



LUISS Guido Carli

Department of Business and Management
Bachelor's Degree in Management and Computer Science
Course of Data Analysis for Business

**Data-Driven Approaches
for Soccer Match Analysis**

Supervisor:
Prof. Francesco Iafrate

Candidate:
Martino Olivieri
248021

Accademic Year 2021/22

Summary

This thesis explores possible solutions to automate match analysis in soccer. It makes use of the spatial-temporal data set published by the *Wyscout* company. The first part will be entirely dedicated to exploring and visualizing the data together with a concrete description of possible applications. While the second part is completely devoted to the definition and implementation of two machine learning models: the first is supervised and is used to get a numerical evaluation of every player. In contrast, the second is a clustering model that detects relevant areas on a soccer pitch. The analysis was entirely done by exploiting the possibilities offered by the R programming language, while the parsing of the JSON files was done using Python; details regarding the script can be found in the appendix.

Contents

Summary	II
List of Figures	V
List of Tables	VI
1 Introduction	1
1.1 Introduction to match analysis powered by data	1
1.2 Collecting data in soccer	2
1.3 Introduction to the structure of the thesis	2
2 Data Pre-Processing	3
2.1 Data description	3
2.2 Exploratory data analysis	6
2.3 Data cleaning	13
2.4 Feature Engineering	14
3 Models	17
3.1 Player rating	17
3.1.1 Team Performance Extraction	17
3.1.2 Feature Weight Extraction	18
3.1.3 Individual Performance Extraction	18
3.1.4 Application	18
3.2 Clustering For Role Detection	21
3.2.1 Introduction To The Purpose Of The Model	21
3.2.2 K-mean Clustering explained	21
3.2.3 Model implementation	21
3.2.4 Application	23
Conclusions	25

A Python and R Code	26
A.1 Python	26
A.1.1 JSON parsing	26
A.2 R Code	30
A.2.1 Data cleaning and Visualization	30
A.2.2 Supervised Model(SVM)	43
A.2.3 Unsupervised Model (K-means)	51
Bibliography	54

List of Figures

2.1	Frequency of each event type	7
2.2	Heat map of the most relevant events type	7
2.3	Distributions of the <i>goal</i> and <i>pass</i> column	8
2.4	Frequency of the red cards in each time frame of the matches .	9
2.5	Position of the goals scored in <i>Roma - Chievo (4-1)</i>	10
2.6	Shots from both teams in <i>Roma - Chievo (4-1)</i>	10
2.7	Possible report of <i>Roma - Chievo (4-1)</i>	11
2.8	Network of passes for A.S Roma in <i>Roma - Chievo Verona 4-1</i> <i>2017/18</i>	12
3.1	Top 8 positive and negative feature weights	19
3.2	Visualization of the results distribution	20
3.3	Silhouette Score	22
3.4	K-means clustering result with K=8	23

List of Tables

2.1	Features selected for the modelling	16
3.1	SVM results	19
3.2	5 of the top player ratings selected by the model	20
3.3	Silhouette results	22
3.4	Interpretation of the clusters	23

Chapter 1

Introduction

1.1 Introduction to match analysis powered by data

As every team sport, soccer is becoming more and more based on strategies. The Oxford Dictionary defines tactics as "a carefully prepared action or strategy to achieve a specified aim." Naturally, the goal of competitive soccer is to win the game. For this reason, the coach's role in a soccer team is evolving, requiring many more skills with respect to the past. Hence new figures, such as a team of match analysts, are now starting to complement the coach. Data is also central in modern soccer, as many decisions are powered by analytics. A quick example could be predicting a player's injuries based on physical and spatial information [1]. These kinds of applications have a significant impact even on the financial part of a soccer company since an injury could cause substantial monetary losses. Data science in soccer is applied mainly in match analysis which is the objective study of all tactical and technical aspects that a soccer player or team performs during a match. The main issue regards collecting valuable data since it requires advanced hardware and software tools that only a few companies can afford. *A.S Roma* is one of them and Its coach, Jose Mourinho, uses spatial-temporal data during and before the match for decision-making. In the past, the analysis of a match was performed by experts. After looking at all the opponent's previous games, they would compile a report containing the adversary's strengths and weaknesses. Nowadays, by exploiting new technologies, some analysis processes could be automated and hidden patterns could be discovered.

1.2 Collecting data in soccer

The previous section highlighted how difficult and expensive it is to collect spatial-temporal data in soccer. "Aside from occasional attempts, soccer statistics have only recently been produced, thanks to sensor technologies that offer high-fidelity data streams taken from every match" [4]. There are three main methods:

- **GPS:** The soccer companies can track their players via *GPS* devices during both games and training.
- **Computer Vision:** The data can be collected using object tracking models.
- **Soccer-log:** Description of the events that occur during a match collected through proprietary tagging software.

The first two are automated and are considerably more complete and precise, while the third is semi-automated. The main difference between the two is that one provides continuous data (both offensive and defensive phases), while the second only describes events regarding the ball possession. Due to the scarce data availability of GPS and Tracking data, this thesis explores the possible application of the soccer-log data-set published by *wyscout*.

1.3 Introduction to the structure of the thesis

This thesis is structured in two main sections that are actually connected to each other. The first phase of the pipeline regards the exploration and visualization of the data together with a concrete description of possible applications for a coach. This section will try to discover patterns in soccer games: starting from the big picture on the whole dataset and then going deep into the performance of the single teams and players. Instead, the second phase is entirely dedicated to machine learning. Both supervised and unsupervised methods will be applied to discover hidden patterns.

Chapter 2

Data Pre-Processing

2.1 Data description

As stated in the last section, the most challenging part is Collecting the data. The Italian Company *Wyscount* published a free-to-use database that comprehends the season 2017/2018 of the top 5 European national soccer competitions: Spanish first division, Italian first division, English first division, German first division, and French first division. The UEFA nation coefficient indicates that these competitions are the most significant in Europe [4]. A total of seven data sets include details on all tournaments, games, teams, players, events, referees, and coaches. The JSON format (JavaScript Object Notation) is used to convey each data set. For the analysis, only the following tables will be used:

Events is the main table used because it describes every event in every match. The following are the Information contained:

- **eventId**: The identifier of the event type;
- **EventName**: The name of the event;
- **subEventId**: The identifier of the sub event;
- **subEventName**: The name of the sub event;
- **tags**: A list of tags related to the event;
- **eventSec**: The match' second during which the event took place;
- **id**: The unique identifier of the event;

- **matchId**: The identifier related to the match to which the event refers;
- **matchPeriod**: In which period the event took place. First or second half;
- **playerId**: The unique identifier of the player that generated the event;
- **positions**: The coordinates of the pitch where the event started and ended;
- **teamId**: The unique identifier of the event to which the player belongs;

Matches contains information related to every match present in the dataset. The following are the columns:

- **competitionId**: Refers to the unique identifier of the competition;
- **date**: The date and time when the match started in implicit form;
- **duration**: The duration of the match. It can be: “Regular”, “Extra Time” or “Penalties”;
- **gameweek**: The number of week from the start of the league;
- **label**: The result of the match in the following form: ”home Team - Away Team, score”;
- **seasonId**: Refers to the season of the match;
- **status**: If the match has been already played or not;
- **venue**: The stadium in which the match was played;
- **winner**: The unique identifier of the team that won the game;
- **wyId**: The unique identifier of the match;
- **teamsData**: The data refereed to both the teams. It contains information about the formations, the bench, substitutions, coaches and score;
 - **hasFormation**: Binary variable: 1 if a formation is present, 0 if not;
 - **score**: Number of goals scored during the match not including penalties;

- **scoreET**: Number of goals scored during the whole match including extra time (not counting penalties);
- **scoreHT**: Number of goals scored during the first half of the match (not counting penalties);
- **side**: Weather the team was home or away;
- **teamId**: The identifier of the team;
- **coachId**: The identifier of the coach;
- **bench**: the list of player present in the bench, together with some basic statistic for every player related to the match;
- **lineup**: The list of players in the line-up together with some statistics related to the match for each player;
- **substitutions**: The list of substitution that took place during the match;

Players contains all the information related to every player present in the data set. The following are the Information contained:

- **birthArea**: Information about the birth area;
- **birthDate**: The birth date of the palyer (YYYY-MM-DD);
- **currentNationalTeamId**: The unique identifier of the national team in which he plays;
- **currentTeamId**: The unique identifier of the team in which he plays;
- **firstName**: The first name of the player;
- **lastName**: The last name of the player;
- **foot**: The preferred foot of the player;
- **height**: He height of the player (in centimeters);
- **middleName**: The middle name (if any) of the player;
- **passportArea**: The geographic area associated with the player's current passport;
- **role**: The main role of the player;

- **weight**: The weight of the player (in kilograms);
- **wyId**: The identifier of the player, assigned by Wyscout;

Tag contains every possible tag present in the "events" data table to describe each observation: the following are the features:

- **Tag**: The unique identifier of the tag;
- **Label**: The actual name of the tag;
- **Description**: The description of the tag;

2.2 Exploratory data analysis

The first step is the *exploratory data analysis (EDA)*. Usually it is the first phase toward a good report. It comprehends all the methods to investigate relationships between variables and provides an overall description of the tables mainly through data visualization. The EDA process will be focused on the *events* table, since it will be the only one used in the modelling section. The analysis could already result in insightful information for a team's coach. So, together with statistical reports, a possible application will be described.

Event Name

The first interesting variable to visualize is *EventName* which contains what event took place from a set of 10 possible levels.

Fig.2.1 shows that *pass* is the most frequent observation. This was an expected result for who knows a bit about soccer, but it is always good to have statistical confirmation.

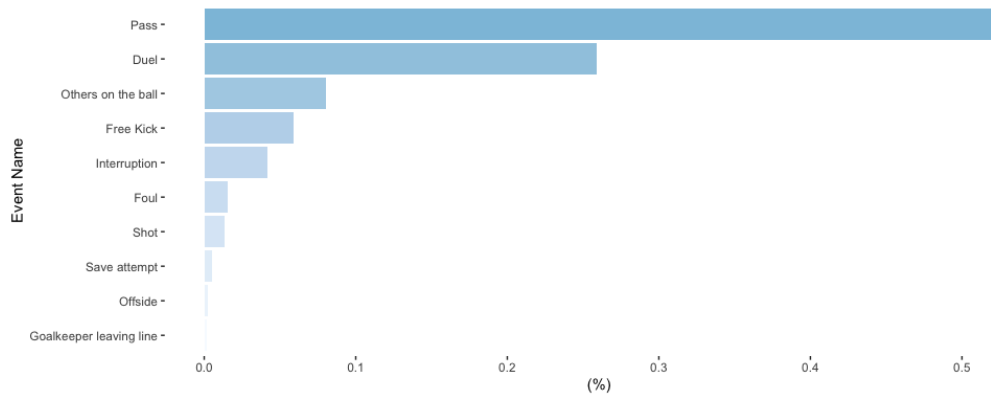


Figure 2.1: Frequency of each event type

Instead the Fig.2.2 shows the heat map of the 4 main events. The *pass* event is very dense in the central area of the football field (fig. a), while the *shots* are mostly concentrated in the penalty area (fig b). For the *Duel* and *Fouls* there is no evident pattern.

