the CRISP-DM model. In image 1.1 you can see the different stages in order to CRISP-DM.

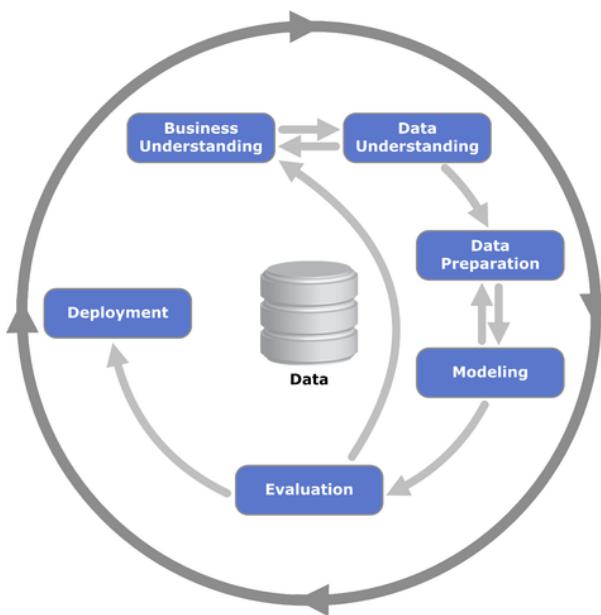


Figure 1.1.: The six stages of CRISP-DM [7]

- Business Understanding

As first step we have to set the goals of the project which are already described in the chapter before. It is necessary to clear what exactly is the required outcome. For this matter in this project we use SCRUM to organize our project. We meet once per week to discuss the achievements from each single person. All 3-4 weeks there is a sprint meeting where we discuss the group achievements and if the project still leads to the right way. In the initial sprint meeting the common parameters of the project are set.

- Data Understanding

Before we can start to create a model or select features we have to understand the Dataset. For this matter it is necessary to understand the structure of the set including the different columns and what they indicate. Additional it is important how many data the set includes and which types of value. It is important to learn also something about soccer and what are important facts about a game.

- Data Preparation

Before it is possible to use the Dataset in a neuronal network it is necessary to prepare the data. For this matter as first step there has to be a feature selection. So one has to decide which columns are important for the later prediction. After we dropped the unnecessary columns it is important to prepare the resulting dataset. The records which have empty attributes have to be deleted or filled with specific values. Additionally the data has to be sorted and aggregated. At the end it has to be split into test and training data.

- Modeling

The resulting dataset from the step before should be ready to be used for different models. In this step we will prepare different models to find a decent one. Additionally it is necessary to define a measurable goal how good the model can get.

- Evaluation

In the end of this semester we will make a evaluation of the accomplished goals and if it is possible to increase the accuracy in the following semester.

- Deployment

As the idea of the project to have a product to show off at the end, this has to be implemented along the way. The finished product should have the ability to predict matches, using the beforehand implemented model(s). It should also list historic games with the prediction of the model, so the user is able to compare prediction versus reality.

## *1. Introduction*

---

The steps ‘Data Preparation’ and ‘Modeling’ will repeat as long as the resulting accuracy is not getting any better or we are satisfied with the resulting one. For all the steps it is necessary to do a lot of research to find the right techniques to do the tasks in a proper way.

To be able deploy a product after the ’Data Preparation’ and ’Modeling’ loop is finished, an architecture and fitting software has to be developed, which makes use of all the previous steps.

## 2. Data Understanding

Data Profiling is the process of examining available data repeatedly throughout a project. At first light profiling assessment was undertaken. This was done using SQL and open source SQLite editors, entity relationship diagram tools and Pandas Profiling. As can be seen in the entity relationship diagram depicted in figure 2.1 the data available for this project is stored in a database consisting of seven tables.

The tables "Country", "League", "Team" and "Player" hold IDs, names and foreign keys to the tables "Team\_Attributes", "Player\_Attributes" and "Match". A player's birthdate, height and weight can be found in the table "Player".

The table "Team\_Attributes" provides data on the individual teams such as playing style, defense class and so on in 12 numerical and 13 categorical columns or variables. However, 66.5 % of the values for "buildUpPlayDribbling" are missing and therefore this variable is useless.

The table "Player\_Attributes" holds 31 numerical and 4 categorical variables and provides player statistics. Only a small percentage of values is missing.

The table "Match" holds information on football matches in Europe from 2008 until 2016 which is made up of 64 numerical and 19 categorical variables including goals scored, ball possession and odds quoted by several bookkeepers. Unfortunately a high percentage of the data on which players played in a match is missing. Moreover, ball possession and shot statistics are in XML format and not available for each match either. The odd variables (from different bookkeepers) are highly correlated.

Together with the product owners, it was decided that this project, or at least its first phase, and the features which have to be developed will be focusing on soccer matches and their outcomes and not on individual team- or player-attributes and statistics. Therefore, the tables "Team\_Attributes" and "Player\_Attributes" are not of interest in this phase. The tables "Country", "League", "Team" and "Player" are not relevant either, since they do not hold any useful information for the first phase.

## 2. Data Understanding

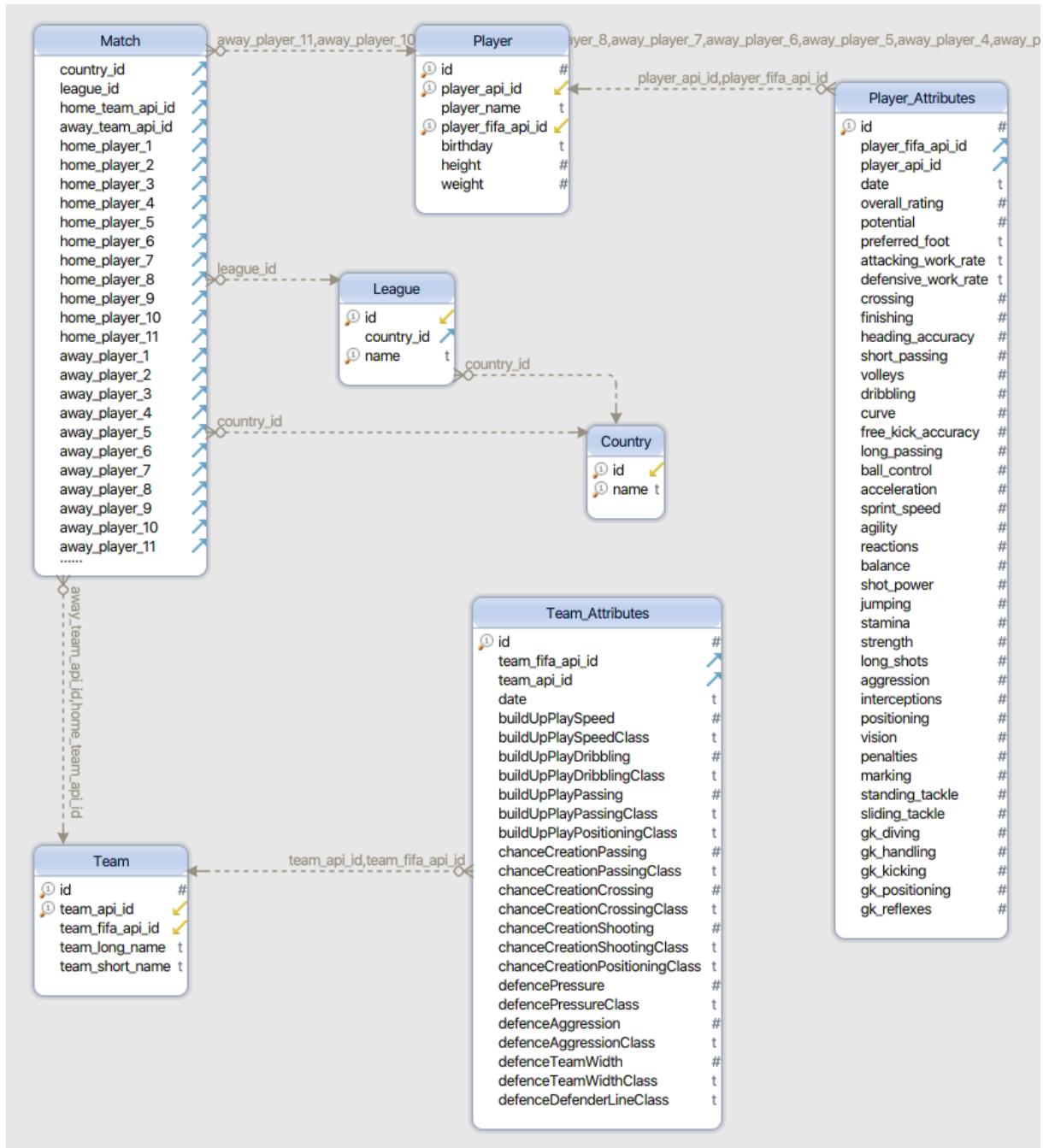


Figure 2.1.: Entity Relationship Diagram European Soccer Database

### 3. Feature Selection and Extraction

As stated in chapter 2 the preproject team agreed with the product owners, to focus on the features of soccer matches and their outcomes in this first phase of the project. Features like the number of scored goals per match/ team or the amount of won, lost or draw matches per team could be calculated directly out of the data at hand. Those features build the base dataset for the analyses we want to pursue. In detail, these are: Odds for the home-team winning, odds for a draw, odds for the away-team winning, amount of wins, draws and losses for the home-team, amount of wins, draws and losses for the away-team, amount of scored and received goals of the home-team and amount of scored and received goals for the away-team. In total there are 13 features in the base dataset. As said before, features about the ball possession or shot-statistics are stored as XML-documents within the database. After consultation with the product owners the team agreed to extract these values out of the XML-documents to get some additional features. The structure of these XML-documents looks like this:

```
1 <shoton>
2   <value>
3     <stats>
4       <shoton>1</shoton>
5     </stats>
6     <event_incident_typefk>137</event_incident_typefk>
7     <elapsed>8</elapsed>
8     <subtype>distance</subtype>
9     <player1>27462</player1>
10    <sortorder>3</sortorder>
11    <team>9912</team>
12    <n>219</n>
13    <type>shoton</type>
14    <id>375900</id>
15  </value>
16 </shoton>
17 <shotoff>
18   <value>
19     <stats>
20       <shotoff>1</shotoff>
21     </stats>
22     <event_incident_typefk>9</event_incident_typefk>
23     <elapsed>9</elapsed>
24     <subtype>distance</subtype>
```

```

25 <player1>30706</player1>
26 <sortorder>1</sortorder>
27 <team>8697</team>
28 <n>220</n>
29 <type>shotoff</type>
30 <id>375902</id>
31 </value>
32 </shotoff>
```

Listing 3.1: XML structure of shot-statistics

```

1 <possession>
2   <value>
3     <comment>51</comment>
4     <event_incident_typefk>352</event_incident_typefk>
5     <elapsed>23</elapsed>
6     <subtype>possession</subtype>
7     <sortorder>1</sortorder>
8     <awaypos>49</awaypos>
9     <homepos>51</homepos>
10    <n>79</n>
11    <type>special</type>
12    <id>462628</id>
13  </value>
14 </possession>
```

Listing 3.2: XML structure of ball-possession-statistics

The shoton and shotoff values can be aggregated per match and team to get features describing the total amount of shots and the amount of shots on the goal per team. Combined with the amount of scored goals per team additional features such as shot accuracy (ratio of shots on the goal to the total amount of shots) and shot efficiency (ratio of scored goals to the total amount of shots) can be calculated per team and match. The additional features after the extraction are: home/away-shots, home/away-shots-on-target, home/away-shot-accuracy, home/away-shot-efficiency. The dataset counts 21 features in total now.

The values for ball-possession are given in percent for each half-time and if necessary for the overtime aswell. This way, the total amount of ball-possession in minutes per team can be summed up quite easily for the whole match.

After the extraction the following features are added to the base dataset: home-possession and away-possession. The features add up to 23 in total now.

# 4. Data Preprocessing

Based on [1] five sliding windows for match data were implemented. The idea of these sliding windows is to deliver match data in historical order and to aggregate data on a team's performance during the last 10 games. First the relevant data was extracted in chronological order with the following SQL statement:

```
1 SELECT id, country_id, league_id, season, date, home_team_api_id,
      away_team_api_id, home_team_goal, away_team_goal, B365H, B365D,
      B365A, shoton, shotoff, possession
2 FROM Match
3 ORDER BY date asc
```

Listing 4.1: SQL code for Sliding Window

Since the bookkeeper odds are highly correlated, as has been pointed out in [chapter 2](#), only Bet365 data (B365H = odds home team wins, B365A = odds away team wins, B365D = odds draw) will be used.

## 4.1. Sliding Windows or Data Options

This part defines five different sliding windows which are used to predict the final results and the scored goals.

