## Way to the Triumph

Our main objective in this paper is to explore the complex interactions among the players on the field that can influence the team's success overall. Besides this, we want to determine the team dynamics and individual players'qualities and how these factors influence the team's performance throughout the season. Finally, we want to present a specific team strategy to improve team performance for the next season.

First of all, we apply the Social Network Theory to evaluate the passing pattern behind each match. We then choose the most representative match of both winning and losing. By replacing node as players and passes as connections, we generate the adjacency matrix that records the passing frequency between each player. We then find the average position of players for both winning and losing games and draw them as nodes on a simulated football field. After this, we complete the edges as the passing frequency that was previously generated by the adjacency matrix. Thus, we generate a Passing Network.

From the passing network, we observe that the passings of the winning game are more spread out and the passings of losing game most occupy a small region on Huskies'own field. After this, in order to understand the importance of a particular node, we define two terms as "Active Passing" and " Passive Receiving", which denote as the passes originated from this player and pass received by this player. Then, we introduce the formula"Degree Centrality" to analyze their contribution to the entire network, which leads to the phenomenon of Dyadic and Triadic Configuration. Understanding the fact of both passing and receiving the ball through these configurations, we then decide to use a new equation on individual performance on the Passing Network.

We then proceed to understand other traits than can quantify a player's quality and team performance. We first collect X and Y coordinates of the shots taken during the entire season to draw a line of which most shots are taken. Behind this line, from the perspective of Huskies, is the danger zone in which Huskies are most likely to make the shot. However, after creating the model, we also realize that the results are not precise enough since it only contains the X coordinates. So, we use the Central Limit theorem to filter out the precise danger zone, with both X and Y coordinates, which we conclude that it is important to send the balls through the passing network and reach the danger zone and make the shot.

We then proceed to define the qualities of a player and formations, so that we can produce the best team strategy based on these two factors to Huskies. We approach this problem from four dimensions: Duel, Pass, Foul and Shot. Moreover, we produce a generic equation for describing these qualities. We then locate the best players by matches and compare them with relative winning rates and team formations. Then, we want to determine which coaches were more frequent to use the formation that generates the most win. Then, we want to further the evaluation of player skills and determine who should be on the best players list. Through this analysis, we also rule out the player who does not perform very well during the season and suggests the coach trade them into other teams.

Through our presentation of Passing Network, numerous traits and the relationship of different aspects of the team, either individually or collectively, we believe that our model and analysis can help Huskies to become better in the next season.

# Contents

# 1 Introduction

## 1.1 Context

Sports analysis is an inevitable step in the competitive sports world, especially for soccer. It can not only improve the player performance, but also maximizes the outcome of the team's quality of play. We seek to use data analytic techniques to explore the complex interactions on the soccer field, hoping to provide valuable suggestions to improve Huskies' overall success.

## 1.2 Assumption

1)Data is all accurate and precisely reported.

2)There is no special condition such as "sick", "emotion" or any other events that are not included in the data frame.

3)The Data we ignore such as Interruptions do not have a significant impact on a player's ability.

4)The Conclusion that we make from this season's data should carry a considerable amount of influence to the next reason.

# 2  Our Model

## 2.1  Passing Network Analysis

### 2.1.1  Model Description

The winning of a game is largely dependent on making effective transitions of balls. By applying knowledge of Social Networking Theory, we create an average passing network model that each player represents a node and each passing as an edge. We filter out the average player coordinates and player connections into a multidimensional adjacent matrix for each match. The passing origin player is denoted as $O$(sub player id), and the passing destination player is denoted as $D$(sub player). Thus, we form a $N \times N$ matrix that rows are passing origin players and columns are passing destination players.

$$\begin{bmatrix} (O_{D_1}, D_{D_1}) & (O_{D_1}, D_{D_2}) & (O_{D_1}, D_{F_3}) & ... & (O_{D_1}, D_N) \\ (O_{D_2}, D_{D_1}) & (O_{D_2}, D_{D_2}) & (O_{D_2}, D_{F_3}) & ... & (O_{D_2}, D_N) \\ ... & ... & ... & ... & ... \\ (O_N, D_{D_1}) & (O_N, D_{D_2}) & (O_N, D_{F_3}) & ... & (O_N, D_N) \end{bmatrix}$$

Since each opponent has its own ID and we do not know the performance of each opponent during the entire season. We decide to use the game that generates the most goal difference to represent our passing network through linear regression with x as match ID and y as goal difference. From this graph below, we decide to choose game No.14 as the most winning representative game and game No.23 as the most losing game, since each generates a goal difference of 4 and -4.
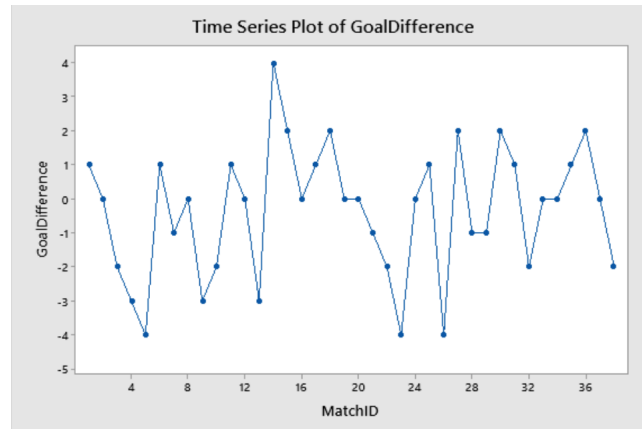


Figure 1: Goal Difference for the entire season

Game 14 started off with the formation 4-4-2, experienced three substitutions from 4-4-2, 4-5-1, and finally ended at 4-6-0. Applying the matrix algorithm above, we generate the following matrix as a table.