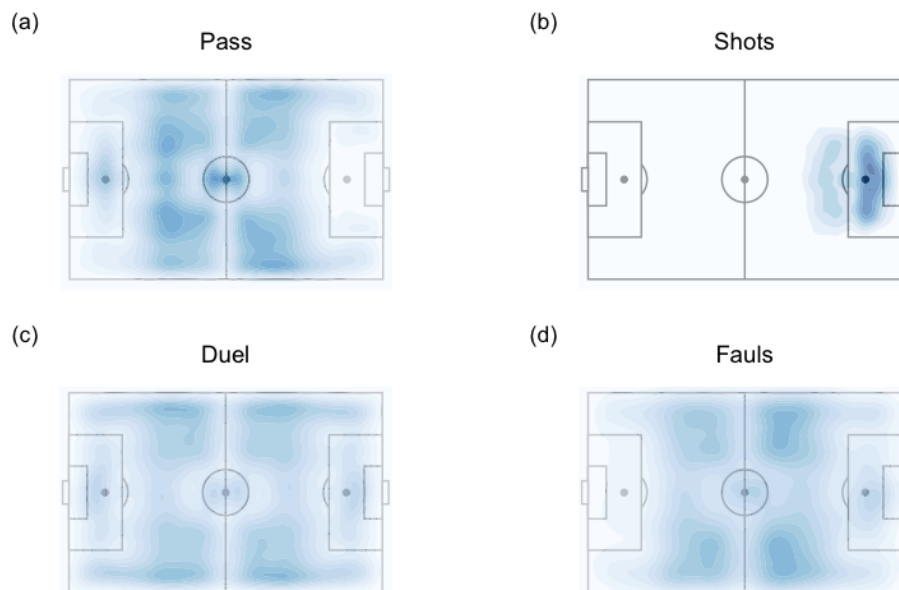


Figure 2.2: Heat map of the most relevant events type

Goal and Pass events

Fig.2.3 shows the statistics of the *goal* and *pass* variables. In particular, fig (a) and (b) demonstrate that in many matches around three goals are

scored while in only a few of them more than 4. Instead, fig (c) and (d) shows the density of respectively accurate and not accurate passes. The mean in each game is 740 for precise passes and 148 for the not accurate pass. Furthermore, it is clear that the two distributions are normal, but on average, more accurate passes are completed with respect to non-accurate. x

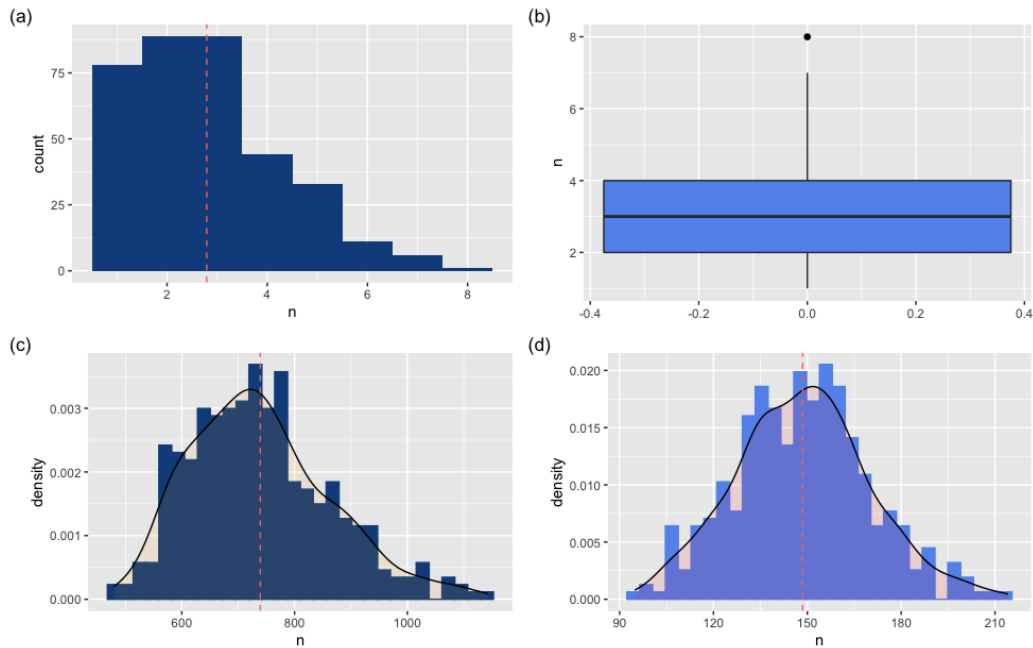


Figure 2.3: Distributions of the *goal* and *pass* column

Distribution along the 90 minutes

The fact that the data are spatial-temporal is essential for the analysis. We could use the *eventMin* column, not just in the ML model, but also to get insightful statistics and visualization. For example, Fig. 2.4(1) shows the frequency of the red cards in each time frame of the matches. The referees seem to give more red cards at the end of the game, maybe because the players are tired or stressed by the result. Instead, Fig.2.4(2) represents the frequency of the yellow cards in each time frame: apart from the start of the game, during which the referees tend to give fewer yellow cards, the frequency looks flat. At the same time, Fig.2.4 (3) regards the frequency of the goals that are well distributed along the 90 minutes.

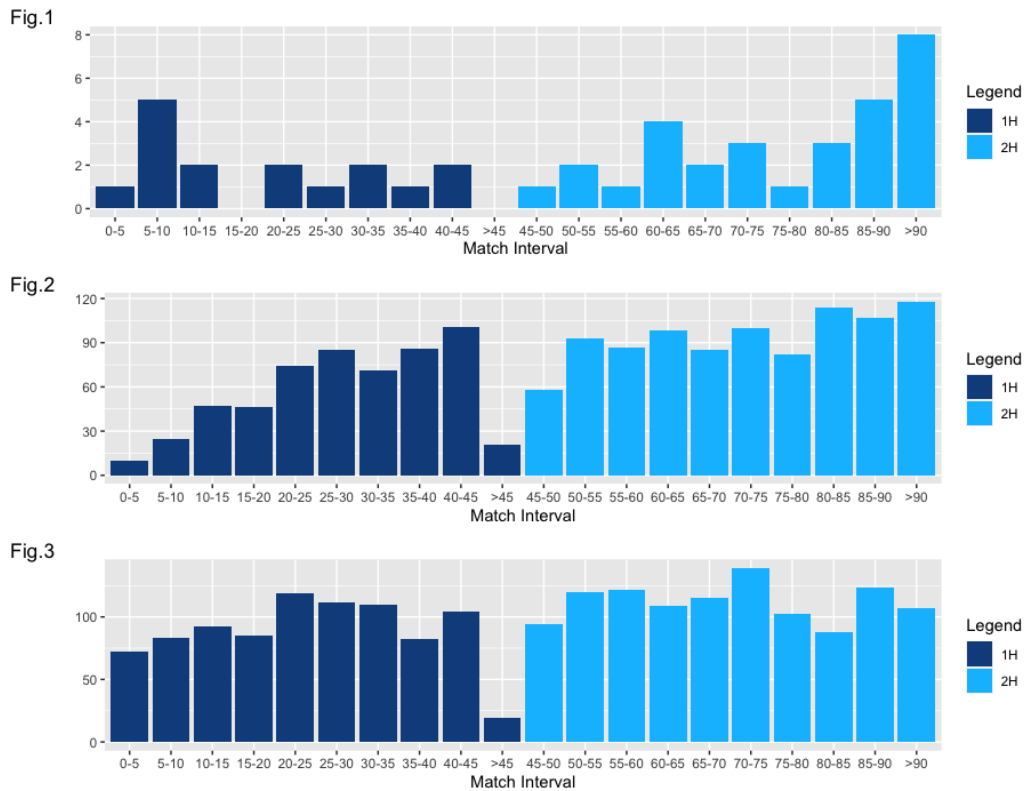
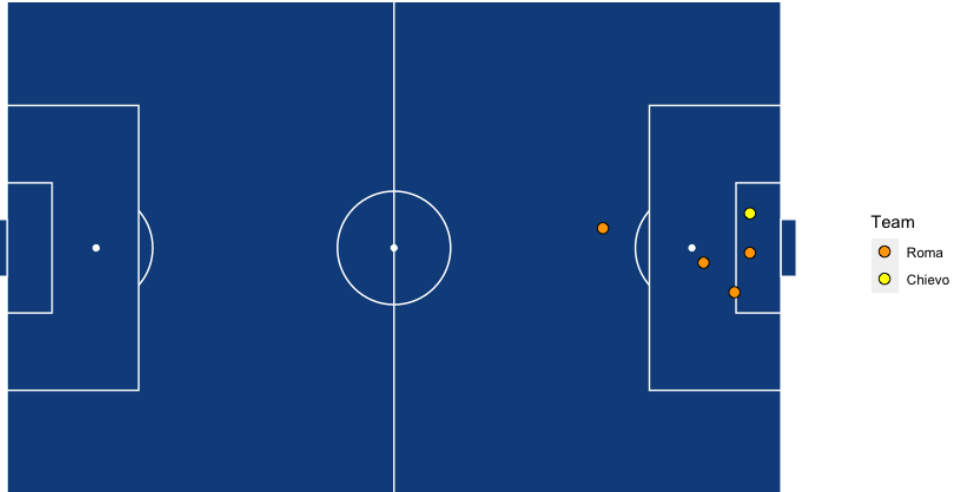
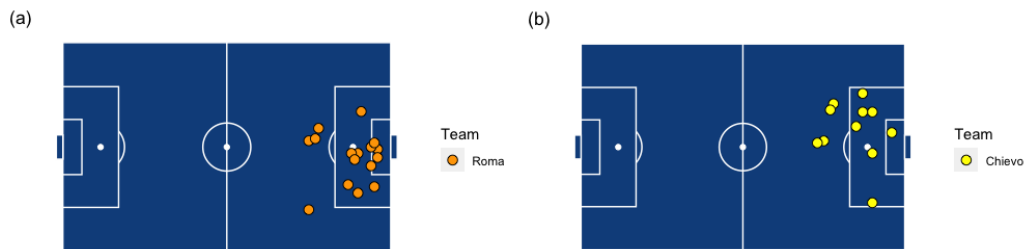


Figure 2.4: Frequency of the red cards in each time frame of the matches

Focus on *Roma - Chievo Verona 4-1 2017/18*

We have a structured agglomerate of data from the season 2017/18. So it would be good to use visualization to understand the whole data set as shown in the last sections. Still, it would be even more insightful to analyze one match in particular from a statistical point of view. For this reason, the match *Roma - Chievo Verona 4-1 2017/18* was selected and analyzed.

Firstly since we have the spatial data about every event, it could be informative to visualize where the goals were scored. Fig.2.5 shows the position in which the players scored a goal during the match. The same visualizations is shown in Fig.2.6, but here the data points represent the shots.

Figure 2.5: Position of the goals scored in *Roma - Chievo (4-1)*Figure 2.6: Shots from both teams in *Roma - Chievo (4-1)*

The same analysis could be enlarged on the whole season of a particular team. In this case, using the *wyscout* dataset, fig Fig.2.7 shows a possible visualization of both *Roma* and *Chievo*. In particular fig (a) shows where *Chievo Verona* scored most of the goals. Fig. (c) and (d) show in which area of the field the teams completed more passes. Instead fig (b) shows where the right back, *Aleksandar Kolarov* played most of the season. All

this information together would be very precious to build a strategy for a specific game.

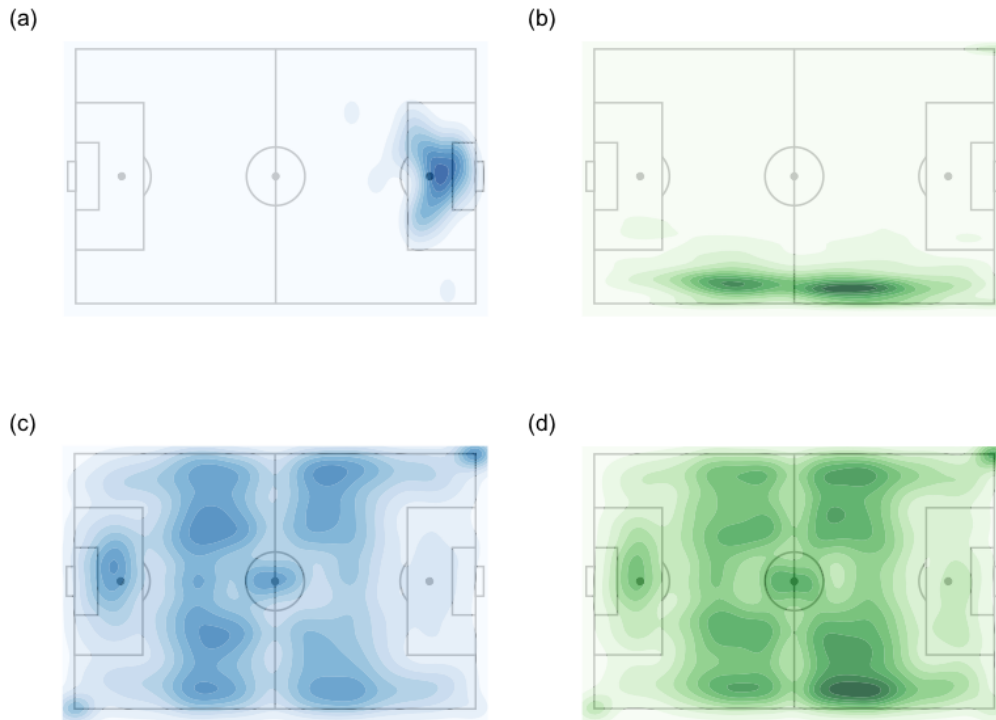


Figure 2.7: Possible report of Roma - Chievo (4-1)

Network Analysis

A more sophisticated visualization concern the analysis of the network of a team. The data structure is not ready for this plot, so some transformation must be applied. To build a network of passes between players of the same team, there should be a sender (who pass the ball) and a receiver (who receives the ball), and both the data should be on the same row. To do so, the following steps are implemented:

- *Apply filter to retrieve only the pass observations;*
- *Place a mark where the ball possession change;*
- *If (the ball possession does not change):*

- Select $idPlayer$ in $row(i+1)$;
- Add it to $row(i)$

- Group by sender and receiver

The resulting table will be composed by three columns (receiver, sender, weight) which are what a network asks as input. Fig. 2.8 shows the graph of *A.S. Roma* in the game *Roma - Chievo Verona 4-1 2017/18* with the following components:

- **Nodes:** Players
- **Edges:** Pass completed
- **Weights:** n of passes completed

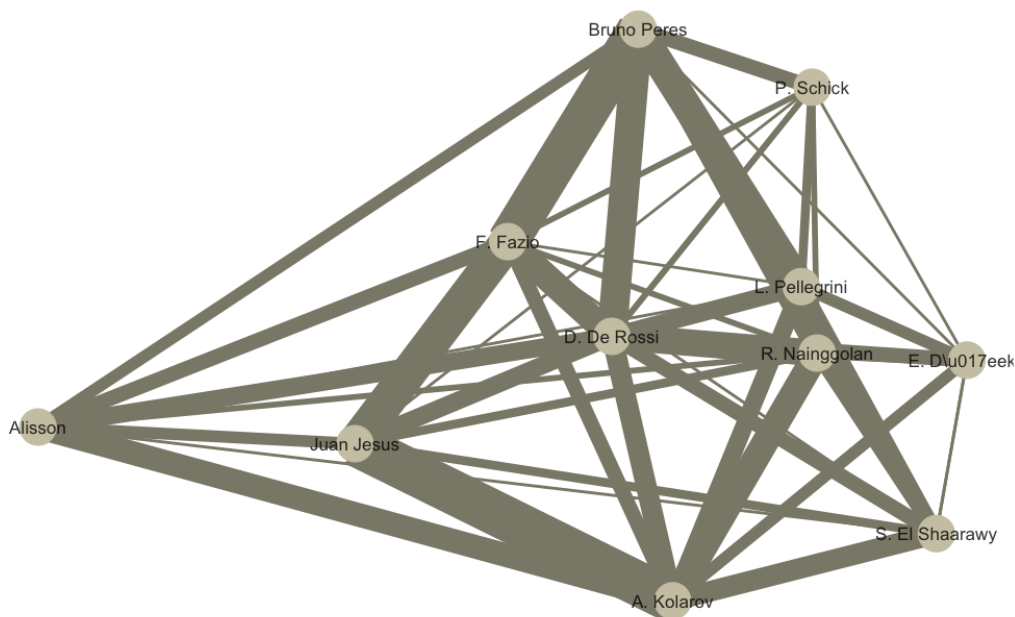


Figure 2.8: Network of passes for A.S Roma in *Roma - Chievo Verona 4-1 2017/18*

For a coach, a set of networks of the last matches played by the team he will face next is fundamental. Maybe a strategy could be to place a player between two highly connected nodes in order to break their ball possession. In this particular case for example, the defenders seems to be highly connected.

Playing with two strikers that press them could be a winning strategy. There exists several possible application of the analysis of a soccer team network. In the paper *Using Network Science to Analyse Football Passing Networks: Dynamics, Space, Time, and the Multi layer Nature of the Game* [5], the authors try to analyze how the network a team changes during a game.

2.3 Data cleaning

Data cleaning is often the most time-consuming part of the process. Data must be processed to be ready to use in the modeling part. Otherwise, the fitting could incur in some problems. The analysis will be focused on the "events" table, but additional features will be taken from the other data-sets. To begin, the first thing to do is to exclude all the observations and the columns that will not be used to make later transformations lighter. For this reason, a good understanding of the topic is needed. To make a quick example, some of outliers should be considered not just from a statistical perspective but also in a concrete sense. Observations such as side fouls should be deleted to make a good clustering model. Firstly only the following observations from the *tag* column were kept (assist, keyPass, interception, red card, yellow card, opportunity, second yellow card, lost, neutral, won, accurate, Feint, dangerous ball lost, counter attack). The above entries were not helpful and insightful for analysis.

The same reasoning was applied to the *eventName* and *subEventName* columns. Some of the entries were filtered out because they were not insightful and valuable for the model. From *eventName* only the following were kept: (Duel, Foul, Free Kick, Others on the ball, Pass, Shot). Instead in the *subEventName* more rows were used (Ground attacking duel, Ground defending duel, Ground loose ball duel, Hand foul, Late card foul, Out of game foul, Protest, Simulation, Violent Foul, Free kick cross, Penalty, Free kick shot, Throw-in, Acceleration, Clearance, Touch, Cross, Hand pass, Head pass, High pass, Launch, Simple pass, Smart pass, Shot, Air duel, Foul, Free Kick, Corner).

The data-set published by *Wyscout* was already cleaned from dirty records. This is because the great work made by their team. Following the tagging, each match goes through a quality control process: The first phase is automatic (an algorithm is used to prevent most operator errors); The second step of quality control is manual and supervised by quality controllers.

2.4 Feature Engineering

The process of extracting features from unprocessed data and converting them into formats appropriate for machine learning models is known as *feature engineering*. The correct features can make modeling less challenging and enable the pipeline to output results of greater quality, making it an essential step in the machine learning process. After the parsing of the *JSON* files, many columns presented an useless format such as nested dictionaries. This kind of format can not be used to fit a model for this reason all the information were extracted and placed in singular new column. This process was applied to the following columns:

- **Events Position:** *format:* xstart: x, ystart: y, both the values were extracted and added to 2 new columns having the names of the keys.
- **Matches Teams data** *format:* dictionary containing the following keys:
 - "coachId":
 - "formation":
 - * "bench": []
 - * "lineup": []
 - * "substitutions": []
 - "hasFormation":
 - "score":
 - "scoreET":
 - "SCoreHT"
 - "scoreP":
 - "side":

Every one of the above information were taken and added the newly created columns having as names the keys.

Event Second

As previously stated, one of the column contains the seconds related to the event. Taking into consideration that each entry includes the seconds passed from the first or second half of the game, the column is entirely converted on a minutes base.

Tags

The most challenging part of the *feature engineering* phase regards the Tag column in the *events* table. This information will be essential in modeling since the ML process will use them as variables. The column has the form of a dictionary that contains the tags. For the modeling stage, it is necessary to have a combination of *eventName*, *subEventName* and *tags*. A total of 67 variables were selected (2.1)

Features	
Duel.Air.duel.accurate	Others.on.the.ball.Acceleration.not.accurate
Duel.Air.duel.not.accurate	Others.on.the.ball.Clearance.accurate
Duel.Ground.attacking.duel.accurate"	Others.on.the.ball.Clearance.not.accurate
Duel.Ground.attacking.duel.not.accurate"	Others.on.the.ball.Touch.assist
Duel.Ground.defending.duel.accurate	Others.on.the.ball.Touch.counter_attack
Duel.Ground.defending.duel.not.accurate	Others.on.the.ball.Touch.dangerous_ball_lost
Duel.Ground.loose.ball.duel.accurate	Others.on.the.ball.Touch.interception
Duel.Ground.loose.ball.duel.not.accurate	Others.on.the.ball.Touch.opportunity
Foul.Hand.foul.red_card	Pass.Cross.accurate
Foul.Hand.foul.yellow_card	Pass.Cross.assist
Foul.Late.card.foul.yellow_card	Pass.Cross.keyPass
Foul.Foul.red_card	Pass.Cross.not.accurate
Foul.Foul.second_yellow_card	Pass.Hand.pass.accurate
Foul.Foul.yellow_card	Pass.Hand.pass.not.accurate
Foul.Out.of.game.foul.red_card	Pass.Head.pass.accurate
Foul.Out.of.game.foul.second_yellow_card	Pass.Head.pass.assist
Foul.Protest.red_card	Pass.Head.pass.keyPass
Foul.Out.of.game.foul.yellow_card	Pass.Head.pass.not.accurate
Foul.Protest.second_yellow_card	Pass.High.pass.accurate
Foul.Protest.yellow_card	Pass.High.pass.assist
Foul.Simulation.yellow_card	Pass.High.pass.keyPass
Foul.Violent.Foul.red_card	Pass.High.pass.not.accurate
Foul.Violent.Foul.yellow_card	Pass.Launch.accurate
Free.Kick.Free.kick.cross.accurate	Pass.Launch.keyPass
Free.Kick.Free.kick.cross.not.accurate	Pass.Launch.not.accurate
Free.Kick.Free.Kick.accurate	Pass.Simple.pass.accurate
Free.Kick.Free.Kick.not.accurate	Pass.Simple.pass.keyPas
Free.Kick.Penalty.not.accurate	Pass.Simple.pass.not.accurate
Free.Kick.Free.kick.shot.not.accurate	Pass.Smart.pass.accurate
Free.Kick.Free.kick.shot.accurate	Pass.Smart.pass.assist
Free.Kick.Throw.in.accurate	Pass.Smart.pass.keyPass
Free.Kick.Throw.in.not.accurate	Pass.Smart.pass.not.accurate
Others.on.the.ball.Acceleration.accurate	Shot.Shot.accurate
Shot.Shot.not.accurate	

Table 2.1: Features selected for the modelling

Chapter 3

Models

3.1 Player rating

This section will reinterpret the *playerRank model* described in the "*PlayeRank: Data-driven Performance Evaluation and Player Ranking in Soccer via a Machine Learning Approach*" [7], together with its implementation.

To be prepared for a game, a numerical valuation of the overall team rating could be precious. For example, the coach could look at each player's ratings and understand who is the most dangerous. Machine learning can help in the process, but the problem is that the *Wyscout* data set doesn't provide a target variable to build a supervised model. There is no objective and numerical evaluation of the players. *Player Rank* tries to find a solution to this problem by ranking the players based on their performances. For this purpose, all the 67 features described in section 2.4 were used. First, a classification model will try to predict the game's outcome. Secondly, the weights estimated by the model will be extracted and used to rank each player based on his performances.

The following sections will describe every step of the *Player Rank* framework.

3.1.1 Team Performance Extraction

The first phase regards the extraction of the performance of each team in each game. To retrieve this information is enough to group by every match. Every observation of the resulting table will be a match represented by the 67 features.

$$X = [x_1, \dots, x_i]$$

$$x_i = \sum_{r=1}^n x_i$$

In the above equation X_i is the team performance vector while x_i is the i variable.

3.1.2 Feature Weight Extraction

The framework's second step is about classifying the match outcome, which is not the primary objective, but it is fundamental in understanding how much each event type influences the result. The likelihood of winning a game is more strongly influenced by some players' actions throughout a match than by others. For example, a red card will probably negatively affect the outcome. The analysis excludes events such as *Goal* because they can make the model biased. Thus from the model is extracted the vector W of feature weights:

$$W = [w_1, \dots, w_i]$$

that represents the importance of each variable on the performance of the team.

3.1.3 Individual Performance Extraction

First, the extraction of each player's (u) performance (r) in each game (m) is needed. It is given by the vector:

$$P_u^m = [x_1, \dots, x_i]$$

in which the x_i represents the i feature. Finally, the vector of weight W will be used to weight each of the x_i :

$$r(u, m) = x \sum_{i=1}^n w_i \cdot x_i$$

The last step is to get the mean over all the m matches for each player u

$$\bar{r}(u, M) = \frac{1}{R} \sum_{i=1}^n r(u, m_i)$$

3.1.4 Application

Using the data set published by *wyscout*, the following section will describe the framework implementation explained in section 3.1. Only the data from the *Italian first division* will be used in this process.

Features Weights Extraction

First, a *Support vector machine* is used as a classification model. It should be trained to predict a binary variable:

- 1: Win
- 0: Lose/Draw

The data set used (760X68) was divided into a training set 0.75 and a test set 0.25. As metric the *AUC* was selected on a 5-fold cross validation. The results are resumed in table 3.1

Metric	AUC	F1	Accuracy	Sensitivity
Result	0.84	0.88	0.85	0.92

Table 3.1: SVM results

Then the weight were extracted from the *SVM*. Fig 3.1 shows the top 8 positive and negative feature weights. It seems that an accurate shot has very positive impact on the outcome of a game while, for example a not accurate simple pass has a negative impact on it.