### 4.1.1. Option 1

For the first sliding window or option, data on shoton, shotoff and possession are not necessary and therefore are dropped. Moreover, rows without odds are dropped as well. The resulting data is saved as a CSV file:

```
1 conn = sqlite3.connect("../data/eusoccerdatabase.sqlite")
2 query = "SELECT id, country_id, league_id, season, date, match_api_id,
      home_team_api_id, away_team_api_id, home_team_goal, away_team_goal, B365H, B365D, B365A, shoton, shotoff, possession FROM Match ORDER BY date asc"
```

#### 4. Data Preprocessing

---

```

3 df = pd.read_sql_query(query, conn)
4 df = df[np.isfinite(df['B365H'])] # drop rows without odds
5 df_noxml=df.drop(['shoton', 'shotoff', 'possession'], axis=1)
6 df_noxml.to_csv("../data/matches_all.csv")

```

Listing 4.2: Python code for matches\_all.csv

This CSV file consists now of 22592 rows each representing one football match including information on how many goals each team scored, what the odds were and the date of the match:

<b>date</b>	<b>match_api_id</b>	<b>home_team_api_id</b>	<b>away_team_api_id</b>	<b>home_team_goal</b>	<b>away_team_goal</b>	<b>B365H</b>	<b>B365D</b>	<b>B365A</b>
2008-08-09 00:00:00	483129	8583	9830	2	1	2.10	3.10	3.75
2008-08-09 00:00:00	483130	9827	7819	2	1	1.57	3.60	6.50
2008-08-09 00:00:00	483131	9746	9831	1	0	2.30	3.00	3.40
2008-08-09 00:00:00	483132	8682	8689	0	1	2.10	3.10	3.80
2008-08-09 00:00:00	483134	9829	9847	1	0	2.40	3.10	3.10

Figure 4.1.: Data for Option 1

This CSV file was then used for the first sliding window. The code for it [2] expects a CSV file and processes it row by row. For every match the outcome of the game based on scored goals is calculated. Moreover, statistics on how each of the two opponents performed in the previous ten games are generated:

<b>odds-home</b>	<b>odds-draw</b>	<b>odds-away</b>	<b>home-wins</b>	<b>home-draws</b>	<b>home-losses</b>	<b>home-goals</b>	<b>home-opposition-goals</b>	<b>away-wins</b>	<b>away-draws</b>	<b>away-losses</b>	<b>away-goals</b>	<b>away-opposition-goals</b>
1.36	4.50	9.00	5	3	2	14	10	2	2	6	8	14
2.25	3.00	3.50	3	4	3	16	12	4	3	3	14	12
2.80	3.00	2.70	5	2	3	10	8	4	4	2	15	13
1.91	3.25	4.33	5	3	2	18	11	3	3	4	10	11
2.10	3.00	4.00	2	4	4	6	9	1	6	3	7	10

Figure 4.2.: Sliding Window Option 1

### 4.1.2. Option 2

This option, which uses a slightly enhanced version [3] of the first sliding window code, adds statistics on goal shots:

home-shots	home-shots_on_target	home-opposition_shots	home-opposition_shots_on_target
137	67	117	53
134	64	151	77
120	58	124	56
177	82	74	37
161	72	74	31

Figure 4.3.: Sliding Window Option 2

Unfortunately due to the poor choice of data origin this option reduces the dataset to 7033 rows. But on the bright side, this option extends the number of features to 21.

### 4.1.3. Option 3

Based on option 2, option 3 additionally calculates shot accuracy statistics with code written and executed in a Jupyter Notebook [8]:

home_shot_accuracy	home_shot_efficiency	home_opposition_shot_accuracy	home_opposition_shot_efficiency
0.489051	0.164179	0.452991	0.301887
0.477612	0.125000	0.509934	0.207792
0.483333	0.172414	0.451613	0.267857
0.463277	0.268293	0.500000	0.324324
0.447205	0.208333	0.418919	0.258065

Figure 4.4.: Sliding Window Option 3

Option 3 has 29 features and 7033 rows. "x\_shot\_accuracy" represents the percentage of the shots which actually hit the goal or goal keeper and "x\_shot\_efficiency" the percentage of the shots which actually were counted as goals.

#### 4.1.4. Option 4

Option 4 again uses an enhanced version [4] of the first sliding window code and extends option 2 with ball possession statistics:

home- possession	home- opposition_shots	home- opposition_shots_on_target	home- opposition_possession
481	117	53	429
540	69	36	369
477	74	31	433
439	133	77	471
435	153	73	475

Figure 4.5.: Sliding Window Option 4

Ball possession is provided in minutes. This option has 6996 rows and 25 features.

#### 4.1.5. Option 5

The last option is based on option 3 and additionally includes ball possession statistics. It has the highest number of features (33) and 6996 rows. The Python code for it can be found in the same Jupyter Notebook as for option 3 [8].

## 4.2. Data Profiling Part 2

Using Pandas Profiling all sliding windows or options were profiled again. Here is and excerpt of the profile for option 3:

## Dataset info

Number of variables	30
Number of observations	7033
Missing cells	0 (0.0%)
Duplicate rows	0 (0.0%)
Total size in memory	1.6 MiB
Average record size in memory	240.0 B

## Variables types

Numeric	27
Categorical	1
Boolean	0
Date	0
URL	0
Text (Unique)	0
Rejected	2
Unsupported	0

### Warnings

away-draws	has 431 (6.1%) zeros	Zeros
away-losses	has 270 (3.8%) zeros	Zeros
away-shots_on_target	is highly correlated with away-shots ( $p = 0.9054704199$ )	Rejected
away-wins	has 148 (2.1%) zeros	Zeros
home-draws	has 422 (6.0%) zeros	Zeros
home-losses	has 271 (3.9%) zeros	Zeros
home-shots_on_target	is highly correlated with home-shots ( $p = 0.9028529298$ )	Rejected
home-wins	has 154 (2.2%) zeros	Zeros

Figure 4.6.: Pandas Profiling Sliding Window 3

All features/columns are numerical and no values (cells) are missing. As expected x-shots and x-shots\_on\_target are highly correlated. Interestingly zero values are marked as warnings. However, in this context they are fine and still usable. Without having domain knowledge and purely relying on tools, one might be inclined to drop these values and lose valuable information. Another interesting information found in Pandas profiles is the total size in memory. 1.8 MiB shows that the data for option 3 is little and definitely small enough to fit into RAM or GPU memory. This information is valuable, as will be shown later.

Since the profiles for the other options are similar to the one described above, they will not be explained.

## 4.3. Prediction of the final result

The prediction of the final result is based on 3 distinguishable classes (H = home team wins, A = away team wins, D = draw).

In the following example, we can see in row 0 that the home team won, which is indicated by a capital "H". The odds that the home team wins were 3.50, 3.30 for a draw and odds of

2.10 that the away team wins. In the last 10 matches the home team won 1 time, finished 3 matches with a draw, lost 6 times and scored 11 goals. Opposing teams were able to win 10 games. Subsequent columns contain the equivalent statistics for the opposing/away team. The second sliding window reduces the total dataset to 7033 rows and generates 21 features.

	result	odds-home	odds-draw	odds-away	home-wins	home-draws	home-losses	home-goals	home-opposition-goals	home-shots	home-shots_on_target	home-opposition_shots	home-opposition_shots_on_target	home-away-wins
0	H	3.50	3.30	2.10	1	3	6	11	16	137	67	117	53	8
1	D	2.50	3.30	2.88	3	1	6	8	16	134	64	151	77	3
2	A	1.91	3.40	4.20	4	2	4	10	15	120	58	124	56	2
3	H	3.25	3.25	2.30	5	2	3	22	12	177	82	74	37	6
4	H	1.20	6.00	19.00	7	2	1	15	8	161	72	74	31	3

Figure 4.7.: Results Sliding Window Option 2

The last part of the data preparation is to normalize the dataset by rescaling it to make all the elements lie between 0 and 1 thus bringing all the values of numeric columns in the dataset to a common scale. Almost all ML algorithm requires data in numerical form. For that, we'll go through the encoding methods that are necessary for the machine to interpret the data(more specifically categorical data) and learn from it.

### 4.3.1. Normalization

Since the input data varies in scale it has to be normalized, except the "result" column which has to be encoded. Normalization is the process which normalizes all input data avoiding negative effects in regards to decisions of a neuron. The following lines of code normalize the data in a dataframe:

```

1 from sklearn import preprocessing
2
3 column_names_to_not_normalize = ['result']
4 column_names_to_normalize = [x for x in list(dataframe) if x not in
5     column_names_to_not_normalize]
6 x = dataframe[column_names_to_normalize].values
7 x_scaled = preprocessing.normalize(x)
8 df_temp = pd.DataFrame(x_scaled, columns=column_names_to_normalize,
9     index = dataframe.index)
10 dataframe[column_names_to_normalize] = df_temp

```

Listing 4.3: Python code for normalization

### 4.3.2. Encoding

The classes H,A,D respectively the outcome of football match which shall be predicted, have to be encoded before they can be used. Using sklearn this can be achieved with only a few lines of code:

```

1 from sklearn import preprocessing
2
3 le = preprocessing.LabelEncoder()
4 le.fit([ "H", "A", "D"])
5 dataframeloc[:, ['result']] = le.transform(dataframe['result'])

```

Listing 4.4: Python code for encoding classes

## 4.4. Prediction of the scored goals

This part differs from the previous one by changing the "result" column to two columns: "home\_team\_goal" and "away\_team\_goal". It means, we are not predicting the winner or draw but we predict the final scored goals. The scored goals are between 0 and 10 for each team.

The following fig. shows the number of goals scored by each team in a game, we can see in row 0 that the home team scored 2 goals and the away team scored 1 goal. The features are based on the last 10 matches as described in the previous sections.

home_team_goal	away_team_goal	odds-home	odds-draw	odds-away	home-wins	home-draws	home-losses	home-goals	home-opposition-goals	home-shots	home-shots_on_target	home-opposition_shots
0	2	1	3.50	3.30	2.10	1	3	6	11	16	137	67
1	2	2	2.50	3.30	2.88	3	1	6	8	16	134	64
2	1	2	1.91	3.40	4.20	4	2	4	10	15	120	58
3	2	1	3.25	3.25	2.30	5	2	3	22	12	177	82
4	3	0	1.20	6.00	19.00	7	2	1	15	8	161	72

Figure 4.8.: Scored Goals Sliding Window Option 2

The next parts describe the standardization and encoding processes.

#### 4.4.1. Normalization

The goal of normalization is to change the values of numeric columns in the dataset to a common scale, without distorting differences in the ranges of values. The normalization is restricted to the input features and not output features.

Since the input data varies in scale it has to be normalized, except the "home\_team\_goal" and "away\_team\_goal" columns which have to be encoded. The following lines of code normalize the data in a dataframe:

```
1 from sklearn import preprocessing
2
3 column_names_to_not_normalize = ['home_team_goal', 'away_team_goal']
4 column_names_to_normalize = [x for x in list(dataframe) if x not in
5     column_names_to_not_normalize]
6 x = dataframe[column_names_to_normalize].values
7 x_scaled = preprocessing.normalize(x)
8 df_temp = pd.DataFrame(x_scaled, columns=column_names_to_normalize,
9     index = dataframe.index)
10 dataframe[column_names_to_normalize] = df_temp
```

Listing 4.5: Scored goals Python code for normalization

#### 4.4.2. Encoding

The goals scored initially are between 0 and 10. To facilitate prediction before starting encoding, we have reduced the maximum number of goals scored in order to reduce the classes. The following figure shows that choosing to have 5 as the maximum number is a good idea because we have a very small number of rows for 6,7,8,9 and 10 goals scored compared to the others.

	home_team_goal	away_team_goal
1	2251	1 2381
2	1732	0 2362
0	1589	2 1401
3	885	3 613
4	378	4 193
5	132	5 52
6	43	6 23
7	13	8 5
8	7	7 2
9	2	9 1
10	1	dtype: int64

Figure 4.9.: The initial number of rows per each scored goal(s)

The second step is to encode 0 goals to -1, 1 goal to -0.6, 2 goals to -0.2, 3 goals to 0.2, 4 goals to 0.6 and 5 or more goals to 1. The following figure Figure 4.10 describes the encoding step.

```
In [8]: def encode(i):
    switcher = {
        0: -1,
        1: -0.6,
        2: -0.2,
        3: 0.2,
        4: 0.6,
        5: 1,
    }
    # 1 be assigned as default value of passed argument (if goals > 5)
    return switcher.get(i, 1)

def decode(i):
    switcher = {
        -1: 0,
        -0.6: 1,
        -0.2: 2,
        0.2: 3,
        0.6: 4,
        1: 5,
    }
    return switcher.get(i, "ERROR! Use Encode Before!")
```

Figure 4.10.: Encode and Decode functions

The encode function can be used with the following lines of code:

```
1 dataframe['home_team_goal'] = dataframe.apply(lambda row: encode(row['
    home_team_goal']), axis=1)
2 dataframe['away_team_goal'] = dataframe.apply(lambda row: encode(row['
    away_team_goal']), axis=1)
```

---

Listing 4.6: Scored goals Python code for encoding classes

After executing the encode function, we have the following result.

```
home_team_goal    away_team_goal
1   -0.6      2251      -0.6      2381
2   -0.2      1732      -1.0      2362
0   -1.0      1589      -0.2      1401
3    0.2       885       0.2       613
4    0.6       378       0.6       193
5    1.0       198       1.0        83
dtype: int64      dtype: int64
```

Figure 4.11.: Encoding results

The decoding function is useful after executing the encoding function in order to have the real scored goals but it keeps predicting 5 as the maximum scored goals (not 10 as it was).