| D\O | D2 | D5 | D6 | D7 | F1 | F2 | G1 | M1 | M10 | M11 | M4 | M6 | M8 | M9 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D2 | 0 | 6 | 8 | 0 | 0 | 13 | 3 | 5 | 0 | 0 | 1 | 2 | 0 | 0 |
| D5 | 5 | 0 | 0 | 0 | 3 | 7 | 3 | 0 | 2 | 1 | 2 | 4 | 3 | 1 |
| D6 | 7 | 1 | 0 | 5 | 5 | 3 | 5 | 5 | 0 | 1 | 1 | 2 | 1 | 0 |
| D7 | 0 | 0 | 4 | 0 | 6 | 4 | 2 | 4 | 0 | 1 | 5 | 0 | 4 | 0 |
| F1 | 0 | 0 | 0 | 1 | 0 | 2 | 0 | 1 | 0 | 0 | 3 | 2 | 3 | 0 |
| F2 | 5 | 9 | 8 | 2 | 4 | 1 | 1 | 6 | 0 | 0 | 4 | 11 | 3 | 0 |
| G1 | 2 | 1 | 0 | 0 | 8 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| M1 | 5 | 5 | 5 | 5 | 1 | 9 | 1 | 0 | 1 | 1 | 7 | 4 | 2 | 0 |
| M10 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| M11 | 0 | 0 | 2 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| M4 | 1 | 1 | 0 | 10 | 1 | 0 | 0 | 3 | 0 | 0 | 1 | 5 | 4 | 0 |
| M6 | 2 | 3 | 0 | 2 | 5 | 6 | 0 | 5 | 0 | 0 | 1 | 0 | 4 | 0 |
| M8 | 0 | 4 | 0 | 3 | 3 | 5 | 0 | 3 | 0 | 0 | 3 | 0 | 0 | 0 |
| M9 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |

Table 1: Match 14

Then, we want to a column summation of Total Passing Players to evaluate the total number of active passes that were generated by this player. This would show us the number of actions that are generated by this player in the network.

**Column summation of total passing per players:**

$$C(D, F, M) = \sum ((D, F, M)_n, (D, F, M)_n)$$

From the row summation of total passing Network, we observe that the three substituted players (M9, M10, M11) have trivial effects on the summation of total passing, which each of them accounted for (2, 3, 6).

We then proceed to do a row summation of this matrix, this would give us the number of passes that are received by this player. It is important to acknowledge this fact as both ends in the passes of soccer are important.

**Row summation of total passing per players:**

$$R(D, F, M) = \sum ((D, F, M)_n, (D, F, M)_n)$$

**Degree Centrality of Active Passing and Passive Receiving:**

Based on the number of passing per players, we decide to use degree centrality to evaluate the importance of a particular node(players) in this passing network. Each passing that originates from

this player will be represented as an active edge connecting to the destination players. We calculate the total sum of the active edge of each node in this network to present the amount of activity that is generated by this node. Then, we calculate the amount of activity that is received by this node to evaluate the flow of activities.

$$T(G) = \sum_{i=1}^{N} R(D, F, M) \qquad C_{(d,p)}(G) = \frac{C(D, F, M)}{T(G)} \qquad C_{(d,r)}(G) = \frac{R(D, F, M)}{T(G)}$$

**Graphical Representation of Match 14 Passing Network Analysis**



Figure 2: Match 14 Passing

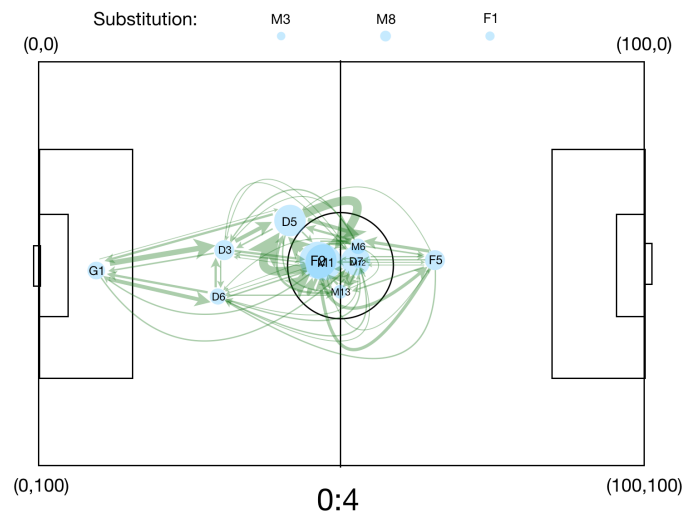**Graphical Representation of Match 23 Passing Network Analysis**



Figure 3: Match 23 Passing

From the Degree Centrality calculation of each player, we concluded that the substitution players M9, M10, M11 are relatively small in their influence, each accounted for 0.5%, 0.8%, and 1.7%, to the entire Passing Network. So, we decided to put them aside and analyze the remaining 97% percent of the network.

We realized that the data set of coordinates is respective to the attacking team perspective. However, since we are only concentrating on the passing network of Huskies, we set up the graph from the Huskies' standpoints, where the left-hand side of Huskies are denoted as (0,0) and the bottom right corner are denoted as (100,100). Then for each player, we put them into the players' average position throughout the match. The width of the connections is proportional to the number of connections between the two players. Lastly, the size of the node is relative to the number of passes they make.

From the match 14 graph, we can see that the F2 is the engine of transitioning goals from players to players as it has the largest size. In the Active Passing Degree Centrality dataset, F2 has 14.4% of active passes in the network, ranking the first place out of 14 players. Moreover, in the Passive Receiving Data Sets, F2 also ranks to be the highest place with 15.5%, which means that 15.5% of the balls are transitioned through F2. And the top four players are among both datasets are Active - (F2(14.4%), F1(10.3%), M1(9.5%), M6(9.5%), Passive - (F2 (15.5%), M1(13.2%), D2(10.9%), D6(10.3%). Moreover, the entire team is spread out on both the X and Y-axis, which control the middle section of the field. This visually tells us that the Huskies have very strong control of the ball transitioning and the team is very confident in making the play.