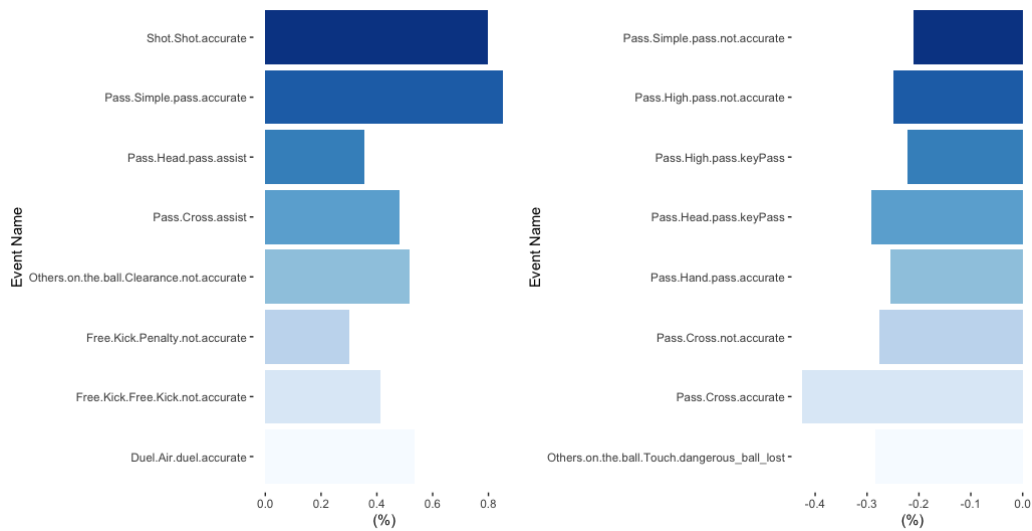


Figure 3.1: Top 8 positive and negative feature weights

Player Ratings

The last step is computing each player's actual rating, as explained in section 3.1.3. The table 3.2 shows some of the best players selected by the model. It looks like the players that played in *Napoli* and *Juventus* are the best ones. In fact, in the *Serie A season 2017/18*, both the teams finished respectively as 2nd and 1st.

Team	Player	Rating
Napoli	Jorginho	0.54
Napoli	Kalidou Koulibaly	0.50
Juventus	Daniele Rugani	0.42
Juventus	Giorgio Chiellini	0.39
Inter	Milan Skriniar	0.38

Table 3.2: 5 of the top player ratings selected by the model

To have a better understanding of the results, fig. 3.2 shows the box plot and distribution of the ratings. The mean is around 0.17, while it looks like only few players are above the 3rd quartile.

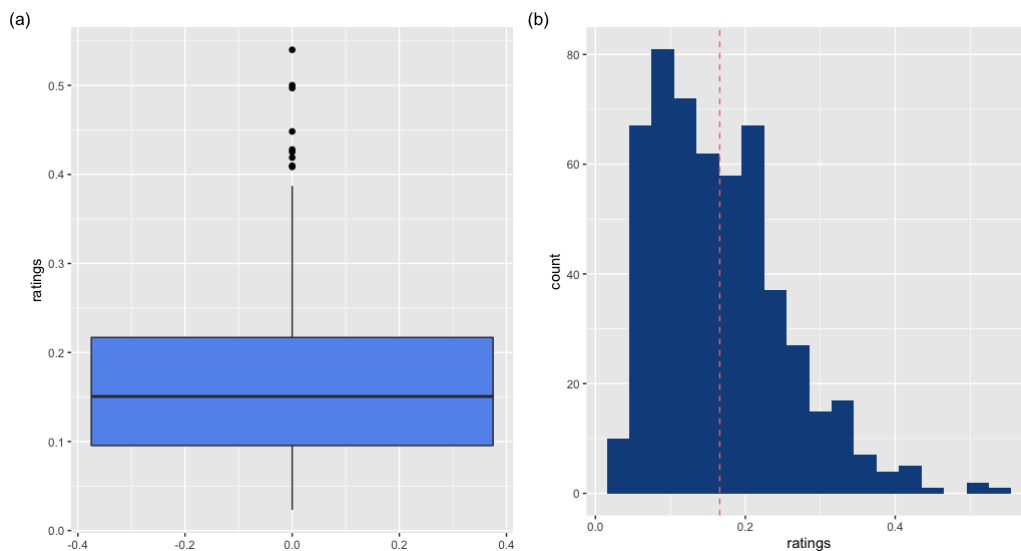


Figure 3.2: Visualization of the results distribution

3.2 Clustering For Role Detection

3.2.1 Introduction To The Purpose Of The Model

The data set published by *wyscout* does not precisely classify each player's role. It only ranks the players according to 3 areas of the field: *defender*, *midfielder*, and *forward*. In soccer, there are many more decomposition of those areas, such as central defender and right/left back. Often, the role given at the start of the game does not reflect the player's actual position during the whole game, which is what matters in building a strategy [8]. To make a quick example, going back to Fig.2.7 (b), *Aleksandar Kolarov* is a right back, so he is a defender, and he is supposed to stay more frequently in the defender's area, but looking at the plot it is clear that he could be classified as right-wing. Thus a coach could inform his left wing to help his left back when defending. For this reason, an unsupervised clustering model is implemented to try to find clear areas in the pitch that reflect the roles.

3.2.2 K-mean Clustering explained

Unsupervised machine learning techniques such as clustering are used to find and group comparable data points in larger datasets without regard to the final result. The chosen model is the Kmean clustering which can be resumed in the following steps.

Lloyd's Algorithm

1. *Specify number of clusters K*
2. *Randomly assign each data point to a cluster*
3. *Compute cluster centroids*
4. *Reassign each point to the closest cluster centroid*
5. *Re-compute cluster centroids*
6. *Repeat 4 and 5 until no improvement is made*

3.2.3 Model implementation

As stated in the section 3.2.1, the objective of the clustering model is to find relevant pitch areas that partially reflect the actual roles of the players. For this purpose, the two columns *xstart* and *ystart* will be used to fit the model. The main objective of the analysis drove this choice. The model does not

need to find the characteristics of each field area in terms of event type (even if it can be an interesting analysis). It needs to look for areas that are strongly correlated with each player's position. The second choice regards the data type in terms of single or grouped observations. The results drove the choice. Using the mean position of each player in each game, the model managed to find more defined areas. The last problem is about the only hyperparameter of the model K . It specifies the number of clusters that should be created. The selection of K can be made using different methods. In this case, the choice was driven both by numerical results and the model's objective. It does not make sense to use a small K . Otherwise, only a few non-informative areas will be clustered. As an objective selection, the silhouette score is used, as shown in Fig 3.3.

Silhouette Score

2	3	4	5	6	7	8	9	10
0.4107	0.3884	0.3728	0.3668	0.3667	0.3765	0.3872	0.3842	0.3701

Table 3.3: Silhouette results

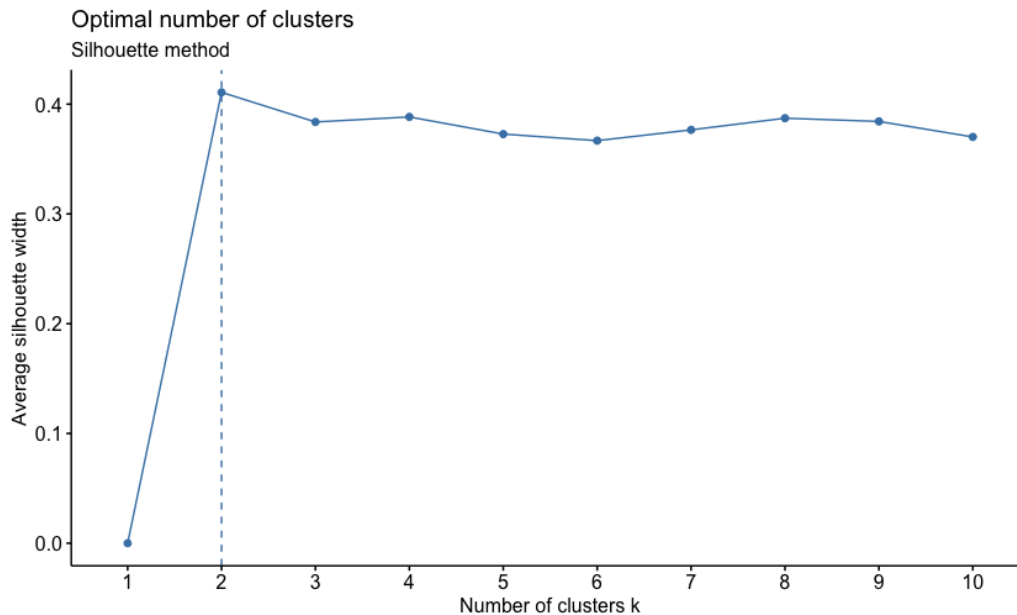


Figure 3.3: Silhouette Score

By looking just at the numbers in table 3.3, it is clear that 2 clusters give the highest score, but it would not make any sense for our analysis to divide

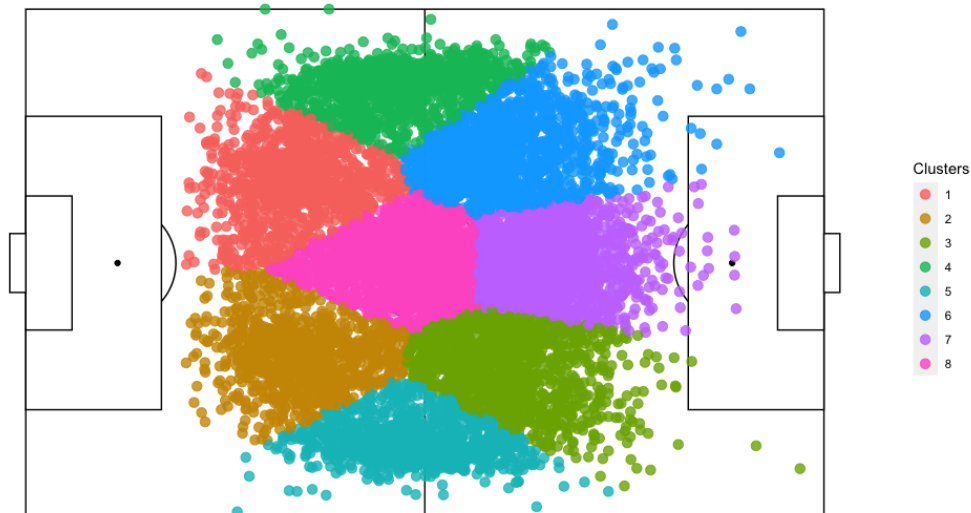
5	6	7	8
Right Back	Left Wing	Striker	Central Midfilder

Table 3.4: Interpretation of the clusters

the field into two sections. Thus by matching the score and the objective, we can choose $k=8$, which is the second highest score.

Fig 3.4 shows the clustering result using a $k=8$. The model clearly managed to delimit eight-pitch areas that could reflect the players' roles. The following table gives a possible interpretation of the 8 clusters.

Cluster	1	2	3	4
Interpretation	Left defender	Right Defender	Right Wing	Left Back

Figure 3.4: K-means clustering result with $K=8$

3.2.4 Application

There are several possible applications of this clustering model. A coach could need a report about the team he will face next week. Precious information would be to classify each player in one of the 8 clusters, so a more targeted

strategy could be implemented. To make an example, it would be insightful to understand where a player plays most of his match so to make, if necessary more density in specific areas. An even more specific application could be to use those 8 clusters to make a network of passes, to understand dangerous areas.

Conclusions

This thesis presented a set of possible solutions for automating several soccer match analysis tasks using a data-driven approach. A static but functional method proved how visualization could help capture insights about every team, from the performance of a single-player to the analysis of the whole team considered as a complex system. Two machine learning models were also used in the process. The first may help capture information about the rating of each player or team, and the second could help get insightful statistics about every player's actual position during a match. In conclusion, a match analysis powered by data is possible. It is clear that patterns exist in how a player or a team plays in a soccer match. But there are limitations. Since soccer is a sport of creativity and unpredictability, it is challenging to make specific analyses. The hope is that this thesis managed to show the possibilities of using data to automate several tasks in the match analysis, highlighting patterns in this kind of data.