## 4.5. Functions and Input Data For Neural Networks

The repetitive data preparation tasks encoding, normalization, splitting data into training and test data as well as getting data and prediction labels have been solved using Python functions. Furthermore, Pandas dataframes cannot be used as input for neural networks written with Tensorflow and Keras. The complete code, which has been partly described above, can be found here [5].

Since the data for all options is very little and fits completely into RAM, as has been described above, no input functions for Tensorflow models, e.g. neural networks, had to be written. An excerpt of how the final input data as numpy array for the neural networks described in chapter 5 looks like can be seen here:

<b>odds-home</b>	<b>odds-draw</b>	<b>odds-away</b>	<b>home-wins</b>	<b>home-draws</b>	<b>home-losses</b>	<b>home-goals</b>	<b>home-opposition-goals</b>
0.049957	0.165301	0.330601	0.183667	0.110200	0.073467	0.514268	0.367334
0.077897	0.103862	0.121172	0.103862	0.138483	0.103862	0.553931	0.415448
0.109311	0.117119	0.105407	0.195198	0.078079	0.117119	0.390396	0.312317
0.068789	0.117049	0.155945	0.180075	0.108045	0.072030	0.648271	0.396166
0.108097	0.154424	0.205899	0.102949	0.205899	0.205899	0.308848	0.463272

Figure 4.12.: Excerpt of Input Data for Neural Networks as Numpy Array

# **5. Modeling**

In this phase, we have successfully completed the data preparation phase and our five sliding windows are ready to use.

Modeling is an iterative process, in which we can apply several modeling techniques to the same problem using the default parameters and then fine-tune them until we satisfy our quality criteria. There is not a single model and a single execution which can satisfactorily answer our questions. For this, we tested several models to find the one that best fits our problem.

This phase comprises tasks such as select modeling techniques, generate test design, build model and assess model.

## **5.1. Predicting the winner**

### **5.1.1. Select Modeling Techniques**

As a first step in modeling, we decided to choose Supervised Machine Learning algorithms, to perform multi-class classification because our objective is to predict the final result of a match between two teams, if there is a win by the home team, a draw or a win by the away team.

Therefore, we have selected Decision Trees and Neural Networks as techniques in order to test its performance and find the most appropriate for our project.

### **5.1.2. Generate Test Design**

This part refers to the generation of a procedure to test the model quality and validity needs, before building our models.

For some modeling techniques, we have divided our dataset into training and test sets, the model is built based on the training set, and its quality is estimated based on the test set, which represents 30% of the dataset.

We also took care to not shuffle the dataset as we need to keep the last 10 matches of the sliding windows dataset in the correct order.

For others, we only used 5% for test set, a small amount because we wanted to keep as much data as possible for training and validation. The remaining 95% are divided into 80% training and 20% validation datasets.

For training the models, we used automated stop as a strategy, after the training loss did not improve more than 0.0001 for 10 consecutive epochs or the model exceeds 1000 training epochs.

To evaluate the models, we used the accuracy results as criteria.

### 5.1.3. Build Models

The aim of this part is to build several models before comparing the results.

Most modeling techniques have a number of parameters that can be adjusted to control the modeling process.

For our Supervised Machine Learning Algorithms, we used scikit-learn API to build a Decision Tree and Multi-Layer Perceptron neural network. We also used the Tensor-Flow/Keras framework to build basic sequential neural networks.

The Decision Tree can be modified by adjusting the depth of the tree. For Neural Networks, we can change the number of hidden layers, the neurons per layer and other parameters.

#### Decision Tree Classifier

A decision tree is a simple classification representation that learns from the data with a set of if-then-else decision rules.

Using the decision tree algorithm, we start at the root of the tree and divide the data on the feature that results in the largest information gain. We can then repeat this procedure until the leaves are pure.

In our project, we set the depth of the tree to four.

The Decision Tree Classifier achieved an accuracy of 52.95% using the first sliding window option as a dataset, as shown in the following Figure.

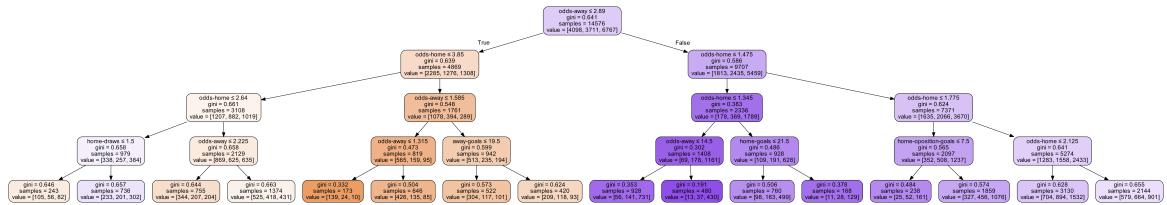


Figure 5.1.: Decision Tree (max\_depth=4)

## Multi-layer Perceptron

Multi-layer Perceptron (MLP) is a supervised learning algorithm consisting of three layers: one input layer, one hidden layer, and one output layer. The units in the hidden layer are fully connected to the input layer, and the output layer is fully connected to the hidden layer.

In our MLP Model, we used the first sliding window with 13 features in the input layer, two hidden layers, 52 neurons in the first one and 32 neurons in the second one. For the output layer, we have 3 neurons.

To be able to solve our problem, we used the sigmoid activation function(logistic) for the hidden layers and the softmax activation function for the output layer. We also used a stochastic gradient descent optimizer as a solver.

As we see in the following plot, the graph of the cost function indicating that the training algorithm converged after the 90th epoch.

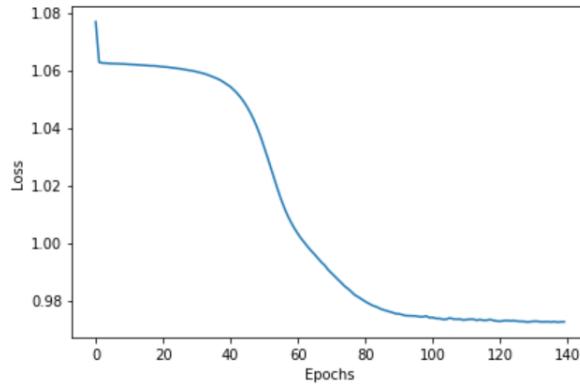


Figure 5.2.: MLP Cost function

The last step is to evaluate the performance of the model by calculating the accuracy of the prediction. We obtained 53.45% for the training dataset and 52.77% for the testing dataset.

## Keras Sequential Neural Network

To build the model, we used Sequential as a model type. Sequential is the easiest way to create a model in Keras. It allows to build a model layer by layer. Each layer has weights that correspond to the layer that follows it.

The chosen layer type is 'dense'. Dense is a standard layer type that works for most cases. In a dense layer, all the nodes in the previous layer connect to the nodes in the current one.

As activation function for the hidden layers 'Rectified Linear Activation' (ReLU) was used and 'Softmax' for the output layer. Softmax sums the output up to 1 so that the output can be interpreted as probabilities. The model will then make its prediction according to the option which has a higher probability.

The first layer needs an input shape. The input shape specifies the number of rows and columns in the input.

The last layer is the output layer. It has three nodes - one for each option: Home Win, Draw or Away Win, which is for our prediction, as shown in the following line of codes.

```
1 model = tf.keras.Sequential([
2     layers.Dense(13, activation='relu', input_shape=(train_X.shape[1],)),
3     layers.Dense(16, activation='relu'),
4     layers.Dense(8, activation='relu'),
5     layers.Dense(3, activation='softmax')
6 ])
```

Listing 5.1: Python code for simple Keras Sequential Model Instantiation

To compile the model, we chose Adam as an optimizer. The Adam Optimizer adjusts the learning rate throughout the training.

The learning rate determines the speed at which the optimal weights for the model are calculated.

For the loss function, we chose *sparse\_categorical\_crossentropy*. It is one of the most common choices for classification. A lower score indicates that the model is performing better.

## 5. Modeling

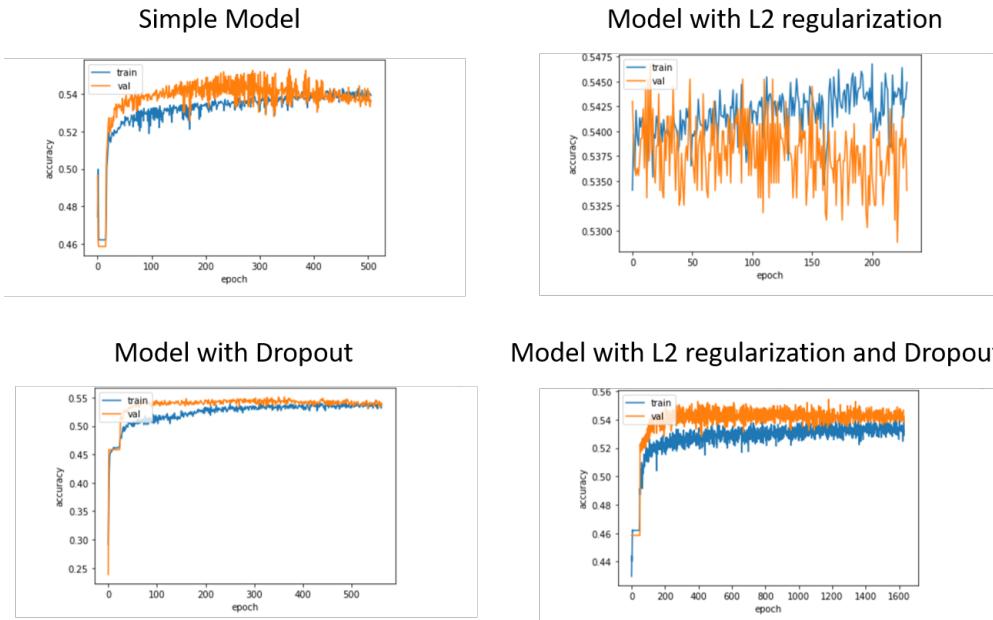
---

Weight regularization is a regularization technique that provides an approach to reduce over-fitting of a deep learning neural network model on training data and to improve the performance of the model on new data.

By default, no regularizer is used in layers. For this, we made some models with the addition of the L2 regularization, which is the sum of the squared weights.

In other models, we have added dropout, which is another regularization technique for neural networks, to avoid over-fitting in neural networks. Additionally we combined both regularization techniques in one model. [Table 5.1](#) shows the test-accuracies of the models in comparison.

The evolution of the loss-function of the generated models over time can be seen in [Figure 5.3](#).



[Figure 5.3](#).: Evolution of the loss function of the models over time

To get a better feeling for the impact the number of hidden layers and the amount of neurons within each hidden layer have on the overall performance of the model, we conducted a series of tests.

We tested each sliding window option with a varying amount of hidden layers (**H1-H4**) and a varying number of neurons per layer (**High, Medium, Low** and **Funnel**).

The architectures of the neural nets can be found in [section A.1](#).

	<b>Normal</b>	<b>L2-Weight Regularization</b>	<b>Dropout</b>	<b>Dropout &amp; Regularization</b>
<b>Model01</b>	0.5220729	0.5191939	0.5191939	0.5268714
<b>Model02</b>	0.5198864	0.52272725	0.5369318	0.50852275
<b>Model03</b>	0.4943182	0.4971591	0.5113636	0.50852275
<b>Model04</b>	0.5057143	0.5114286	0.49142858	0.4857143
<b>Model05</b>	0.5	0.49714285	0.4942857	0.50285715

Table 5.1.: Test-accuracy of various models with different parameters

<b>Model01</b>	<b>H</b>	<b>M</b>	<b>L</b>	<b>F</b>
<b>H1</b>	0.537428	0.53358924	0.5460653	-
<b>H2</b>	0.53262955	0.537428	0.53454894	-
<b>H3</b>	0.5422265	0.537428	0.47408828	0.5393474
<b>H4</b>	0.5431862	0.5393474	0.537428	0.537428

Table 5.2.: Test-accuracies for variation of hidden layers and neurons for sliding window option 1

The test-accuracies for the models based on sliding window option 1 can be seen in Table 5.2.

Table 5.3 shows the test-accuracies for the models based on sliding window option 2.