However, on the contrary, match 23 Passing Network shows us a very different graph. The passing zone of the entire team is more condense as the average positions of multiple team members are squeezed in the central regions. The passes are shorter compared to Match 14 and most of the Huskies' passes are on their own half of the field.
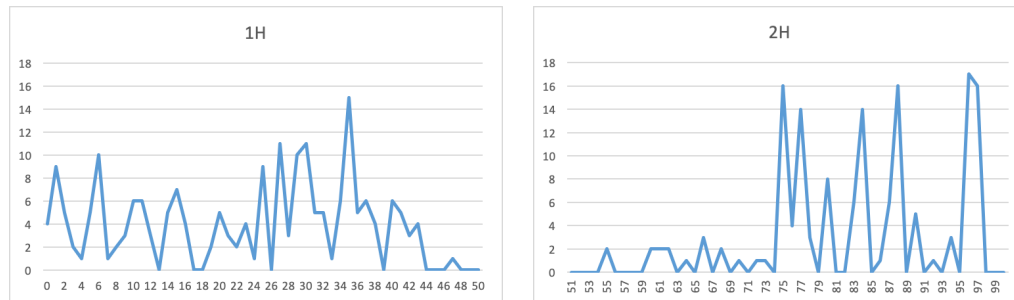
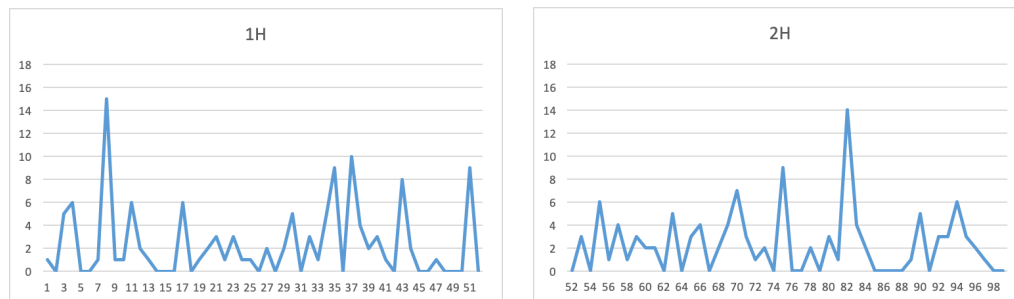Figure 4: Match 14: Minute to minute passing



Figure 5: Match 23: Minute to minute passing

Moreover, the total passes throughout the game are also less than Match 14. The above graph shows the minute to minute passing between Match 14 and Match 23 with the ratio of 347 to 236 From here, we realize the correlation between effective passing and winning, and how passing can be a vital indicator of team member collaboration, winning and personal strength.

## 2.2   Passing Analysis Throughout the Entire Season

We collect the passing events for the entire season and decide to allocate the passing events for the wins and not wins(losses and ties). Then for each dataset, we decide to average them out to see the passing pattern for this entire season. And we find a significant difference on wins and non-winning games
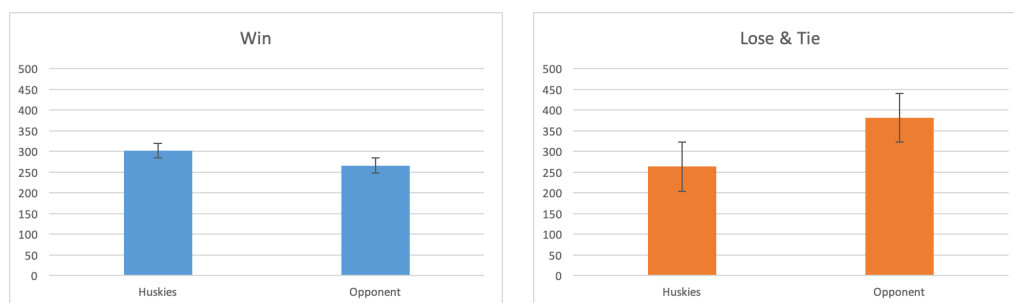


Figure 6: Season Passing Stats

The number of passes on the win matches outnumbers the passes on the lost matches by 50. Moreover,

the number of passes for Huskies during the losing matches is significantly lower than the opponent, which usually implies a lesser control of the ball than the opponent does and less of the time which the ball is on the opponent's field.

## 2.3   Dyadic and Triadic configurations

From the Graph above for Match 14, we observe that the team collaboration between each member has become a vital factor in constituting the win. For example, the interaction between F2 and M6 forms a strong pairwise connection that ensures the safety of the ball control when the attack is potentially facing a group of strong opponents that try to surround the attacking player. Moreover, the Triadic configuration between player F1, M6 and D7 helps to provide support when F1 tries to penetrate the opponent's defense. The relationships between players help to generate more energy flow through the network.
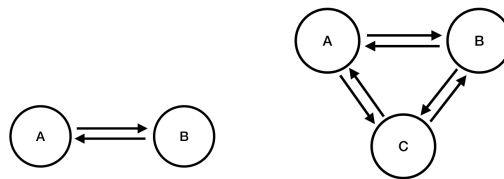


Figure 7: Dyadic and Triadic Configurations

| A | B | Frequency |
|---|---|---|
| M1 | M3 | 65 |
| D4 | M3 | 33 |
| D5 | F2 | 28 |
| D5 | M1 | 26 |

Table 2: Dyadic

| A | B | C | Frequency |
|---|---|---|---|
| F2 | M1 | M3 | 10 |
| D5 | F2 | M1 | 9 |
| D1 | M1 | M3 | 8 |
| F2 | M1 | M4 | 6 |
| D6 | M1 | M3 | 6 |
| D2 | M1 | M3 | 6 |

Table 3: Triadic

Has formed a central structure along the middle line of the field. It represents the front base for the attacking team members to receive support and the first defense line that protects the opponent team to move further into the Huskies field. Another important configuration during Match 14 was the interaction between F2, M1, and M4. Moreover, these two configurations overlap with F2 and M1, which further connect the entire factor. So, in order to optimize the passing network. It is vitally important to increase these interactions between players. Our advice to the coach of Huskies is that they should practice more pairwise or group of three team plays to increase the efficiency of transitioning the balls to the opponent's field.

## 2.4 Composite Degree Centrality Assessment

From the examples and data provided above, we realize the importance of passing from multiple dimensions. If we see passing as an edge in a group, then it is important to acknowledge that both ends of the edge are equally important. Thus, we divide passing to two ends as active passing and passive receiving which both should count as an individual's playing ability into making contributions to the entire passing network. Thus we develop the Composite Degree Centrality Model into evaluations of team members.

For Individual Match:

$$\frac{C_{(d,p)}(G) \times C_{(d,r)}(G)}{\sum \left( C_{(d,p)}(G) \times C_{(d,r)}(G) \right)}$$

For the Entire Season:

$$\sum \frac{C_{(d,p)}(G) \times C_{(d,r)}(G)}{\sum \left( C_{(d,p)}(G) \times C_{(d,r)}(G) \right)}$$

Thus, we can determine the individual contribution of players into constructing the passing network.

