
Big data bowl 2022

Massimo Hong
Student ID: 2020280082
Lohan Meunier
Student ID: 2021280356

Abstract

1 NFL is one of the biggest and most prestigious sport league in the United States,
2 with the highest average attendance of any sport in the world.
3 The kaggle competitions has a list of compulsory, but not limited, tasks to complete.
4 In our project we will also gather and filter data from the seasons 2018-2020,
5 analyse the various statistics and feed the obtained parameters to the neural network
6 in order to get the win probability of a team against a rival. In the end, we will give
7 the same inputs to a logistic regression model and compare the respective outputs.
8 Other than win probability, it could also be possible to predict whether a certain
9 play or strategy will be successfull or not.
10 A neural network *should* be more accurate than a logistic regression, that is essen-
11 tially just a "subset" of the former.

12 1 Introduction to NFL

13 The National Football League (NFL) is a professional American football league consisting of 32
14 teams, divided equally between the National Football Conference (NFC) and the American Football
15 Conference (AFC). The NFL is one of the four major North American professional sports leagues,
16 the highest professional level of American football in the world.

17 The NFL's eighteen-week regular season runs from early September to early January, with each team
18 playing seventeen games and having one bye week. Following the conclusion of the regular season,
19 seven teams from each conference (four division winners and three wild card teams) advance to the
20 playoffs, a single-elimination tournament culminating in the Super Bowl, which is usually held on
21 the first Sunday in February and is played between the champions of the NFC and AFC[1].

22 2 Dataset

23 The dataset is divided into multiple csv files:

- 24 • game.csv, contains teams playing involved in each game.
- 25 • plays.csv, contains information regarding the plays of a game.
- 26 • players.csv, contains information regarding the players.
- 27 • tracking.csv, contains player tracking for each season.
- 28 • PFFScoutingData.csv, contains play level scouting for each game.

29 3 Data analysis

30 We can use the gathered information to determine factors like:

- 31 • what play causes a team to win or lose the most points.
- 32 • the most important attacking player on the team.
- 33 • if the team has an advantage against the rival team based on the players, or the overall
- 34 playstyle (e.g. difference in height and weight can mean a lot in this sport).
- 35 • what could be the contribution of a certain player to the team.

36 All these factors can be used to better predict the outcome of a certain play, strategy, or the match
37 itself.

38 4 Required tasks of the competition

- 39 • Create a new special teams metric.
- 40 • Quantify special teams strategy. Special teams' coaches are among the most creative and
- 41 innovative in the league. Compare and contrast how each team game plans. Which strategies
- 42 yield the best results? What are other strategies that could be adopted[2]?
- 43 • Rank special teams players. Each team employs a variety of players (including longsnappers,
- 44 kickers, punters, and other utility special teams players). How do they stack up with respect
- 45 to one another[2]?

46 5 Which neural network to use

47 We will divide the data into training set and test set. The training set will be approximately 2/3 of the
48 total data set.

49 A possible choice for the neural network is the Probabilistic neural network (PNN), as these kind of
50 networks usually generate an accurate target probability score.

51 Another option is to use a classification network to predict whehter a play was successfull or not.

52 6 Regression Model

53 It is more natural to use a logistic regression approach rather than a linear model, even though the
54 output is not exactly binary. A few tweaks might be necessary in order to make it work correctly, such
55 as using odds instead of probabilities, for an easier interpretation. The advantages of a regression
56 model is its easy implementation and ready to use methods offered by the python *sklearn* library.

57 7 Conclusion

58 After getting the results from the neural network and the regression model, we can proceed to compare
59 them. If everything is implemented correctly, the neural network should have a higher accuracy score
60 than the regression model.

61 References

- 62 [1] Wikipedia. *National Football League*.
- 63 [2] Kaggle. *NFL Big Data Bowl 2022*.