Table 5.4 shows the test-accuracies for the models based on sliding window option 3.

Table 5.5 shows the test-accuracies for the models based on sliding window option 4.

Table 5.6 shows the test-accuracies for the models based on sliding window option 5.

<b>Model02</b>	<b>H</b>	<b>M</b>	<b>L</b>	<b>F</b>
<b>H1</b>	0.52840906	0.53409094	0.52840906	-
<b>H2</b>	0.5198864	0.5255682	0.53125	-
<b>H3</b>	0.5255682	0.53125	0.5369318	0.5369318
<b>H4</b>	0.5198864	0.5198864	0.53977275	0.5255682

Table 5.3.: Test-accuracies for variation of hidden layers and neurons for sliding window option 2

<b>Model03</b>	<b>H</b>	<b>M</b>	<b>L</b>	<b>F</b>
<b>H1</b>	0.53977275	0.5426136	0.5511364	-
<b>H2</b>	0.5625	0.5568182	0.54545456	-
<b>H3</b>	0.53977275	0.5369318	0.5625	0.53977275
<b>H4</b>	0.5625	0.5426136	0.54545456	0.5426136

Table 5.4.: Test-accuracies for variation of hidden layers and neurons for sliding window option 3

<b>Model04</b>	<b>H</b>	<b>M</b>	<b>L</b>	<b>F</b>
<b>H1</b>	0.5342857	0.5342857	0.5342857	-
<b>H2</b>	0.5257143	0.5371429	0.5342857	-
<b>H3</b>	0.5314286	0.5342857	0.5314286	0.52
<b>H4</b>	0.5285714	0.5285714	0.5285714	0.5228571

Table 5.5.: Test-accuracies for variation of hidden layers and neurons for sliding window option 4

<b>Model05</b>	<b>H</b>	<b>M</b>	<b>L</b>	<b>F</b>
<b>H1</b>	0.5342857	0.52	0.5314286	-
<b>H2</b>	0.5228571	0.5314286	0.5371429	-
<b>H3</b>	0.5342857	0.5228571	0.5314286	0.5371429
<b>H4</b>	0.5228571	0.5257143	0.5314286	0.52

Table 5.6.: Test-accuracies for variation of hidden layers and neurons for sliding window option 5

## 5.2. Predicting the final goals

### 5.2.1. Regression

Regression is used in a problem when the output variable is a real or continuous value, such as "number of goals" in our project. We are using machine learning regression algorithms i.e. multi-output regressor, decision tree regressor, random forest regressor, mlp regressors and Deep Learning multi perceptron neural networks with Keras to perform regression in our project to predict the number of goals scored per each team. Keras is a deep learning library that wraps the efficient numerical libraries Theano and TensorFlow. Main steps involved in the regression are-load a CSV dataset and make it available to Keras, preprocessing, create a neural network models with Keras for a regression problem, apply quality criteria to the models, tune the network topology of models with Keras and selection of the best model among all for making predictions of new data.

### 5.2.2. Build Models

First, we have five different machine learning regression algorithms to build models and then compared their results in term of test accuracy.

We have also used plain multiperceptron model as deep learning artificial neural network algorithm with keras . We can create Keras models and evaluate them with scikit-learn because scikit-learn is good at evaluating models and allow us to use powerful data preparation and model evaluation schemes with very few lines of code.

#### Multi Output Regressors

First one is Multioutput regression which consists of fitting one regressor per target. This is a simple strategy for extending regressors that do not natively support multi-target regression. Multioutput regressor can be created using Gradient Boosting Regressor and support vector machines to make predictions on test data.

Gradient Booster works on boosting or improving the week learner predictions and involve a loss function to be optimized, a week learner to make predictions and an additive model to add weak learners to minimize the loss function.

The final training accuracy Final is 20.59%, training accuracy home team goals is 17.10% and training accuracy away team goals: 24.08%

The final test accuracy for gradient booster is 18.89% , test accuracy for home team

goals is 12.93% and test accuracy for away team goals is 24.86%. Also, the coefficient of determination R square value of the prediction is calculated to be 0.1612.

Second multioutput regressor used is support vector regressor which uses the same principle as for classification and it relies on kernel functions. The final training accuracy Final is 16.22%, training accuracy home team goals is 13% and training accuracy away team goals: 19.43%

The final test accuracy for gradient booster is 16.26% , test accuracy for home team goals is 10.51% and test accuracy for away team goals is 22.02%. Also, the coefficient of determination R square value of the prediction is calculated to be 0.0882

### **5.2.3. Decision Tree Regressor**

The decision trees is used to predict simultaneously the noisy x and y observations of a circle given a single underlying feature. As a result, it learns local linear regressions approximating the circle.

The final training accuracy Final is 23.18%, training accuracy home team goals is 32.53% and training accuracy away team goals: 13.83%

The final test accuracy for gradient booster is 22.37% , test accuracy for home team goals is 30.82% and test accuracy for away team goals is 13.83%. Also, the coefficient of determination R square of the prediction is calculated to be 0.0898

### **5.2.4. Random Forest Regressor**

A random forest is an ensemble model that consists of many decision trees. Predictions are made by averaging the predictions of each decision tree. It is a forest that is a collection of trees. This makes random forests a strong modeling technique that's much more powerful than a single decision tree.

The final training accuracy Final is 45.47%, training accuracy home team goals is 43.53% and training accuracy away team goals: 47.42%

The final test accuracy for gradient booster is 16.12% , test accuracy for home team goals is 14.06% and test accuracy for away team goals is 18.18%. Also, the coefficient of determination R square of the prediction is calculated to be 0.0845

### 5.2.5. MLP Regressor

It is Multi-layer Perceptron regressor. This model optimizes the squared-loss using Adam or stochastic gradient descent. Both optimizers are used to update weights for better predictions. Adam is an optimization algorithm that can be used instead of the classical stochastic gradient descent procedure to update network weights iterative based in training data. /newline

#### Results using Adam Optimizer

The final training accuracy Final is 18.53%, training accuracy home team goals is 14.02% and training accuracy away team goals: 23.04%

The final test accuracy for gradient booster is 18.82% , test accuracy for home team goals is 11.36% and test accuracy for away team goals is 26.28%. Also, the coefficient of determination R square of the prediction is calculated to be 0.1538

#### Results using Stochastic Gradient Descent Optimizer

The final training accuracy Final is 15.01%, training accuracy home team goals is 0.14% and training accuracy away team goals: 29.88%

The final test accuracy for gradient booster is 15.20% , test accuracy for home team goals is 0.14% and test accuracy for away team goals is 30.26%. Also, the coefficient of determination R square of the prediction is calculated to be negative 0.0015

The r square value of the model must be in between 0 and 1 and as close to 1 is the better. As observed, none of the regressor model using machine learning algorithms has provided us the satisfactory accuracy to use it in backend of our website. Hence, we have decided to use deep learning with tensorflow and keras to perform regression.

### 5.2.6. Regression using Keras Deep Learning

This part describes the usage of deep learning techniques and keras deep learning library that wraps the efficient numerical libraries Theano and TensorFlow to do regression or we can say to predict the number of goals for each teams. It includes the following common main steps- Import the dataset, preprocessing of dataset, baseline model building, making predictions, evaluate models on the basis of quality criteria and parameter tuning.

## Develop a Baseline Neural Network Model

In this section we will create a baseline neural network model for the regression problem. We can create Keras models and evaluate them with scikit-learn which allow us to excels at evaluating models and will allow us to use powerful data preparation and model evaluation schemes with very few lines of code.

Below we define the function to create the baseline model to be evaluated. It is a simple model that has a densely connected hidden layer with the higher number of neurons as input attributes (30). The network uses good practices such as the rectifier activation function for the hidden layer. No activation function is used for the output layers because it is a regression problem and we are interested in predicting numerical values directly without transform.

Also, the mean squared error and mean absolute errors are our loss functions which is an estimate of how accurate the neural network is in predicting the test data. We can also used 20 percent of training data for validation and 80% of the training data is used to test the model, while the remaining 20% is used for testing purposes.

The baseline artificial neural network model for regression in shown in the following figure Figure 5.4

```
def build_model():
    model= tf.keras.models.Sequential()
    model.add(tf.keras.layers.Dense(units=30, activation='relu', input_shape=(train_X02.shape[1],)))
    model.add(tf.keras.layers.Dense(units=20, activation='relu'))
    model.add(tf.keras.layers.Dense(units=10, activation='relu'))
    model.add(tf.keras.layers.Dense(units=2))

    optimizer = tf.keras.optimizers.RMSprop(0.001)

    model.compile(loss='mse',
                  optimizer=optimizer,
                  metrics=['mae', 'mse', 'accuracy'])
    return model
```

Figure 5.4.: Baseline model for regression using deep learning

## Train Model

The model has trained for 1000 epochs by keeping validation split of 0.2 and the training and validation accuracy is recorded in the history object. we have observed that the more epochs are run, the lower our MSE and MAE become, indicating improvement in accuracy across each iteration of our model.

Keras is calculating both the training loss and validation loss, i.e. the deviation between the

predicted  $y$  and actual  $y$  as measured by the mean squared error. Let's see our respective losses plot using graph which compare the validation loss and training loss in following figure [Figure 5.5](#)

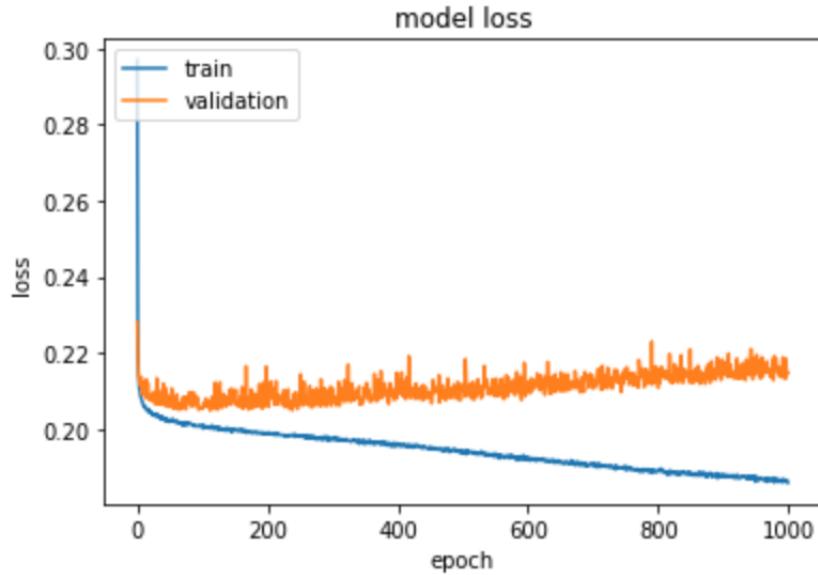


Figure 5.5.: Display of Validation and Training Loss

### **Make Predictions and Evaluate Models using Quality Criteria**

We have predicted the number of goals using data in the testing set, however the predictions we have received are in the encoded form which we have decoded using our decode functions. The predictions made on test data are shown in the following figure [Figure 5.6](#)

-	QM-HG-Train	QM-AG-Train	FinalQM-Train	QM-HG-Test	QM-AG-Test	FinalQM-Test
Model	0.83	0.84	0.83	0.81	0.83	0.82

Table 5.7.: Quality Model for Away Team and Home Team

```
In [35]: y_train_pred
Out[35]: array([[-0.49964523, -0.48719516],
               [-0.4722548 , -0.6432786 ],
               [-0.48331586, -0.5827433 ],
               ...,
               [-0.43938887, -0.5005152 ],
               [-0.53124684, -0.49210525],
               [-0.14386953, -0.6515794 ]], dtype=float32)

In [36]: train_y02
Out[36]: array([[-0.2, -0.6],
               [-0.2, -0.2],
               [-0.6, -0.2],
               ...,
               [-0.2,  0.6],
               [-1. ,  0.2],
               [-0.2, -1. ]])
```

Figure 5.6.: Predictions on Test Dataset

Finally, the evaluation of models has been performed using Quality criteria explained in the Chapter 6 quality criteria section of the report. Then we calculate the final quality model for both teams away team and home team respectively. The results of the quality models are shown in the following table where QM is abbreviation for Quality Model, HG is abbreviation for HomeTeam Goals and AG is abbreviation for AwayTeam Goals

## 5.2.7. Multi Class Classification

Another idea for predicting the final goals is to use multi class classification. For this matter a common classification model can be used.

```
1 inputs = tf.keras.layers.Input(shape=(21,))
2 layer1 = layers.Dense(10, activation="relu", name="layer1")
3 layer2 = layers.Dense(20, activation="relu", name="layer2")
4 layer3 = layers.Dense(6, activation="softmax", name="layer3")
5 model = layer3(layer2(layer1(inputs)))
```

```

6
7 home-team_1 = model
8 away-team_2 = model
9 model = tf.keras.models.Model(
10     inputs=inputs, outputs=[home-team_1, away-team_2])
11 model.compile(optimizer="Adam", loss='sparse_categorical_crossentropy',
12 metrics=[ "acc"])
12 history = model.fit(x, (y_home-team, y_away-team), epochs=100)

```

Listing 5.2: Python code for multi class classification

As the reader can see in listing 5.2 the used model has one input layer with the shape of 21 neurons, because there are 21 features. There are 2 hidden layers with 10 and 20 neurons and one output layer with 6 neurons. The 6 neurons are because we are only predicting goals from 0 to 5 and if the goals are greater than 5 (which does not happen very often, as it is described in the Data Preprocessing part) we count them as 5 goals. In line 5 it is shown how the model is created from the defined layers before. This model is used for both predictions (line 7 and 8), the home team score prediction and the away team score prediction. In line 11 is the model compilation with the two output vectors and in line 12 the training with these vectors. So instead of applying the model on only one output vector it is just used for two of them.

```

1 model.evaluate(test02,( testLabels_hometeam, testLabels_awayTeam))

```

Listing 5.3: Python code for multi class classification model evaluation

As the reader can see in listing 5.3 the evaluation has to be adjusted in the same way. It is necessary to provide the model with two evaluation vectors as well. The result is a array with 5 values, the whole loss of the model, the loss per team and the accuracy per team. With this model we achieve an accuracy of 30.53% for the home team and 37.64% for the away team.