# 3  Danger Zone

## 3.1  Objective

To use the **Central Limit Theorem(CLT)** to find the "Danger Zone" that reflects the most frequent region of producing shots in the future.

## 3.2  Initial Model

Since we looped over the shots produced by Huskies and collected all the possible attempts, we decided to use median ($x = 88$) to show the "danger line" that is most likely to produce shots, and intuitively, any area closer to the opponent's goal can be regarded as the "danger zone":
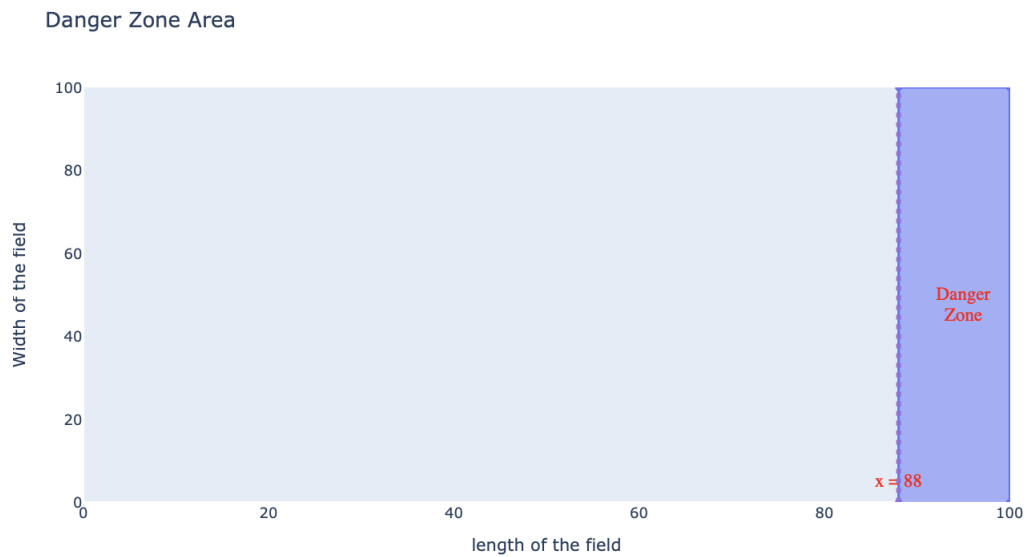


Figure 8: Danger Zone

## 3.3  New Model

As we observe that Figure 2 has a high standard error of $(8.989, 12.404)$ for the shot collection, we decide to use the Central Limit Theorem to find the most-likely-shot zone.

$$\mu = \frac{1}{N} \sum_{i=1}^{N} x_i \qquad\qquad \sigma = \sqrt{\frac{\sum_{i=1}^{N} (x_i - \mu)}{N}}$$

$\mu$ : the total mean of location of shots from winning games (the population mean)
$x_i$: each specific location of shots
$N$: total number of shots
$\sigma$: standard deviation of the total winning shots

By collecting all the useful shot coordinates of Huskies through last season together, we chose to pick up 119 choices as a sample size each time, which is the total number of shots Huskies made in the winning games. Then, based on the law of large numbers, we decided to do 100000 times to minimize the close heuristic effect which we want to predict the future result instead of this season. The final distribution graph we produced based on the bootstrapping result is undoubtedly similar to the normal distribution:



Figure 9: Distribution of $X$(left) and $Y$(right) coordinates

As the graph shows, we can finally specify the "danger zone" as a rectangular with the coordinates: (82.62, 47.12), (82.62, 55.18), (88.27, 47.12), (88.27 55.18). Correspondingly, the standard error right now is (0.83,1.14), which is conspicuously better than the original data we extracted from this season of Huskies.
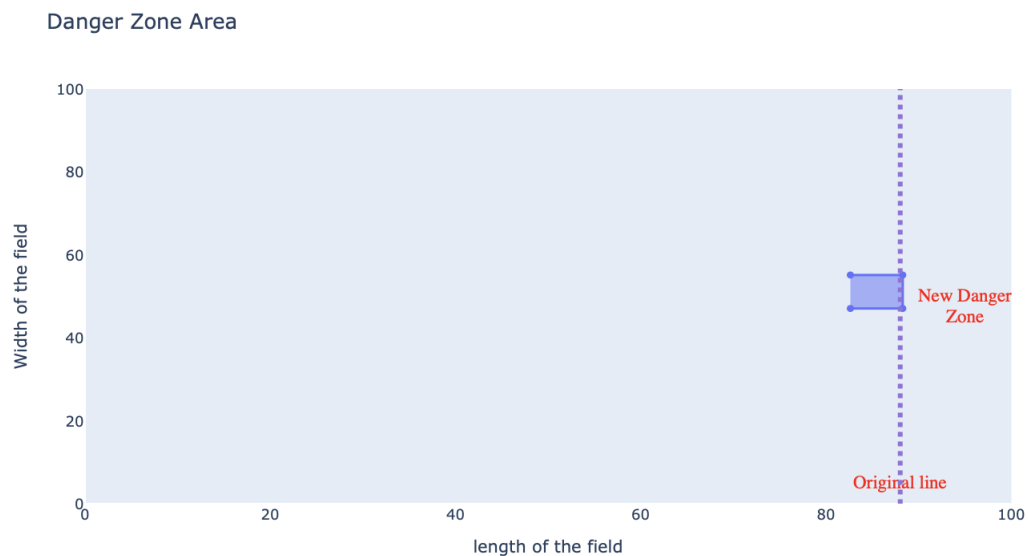


Figure 10: Final Danger Zone

Now, we can say that whenever and whatever the Huskies reach out to this region, then it most likely to be counted as a successful play to produce shots. We can use this region to judge the team performance for each game and the ability of coaches to create the shooting chance for players.

# 4  Player Abilities

## 4.1  Objectives

Through finding the best way to define the ability of any soccer player, we would like to make the following achievements:
1) Locating the best players among Huskies with the comparatively higher overall skills and analyze the relationship between the number of appearances of them with the winning rate and team formations.
2) Deciding which coaches more frequently use the best team formation with the best players in this season and which were not, then we can make the improvement around this point.
3) Find the bottom ability players, so that we can adjust their positions with the team formation or even trade off.