Appendix A

Python and R Code

A.1 Python

A.1.1 JSON parsing

The following part of the python code handle all the parsing of the JSON files together with some pre-processing and feature extraction.

```
1 import json as js
2 import pandas as pd
3 from utilities import *
4
5
6 #TAGs
7 tags2name='CSV/tags2name.csv'
8
9
10 with open(tags2name) as train_file:
11     tags2name_df = pd.read_csv(train_file)
12
13
14 tags2name_df.drop("Description", axis=1, inplace=True)
15
16 tags2name_df_pr=tags2name_df.loc[(tags2name_df['Label']=='assist')
17 |
18     (tags2name_df['Label']=='keyPass') |
19     (tags2name_df['Label']=='interception') |
20     (tags2name_df['Label']=='opportunity') |
21     (tags2name_df['Label']=='red_card') |
22     (tags2name_df['Label']=='yellow_card') |
```

```

    (tags2name_df['Label']=='second_yellow_card')|
19 (tags2name_df['Label']=='lost') |
    (tags2name_df['Label']=='neutral') |
    (tags2name_df['Label']=='won') |
20 (tags2name_df['Label']=='accurate') |
    (tags2name_df['Label']=='not accurate') |
    (tags2name_df['Label']=='Feint') |
21 (tags2name_df['Label']=='dangerous_ball_lost')|
22 (tags2name_df['Label']=='counter_attack']]
23
24
25 tags2name_dict_2=dict(tags2name_df.values)
26
27
28 tags2name_dict_PR=dict(tags2name_df_pr.values)
29
30 #TEAMS
31 teams = 'teams.json'
32
33 with open(teams) as train_file:
34     teams_df = pd.read_json(train_file)
35
36
37
38 teams_dict=dict(teams_df.iloc[:,[2,1]].values)
39
40
41 teams_df.to_csv('/Users/martino/Desktop/dataset
    tesi/CSV/teams.csv')
42
43
44 #PLAYERS
45 players = 'players.json'
46
47 with open(players) as train_file:
48     players_df = pd.read_json(train_file)
49
50 players_dict=dict_players(players_df)
51
52 players_df.to_csv('/Users/martino/Desktop/dataset
    tesi/CSV/players.csv')
53
54
```

```

55 #Next thing to do is to parse the json file of the Events.
56 #that will be the core of the analysis
57
58 #the following functions from the utilities will be applied :
59     #cordinates: to extract the x y cortidates from the positions
        column
60     #get_names_players: to extract the full name of every player
61     #get_names_teams: to extract the name of every name
62     #get_tags: to expand the tags column since every events could
        have more than one tag
63     #one hot encode: for every tag so to have a matrix (0,1)
        tag,event
64
65 #Insted for the DF used to dit the models a combination of event
        name and tag is used to for the matrix
66     #filter out some of the useless avents
67     #get_tags_for_player_rank: form the matrix
68
69
70 #-----Italy-----
71 events_Italy='events/events_Italy.json'
72
73
74 with open(events_Italy) as train_file:
75     events_Italy_df = pd.read_json(train_file)
76
77
78 cordinates(events_Italy_df)
79
80 get_names_players(events_Italy_df,players_dict)
81
82 get_names_teams(events_Italy_df,teams_dict)
83
84 events_Italy_df_DV=get_tags(events_Italy_df,tags2name_dict_2)
85
86 events_Italy_df_DV = pd.DataFrame(events_Italy_df_DV,columns
        =['id','tags'])
87 events_Italy_df_DV= pd.crosstab(events_Italy_df_DV['id'],
        events_Italy_df_DV['tags'])
88 events_Italy_df_DV=events_Italy_df_DV.reset_index()
89 events_Italy_df_DV = pd.merge(events_Italy_df, events_Italy_df_DV,
        how='inner')
90 events_Italy_df_DV.to_csv('/Users/martino/Desktop/dataset

```

```

tesì/CSV/CSV_events/events_Italy_DV.csv')
91
92 events_Italy_df=
    events_Italy_df.loc[(events_Italy_df['eventName']=='Duel') |
93         (events_Italy_df['eventName']=='Foul') |
94         (events_Italy_df['eventName']=='Free Kick') |
95         (events_Italy_df['eventName']=='Others on the ball')|
96         (events_Italy_df['eventName']=='Pass')|
97         (events_Italy_df['eventName']=='Shot')]
98
99 events_Italy_df=
    events_Italy_df.loc[(events_Italy_df['subEventName']=='Ground
    attacking duel') |
100    (events_Italy_df['subEventName']=='Ground defending duel')|
101    (events_Italy_df['subEventName']=='Ground loose ball duel')|
102    (events_Italy_df['subEventName']=='Hand foul')|
103    (events_Italy_df['subEventName']=='Late card foul')|
104    (events_Italy_df['subEventName']=='Out of game foul')|
105    (events_Italy_df['subEventName']=='Protest')|
106    (events_Italy_df['subEventName']=='Simulation')|
107    (events_Italy_df['subEventName']=='Violent Foul')|
108    (events_Italy_df['subEventName']=='Free kick cross')|
109    (events_Italy_df['subEventName']=='Penalty')|
110    (events_Italy_df['subEventName']=='Free kick shot')|
111    (events_Italy_df['subEventName']=='Throw in')|
112    (events_Italy_df['subEventName']=='Acceleration')|
113    (events_Italy_df['subEventName']=='Clearance')|
114    (events_Italy_df['subEventName']=='Touch')|
115    (events_Italy_df['subEventName']=='Cross')|
116    (events_Italy_df['subEventName']=='Hand pass')|
117    (events_Italy_df['subEventName']=='Head pass')|
118    (events_Italy_df['subEventName']=='High pass')|
119    (events_Italy_df['subEventName']=='Launch')|
120    (events_Italy_df['subEventName']=='Simple pass')|
121    (events_Italy_df['subEventName']=='Smart pass')|
122    (events_Italy_df['subEventName']=='Shot')|
123    (events_Italy_df['subEventName']=='Air duel')|
124    (events_Italy_df['subEventName']=='Foul')|
125    (events_Italy_df['subEventName']=='Free Kick')|
126    (events_Italy_df['subEventName']=='Corner')]
127
128 events_Italy_rank_player_df=get_tags_for_player_rank(events_Italy_df
129 ,tags2name_dict_PR)

```

```
130
131 events_Italy_rank_player_df =
    pd.DataFrame(events_Italy_rank_player_df, columns = ['id', 'tags'])
132
133 events_Italy_rank_player_df=
    pd.crosstab(events_Italy_rank_player_df['id'],
    events_Italy_rank_player_df['tags'])
134
135 events_Italy_rank_player_df=events_Italy_rank_player_df.reset_index()
136
137 events_Italy_rank_player_df = pd.merge(events_Italy_df,
    events_Italy_rank_player_df, how='inner')
138
139 events_Italy_rank_player_df.to_csv('/Users/martino/Desktop/dataset
    tesi/CSV/CSV_events/events_Italy_PR.csv')
140
141
142 matches_Italy='matches/matches_Italy.json'
143
144
145 with open(matches_Italy) as train_file:
146     matches_Italy_df = pd.read_json(train_file)
147
148 get_info_team(matches_Italy_df)
149
150 matches_Italy_df.to_csv('/Users/martino/Desktop/dataset
    tesi/CSV/CSV_matches/matches_Italy.csv')
```

A.2 R Code

A.2.1 Data cleaning and Visualization

The following code regards every aspects of the cleaning and visualization of the data

```
1 require(tidyverse)
2 require(dplyr)
3 library(ggplot2)
4 library(ggsoccer)
5 library(raster)
6 require(ggpubr)
7 require(forcats)
```

```

8 library(soccermatics)
9
10 #import the data
11 events_italy=read.csv(file='CSV/CSV_events/events_Italy.csv',
12                       header=T,
13                       sep=',',
14                       stringsAsFactors=T,
15                       dec='.')
16
17 #take a look at the data
18 str(events_italy)
19 summary(events_italy)
20
21 #-----Data cleaning-----
22 #eventSec column
23 #delete useless columns (positions,X,tags,NA.)
24 events_italy=events_italy%>%dplyr::select(-c(positions,X,tags,NA.))
25
26
27 #convert sec into minutes taking into account the division of
    first and second half
28
29 secondH <- events_italy%>% dplyr:: filter (matchPeriod=='2H')
30 secondH$eventMin <- ((secondH$eventSec)/60)+45
31
32 #discretization of 'eventSec' variable for visualization
33
34 secondH$timeframe <- cut(secondH$eventMin,
35                          breaks =
36                            c(45,50,55,60,65,70,75,80,85,90,120),
37                            labels=c('45-50', '50-55', '55-60', '60-65',
38                                      '65-70', '70-75', '75-80', '80-85',
39                                      '85-90', '>90'))
40
41 firstH <- events_italy%>% dplyr:: filter (matchPeriod=='1H')
42 firstH$eventMin <- ((firstH$eventSec)/60)
43
44 firstH$timeframe <- cut(firstH$eventMin,
45                          breaks = c(0,5, 10, 15, 20, 25
46                                      ,30,35,40,46,60),
47                          labels=c('0-5', '5-10', '10-15', '15-20',
48                                    '20-25', '25-30',
49                                    '30-35', '35-40', '40-45',
50                                    '>45'))

```



```

44
45 events_italy <- rbind(firstH, secondH)
46
47 #make 2 copies for later use
48 events_italy_PR <- events_italy
49
50 events_italy_RD <- events_italy
51
52 #-----Data visualization-----
53 events_italy_DV=read.csv(file='CSV/CSV_events/events_Italy_DV.csv',
54                          header=T,
55                          sep=',',
56                          stringsAsFactors=T,
57                          dec='.')
58 #eventSec column
59 #delete useless columns (positions,X,tags,NA.)
60 events_italy_DV=events_italy_DV%>%dplyr::select(-c(positions,X,tags,NA.))
61
62 #convert sec into minutes taking into account the division of
63   first and second half
64 secondH <- events_italy_DV%>% dplyr:: filter (matchPeriod=='2H')
65 secondH$eventMin <- ((secondH$eventSec)/60)+45
66
67 #discretization of 'eventSec' variable for visualization
68
69 secondH$timeframe <- cut(secondH$eventMin,
70                          breaks =
71                            c(45,50,55,60,65,70,75,80,85,90,120),
72                            labels=c('45-50', '50-55', '55-60', '60-65',
73                                      '65-70', '70-75', '75-80', '80-85',
74                                      '85-90', '>90')
75 )
76
77 firstH <- events_italy_DV%>% dplyr:: filter (matchPeriod=='1H')
78 firstH$eventMin <- ((firstH$eventSec)/60)
79
80 firstH$timeframe <- cut(firstH$eventMin,
81                          breaks = c(0,5, 10, 15, 20, 25
82                                      ,30,35,40,46,60),
83                          labels=c('0-5', '5-10', '10-15', '15-20',
84                                      '20-25', '25-30', '30-35', '35-40', '40-45',
85                                      '>45')
86 )

```

```
81
82 events_italy_DV <- rbind(firstH, secondH)
83 #bar chart for frequency along the game of yellow and red cards
84   together with goals
85
86 library(RColorBrewer)
87
88 mycolors <- colorRampPalette(brewer.pal(8, "Blues"))(20)
89
90 red_card_plot <- ggplot(events_italy_DV, aes(x=timeframe ,
91   y=red_card, fill=matchPeriod ))+
92   geom_bar(stat = 'identity')+
93   labs(
94     x = "Match Interval", y = "",
95     tag = "Fig.1",
96     color = "Match Period")+
97   scale_fill_manual("Legend", values = c("1H" = "dodgerblue4", "2H"
98     = "deepskyblue"))+
99   theme(plot.title = element_text(size=18))
100
101 yellow_card_plot <- ggplot(events_italy_DV, aes(x=timeframe,
102   y=yellow_card, fill=matchPeriod))+
103   geom_bar(stat = 'identity')+
104   labs(
105     x = "Match Interval", y = "",
106     tag = "Fig.2",
107     color = "Match Period")+
108   scale_fill_manual("Legend", values = c("1H" = "dodgerblue4", "2H"
109     = "deepskyblue"))+
110   theme(plot.title = element_text(size=18))
111
112 Goal_plot <- ggplot(events_italy_DV, aes(x=timeframe, y=Goal,
113   fill=matchPeriod))+
114   geom_bar(stat = 'identity')+
115   labs(
116     x = "Match Interval", y = "",
117     tag = "Fig.3",
118     color = "Match Period")+
119   scale_fill_manual("Legend", values = c("1H" = "dodgerblue4", "2H"
120     = "deepskyblue"))+
121   theme(plot.title = element_text(size=18))
```

```

117 figure2 <- ggarrange(red_card_plot,yellow_card_plot,Goal_plot,
118   ncol = 1, nrow = 3)
118
119
120 #Bar chart for event frequency
121
122 events_italy_DV %>%
123   count(eventName) %>%
124   mutate(perc = n / nrow(events_italy_DV)) -> tips2
125
126 tips2$eventName <- factor(tips2$eventName,
127   levels = c('Goalkeeper leaving
128     line','Offside','Save attempt','Shot',
129     'Foul','Interruption','Free
130     Kick','Others on the
131     ball','Duel','Pass'))
132
133 levels(tips2$eventName)
134 levels(events_italy_DV$eventName)
135
136 ggplot(tips2, aes(x = eventName, y = perc, fill=eventName)) +
137   geom_bar(stat = "identity")+coord_flip()+
138   scale_fill_manual(values = mycolors)+
139   theme(panel.background = element_rect(fill =
140     "white"),legend.position="none",plot.title =
141     element_text(size=18))+
142   labs(x = "Event Name", y = "(%)")
143
144 # Pitch heatmap for type of event
145
146 pass_heatmap<- ggplot() +
147   annotate_pitch( fill= 'white', colour = "black") +
148   theme_pitch() +
149   theme(panel.background = element_rect(fill =
150     "White"),legend.position="none")+
151   geom_density_2d_filled(data =
152     events_italy_DV%>%dplyr::filter(eventName=='Pass'),
153   aes(x = x_start, y = y_start ),alpha=0.8
154     )+scale_fill_manual(values = mycolors)+
155   labs(title = "Pass",
156     x = "", y = "",
157     tag = "(a)")+
158   theme(plot.title = element_text(size = 14,hjust = 0.5))