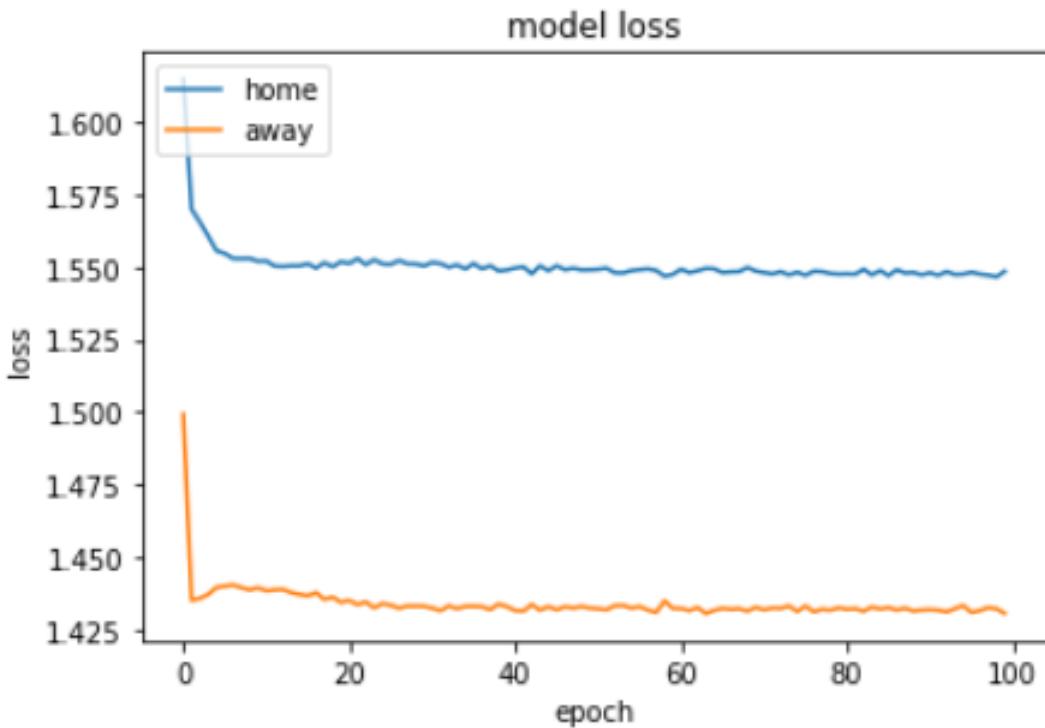


Figure 5.7.: Multi Class Classification loss graph

In figure 5.7 the reader is able to see, that the loss is higher for the home team than the away team. In the beginning the value is decreasing very much and after approximative epoch 10 the value is quite stable. The same shows the accuracy graph in figure 5.8. The accuracy for the away team is better than the one for the home team. And it is only increasing in the beginning. Afterwards it is a little bit noisy but the accuracy is not increasing significant in any epoch.

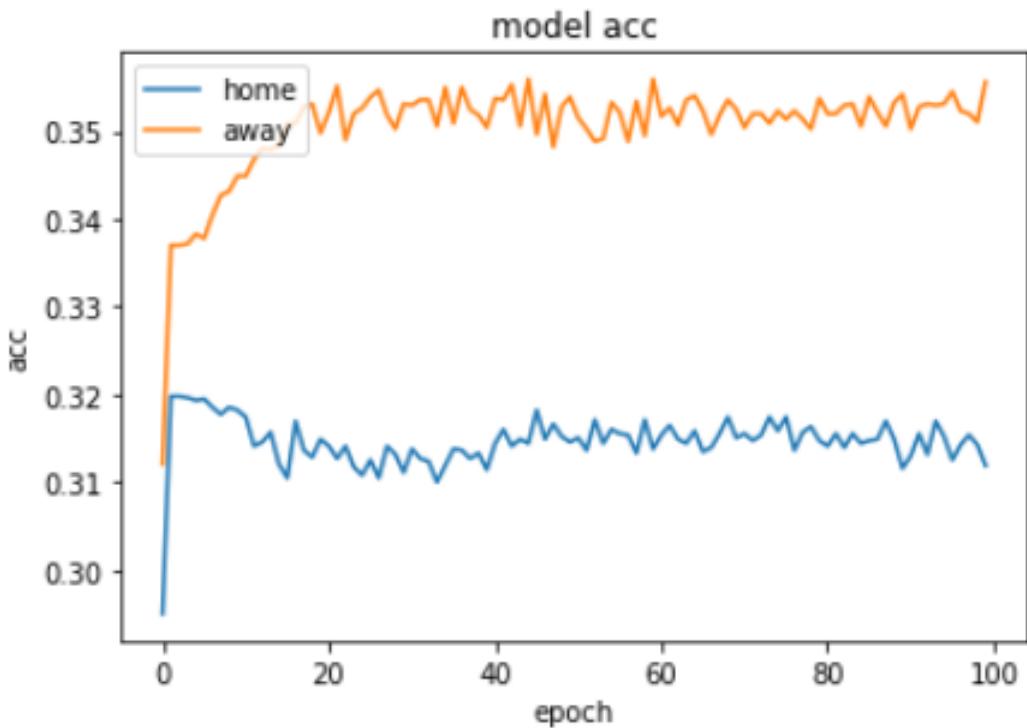


Figure 5.8.: Multi Class Classification accuracy graph

After using the own built quality criteria this approach reaches 80% accuracy. In this point there was no comparison with different models done. This will happen in the next semester. Maybe more hidden layers or more neurons can improve the outcome. Another idea for improving the result is to try the different sliding windows approaches. For this test sliding window 02 was used.

# 6. Evaluation

In this phase of the CRISP-DM Process we are going to evaluate the models and the features, we previously generated.

## 6.1. Prediction of the final result

In this part, the evaluation criteria for the models and in the same way for the features which are used, is based on the test-accuracy. The highest accuracy we reached jet, have been 56.25% with the Sliding Window Option 1, which means with the highest amount of features. The more features we have the less training samples are there for our training, so it could be better to train with more samples and less features. But because of the reason that every models are in the same range, which means there is no model which has a much smaller accuracy, it is not possible to tell the one with the highest accuracy the best model or has the best feature selection in every case. If you are using different data or a higher or smaller amount of features it could be possible to reach a better accuracy with another model than the one which reached the highest accuracy in the actual case. The first ranked model is using Keras Sequential Neural Network, as shown in [Table 6.1](#):

Rank	Model	Sliding Window Option	Parameters	Test-Accuracy
1	Keras Sequential Neural Network	3	Hidden layers: 2 (21, 21 neurons)	0.5625
2	Multi-layer Perceptron	1	Hidden layers: 2 (52, 32 neurons)	0.5345
3	Decision Tree	1	Depth = 4	0.5295

Table 6.1.: Comparison of classifiers

As you can see in the table Table 6.1, too, you can not even tell, that the more features you use, the better the outcome will be. For example with the Multi-layer Perceptron you get a worse accuracy with more features. But for the actual situation we recommend to use the first ranked model for predicting the outcome for football games, but we will maybe use a different model with other features in the future.

## 6.2. Predicting the final goals

In this part, we tested different evaluation criteria in order to compare the models. First, we created our own quality criteria to ensure that we can compare regression results and classification results. Second, we used test-accuracy.

### 6.2.1. Quality Criteria

Our quality criteria is divided into two parts: Rounding the predicted values and calculating the score.

#### Round the predicted values

The main idea of creating our own rounding function is to assign the predicted values of the model to the nearest real class. For example, if we predict 0.54, it will be rounded to 0.6 because it is between 0.33 and 0.67. 0.60 is the encoded value to score 4 goals, as shown in the decoded value column. The following table shows the intervals of all possible values that we can predict and the associated decoded value.

Predicted value Range	Round predicted value	Decoded value
[-1, -0.67[	-1	0
] -0.67, -0.33]	-0.60	1
] -0.33, 0]	-0.20	2
] 0, 0.33]	0.20	3
] 0.33, 0.67]	0.60	4
] 0.67, 1]	1	5

Table 6.2.: Round and Decode predicted values

The decoding function is used both in the predicted values and in the original values before comparing them and calculating the scores. We have defined this function to facilitate the comparison process.

### **Calculate the scores**

To calculate the score, we first calculate the absolute value between the original value and that predicted between the goals scored by the home team and the away team. Then, we encode this calculated value which is between 0 and 5 to a value between 0 and 1. 0 goal corresponds to a score equal to 1 and 5 goals corresponds to a score equal to 0. The following table shows all the possible values.

Abs (yoriginal - ypred)	Degree difference
0	1
1	0.8
2	0.6
3	0.4
4	0.2
5	0

Table 6.3.: Calculate the degree difference

After calculating the difference in degree for goals scored for home and away teams, the score is the average of the two columns. The final result will be the average of all the rows in the table which is between 0 and 1.

#### **6.2.2. Comparing the regressors**

This section highlights the different regressor models created for two different datasets with different features and then these models are compared to find the best model to make predictions on the new data or test data, we have created three different models using tried trial and error to set parameters and hiddenlayers units.

In classification tasks, it is easy to calculate sensitivity or specificity of classifier because output is always binary correct classification, incorrect classification. So we can count good/bad answers and based on the confusion matrix calculate some measurements. But in regression tasks, the output is a number. So we can't just say is it correct or incorrect but instead we should measure "how far from true solution we are" by either calculating

## 6. Evaluation

---

Dataset Sliding02 Models	QM-HG-Train	QM-AG-Train	FinalQM- Train	QM-HG-Test	QM-AG-Test	FinalQM- Test
Model1	0.83	0.84	0.83	0.81	0.83	0.82
Model2	0.82	0.84	0.83	0.81	0.84	0.83
Model3	0.82	0.84	0.83	0.81	0.84	0.83

Table 6.4.: Quality Model for models with different hidden units for dataset sliding02

Dataset Sliding03 Models	QM-HG-Train	QM-AG-Train	FinalQM- Train	QM-HG-Test	QM-AG-Test	FinalQM- Test
Model1	0.83	0.84	0.84	0.81	0.83	0.82
Model2	0.80	0.83	0.82	0.80	0.84	0.82
Model3	0.83	0.84	0.83	0.81	0.84	0.82

Table 6.5.: Quality Model for models with different hidden units for dataset sliding03

coefficient of determination R square or by focusing on minimizing the mean squared error also known as loss. That is why in our regression models, we have calculated loss for validation as well as for training. For all the three models, input units have been taken equal to features that is 21 units and as we are predicting two outputs- Home Team Goals and Away Team Goals, the output layers are having 2 units or neurons. For the dataset sliding02, we have used 3 hidden layers from which first hidden layer has 30 units which makes the neural network densely connected with 20 units in the second hidden layer and 10 in the third hidden layer. While training the model, we have observed the validation accuracy of 68 percent and Training accuracy of 70 percent. For second model, we have used three hidden layers with 21 units in first hidden layer making the artificial neural network as fullyconnected, 10 units in second hidden layer and 5 in third hidden layer which give us validation accuracy of 66 percent and training accuracy of 72 percent. And the third model have 4 hidden layers with 21 units in first hidden layer layer, 14 units in second hidden layer, 12 units in third hidden layer and 10 units in last hidden layer which give us validation accuracy of 56 percent and training accuracy of 69 percent.