## 4.2  Team Model

We first need to analyze which kinds of data can be used to define the ability of players from the database we have. After researching on the relevant articles, we finally came across one paper called Characteristics That Make A Good Soccer Player. It describes the ability of any soccer player as team work, intelligence, attitude and concentration. Since we only have limited data set and simple code name, we decided to use the Duel, Pass, Foul, Shot and Interruption to define the players' skill sets. However, as we looped over the whole database, there only occurs four times interruptions, but for the other aspect we can receive a pretty decent amount of number to analyze, thus we decided to drop Interruption out of consideration. In order to reflect the accurate skill level from each aspect for each player, we would like to create a simple equation as:

$$E_t = E_d + E_p + E_f + E_s$$

Here, we define $E_t$, $E_d$, $E_p$, $E_s$ respectively as the total efficiency for each player, the efficiency of Duel, the efficiency of Pass and the efficiency of Shout.

In order to reflect the efficiency of each component from database, we think that the most efficient way to achieve the result is to use this formula: $E(s) = \frac{N}{P}$, where x stands for the total number of $s$ and $P$ stands for the total number of game the player participated. We decide to regard the efficiency as the amount of the specific action he can make per game, by calculating the total number of a specific aspect from the player ($N$) and divide it with the number of played games ($P$), since each player had played different times last season as shown below:
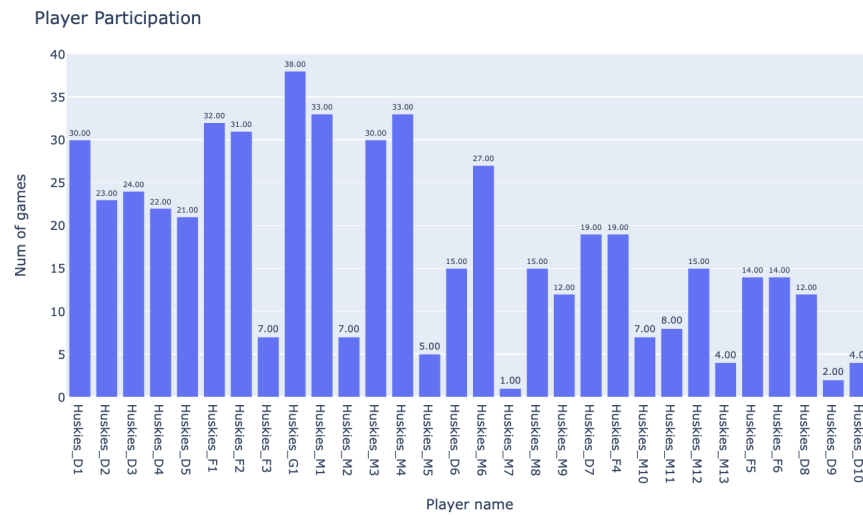
Figure 11: Player Participation

Then following the calculation,we encounter this problem of raw data unbalancing as they have different parameters.
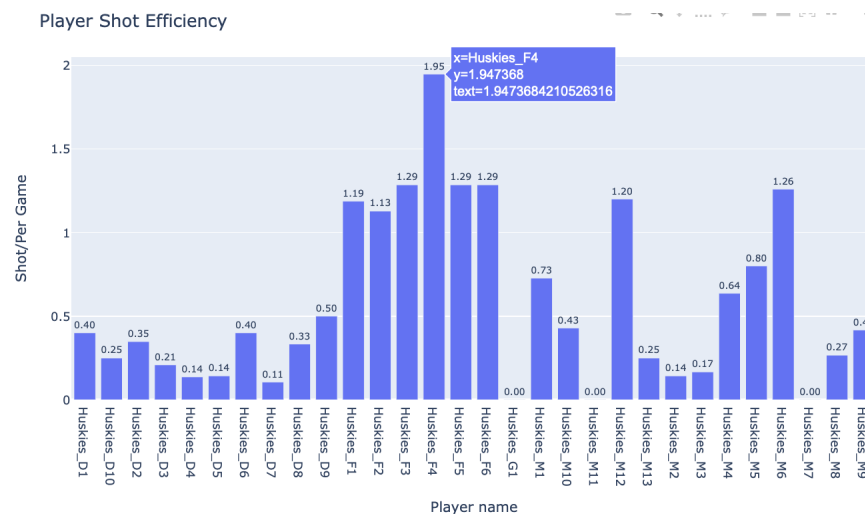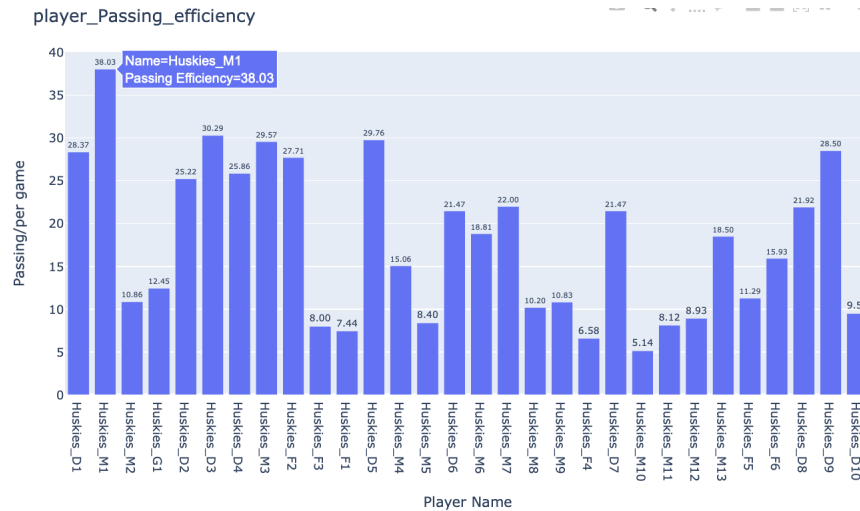


Figure 12: Player Shot Efficiency

Figure 13: Player Passing Efficiency

As the graphs show, the efficiency of passing is much higher than the efficiencies of Shot, and indeed Duel and Passing are much higher than the rest, so we can not directly use the raw score, thus instead, we create this curving function:

$$f(x) = \frac{100x}{H}$$

(Function represents that the highest one ($H$) in each aspect will become 100% and the other ones are computed as the percentage of H they attained) By applying the formula, the data becomes more representative as:



Figure 14: Player Pass Efficiency

We approach the same method for Foul Efficiency, Duel Efficiency and Shot Efficiency (see in Appendix).