```

```

151
152 shot_heatmap<- ggplot() +
153   annotate_pitch( fill= 'white', colour = "black") +
154   theme_pitch() +
155   theme(panel.background = element_rect(fill =
156     "White"),legend.position="none")+
157   geom_density_2d_filled(data =
158     events_italy_DV%>%dplyr::filter(eventName=='Shot'),
159     aes(x = x_start, y = y_start ),alpha=0.6 )+
160   scale_fill_brewer(palette="Blues")+
161   labs(title = "Shots",
162     x = "", y = "",
163     tag = "(b)")+
164   theme(plot.title = element_text(size = 14,hjust = 0.5))
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999

```

```

188
189
190 figure <- ggarrange(pass_heatmap, shot_heatmap, duel_heatmap,
191   fowl_heatmap,
192   ncol = 2, nrow = 2)
193 #focus on a single game
194
195 focus_DV <- events_italy_DV%>%dplyr::filter(matchId==2576300)
196
197 focus_goal <- focus_DV%>%dplyr::filter(Goal==1 & eventName=='Shot')
198
199 #Goals
200 focus_goal_plot<- ggplot() +
201   annotate_pitch( fill= 'dodgerblue4', colour = "white") +
202   theme_pitch() +
203   geom_point(data = focus_goal,aes(x = x_start, y = y_start, fill =
204     t_name),
205     shape = 21,
206     size = 3)+
207   labs(
208     x = "", y = "")+
209   theme(plot.title = element_text(size = 14,hjust = 0.5),
210     panel.background = element_rect(fill = "White"))+
211   scale_fill_manual("Team", values = c( "Roma" =
212     "orange","Chievo"="yellow"))
213
214 focus_shot_r <- focus_DV%>%dplyr::filter(eventName=='Shot' &
215   t_name=='Roma')
216
217 #shots Roma
218 focus_shot_plot_r<- ggplot() +
219   annotate_pitch( fill= 'dodgerblue4', colour = "white") +
220   theme_pitch() +
221   geom_point(data = focus_shot_r,aes(x = x_start, y = y_start, fill
222     = t_name),shape = 21,size = 3)+
223   labs( x = "", y = "",
224     tag = "(a)")+
225   theme(plot.title = element_text(size = 14,hjust = 0.5),
226     panel.background = element_rect(fill = "White"))+
227   scale_fill_manual("Team", values = c( "Roma" = "orange"))
228
229 focus_shot_c <- focus_DV%>%dplyr::filter(eventName=='Shot' &

```

```

    t_name=='Chievo')
226
227 #shots Chievo
228 focus_shot_plot_c<- ggplot() +
229   annotate_pitch( fill= 'dodgerblue4', colour = "white") +
230   theme_pitch() +
231   geom_point(data = focus_shot_c,aes(x = x_start, y = y_start, fill
    = t_name),
232             shape = 21,
233             size = 3)+
234   labs(
235     x = "", y = "",
236     tag = "(b)")+
237   theme(plot.title = element_text(size = 14,hjust = 0.5),
238         panel.background = element_rect(fill = "White"))+
239   scale_fill_manual("Team", values = c( "Chievo" = "yellow"))
240
241 figure3 <- ggarrange(focus_shot_plot_r, focus_shot_plot_c,
242                     ncol = 2, nrow = 1)
243
244 focus_chievo <- events_italy_DV%>%dplyr::filter(teamId==3165 &
    Goal==1 & eventName!='Save attempt')
245
246
247 #chievo heat map shots
248 mycolors2 <- colorRampPalette(brewer.pal(8, "Blues"))(12)
249
250 focus_goal_chievo_heatmap<- ggplot() +
251   annotate_pitch( fill= 'white', colour = "black") +
252   theme_pitch() +
253   scale_fill_manual(values = mycolors2)+
254   theme(panel.background = element_rect(fill =
    "White"),legend.position="none")+
255   geom_density_2d_filled(data =
    events_italy_DV%>%dplyr::filter(teamId==3165 & Goal==1&
    eventName!='Save attempt'),
256                          aes(x = x_start, y = y_start ),alpha=0.8 )+
257   labs( x = "", y = "",
258         tag = "(a)")+
259   theme(plot.title = element_text(size = 14,hjust = 0.5))
260
261 mycolors3 <- colorRampPalette(brewer.pal(8, "Greens"))(11)
262

```

```

263 #kolarov heatmap
264 focus_kolarov_heatmap<- ggplot() +
265   annotate_pitch( fill= 'white', colour = "black") +
266   theme_pitch() +
267   scale_fill_manual(values = mycolors3)+
268   theme(panel.background = element_rect(fill =
269     "White"),legend.position="none")+
270   geom_density_2d_filled(data =
271     events_italy_DV%>%dplyr::filter(p_name=='Aleksandar Kolarov'),
272     aes(x = x_start, y = y_start ),alpha=0.8 )+
273   labs(
274     x = "", y = "",
275     tag = "(b)")+
276   theme(plot.title = element_text(size = 14,hjust = 0.5))
277
278 #pass chievo heatmap
279 focus_pass_chievo_heatmap<- ggplot() +
280   annotate_pitch( fill= 'white', colour = "black") +
281   theme_pitch() +
282   scale_fill_manual(values = mycolors2)+
283   theme(panel.background = element_rect(fill =
284     "White"),legend.position="none")+
285   geom_density_2d_filled(data =
286     events_italy_DV%>%dplyr::filter(teamId==3165 &
287     eventName!='pass'), aes(x = x_start, y = y_start ),alpha=0.8
288     )+
289   labs( x = "", y = "",
290     tag = "(c)")+
291   theme(plot.title = element_text(size = 14,hjust = 0.5))
292
293 focus_pass_Roma_heatmap<- ggplot() +
294   annotate_pitch( fill= 'white', colour = "black") +
295   theme_pitch() +
296   scale_fill_manual(values = mycolors3)+
297   theme(panel.background = element_rect(fill =
298     "White"),legend.position="none")+
299   geom_density_2d_filled(data =
300     events_italy_DV%>%dplyr::filter(teamId==3158 &
301     eventName!='pass'),
302     aes(x = x_start, y = y_start ),alpha=0.8 )+
303   labs(
304     x = "", y = "",
305     tag = "(d)")+

```

```
297   theme(plot.title = element_text(size = 14,hjust = 0.5))
298
299
300 figure5 <- ggarrange(focus_goal_chievo_heatmap,
301   focus_kolarov_heatmap, focus_pass_chievo_heatmap,
302   focus_pass_Roma_heatmap,
303   ncol = 2, nrow = 2)
304
305 #Network analysis
306
307 match_italy_NA=read.csv(file='CSV/CSV_matches/matches_Italy.csv',
308   header=T,
309   sep=',',
310   stringsAsFactors=T,
311   dec='.')
312
313 events_italy_NA <-
314   focus_DV%>%dplyr::select(c(eventName,playerId,matchId,
315   teamId,eventMin,x_start,y_start,x_end,y_end))
316
317 events_italy_NA_try <-
318   events_italy_NA%>%dplyr::filter(eventName=='Pass')
319
320 #build the table
321
322 i=1
323 r=2
324 possesso <- data.frame(NA_col = rep(NA, nrow(events_italy_NA_try)))
325
326 for(elm in events_italy_NA_try$teamId){
327
328   if(r<nrow(events_italy_NA_try)+1){
329     if(elm==events_italy_NA_try$teamId[r]){
330
331       possesso[ i,] <- NA
332
333     }
334     else{
335       possesso[ i, ] <- 'interrupt'
336     }
337   }
338   i <- i+1
339 }
```



```

336   r <- r+1
337   }
338 }
339
340 events_italy_NA_try$possesso <- possesso$NA_col
341
342 #find the receiver
343 i=1
344 r=2
345 ricevitori <- data.frame(NA_col = rep(NA,
    nrow(events_italy_NA_try)))
346
347 for (elm in events_italy_NA_try$possesso ) {
348   if( is.na(elm) ){
349
350     ricevitori[ i,] <- events_italy_NA_try$playerId[r]
351   }
352   else{
353
354     ricevitori[ i,] <- NA
355   }
356   i <- i+1
357   r <- r+1
358 }
359
360 events_italy_NA_try$receiver <- ricevitori$NA_col
361
362 #find the weights
363
364 events_italy_NA_2 <- events_italy_NA_try%>%
365 dplyr::select(c(playerId,receiver,teamId))%>%
366 drop_na()%>%
367   dplyr::filter(teamId==3158)
368
369 events_italy_NA_2 <-data.frame(
    events_italy_NA_2%>%dplyr::group_by(playerId,receiver)%>%
370 dplyr::summarise(weights=n())%>%
371   rename(from=playerId,to=receiver))
372
373 #let's try with a match==2575959
374 edges
    <-events_italy_NA_2%>%dplyr::filter(!(from=='340019' |from=='25405' |
375 from=='92966' |to=='340019' |to=='25405' |to=='92966'))

```

```

376
377 edges <- edges %>%
378   filter(!(from == to))
379
380 edges$weights <- (edges$weights - min(edges$weights)) /
   (max(edges$weights) - min(edges$weights))
381
382 str(edges)
383
384 nodes <- data.frame(events_italy_NA_2%>%
385   dplyr::group_by(from)%>%
386   dplyr::summarise(x=n()))%>%
387   dplyr::select(-c(x))%>%dplyr::rename(id=from)%>%
388   filter(!(id=='340019'|id=='25405'|id=='92966'))
389
390 players_NA=read.csv(file='CSV/players.csv',
391                    header=T,
392                    sep=',',
393                    stringsAsFactors=T,
394                    dec='.')
395
396 players_NA <- players_NA%>%dplyr::select(wyId,
   shortName)%>%rename(id=wyId)
397
398
399 nodes_names<-merge(x=nodes,y=players_NA,by="id",all.x=TRUE)
400
401 #build the network
402
403 library(network)
404
405 routes_network <- network(edges, vertex.attr = nodes,
406   matrix.type = "edgelist", ignore.eval = FALSE, loops = TRUE)
407
408 library(GGally)
409 library(network)
410 library(sna)
411 library(ggplot2)
412
413 position_NA <-
   focus_DV%>%dplyr::select(c(playerId,x_start,y_start))%>%
414   dplyr::group_by(playerId)%>%
415   dplyr::summarise(x_mean=mean(x_start),y_mean=mean(y_start))%>%

```

```

416   dplyr::rename(id=playerId)
417
418   nodes_names_pos<-merge(x=nodes_names,y=position_NA,by="id",all.x=TRUE)
419
420   formation_NA = as.matrix(data.frame(x = nodes_names_pos$x_mean , y
    = nodes_names_pos$y_mean ))
421
422   roma_network<- ggnet2(routes_network,mode=formation_NA ,size =
    12,edge.size = "weights", label =nodes_names$shortName,color =
    "cornsilk3",
423     label.alpha = 0.8,label.color =
    "black",edge.color='cornsilk4')+
424   theme(panel.background = element_rect(fill = "white"))
425
426   #Distributions variables
427   goals <- events_italy_DV %>% dplyr::filter(!(eventName=='Save
    attempt'))%>%
428     dplyr::filter(Goal==1)%>%dplyr::group_by(matchId)%>%dplyr::summarise(n=n())
429
430   hist_goals <- ggplot(goals,aes(x=n))+
431     geom_histogram(binwidth=1,fill='dodgerblue4')+
432     geom_vline(data=goals, aes(xintercept=mean(goals$n), color="red"),
433       linetype="dashed")+
434     theme(legend.position="none")+
435     labs(tag = "(a)")
436
437
438   box_plot_goals <-
    ggplot(goals,aes(y=n))+geom_boxplot(outlier.colour="black",
    outlier.shape=16,outlier.size=2,notch=FALSE,
    fill='cornflowerblue')+
439   labs(tag = "(b)")
440
441   pass_hist <- events_italy_DV%>%dplyr::filter(eventName=='Pass'&
    accurate==1)%>%
442     dplyr::group_by(matchId)%>%dplyr::summarise(n=n())
443
444   hist_pass <- ggplot(pass_hist,aes(x=n))+
445     geom_histogram(aes(y=..density..), colour="dodgerblue4",
    fill="dodgerblue4")+
446     geom_density(alpha=.1, fill="orange") +
447     geom_vline(data=goals, aes(xintercept=mean(pass_hist$n),
    color="red"),