Similarly, similar experiments have been made with dataset sliding03 as well. For all the models we have also evaluated the quality for each team and for both training dataset and test dataset using the quality criteria defined in the above section. The evaluated qualities for all three models for dataset sliding02 and sliding03 have displayed in the respective tables [Table 6.4](#) and [Table 6.5](#) where QM is abbreviation of Quality Model, HG is abbreviation of Home goals and AG is abbreviation of Away goals .

We have used first model from dataset sliding02 as our best model as the validation and training accuracy of this model is better than the other models with a percentage of 68 percent and 72 percent respectively, ideally the validation accuracy must be slightly better than the training accuracy to prevent overfitting of the model. Because if the training loss is higher than validation loss then there are chances of overfitting and if the training loss is very less than validation loss then there are chances of underfitting. But still the

validation accuracy is quite lower than the training accuracy.

Also, as our approach is to minimize the mse that is mean squared error in order to get predictions on test data as close to actual data, this model provides us with minimum training mse of 0.19 and validation mse of 0.21. Hence, we are considering it as our best model.

Using quality criteria defined, we have observed a final quality of performance for model 1 for Training data is 83 percent and for test data it is 82 percent. One important point is that the quality evaluation for all the models are showing almost similar values to each other which is not an ideal solution, so there are lot of scope of improvement in the project.

### **6.2.3. Comparing Regression with Multi Class Classification**

As it is described before the Multi Class Classification has still many possibilities to improve, so in this point we only compare the actual model with the best regression model. For this we will use our own Quality Criteria, but as it is described in the part ‘Quality Criteria’, this has some disadvantages and has to be improved in the future. Because of this, the test accuracy is also used for the comparison. For the Multi Class Classification approach we are reaching a test accuracy of 30.54% for the away team and 37.64% for the home team, which makes an average of 34.09%. With the regression model we are reaching a test accuracy of [HEMI U HAVE TO FILL IN]. Using our quality criteria we are reaching and average of 80.38% (77.30% for the home team and 83.47% for the away team) with the Multi Class Classification and [HEMI U HAVE TO FILL IN] with the regression approach. Because of the better results with the regression model the team decided to use this model for the homepage. After improving the regression and the Multi Class Classification this could change.

preferably fungible

# 7. Deployment

This chapter is about the deployment step. Details about the backend are described first, then the API as interface between the latter and the frontend gets explained and afterwards the implementation of the frontend will be worked through.

## 7.1. Backend

Here, the overall architecture, structure and functionality of the backend is explained. Afterwards each of its component is described in detail. At the end, the installation and usage of the backend are shown.

### 7.1.1. Overview

The backend has the job to integrate several independent functionalities with each other and by that being able to serve the frontend. As each of the functionalities is independent of the others, the main idea is to keep them modularized. This not only increases maintainability but also allows for independent development and testing of each module. These modules and their interactions with each other are shown in [Figure 7.1](#) and further described in the following.

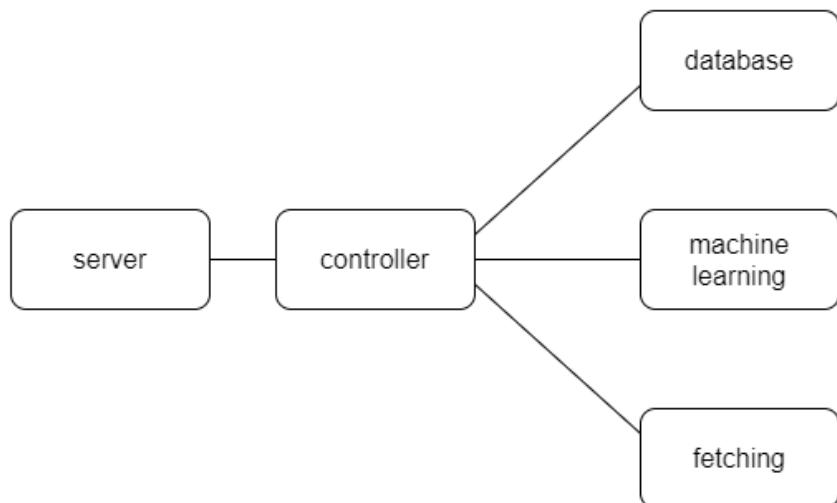


Figure 7.1.: Architecture of the backend

### 7.1.2. Server

The server module consists only of one equally named file and is the entry point of the backend. It's job is to hold a python-specific web server named 'flask'. The only thing except for the flask-part in this module is the access to the controller.

So in this file, every API-route is defined and for each route, an accordingly named function within the controller is called. The return value of it is then returned to the HTTP client, which sent the request.

### 7.1.3. Controller

Like the server module, the controller also consists of only one file. Yet it is the heart of the backend and orchestrates every action. The server requests answers from the controller according to the accessed route and the controller pulls all the triggers and data moving to create the desired answer. The file also contains all required global variables for every model - be it database-filename (more on that later) or model names. Therefore it is the single point of truth for backend variables. It contains one method for each route of the server plus one to predict a single match, as this requires playing together of the database and machine learning module.

To get a list of matches, be it with the results of classification or regression, a method from the database-module is called and before returning the list of matches, each match gets slightly modified. First, the date, which is stored as string in the database (due to the lack of a date-type in sqlite) is formatted to match the API-format (described in the next chapter). Afterwards, the actual result in the API-conform is created out of the two distinct goal-values.

Fetching new games and inserting them into the database requires a loop over all matches returned by the fetching module, preparing data for prediction, predict both the classification and regression result and writing each match with all the values into the database. This ensures, that the database always contains predictions for each match.

Retraining the classification and regression model can be done independently, but the workflow is very similar. First every match is retrieved from the database, then each match gets enriched by the additional data required for prediction, then the machine learning model receives this new list, executes the training and returns the new model. The controller then again calls the machine learning module and saves the new model to disk. Afterwards every match is repredicted and updated in the database.

### 7.1.4. Database

This module has its own folder, because it requires several files. First of all there is a python file, which contains all database-specific methods. Out of simplicity and experience, the chosen database-engine was sqlite, whose file is also stored in this folder. Each of the methods gets the filename of the database passed as argument. They contain an enormous amount of really long comments, so future developers have it easier to understand not only the database schematic, but also all of the data creation. There are functions for creating the database file (if it not already exists), add a new match to it, update either the prediction of classification or regression. While these were methods to add data to or update data in the database there are also three functions for getting data from the database. First, there is the option to get a single match, which is only used for testing purposes. Next comes a high-level function that returns the cumulated stats of the last ten matches of a specific team. It takes into account both playing as hometeam and as awayteam. Last but not least there is a function to receive all matches in the database, sorted from newest to oldest. The file contains one additional helper method which is for retrieving data from the database as type dictionary instead of array for each entry.

### 7.1.5. Fetching

The fetching module also has its own folder, but it would not explicitly required to have so. This is due to a earlier development time, when the required data for fetching was downloaded as file and not directly into a variable. Later on it was decided to leave it like this, as through the folder structure it gets clear, that database, fetching and machine learning are quite different from the server and controller. The module contains only one function, which is for fetching a complete season. Which season it is, is announced by an argument placed by the controller. There are two datasources that have to be fetched, their results have to be merged together and then get returned to the controller. The first datasource is a csv-file which is downloaded from <http://www.football-data.co.uk/mmz4281/<season>/D1.csv> (`<season>` should match the desired season in the format of f.e. '`19-20`') and contains every finished match of the season with a lot of additional information needed for the machine learning part like odds, shots per team, etc. The other datasource is the website [https://www.worldfootball.net/all\\_matches/bundesliga-20<XX>-20<YY>/](https://www.worldfootball.net/all_matches/bundesliga-20<XX>-20<YY>/) where `<XX>` and `<YY>` also equal to f.e. '`19`' and '`20`' respectively. This additional datasource is required to retrieve also upcoming matches which were not contained in the first datasource. Additionally, this second datasource contains a better readable teamname instead of abbreviations. After

both matchlists are merged, the resulting, more complete list is returned back to the controller.

Currently, the upcoming matches are not added to the returned list. This is due to the fact, that the machine learning requires f.e. the betting odds as input and this information is not freely available.

### **7.1.6. Machine learning**

Like the two previous submodules, the machine learning module has its own folder. This is necessary because this module also contains the models for classification and regression. The single file with the multiple methods contain a simplified version of the same code used in the jupyter notebooks for development, which is why they aren't explained in detail. Exceptions are a function to execute one of the two models on a match, a simplified retraining method and methods for preparation of data for predictions. For the latter the high-level function receives data for a specific match, which should be predicted and a list of the last 10 matches for each team. It then uses a lower-level function to calculate the total amounts of f.e. shots, shots on target etc. and pulls together the calculated values to return to the controller. Because the classification-training requires slightly differently prepared data as the regression-training, there are two closely related but still independent methods for each of the data preparations.

### **7.1.7. Installation and usage**

To run the backend python is required. The software was developed under version 3.8.2, but should support newer version as well. To install all of the dependencies, there is a file named `requirements.txt` which contains all the additional software with their respective versions. To install them, simply run `pip install -r requirements.txt`. To start the backend, just run `py server.py` and wait until the message with the local port appears (approx. 2s).

## **7.2. API**

As the back- & frontend were developed independent, an API as interface between them was developed. All of the below calls are HTTP GET methods.

In total there are two ways a call is answered by the backend. The first one is the

return of a JSON-object, that the frontend has to process. The second one is a simple '`ok`' string for actions that the backend has successfully finished. While the latter is self-explaining, the former had to be agreed upon for both front- and backend.

### 7.2.1. /soccerGamesClassification

As shown in [7.1](#), this route returns a JSON-object named 'SoccerMatches' containing all soccer matches, with their respective dates, names of both teams, the predicted result and the actual result (if available). Every of those attributes is of type string. For the teamnames this definitely makes sense, but as JSON lacks a type for date, it was also chosen to be a string in the format of '`DD/MM/YYYY`'. Again, both front- and backend had to make sure to use the same formatting for this string. The predicted Result is actually a single character (namely '`D`' for draw, '`H`' for home team win and '`A`' for away team win), but again, JSON doesn't have a dedicated type for it. For the actual result of the match it would also be possible to create two integers, one goal-counter for each team, but this way the backend would just concatenate both numbers with a colon in between and pass this to the frontend.

```

1 {
2   "SoccerMatches": [
3     {
4       "date": <string>,
5       "homeTeam": <string>,
6       "awayTeam": <string>,
7       "predictedResult": <string>,
8       "actualResult": <string>
9     }, ...
10   ]
11 }
```

Listing 7.1: JSON structure of API call on /soccerGamesClassification

### 7.2.2. /soccerGamesRegression

This route returns the exact same JSON-object as the previous one and as already shown in [7.1](#). The only difference is the field content of the predicted result: Instead of the

previous single character, this time the string has the same format as the actual result with the structure of <predicted home team goals>:<predicted away team goals>.

### 7.2.3. /retrainClassification

As the backend is able to create a newly trained model, this feature has to be accessed from the frontend. It does so, by sending a GET request to this route and waiting for it to return 'ok'. When this happens, the backend has finished creating a new model to use for classification.

### 7.2.4. /retrainRegression

Closely related to the previous /retrainClassification route, the only difference is the model, which is newly created by the backend; Instead of the model for classification, this one overwrites the model used for regression.

### 7.2.5. /fetchNewMatches

With a GET call on this route, the backend is commanded to download a new list of matches, both historic and upcoming ones. After the download, it has to integrate them into the database for further usage.

After reading about all the routes, it could be argued, that there is no route for (re-)execution of the prediction for all matches in the database. This is possible due to the fact, that whenever the backend either fetches new matches or retrains a model, it (re-)executes both models on every match in the database. This may not be as efficient as a dedicated route for each model, but it is just as effective. Additionally, this prevents locks, as everything is ran sequentially. One more reason is the inefficiency introduced by working through all matches in the database two times - once for each model.

## 7.3. Frontend

-> TODO Khaled

## **8. Conclusion**

We reached all main goals for the project. The best model has a decent accuracy and we have learned many things about machine learning. Through the team work with SCRUM we were able to delegate the tasks in a way, that we get the best end solution. Additionally we learned how to apply SCRUM for team work, which will help us in our further working life. All steps were necessary for the final outcome, this means we did not got stuck in a wrong direction. The project is proper documented, that a future team is able to continue with our work on the project. The biggest issues we had especially in the beginning of the project, because we did not really know how we should reach our goals or what goals we really had. For the future it would be awesome if the communication between the product owners and the students would be better. Not only in how we should start and what are the main goals, but also how the grade will be exactly built and how we have to structure the report. It would be a good way to create a one- or two-paper with all the common guidelines that the students know exactly what the product owners are require and how they will evaluate the outcome. We are not in a stage that we could actually really earn very much money with our model, which means, that we are only hardly over a 50% accuracy. We are still have the goal to improve our model to reach higher accuracies. For the next semester the main part will be improving the model through additional features and through changing the model.