Also, since we calculate the foul efficiency as how many fouls a player can produce per game, so should be counted as a negative factor. Therefore, the new equation should be:

$$E_t = E_d + E_p - E_f + E_s$$

Thus, we can right now locate the best players and worst players for each position since we have their specific efficiency percentage of different divisions.

## 4.3  Division Model

After reranking all the players, we find these following results:

*We assume that Goalkeeper should not be ranked since there is no substitution and intuitively he can not make many contributions on these aspects.

| Type\Ranking | Top1 | Top2 | Top3 | Top4 | Worst1 |
|---|---|---|---|---|---|
| Forward | 1.66 F6 | 1.20 F2 | 1.12 F3 | 1.10 F4 | 0.81 F5 |
| Midfielder | 1.44 M6 | 1.33 M1 | 1.08 M12 | 1.00 M7 | 0.31 M10 |
| Defender | 1.29 D1 | 1.26 D3 | 1.23 D8 | 1.17 D9 | 0.21 D10 |

Table 4: Total Efficiency

From the table above, we can see that F6, M6, and M1 are the most valuable players from all dimensions. However, to make the ranking fairer for each position, we decide to rank again with the specific traits that are more relevant to them. For example, forwards should be hinged on the Shot and Duel Efficiency, defenders and the midfielders should be analyzed from Pass, Duel, and Foul Efficiency.

$$E_t = E_d + E_s$$

| Type\Ranking | Top1 | Top2 | Top3 | Top4 | Top5 | Worst1 |
|---|---|---|---|---|---|---|
| Forward | 1.82 F4 | 1.61 F1 | 1.52 F5 | 1.44 F6 | 1.14 F2 | 0.99 F3 |

Table 5: Forward

However, from the **former Passing Network**, we find that F2 is more inclined to pass so we can not simply judge him based on this standard. Yet, F3 is relatively below others' ability.

$$E_t = E_d + E_p - E_f$$

The ranking is:

| Type\Ranking | Top1 | Top2 | Top3 | Top4 | Top5 | ... | Worst1 | Worst2 |
|---|---|---|---|---|---|---|---|---|
| Midfielder | 1.01 M7 | 0.96 M1 | 0.80 M6 | 0.67 M9 | 0.63 M3 | ... | 0.12 M2 | 0.09 M10 |

Table 6: Midfielder

We can now put M7, M1, M6 into the final top Midfielder list, and vice versa for M10, M2.

$$E_t = E_d + E_p - E_f$$

The ranking is:

| Type\Ranking | Top1 | Top2 | Top3 | Top4 | Top5 | ... | Worst1 | Worst2 |
|---|---|---|---|---|---|---|---|---|
| Defender | 1.16 D3 | 1.08 D1 | 1.06 D8 | 0.92 D9 | 0.89 D4 | ... | 0.62 D2 | 0.08 D10 |

Table 7: Defender

From the ranking above, D3, D1, D8 are ranked into the final top Defender list and add D2, D10 into the worst list.

So far, the top player list we have are: F4, F1, F5, F2, M7, M1, M6, D3, D1, D8; and the worst player are: F3, M2, M10, D2, D10.To double-check their team collaboration with the whole team, we decide to combine with the team configurations.

According to Table 2 and 3, we can conclude that M1, D3, D1 are the most important core players for Huskies since plays almost all include him, but M3 has a high chance to cause the foul and lower chance of duel so his ability score is not conspicuous from the rankings.

By locating all the games that the top players were included, we find F2 played all the 10 winning games and no matter which coach could lead the team, they won nonetheless, so we can conclude that F2 should also be in the core list. However, for M1, D3, D1, among the games they all attended, they won 6 times, tied 7 times and only lost 5 times, which occupies 34% of the winning rate, there seems to be no difference, so we can conclude those core players are not the significant determinator for winning. However, for D3, all the games he has not played, the Huskies only won 3 times. For the worst players, M2 has the highest amount of lost games (6), and then F3, M10, and D10, but the game D2 played has the highest amount of winning games (9), so we can then delete him from the worst player list.

Then from all the analyses, we can now decide the players to be traded in the future: M2, F3, M10; and F2 D3 are the player most important for the team.

By combining the frequency of team formation with the core player in the games, we find that under 4-2-2 and 4-3-3 formations, they can maximize the ability of core players and they won most times. Therefore, in the future, this can be counted as criteria to follow if we can reasonably assume core players will not leave the team. Also, even though coach 3 led most of the games in the past season,

and there is a high probability that coach 1and coach 2 will use different strategy, but so far, we can see from the table, coach 2 did not use the most winning formation frequently and the only two he coached are all with the formation as 4-2-2. Therefore, we can make the conclusion that coach 2 may need to try more winning formation (4-2-2,4-3-3) in the future if possible. Here is a graph shown the formation frequency of Coach 2:
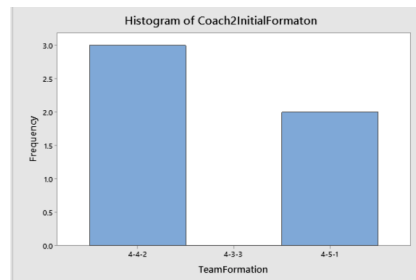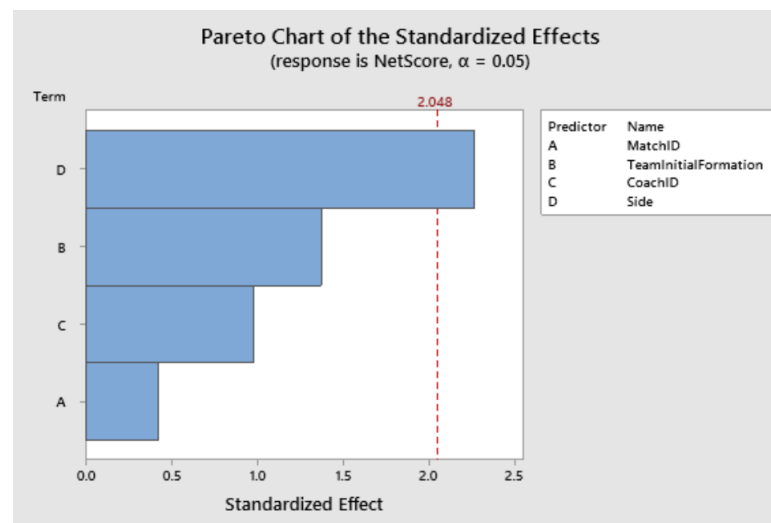
Figure 15: Coach2 Stats

# 5 Further Discussion

Figure 16: Regression of outcomes

Since after analyzing the relationship between winning games and the other aspects (formation, coach, side), we find that coaches do not have a dominant influence on the outcomes of games, so besides the advice of changing formation, it is not necessary to make further discussion at this moment. However, whether playing at home or away is the biggest determinator among all the aspects, but the schedule normally can not arrange by the team itself, so we do not have the ability to make improvements in this aspect.