```

```

448         linetype="dashed")+
449     theme(legend.position="none")+
450     labs(tag = "(c)")
451
452 pass_hist_not_acc <-
453     events_italy_DV%>%dplyr::filter(eventName=='Pass'&
454     not.accurate==1)%>%
455     dplyr::group_by(matchId)%>%dplyr::summarise(n=n())
456
457 hist_pass_not_acc <- ggplot(pass_hist_not_acc,aes(x=n))+
458     geom_histogram(aes(y=..density..), colour="cornflowerblue",
459     fill="cornflowerblue")+
460     geom_density(alpha=.1, fill="red") +
461     geom_vline(data=goals, aes(xintercept=mean(pass_hist_not_acc$n),
462     color="orange"),
463     linetype="dashed")+
464     theme(legend.position="none")+
465     labs(tag = "(d)")
466
467 figure <- ggarrange(hist_goals, box_plot_goals, hist_pass,
468     hist_pass_not_acc,
469     ncol = 2, nrow = 2)
470
471 mean(pass_hist_not_acc$n)

```

A.2.2 Supervised Model(SVM)

The following section shows the code of modelling phase regarding the classification task.

```

1
2
3 #take a look at the variables and choose which one to exclude
4 str(events_italy_PR)
5 summary(events_italy_PR)
6
7 #exclude the useless variable for this analysis (we only need the
8   IDs of plays,
9 #teams and match together with the events
10 events_italy_PR=events_italy_PR%>%dplyr::select(-c(eventId,subEventName,
11     eventName, matchPeriod,
12     eventSec,subEventId,id,x_start,y_start,x_end,y_end,timeframe))

```

```

10
11
12 #group by matchId and playerID in order to get the sum of each
    event for every player in every match
13 events_italy_PR_players=events_italy_PR%>%dplyr::select(-c(p_name,t_name))
14
15 events_italy_PR_players <- events_italy_PR_players%>%
    group_by(matchId,playerId,teamId)%>%summarise_all(list(sum))
16
17 #by analysing the staitstics of each column we can exclude the some
    of them which are useless
18 #summary(events_italy_PR_players)
19 #str(events_italy_PR_players)
20
21 events_italy_PR_players=events_italy_PR_players%>%
22   dplyr::select(c(matchId,playerId,teamId,
23   Duel.Air.duel.accurate,Duel.Air.duel.not.accurate,
24   Duel.Ground.attacking.duel.accurate,
25   Duel.Ground.attacking.duel.not.accurate,
26   Duel.Ground.defending.duel.accurate,
27   Duel.Ground.defending.duel.not.accurate,
28   Duel.Ground.loose.ball.duel.accurate,
29   Duel.Ground.loose.ball.duel.not.accurate,
30   Foul.Hand.foul.red_card,Foul.Hand.foul.yellow_card,
31   Foul.Late.card.foul.yellow_card,Foul.Foul.red_card,
    Foul.Foul.second_yellow_card,
32   Foul.Foul.yellow_card,Foul.Out.of.game.foul.red_card ,
    Foul.Out.of.game.foul.second_yellow_card,
33   Foul.Out.of.game.foul.yellow_card,Foul.Protest.red_card
    ,Foul.Protest.second_yellow_card,
34   Foul.Protest.yellow_card,
    Foul.Simulation.yellow_card,Foul.Violent.Foul.red_card,
35   Foul.Violent.Foul.yellow_card,Free.Kick.Free.kick.cross.accurate,
    Free.Kick.Free.kick.cross.not.accurate,
36   Free.Kick.Free.Kick.accurate,Free.Kick.Free.Kick.not.accurate,
37   Free.Kick.Penalty.not.accurate,
38   Free.Kick.Free.kick.shot.not.accurate,Free.Kick.Free.kick.shot.accurate,Free.Kick.T
39   Free.Kick.Throw.in.not.accurate,
40   Others.on.the.ball.Acceleration.accurate,
41   Others.on.the.ball.Acceleration.not.accurate,
    Others.on.the.ball.Clearance.accurate,
42   Others.on.the.ball.Clearance.not.accurate,
43   Others.on.the.ball.Touch.assist,

```

```

44 Others.on.the.ball.Touch.counter_attack,
45 Others.on.the.ball.Touch.dangerous_ball_lost,
46 Others.on.the.ball.Touch.interception,
47 Others.on.the.ball.Touch.opportunity,Pass.Cross.accurate,
48 Pass.Cross.assist,Pass.Cross.keyPass,
    Pass.Cross.not.accurate,Pass.Hand.pass.accurate,
49 Pass.Hand.pass.not.accurate,
50 Pass.Head.pass.accurate,Pass.Head.pass.assist,
51 Pass.Head.pass.keyPass ,Pass.Head.pass.not.accurate,
52 Pass.High.pass.accurate,Pass.High.pass.assist,
53 Pass.High.pass.keyPass,Pass.High.pass.not.accurate,
54 Pass.Launch.accurate,Pass.Launch.keyPass,Pass.Launch.not.accurate,
55 Pass.Simple.pass.accurate,
56 Pass.Simple.pass.keyPass,
57 Pass.Simple.pass.not.accurate,Pass.Smart.pass.accurate,
    Pass.Smart.pass.assist,
58 Pass.Smart.pass.keyPass, Pass.Smart.pass.not.accurate,
59 Shot.Shot.accurate,Shot.Shot.not.accurate))
60
61 col_names <- colnames(events_italy_PR_players )
62 #same thing for the team dataset
63 #group by matchId and teamID in order to get the sum of each event
    for every team in every match
64 events_italy_PR_teams=events_italy_PR%>%dplyr::select(-c(playerId,p_name,t_name))
65
66 events_italy_PR_teams <-
    events_italy_PR_teams%>%group_by(matchId,teamId)%>%summarise_all(list(sum))
67
68 #get the binary variable for winner
69 match_italy_PR=read.csv(file='CSV/CSV_matches/matches_Italy.csv',
70                          header=T,
71                          sep=',',
72                          stringsAsFactors=T,
73                          dec='.')
74
75
76 match_italy_PR <- match_italy_PR%>%dplyr:: select(c(wyId ,
    winner))%>%dplyr:: rename(matchId=wyId)
77
78
79 events_italy_PR_teams <-
    merge(events_italy_PR_teams,match_italy_PR, by='matchId'
    ,all.x=TRUE)

```

```

80
81 events_italy_PR_teams$outcome <-
      ifelse(events_italy_PR_teams$winner==events_italy_PR_teams$teamId,1,0)
82
83 #by analysing the staitistics of each column we can exclude the some
      of them which are useless
84 events_italy_PR_teams=events_italy_PR_teams%>%
85   dplyr::select(c(matchId,teamId,outcome,Duel.Air.duel.accurate,
86   Duel.Air.duel.not.accurate,Duel.Ground.attacking.duel.accurate,
87   Duel.Ground.attacking.duel.not.accurate,
88   Duel.Ground.defending.duel.accurate,
89   Duel.Ground.defending.duel.not.accurate,
90   Duel.Ground.loose.ball.duel.accurate,
91   Duel.Ground.loose.ball.duel.not.accurate,
92   Foul.Hand.foul.red_card,
93   Foul.Hand.foul.yellow_card,
94   Foul.Late.card.foul.yellow_card,Foul.Foul.red_card,
      Foul.Foul.second_yellow_card,
95   Foul.Foul.yellow_card,Foul.Out.of.game.foul.red_card ,
      Foul.Out.of.game.foul.second_yellow_card,
96   Foul.Out.of.game.foul.yellow_card,Foul.Protest.red_card
      ,Foul.Protest.second_yellow_card,
97   Foul.Protest.yellow_card,
      Foul.Simulation.yellow_card,Foul.Violent.Foul.red_card,
98   Foul.Violent.Foul.yellow_card,Free.Kick.Free.kick.cross.accurate,
      Free.Kick.Free.kick.cross.not.accurate,
99   Free.Kick.Free.Kick.accurate,Free.Kick.Free.Kick.not.accurate,
100   Free.Kick.Penalty.not.accurate,Free.Kick.Free.kick.shot.not.accurate,
101   Free.Kick.Free.kick.shot.accurate,
102   Free.Kick.Throw.in.accurate,Free.Kick.Throw.in.not.accurate,
103   Others.on.the.ball.Acceleration.accurate,
104   Others.on.the.ball.Acceleration.not.accurate,
      Others.on.the.ball.Clearance.accurate,
105   Others.on.the.ball.Clearance.not.accurate,
106   Others.on.the.ball.Touch.assist,Others.on.the.ball.Touch.counter_attack,
107   Others.on.the.ball.Touch.dangerous_ball_lost,
108   Others.on.the.ball.Touch.interception,
109   Others.on.the.ball.Touch.opportunity,Pass.Cross.accurate,
110   Pass.Cross.assist,Pass.Cross.keyPass,
      Pass.Cross.not.accurate,Pass.Hand.pass.accurate,
111   Pass.Hand.pass.not.accurate,
112   Pass.Head.pass.accurate,Pass.Head.pass.assist,
113   Pass.Head.pass.keyPass ,Pass.Head.pass.not.accurate,

```

```

114   Pass.High.pass.accurate,Pass.High.pass.assist,
115   Pass.High.pass.keyPass,Pass.High.pass.not.accurate,
116   Pass.Launch.accurate,Pass.Launch.keyPass,
117   Pass.Launch.not.accurate,Pass.Simple.pass.accurate,
118   Pass.Simple.pass.keyPass,
119   Pass.Simple.pass.not.accurate,Pass.Smart.pass.accurate,
        Pass.Smart.pass.assist,
120   Pass.Smart.pass.keyPass, Pass.Smart.pass.not.accurate,
121   Shot.Shot.accurate,Shot.Shot.not.accurate))
122
123 #delete players not present in the player_df
124
125 players_PR=read.csv(file='CSV/players.csv',
126                    header=T,
127                    sep=',',
128                    stringsAsFactors=T,
129                    dec='.')
130
131 players_PR_prov <-
        players_PR%>%dplyr::select(c(firstName,lastName,wyId))%>%rename(playerId=wyId)
132
133 players_PR <- players_PR%>%dplyr::select(c(wyId))%>%dplyr::
        rename(playerId=wyId)
134
135 events_italy_PR_players <-
        merge(events_italy_PR_players,players_PR, by='playerId' )
136
137 #-----Feature weights extraction----
138 #standardize the data
139
140 df_standard <-
        events_italy_PR_teams%>%dplyr::select(-c(teamId,matchId,outcome))
141
142 df_standard <- data.frame(scale(df_standard))
143
144 df_standard$outcome <- as.factor(events_italy_PR_teams$outcome)
145
146 #----- Fit the SVM linear classifier-----
147
148 #train and test set
149
150 library(caTools)
151