# Bibliography

- [1] Andrew Carter. *Deep Neural Network (DNN) Football/Soccer Predictor*. URL: <https://github.com/AndrewCarterUK/football-predictor>.
- [2] Andrew Carter; Sergej Dechant. *Sliding Window 01*. 2019. URL: [https://github.com/thu-soccer/project/blob/master/code/process\\_csv\\_data.py](https://github.com/thu-soccer/project/blob/master/code/process_csv_data.py).
- [3] Andrew Carter; Sergej Dechant. *Sliding Window 02*. 2019. URL: [https://github.com/thu-soccer/project/blob/master/code/process\\_csv\\_data\\_shots.py](https://github.com/thu-soccer/project/blob/master/code/process_csv_data_shots.py).
- [4] Andrew Carter; Sergej Dechant. *Sliding Window 03*. 2019. URL: [https://github.com/thu-soccer/project/blob/master/code/process\\_csv\\_data\\_shots\\_possession.py](https://github.com/thu-soccer/project/blob/master/code/process_csv_data_shots_possession.py).
- [5] Sergej Dechant. *Google Colab Notebook for Data Preparation and Tensorflow Neural Networks*. 2019. URL: [https://github.com/thu-soccer/project/blob/master/colab/colab\\_nn.ipynb](https://github.com/thu-soccer/project/blob/master/colab/colab_nn.ipynb).
- [6] Kaggle Inc. *European Soccer Database*. 2019. URL: <https://www.kaggle.com/hugomathien/soccer>.
- [7] Kenneth Jensen. *IBM SPSS Modeler CRISP-DM Guide*. 2012. URL: <ftp://public.dhe.ibm.com/software/analytics/spss/documentation/modeler/18.0/en/ModelerCRISPDM.pdf>.
- [8] Sergej Dechant; Martin Schmid. *Jupyter Notebook Sliding Windows*. 2019. URL: [https://github.com/thu-soccer/project/blob/master/code/sliding\\_window.ipynb](https://github.com/thu-soccer/project/blob/master/code/sliding_window.ipynb).

# A. Appendix

## A.1. Architectures for various neural nets

Table A.1.: Test Variation of Hidden-Layers and Neurons for Neural Nets

Name	Input	HLayer1	HLayer2	HLayer3	HLayer4	Output
model01_H1_H	13	13	-	-	-	3
model02_H1_H	21	21	-	-	-	3
model03_H1_H	29	29	-	-	-	3
model04_H1_H	25	25	-	-	-	3
model05_H1_H	33	33	-	-	-	3
model01_H1_M	13	9	-	-	-	3
model02_H1_M	21	12	-	-	-	3
model03_H1_M	29	14	-	-	-	3
model04_H1_M	25	13	-	-	-	3
model05_H1_M	33	16	-	-	-	3
model01_H1_L	13	4	-	-	-	3
model02_H1_L	21	5	-	-	-	3
model03_H1_L	29	7	-	-	-	3
model04_H1_L	25	6	-	-	-	3
model05_H1_L	33	8	-	-	-	3
model01_H2_H	13	13	13	-	-	3
model02_H2_H	21	21	21	-	-	3

continued on next page

### A. Appendix

---

Name	Input	HLayer1	HLayer2	HLayer3	HLayer4	Output
model03_H2_H	29	29	29	-	-	3
model04_H2_H	25	25	25	-	-	3
model05_H2_H	33	33	33	-	-	3
model01_H2_M	13	9	9	-	-	3
model02_H2_M	21	12	12	-	-	3
model03_H2_M	29	14	14	-	-	3
model04_H2_M	25	13	13	-	-	3
model05_H2_M	33	16	16	-	-	3
model01_H2_L	13	4	4	-	-	3
model02_H2_L	21	5	5	-	-	3
model03_H2_L	29	7	7	-	-	3
model04_H2_L	25	6	6	-	-	3
model05_H2_L	33	8	8	-	-	3
model01_H3_H	13	13	13	13	-	3
model02_H3_H	21	21	21	21	-	3
model03_H3_H	29	29	29	29	-	3
model04_H3_H	25	25	25	25	-	3
model05_H3_H	33	33	33	33	-	3
model01_H3_M	13	9	9	9	-	3
model02_H3_M	21	12	12	12	-	3
model03_H3_M	29	14	14	14	-	3
model04_H3_M	25	13	13	13	-	3
model05_H3_M	33	16	16	16	-	3

continued on next page

### A. Appendix

---

Name	Input	HLayer1	HLayer2	HLayer3	HLayer4	Output
model01_H3_L	13	4	4	4	-	3
model02_H3_L	21	5	5	5	-	3
model03_H3_L	29	7	7	7	-	3
model04_H3_L	25	6	6	6	-	3
model05_H3_L	33	8	8	8	-	3
model01_H3_F	13	13	10	7	5	3
model02_H3_F	21	18	13	9	5	3
model03_H3_F	29	22	16	11	6	3
model04_H3_F	25	20	15	11	6	3
model05_H3_F	33	25	19	12	6	3
model01_H4_H	13	13	13	13	13	3
model02_H4_H	21	21	21	21	21	3
model03_H4_H	29	29	29	29	29	3
model04_H4_H	25	25	25	25	25	3
model05_H4_H	33	33	33	33	33	3
model01_H4_M	13	9	9	9	9	3
model02_H4_M	21	12	12	12	12	3
model03_H4_M	29	14	14	14	14	3
model04_H4_M	25	13	13	13	13	3
model05_H4_M	33	16	16	16	16	3
model01_H4_L	13	4	4	4	4	3
model02_H4_L	21	5	5	5	5	3
model03_H4_L	29	7	7	7	7	3

continued on next page

---

*A. Appendix*

Name	Input	HLayer1	HLayer2	HLayer3	HLayer4	Output
model04_H4_L	25	6	6	6	6	3
model05_H4_L	33	8	8	8	8	3
model01_H4_F	13	13	10	7	5	3
model02_H4_F	21	18	13	9	5	3
model03_H4_F	29	22	16	11	6	3
model04_H4_F	25	20	15	11	6	3
model05_H4_F	33	25	19	12	6	3

end of table

## **A.2. Daily Scrum Logs first therm**



<b>Titel</b>	<b>Daily Scrum</b>
<b>Date</b>	24.10.2019
<b>Sprint</b>	01
<b>Participants</b>	Sergej, Martin, Lisa, Khaled

## 1. What have you done?

- **Lisa:** Installed TF, Articles about TF, Exercises in TF, Data Profiling, Python ML
- **Khaled:** Keras & TF Anaconda Tutorials, Book about NN
- **Sergej:** Nielssens Book NN, Installed TF & Keras, Tutorial for TF, Look at Soccer DB SQLite, Reading on Data Science, Research on Data Science Maths, Basic research on Linear Regression
- **Martin:** Keras & TF + Installation

## 2. What will you do?

- **Lisa:** Installing Keras, Do Examples, Research about ML
- **Khaled:** Neuroal Network, Data Profiling
- **Sergej:** Get used to TF for SQL based data
- **Martin:** Data Profiling, getting used to TF/ Keras

## 3. Issues?

- **Lisa:** -
- **Khaled:** -
- **Sergej:** -
- **Martin:** -



<b>Titel</b>	<b>Daily Scrum</b>
<b>Date</b>	31.10.2019
<b>Sprint</b>	1
<b>Participants</b>	Khaled, Lisa, Sergej, MArtin

## 1. What have you done?

- **Lisa:** Read Articles about keras + videos, installation, Example for Keras and TF
- **Khaled:** data profiling, started data pre-processing (missing data), started using sci-kit learn library, Research Neuronal Network
- **Sergej:** reading on classification, linear regression, handwriting recognition, data profiling, made some examples with TF
- **Martin:** research for keras & tf, data profiling

## 2. What will you do?

- **Lisa:** Do more examples, do some tests with our database
- **Khaled:** continue data processing
- **Sergej:** data profiling, research
- **Martin:** data profiling & preprocessing

## 3. Issues?

- **Lisa:** -
- **Khaled:** -
- **Sergej:** -
- **Martin:** -



<b>Titel</b>	<b>Daily Scrum</b>
<b>Date</b>	07.11.2019
<b>Sprint</b>	Sprint 02
<b>Participants</b>	Khaled, Sergej, Martin

## 1. What have you done?

- **Lisa:** Analyzed the Dataset
- **Khaled:** Started analysing the data
- **Sergej:** Spent time on data profiling, literature research on dp, slice & dice (pivot tables)
- **Martin:** Research on feature selection, data profiling

## 2. What will you do?

- **Lisa:** feature selection
- **Khaled:** Starting to use both teams as input for analytics
- **Sergej:** write some pages for the report on data profiling
- **Martin:** research on feature selection, data profiling

## 3. Issues?

- **Lisa:-**
- **Khaled:** Maybe using google collab?
- **Sergej:** -
- **Martin:** -



Titel	Daily Scrum
Date	14.11.2019
Sprint	Sprint 02
Participants	Khaled, Sergej, Martin, Lisa

## 1. What have you done?

- **Lisa:** research on feature selection, found tutorial and diagramm
- **Khaled:** continued data pre-processing, mapping feature for season
- **Sergej:** data profiling + documentation
- **Martin:** extracted reduced dataset from Bundesliga seasons 2014/15 and 2015/16, SQL, research

## 2. What will you do?

- **Lisa:** do the tutorial, learn more about fs
- **Khaled:** find good features for hidden layers
- **Sergej:** data profiling + documentation
- **Martin:** research on feature selection, normalization, prediction, NN

## 3. Issues?

- **Lisa:** database is very large
- **Khaled:** database is very large
- **Sergej:** -
- **Martin:** -



Titel	Daily Scrum
Date	21.11.2019
Sprint	Sprint 02
Participants	Khaled, Sergej, Martin, Lisa

## 1. What have you done?

- **Lisa:** did the tutorial, analyzed xml-columns, research feature selection
- **Khaled:** handled high correlation features, handled NaN/null-values, provided model
- **Sergej:** started report LaTeX, research and trials of feature selection
- **Martin:** documentation of data extraction, research feature selection and normalization

## 2. What will you do?

- **Lisa:** evaluate the model, find some new features, apply possession-feature
- **Khaled:** include possession feature
- **Sergej:** maybe try different model
- **Martin:** normalization, feature selection

## 3. Issues?

- **Lisa:** -
- **Khaled:** -
- **Sergej:** waiting for the VM
- **Martin:** -



Titel	Daily Scrum
Date	28.11.2019
Sprint	03
Participants	Khaled, Lisa, Sergej, Martin

## 1. What have you done?

- **Lisa:** implemented time window, extracted average amount of points feature
- **Khaled:** average goal feature, online course for model implementation, looked on possession feature
- **Sergej:** reading on feature extraction and how to adjust a model in the right way
- **Martin:** extract feature avg. earned points, Reading raschka ebook python ml

## 2. What will you do?

- **Lisa:** trying to make a model, find extract the features differently
- **Khaled:** continue research
- **Sergej:** implement additional sliding window approach
- **Martin:** Research, continue the Raschka book

## 3. Issues?

- **Lisa:** -
- **Khaled:** -
- **Sergej:** -
- **Martin:** -



<b>Titel</b>	<b>Daily Scrum</b>
<b>Date</b>	05.12.2019
<b>Sprint</b>	03
<b>Participants</b>	Khaled, Lisa, Sergej, Martin

## 1. What have you done?

- **Lisa:** tried to build a TF Model, read Raschka
- **Khaled:** finished certificate about deep learning with TF by IBM,
- **Sergej:** implemented sliding window, tried to build a model with TF,
- **Martin:** Normalization/Standardization Raschka

## 2. What will you do?

- **Lisa:** Look at sliding window from Sergej, feature extraction
- **Khaled:** try logistic regression and learned from certificate
- **Sergej:** TF Serving, use sliding window, try different classifiers
- **Martin:** extract the XML-feature shot-on-goal

## 3. Issues?

- **Lisa:** -
- **Khaled:** -
- **Sergej:** -
- **Martin:** -



<b>Titel</b>	<b>Daily Scrum</b>
<b>Date</b>	12.12.2019
<b>Sprint</b>	03
<b>Participants</b>	Khaled, Lisa, Sergej, Martin

## 1. What have you done?

- **Lisa:** built model with scikit learn, accuracy of 52%, understanding of sergejs sliding window, setup old LaTeX template
- **Khaled:** worked on logistic regression with TF
- **Sergej:** spent more time on sliding window, played around with TF
- **Martin:** feature extraction from xml columns

## 2. What will you do?

- **Lisa:** start with the report and more feature selecting, try to build a model
- **Khaled:** analyzing the data
- **Sergej:** get understanding of how to use the model
- **Martin:** feature selection

## 3. Issues?

- **Lisa:** -
- **Khaled:** -
- **Sergej:** -
- **Martin:** -



Titel	Daily Scrum
Date	19.12.2019
Sprint	04
Participants	Khaled, Lisa, Sergej, Martin