# 6 Evaluation

## 6.1 Strength

Before creating each model, we carefully cleared the data set to efficiently minimize the noises and outliers from the miscellaneous data, then we save plenty of time on finding the useful information.

Our models are multi-dimensional. Our models not only represent the different layers from the shallow to the bottom to prove our final thesis, but we also analyze the data from different perspectives, such as the Player Efficiency and Coach Formation Frequency and so on. Therefore, we can depict a comprehensive picture consisting of various aspects.

Last but not least, the models we have can be kept for future use. Since "danger zone" and the diagrams of different types of configurations can all reflect the team's characteristic; and this team trait normally is consistent for several seasons, then the Huskies can make a constructional plan based on the drawbacks we analyze from the features.

## 6.2 Weakness

There exist flux and unpredictable differences between our models and the real-world situation. Since from the databases we can analyze, they are not included the scoring time and the corresponding player, so we cannot take the efficiency of shot and cut-in for specific players into consideration, and we cannot make a precise conclusion about the efficiency of attacking for each game. Then, when we evaluate the attribute of each player, we lack enough information.

Some information is arbitrary. As for the coaches most frequent formations, we cannot directly find the formation coach used to produce a positive effect, since there is no scoring time for both sides and there is no information of weathers, so we cannot confidently say that certain type of team formation can make a great contribution for different coaches. Therefore, we can only evaluate the starting team formation for each coach, and make a consequent analysis.

## 6.3 Future Improvement

First, with more available time, we can enhance our description of the "danger zone". If we can convert the cartesian coordinates into real field size, then we can give a more visual-friendly diagram to reflect the real world condition. Also, we can further discuss which certain types of plays (strategies) can attain the scoring zone more easily.

Second, we assume no injury for all players throughout the season, but if we can then analyze the dynamic changes of each player's statistics, the assumption is possible to become a data-supported fact. Also, we only try a little on this part because of the limitation of time, such as:
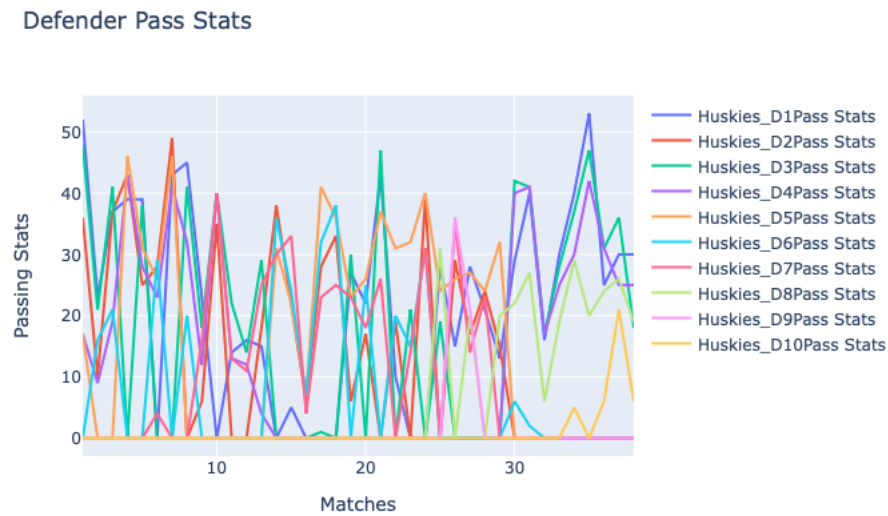
Figure 17: Defender Pass Stats

Third, we can make each graph more precisely. As for the graph of coaches' formation frequency and passing dynamical pattern, sometimes we do not have the legend to decorate the graph, so they are not informative if we do not make additional explanations.

# References

[1]  Alto, Valentina. "Visualizing the Central Limit Theorem with Python." *Medium*, Towards Data Science, 1 Aug. 2019, towardsdatascience.com/visualizing-the-central-limit-theorem-with-python-e89d2ce41788.

[2]  "Characteristics That Make A Good Soccer Player." *Puget Sounds Slammers*, 28 Nov. 2017, pugetsoundslammers.org/characteristics-that-make-a-good-soccer-player/.

[3]  Claywell, Charlie R. "What Is Social Network Theory?" *LoveToKnow*, LoveToKnow Corp, socialnetworking.lovetoknow.com/What-is-Social-Network-Theory.

[4]  Golbeck, Jennifer. "Analyzing Networks." *Introduction to Social Media Investigation*, Syngress, 20 Mar. 2015, www.sciencedirect.com/science/article/pii/B9780128016565000214.

[5]  Richeson, Dave. "How to Curve an Exam and Assign Grades." *David Richeson: Division by Zero*, 23 Mar. 2011, divisbyzero.com/2008/12/22/how-to-curve-an-exam-and-assign-grades/.