```



```
152 set.seed(123)
153 split = sample.split(df_standard$outcome, SplitRatio = 0.75)
154
155 df_standard <- df_standard %>% mutate(outcome =
      factor(outcome, labels = make.names(levels(outcome))))
156
157 training_set = subset(df_standard, split == TRUE)
158 test_set = subset(df_standard, split == FALSE)
159
160 library(caret)
161
162 # Define fitControl
163 fitControl <- trainControl(
164   method = "cv",
165   number = 5,
166   classProbs = TRUE,
167   summaryFunction = twoClassSummary
168 )
169
170 # set random seed and run the model
171
172 set.seed(321)
173 svmFit1 <- train(x = training_set[-68] ,
174                 y=training_set$outcome,
175                 method = "svmLinear",
176                 trControl = fitControl,
177                 preProc = c("center","scale"),
178                 metric="ROC" )
179
180 y_pred_2 = predict(svmFit1, newdata = test_set[-68])
181
182 confusionMatrix(y_pred_2, test_set[, 68])
183 #Sensitivity : 0.9224
184 #Accuracy : 0.8474
185 #roc:0.84
186
187 library("MLmetrics")
188
189 f1_svm<- F1_Score(y_pred_2, test_set[, 68], positive =
      NULL)#0.8806584
190
191 #get weights
192
```

```
193 coefs <- svmFit1$finalModel@coef[[1]]
194
195 mat <- svmFit1$finalModel@xmatrix[[1]]
196
197 weights2 <- t(coefs %*% mat)
198
199 col_names_weights <- colnames(training_set)[-68]
200
201 #visualize weights
202
203 #positive weights
204 weights_dv <- data.frame(col_names_weights,weights2)
205
206 weights_dv <-data.frame(weights_dv[order(-weights2),])[1:8,]
207
208 weights_pos<- ggplot(weights_dv, aes(x = col_names_weights, y =
  weights2, fill=col_names_weights))+
209   scale_fill_brewer(palette="Blues") +
210   geom_bar(stat = "identity")+coord_flip()+
211   theme(panel.background = element_rect(fill =
  "white"),legend.position="none",plot.title =
  element_text(size=18))+
212   labs(x = "Event Name", y = "(%)")
213
214 #negative weight
215 weights_dv_2 <- data.frame(col_names_weights,weights2)
216
217 weights_dv_2 <-data.frame(weights_dv_2[order(weights2),])[1:8,]
218
219 weights_neg <- ggplot(weights_dv_2, aes(x = col_names_weights, y =
  weights2, fill=col_names_weights))+
220   scale_fill_brewer(palette="Blues") +
221   geom_bar(stat = "identity")+coord_flip()+
222   theme(panel.background = element_rect(fill =
  "white"),legend.position="none",plot.title =
  element_text(size=18))+
223   labs(x = "Event Name", y = "(%)")
224
225 figure6 <- ggarrange(weights_pos,weights_neg, ncol = 2, nrow = 1)
226
227
228 #---- Player Rating Phase -----
229
```

```

230 #dot product vector of players performances * vector of feature
      weights
231 events_italy_PR_players_rating <-
      events_italy_PR_players%>%dplyr::select(-c(playerId,matchId,teamId))
232
233
234 events_italy_PR_players_rating_M <-
      data.matrix(events_italy_PR_players_rating)
235
236 weighted_events_italy_PR_players <-
      events_italy_PR_players_rating_M %*%weights2
237
238 range01 <- function(x){(x-min(x))/(max(x)-min(x))}
239
240 weighted_events_italy_PR_players_st <-
      data.frame(perfomance_for_match=range01(weighted_events_italy_PR_players)
      )
241
242 weighted_events_italy_PR_players_st$matchId <-
      events_italy_PR_players$matchId
243
244 weighted_events_italy_PR_players_st$playerId <-
      events_italy_PR_players$playerId
245
246 library(qcc)
247 events_italy_PR_prov <-
      events_italy_PR%>%dplyr::select(c(playerId,p_name))%>%distinct()
248
249 player_ratings <- weighted_events_italy_PR_players_st%>%
250   dplyr::select(c(playerId,perfomance_for_match))%>%
      dplyr::group_by(playerId)%>%
251   dplyr::summarise(ratings=mean(perfomance_for_match))
252
253 player_ratings <-
      merge(x=player_ratings,y=events_italy_PR_prov,by='playerId',all.x=TRUE
      )
254
255 #visualize the results
256
257 ratings_plot <- ggplot(data=player_ratings, aes(y=ratings))+
258   geom_boxplot(outlier.colour="black", outlier.shape=16,
259     outlier.size=2, notch=FALSE,
260     fill='cornflowerblue')+theme(legend.position="none")+

```

```

261   labs(tag = "(a)")
262
263   hist_goals <- ggplot(player_ratings,aes(x=ratings))+
264     geom_histogram(binwidth=0.03,fill='dodgerblue4')+
265     geom_vline(data=player_ratings,
266               aes(xintercept=mean(player_ratings$ratings), color="red"),
267                 linetype="dashed")+
268     theme(legend.position="none")+
269     labs(tag = "(b)")
270
271   figure7 <- ggarrange(ratings_plot,hist_goals, ncol = 2, nrow = 1)

```

A.2.3 Unsupervised Model (K-means)

```

1
2   #-----Clustering For Role detection-----
3
4   library(plotly)
5   library(factoextra)
6   library(NbClust)
7   library(class)
8   library(MASS)
9   require(pROC)
10  library(ISLR)
11  library(RColorBrewer)
12  library(wesanderson)
13
14  #data cleaning
15  events_italy_RD=events_italy_RD%>%
16    dplyr::select(c(matchId,teamId,playerId,x_start,
17                  y_start,Duel.Air.duel.accurate,Duel.Air.duel.not.accurate,
18                  Duel.Ground.attacking.duel.accurate,
19                  Duel.Ground.attacking.duel.not.accurate,
20                  Duel.Ground.defending.duel.accurate,
21                  Duel.Ground.defending.duel.not.accurate,
22                  Duel.Ground.loose.ball.duel.accurate,
23                  Duel.Ground.loose.ball.duel.not.accurate,
24                  Foul.Hand.foul.red_card,Foul.Hand.foul.yellow_card,
25                  Foul.Late.card.foul.yellow_card,Foul.Foul.red_card,
26                  Foul.Foul.second_yellow_card,

```

```

25     Foul.Foul.yellow_card,Foul.Out.of.game.foul.red_card ,
        Foul.Out.of.game.foul.second_yellow_card,
26     Foul.Out.of.game.foul.yellow_card,Foul.Protest.red_card
        ,Foul.Protest.second_yellow_card,
27     Foul.Protest.yellow_card,
        Foul.Simulation.yellow_card,Foul.Violent.Foul.red_card,
28     Foul.Violent.Foul.yellow_card,Free.Kick.Free.kick.cross.accurate,
        Free.Kick.Free.kick.cross.not.accurate,
29     Free.Kick.Free.Kick.accurate,Free.Kick.Free.Kick.not.accurate,
30     Free.Kick.Penalty.not.accurate,
31     Free.Kick.Free.kick.shot.not.accurate,Free.Kick.Free.kick.shot.accurate,
32     Free.Kick.Throw.in.accurate,
33     Free.Kick.Throw.in.not.accurate,Others.on.the.ball.Acceleration.accurate,
34     Others.on.the.ball.Acceleration.not.accurate,
        Others.on.the.ball.Clearance.accurate,
35     Others.on.the.ball.Clearance.not.accurate,
36     Others.on.the.ball.Touch.assist,Others.on.the.ball.Touch.counter_attack,
37     Others.on.the.ball.Touch.dangerous_ball_lost,
38     Others.on.the.ball.Touch.interception,
39     Others.on.the.ball.Touch.opportunity,Pass.Cross.accurate,
40     Pass.Cross.assist,Pass.Cross.keyPass,
        Pass.Cross.not.accurate,Pass.Hand.pass.accurate,
41     Pass.Hand.pass.not.accurate,
42     Pass.Head.pass.accurate,Pass.Head.pass.assist,
43     Pass.Head.pass.keyPass ,Pass.Head.pass.not.accurate,
44     Pass.High.pass.accurate,Pass.High.pass.assist,
45     Pass.High.pass.keyPass,Pass.High.pass.not.accurate,
46     Pass.Launch.accurate,Pass.Launch.keyPass,
47     Pass.Launch.not.accurate,Pass.Simple.pass.accurate,
48     Pass.Simple.pass.keyPass,
49     Pass.Simple.pass.not.accurate,Pass.Smart.pass.accurate,
        Pass.Smart.pass.assist,
50     Pass.Smart.pass.keyPass,
        Pass.Smart.pass.not.accurate,Shot.Shot.accurate,Shot.Shot.not.accurate))
51
52
53 #as_factor
54 events_italy_RD$playerId <- as.factor(events_italy_RD$playerId)
55
56 events_italy_RD$teamId <- as.factor(events_italy_RD$teamId)
57
58
59 #take into account only the usefull type of events

```

```
60
61 coordi=events_italy_RD %>%
62   dplyr::select(c(x_start, y_start, playerId,
63     matchId))%>%dplyr::group_by(matchId,playerId)%>%
64   dplyr::summarise(avg_pos_X = mean(x_start), avg_pos_Y=
65     mean(y_start), .groups = 'drop')%>%
66   dplyr::filter(avg_pos_X>20)%>%dplyr::select(-c(playerId,matchId))
67
68 #k=8
69 # Silhouette method
70 #Measures the quality of a clustering and determines how well each
71 #point lies within its cluster.
72 silhouette <- fviz_nbclust(coordi, kmeans, method = "silhouette") +
73   labs(subtitle = "Silhouette method")
74
75 silhouette_score <- data.frame(silhouette[["data"]][["y"]])
76
77 km_1 = kmeans(coordi, 8, nstart = 1, iter.max = 1e2)
78
79 coordi$Clusters <- as.factor( km_1$cluster)
80
81 par(mfrow=c(1,1))
82
83 p <- ggplot(coordi, aes(x=avg_pos_X, y=avg_pos_Y)) +
84   annotate_pitch(fill = "white", colour = "black") +
85   theme_pitch()+
86   geom_point(aes(col =Clusters),alpha=0.8 ,size =3) +
87   scale_fill_manual("Legend")+
88   coord_fixed()
89
90 results <- data.frame(coordi$avg_pos_X,coordi$avg_pos_Y)
91 results <- data.frame(NA_col = rep(NA, nrow(coordi)))
```

Bibliography

- [1] Rossi A, Pappalardo L, Cintia P, Iaia FM, Fernández J, Medina D (2018) Effective injury forecasting in soccer with GPS training data and machine learning. *PLoS ONE* 13(7): e0201264. <https://doi.org/10.1371/journal.pone.0201264>
- [2] Prakoso, Muhammad Lumintuarso, Ria. (2021). Analysis of the Use of Statistical Data in the Formulation of Strategies, Tactics and Evaluation of Football Matches. 10.2991/ahsr.k.210707.009.
- [3] Ballesta, Christian Garcia-Romero, Jerónimo José Carlos, Fernández-García Alvero Cruz, Jose Ramon. (2015). Current methods of soccer match analysis. *Revista Internacional de Medicina y Ciencias de la Actividad Fisica y del Deporte*. 15. 785-803.
- [4] Pappalardo, L., Cintia, P., Rossi, A. et al. A public data set of spatio-temporal match events in soccer competitions. *Sci Data* 6, 236 (2019). <https://doi.org/10.1038/s41597-019-0247-7>
- [5] Buldú JM, Busquets J, Martínez JH, Herrera-Diestra JL, Echegoyen I, Galeano J and Luque J (2018) Using Network Science to Analyse Football Passing Networks: Dynamics, Space, Time, and the Multilayer Nature of the Game. *Front. Psychol.* 9:1900. doi: 10.3389/fpsyg.2018.01900
- [6] Cintia, P., Rinzivillo, S. Pappalardo, L. Network-based Measures for Predicting the Outcomes of Football Games. *Proceedings of the 2nd Workshop on Machine Learning and Data Mining for Sports Analytics (MLSA)*, 46–54 (2015).
- [7] Luca Pappalardo, Paolo Cintia, Paolo Ferragina, Emanuele Massucco, Dino Pedreschi, and Fosca Giannotti. 2019. PlayeRank: Data-driven Performance Evaluation and Player Ranking in Soccer via a Machine Learning Approach. *ACM Trans. Intell. Syst. Technol.* 10, 5, Article 59 (September 2019), 27 pages. <https://doi.org/10.1145/3343172>

-
- [8] Jonny Whitmore, Thomas Seidl; *Shape Analysis: Automatically Detecting Formations*; Stats Perform; <https://www.statsperform.com/resource/shape-analysis-automatically-detecting-formations/>
- [9] Jonathan Whitmore; *Introducing Movement Chain*; Stats Perform; <https://www.statsperform.com/resource/introducing-movement-chains/>.
- [10] Rein, R., Memmert, D. Big data and tactical analysis in elite soccer: future challenges and opportunities for sports science. SpringerPlus 5, 1410 (2016). <https://doi.org/10.1186/s40064-016-3108-2>
- [11] Schumaker, Robert P. et al. "Sports Data Mining." (2010).
- [12] Taki, Tsuyoshi and Jun-ichi Hasegawa. "Visualization of dominant region in team games and its application to teamwork analysis." Proceedings Computer Graphics International 2000 (2000): 227-235.
- [13] Narizuka, T., Yamazaki, Y. (2019). Clustering algorithm for formations in football games. Scientific Reports, 9.