## 1. What have you done?

- **Lisa:** made a NN with Sklearn with 53% acc, tried model with TF, continued the Report
- **Khaled:** Did some Decision Trees, worked on architecture
- **Sergej:** fixed an error in the data for chronological order, added feature to martins code, tried KNN found optimal K with 51%
- **Martin:** started extraction of ball possession feature

## 2. What will you do?

- **Lisa:** take historical ordered dataset and try Sklearn and TF model again with new data, continue report
- **Khaled:** keep working on the classifiers and on the architecture
- **Sergej:** design a couple of NN, use Schwerins list of most common used NN's
- **Martin:** extract ball possession in minutes per match, try Schwerins list of common used NN's



### 3. Issues?

- Lisa: -
- Khaled: -
- Sergej: -
- Martin: -



<b>Titel</b>	<b>Daily Scrum</b>
<b>Date</b>	09.01.2020
<b>Sprint</b>	04
<b>Participants</b>	Khaled, Lisa, Sergej, Martin

## 1. What have you done?

- **Lisa:** Finished reading raschka ML book, continued with the report, tried the sklearn model and decision tree with historical ordered data,
- **Khaled:** Finished DecissionTrees for 2 new options, research of multy class classifications NN, coding review
- **Sergej:** wrote some code, built some models, researched and implemented keras + TF + google colab
- **Martin:** comparison of different amount of hidden layers and number of neurons per layer, extract ball possession in minutes per match

## 2. What will you do?

- **Lisa:** finishing the report and create the presentation
- **Khaled:** finishing the report and create the presentation
- **Sergej:** finishing the report and create the presentation
- **Martin:** finishing the report and create the presentation

## 3. Issues?

- **Lisa:** -
- **Khaled:** -
- **Sergej:** -
- **Martin:** -

### **A.3. Daily Scrum Logs second therm**



<b>Titel</b>	<b>Daily Scrum</b>
<b>Date</b>	18.03.2020
<b>Sprint</b>	05
<b>Participants</b>	Khaled, Lisa

## 1. What have you done?

- **Lisa:** presentation for the future students (5h)
- **Khaled:** set up the virtual machine (5h), add slides for presentation for the future students (2h)

## 2. What will you do?

- **Lisa:** test different test and learning data, work with google colab
- **Khaled:** compare models with more statistical tools

## 3. Issues?

-



Titel	Daily Scrum
Date	26.03.2020
Sprint	05
Participants	Khaled, Lisa, Hemlata

## 1. What have you done?

- **Lisa:** tested smaller test- data set sizes (10h)
- **Khaled:** extract best 15 models (7h), made test with drop out and weight regularization(3h30)

## 2. What will you do?

- **Lisa:** test bigger test data size 10,15, 30..., research statistical tools
- **Khaled:** research statistical tools, and try to do something

## 3. Issues?

-



Titel	Daily Scrum
Date	11.04.2020
Sprint	06
Participants	Khaled, Lisa, Hemlata, Till

## 1. What have you done?

- **Lisa:** tested bigger test data sizes -> 10,15%, research statistical tools (12h)
- **Khaled:** Research for Technology for the Frontend (3h), set up working environment -> frontend (3h), implement basic website(4h)
- **Hemlata:** Research for ML and DL(3h), worked through already existing Notebooks(3h), Decision Trees, Reading Raschka, Research for Model Improvement, read report(6h)
- **Till:** Looked through Notebooks (8h)

## 2. What will you do?

- **Lisa:** research statistical tools, confirm with Khaled, try some statistical tools
- **Khaled:** Research on the API, Backend...
- **Hemlata:** Work through the remaining Notebooks, Reading Raschka
- **Till:** Look into the slidingWindow.csv Files, look through more Notebooks, evaluate alternative Options, read the report

## 3. Issues?

-



Titel	Daily Scrum
Date	18.04.2020
Sprint	06
Participants	Khaled, Lisa, Hemlata, Till

## 1. What have you done?

- **Lisa:** research statistical tools + first implementation (10h)
- **Khaled:** Research Web API (2h), Define Web architecture (backend, frontend)(2h), Define RESTful API(3h), Implement simple backend with RESTful API using JSON format (3h30)
- **Hemlata:** read report(3h), learned existing code(6h), understanding tf and keras(2h)
- **Till:** read report (3h), started with raschka (3h), looked and improved khalels backend api (4h)

## 2. What will you do?

- **Lisa:** implementation statistical tools, find best models
- **Khaled:** Implement RESTful API frontend (dynamic)
- **Hemlata:** try to understand tf and keras, comparing models and find the best one
- **Till:** learn keras and tf, working with Khaled for the api & backend

## 3. Issues?

-



Titel	Daily Scrum
Date	23.04.2020
Sprint	06
Participants	Khaled, Lisa, Hemlata, Till

## 1. What have you done?

- **Lisa:** implemented all statistical tools for all models and made a excel sheet with outcome, make a notebook with best model (15h)
- **Khaled:** Implement RESTful API frontend (6h), implement functionalities frontend (dynamic)(5h)
- **Hemlata:** read report(2h), learned existing code(3h), understanding tf and keras(6h)
- **Till:** read report (2h), continued with raschka (6h)

## 2. What will you do?

- **Lisa:** research regression, implementing simple model using regression, find best model
- **Khaled:** research regression, implementing simple model using regression
- **Hemlata:** try to understand tf and keras, comparing models and find the best one
- **Till:** implement backend, continue with raschka

## 3. Issues?

-



Titel	Daily Scrum
Date	01.05.2020
Sprint	07
Participants	Khaled, Lisa, Hemlata, Till

## 1. What have you done?

- **Lisa:** found best Model according to accuracy -> model04\_H2\_H (2h) , research regression, research on existing model (9h)
- **Khaled:** data preparation for regression implementation (3h), created sliding windows(5h), research regression (3h)
- **Hemlata:** try understanding keras model, learning about libraries, reading raschka book, tried some changes in already existing code to improve understanding (11h)
- **Till:** started with backend (12h)

## 2. What will you do?

- **Lisa:** research regression, first implementations
- **Khaled:** research regression, first implementations
- **Hemlata:** trying to implement a model by herself, research regression
- **Till:** resume with backend

## 3. Issues?

-



Titel	Daily Scrum
Date	07.05.2020
Sprint	07
Participants	Khaled, Lisa, Hemlata, Till

## 1. What have you done?

- **Lisa:** research regression (5h) , research on existing model -> only make sense for images (2h) , first tries on implementations (5h)
- **Khaled:** research regression (5h), started first implementations (3h), but got bad result
- **Hemlata:** research multi target regression(3h), started with decision tree with multi target output(4h), tried writing models by myself(3h)
- **Till:** mostly finished backend (11h), find out to interpret output (1h)

## 2. What will you do?

- **Lisa:** research regression, more implementations
- **Khaled:** research regression, more implementations
- **Hemlata:** trying to implement a model by herself, research regression, decision tree regression implementation
- **Till:** resume with backend, improve knowledge about existing models, research on fining bookie info etc...

## 3. Issues?

-



Titel	Daily Scrum
Date	14.05.2020
Sprint	07
Participants	Khaled, Lisa, Hemlata, Till

## 1. What have you done?

- **Lisa:** research regression (2h) , implementation of a 10 class classification model best acc -> about 37% (5h), decision tree regression (5h)
- **Khaled:** tested different classifier Multi-Target Regressors: Multi-Output Regressor, Decision Tree Regressor, Random Forest Regressor, MLP Regressor(6h), compare regressors(3h), started with presentation (30min)
- **Hemlata:** implement three regressions, decision tree, random forest regressor, regression through svr(11h)
- **Till:** research for data source <https://www.football-data.co.uk/germany.php> (2h), work on backend (6h)

## 2. What will you do?

- **Lisa:** implementation regression, research regression
- **Khaled:** research regression, continue to improve -> hints from profs
- **Hemlata:** research regression, continue to improve
- **Till:** database, integration of football data

## 3. Issues?

-



Titel	Daily Scrum
Date	21.05.2020
Sprint	08
Participants	Khaled, Lisa, Hemlata, Till

## 1. What have you done?

- **Lisa:** research in multi task classification (4h), research in independent multi task learning (5h)
- **Khaled:** data preparation for reducing classes (2h), better solution for the round of predicting values (2h30), started working for quality criteria function(3h30)
- **Hemlata:** research on all kind of regression algorithms(9h)
- **Till:** worked on database integration (8h), research for data source (3h)

## 2. What will you do?

- **Lisa:** more research custom output layer, research in multi task classification, research in independent multi task learning
- **Khaled:** develop quality criteria
- **Hemlata:** regression implementation with keras
- **Till:** implement data fetching

## 3. Issues?

-



Titel	Daily Scrum
Date	28.05.2020
Sprint	08
Participants	Khaled, Lisa, Hemlata, Till

## 1. What have you done?

- **Lisa:** research in multi task classification, research in independent multi task learning, research custom output layers, tried to implement custom output layer (9h)
- **Khaled:** developed quality criteria (6h), tested on scikit learn regressor(2h)
- **Hemlata:** implementation regression with keras using neural network(7h), different kind of loss(2h), different matrix(3h)
- **Till:** worked on data fetching (10h)

## 2. What will you do?

- **Lisa:** found multi task classification for pictures with cnn, try to adapt it to our problem
- **Khaled:** look through hemlatas code, test other regressors, take the best out of both codes
- **Hemlata:** try to include quality criteria from Khaled, take the best out of both codes
- **Till:** work on data fetching

## 3. Issues?

- we are all very stressed from seminar paper



Titel	Daily Scrum
Date	04.06.2020
Sprint	08
Participants	Khaled, Lisa, Hemlata, Till

## 1. What have you done?

- **Lisa:** implementation in independent multi task learning (8h)
- **Khaled:** worked through Hemlatas code (1h), implemented artificial NN with Tensorflow (5h), regression tensorflow research (3h)
- **Hemlata:** improved regression model, model comparing(9h)
- **Till:** finished data fetching (8h)

## 2. What will you do?

- **Lisa:** multi task learning implementation -> solving existing problems
- **Khaled:** help Lisa and Hemlata 😊
- **Hemlata:** try to include quality criteria from Khaled, take the best out of both codes
- **Till:** finish back end

## 3. Issues?

- we are all very stressed from seminar paper



Titel	Daily Scrum
Date	12.06.2020
Sprint	09
Participants	Khaled, Lisa, Hemlata, Till

## 1. What have you done?

- **Lisa:** more implementation in multi task learning (7h)
- **Khaled:** defined api for frontend and backend for goals prediction(3h), adjusted frontend for final results (6h)
- **Hemlata:** improved quality criteria(4h), research ML(3h), Test different parameters(4h)
- **Till:** defined api for frontend and backend for goals prediction (3h), changed sorting of matches (1h)

## 2. What will you do?

- **Lisa:** more implementation in multi task learning
- **Khaled:** add new webpage for goals, improve the frontend
- **Hemlata:** Continue testing different parameters for the model and csv files (options) for regression
- **Till:** backend for upcoming goal prediction, API calls for fetching new games, get classification results, get regression results and for retrain both models, adjusted JSON of new API

## 3. Issues?

-



<b>Titel</b>	<b>Daily Scrum</b>
<b>Date</b>	18.06.2020
<b>Sprint</b>	09
<b>Participants</b>	Khaled, Lisa, Hemlata, Till

## 1. What have you done?

- **Lisa:** Implemented multitask learning -> 35% acc (10h)
- **Khaled:** add new webpage for goals(5h), improve the frontend(4h), add fetching and retrain button(1h)
- **Hemlata:** continued testing parameters(4h), regression on sliding window 3 -> 84% acc (7h)
- **Till:** backend for upcoming goal prediction (6h), API calls for fetching new games (1h), get classification results (2h), get regression results (1h) and for retrain both models (1h) , adjusted JSON of new API (1h)

## 2. What will you do?

- **Lisa:** Improve multi task classification model
- **Khaled:** Improve first classification model
- **Hemlata:** Find different Model Accuracies,
- **Till:** Retrain Functionality

## 3. Issues?

-



<b>Titel</b>	<b>Daily Scrum</b>
<b>Date</b>	25.06.2020
<b>Sprint</b>	09
<b>Participants</b>	Khaled, Lisa, Hemlata, Till

## 1. What have you done?

- **Lisa:** Implemented multitask learning with data preparation from regression part and the quality criteria -> 81% (8h)
- **Khaled:** check the classification for home, away and draw -> no solution found(12h)
- **Hemlata:** created tables for both datasets, and different parameters -> different acc -> found best one(7h)
- **Till:** added retrain functionality for classification and regression (12h), tests of frontend and backend (1h)

## 2. What will you do?

- **Lisa:** report
- **Khaled:** report
- **Hemlata:** report
- **Till:** report

## 3. Issues?

-