# 7    Appendix



Figure 18: Player Foul Efficiency



Figure 19: Player Duel Efficiency

Figure 20: Play Foul Efficiency



Figure 21: Player Shot Efficiency
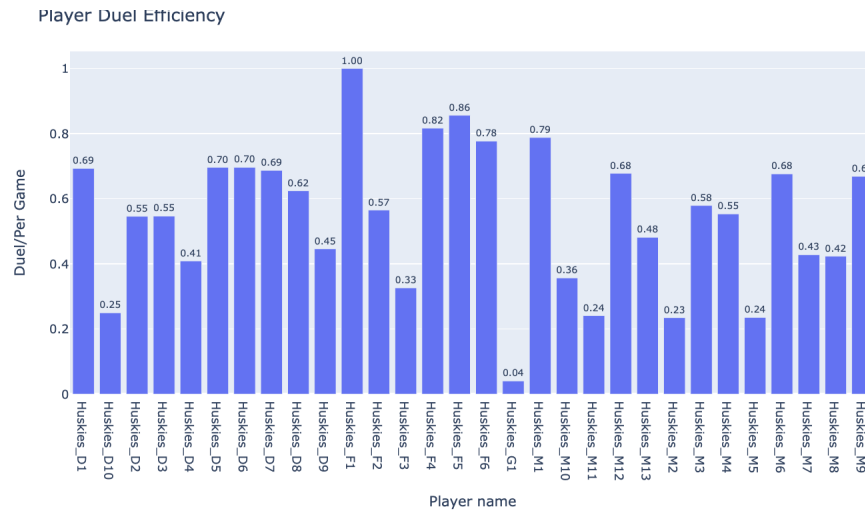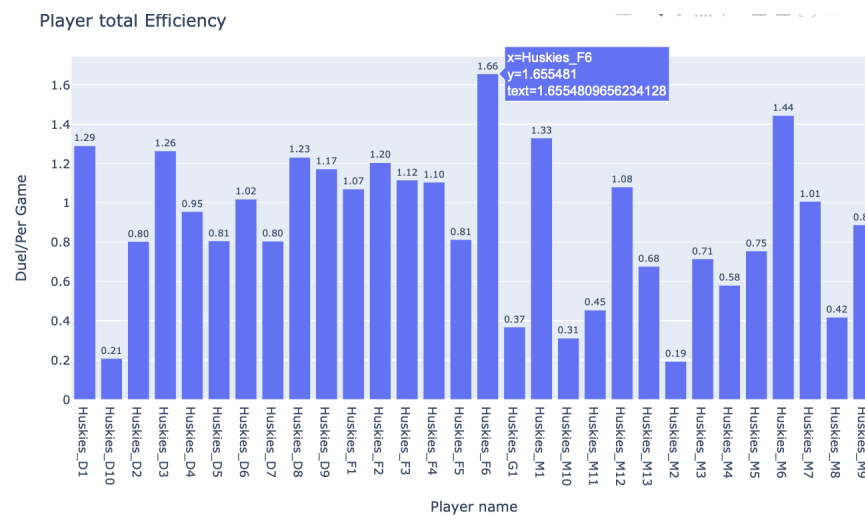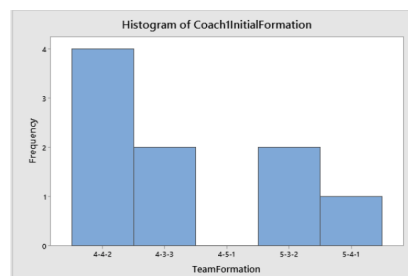
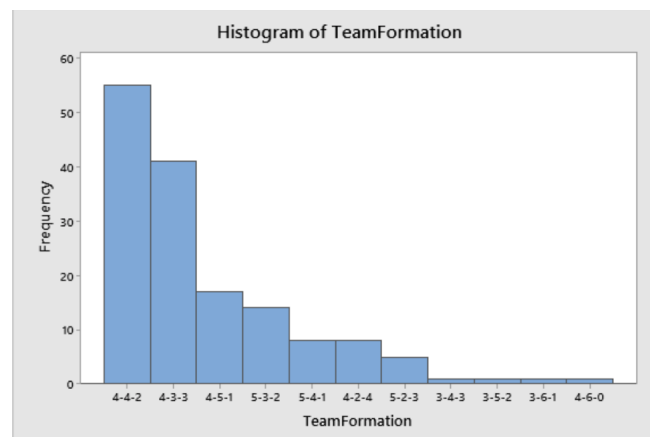Figure 22: Player Doul Efficiency



Figure 23: Player Total Efficiency


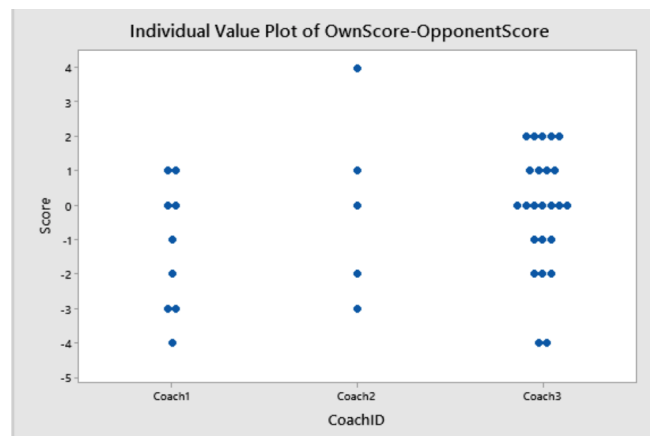
Figure 24: Coach1 Initial Formation
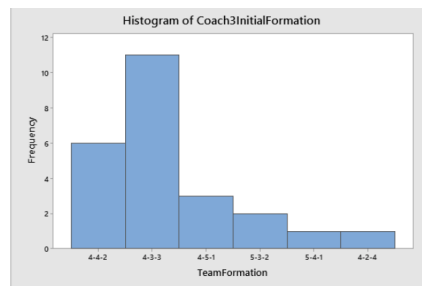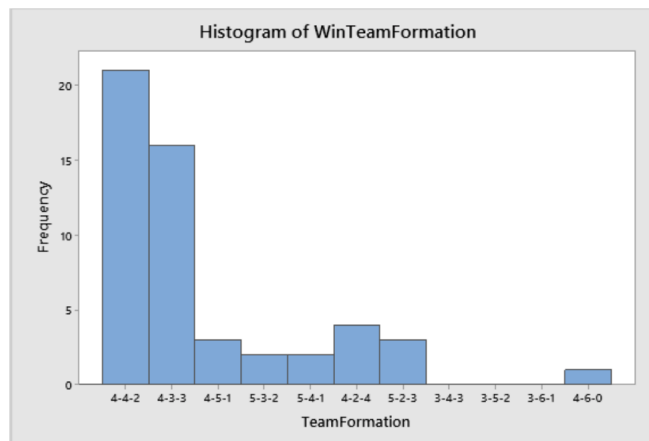
Figure 25: Coach3 Initial Formation



Figure 26: Coach Outcomes



Figure 27: Team Formation

Figure 28: Win Team Formation