

SCHOOL OF COMPUTING, ENGINEERING AND DIGITAL TECHNOLOGIES

MSC. DATA SCIENCE WITH ADVANCED PRACTICE

BAAH NKANSAH KWAKU

C2082617



BIG DATA AND BUSINESS INTELLIGENCE

CIS 4008-N

**ENGLISH PREMIER LEAGUE 2020/21 ANALYSIS
SECTION 1**

BUSINESS INTELLIGENCE REPORT

JANUARY 10, 2023

Contents

Executive Summary	2
Data source description and BI Requirements	4
Introduction	4
Data Source and Description	4
Rationale for choosing this dataset	6
BI Requirements\Questions	7
Finding based on analysis and evaluation	8
Key Findings	12
Recommendations	13
Data Pre-Processing or Data Cleansing	16
Loading the data	16
Data cleaning and Pre-processing	19
Removing Errors/Nulls/Columns	21
Creating and renaming columns	25
BI Data Modelling via Star Schema - Facts and Dimensions	26
Creating Dimension Tables	26
Adding index columns to fact table.	32
Duplicating and splitting columns.	36
Creating relationships.	39
DAX and M Language	43
Dashboard	46
Dashboard 1: Homepage	46
Dashboard 2: Club Performance Analysis.	47
Dashboard 3: Player Profile Analysis.	48
Dashboard 4: Positions Analysis	49
Dashboard 5: Nationality Analysis	50
Dashboard 6: AI Insights	51
REFERENCES	54

Executive Summary

Following the establishment and inauguration of the premier league in the 1992/93 season, the English premier league and the clubs who have participated in the league have attained significant and remarkable success over the years while establishing a formidable global reputation for the best standard and highly entertaining style of football. Conclusively, it could be said that the English premier league 2020/21 season was a defensive season more than an offensive and attacking one and this can be justified by making a comparative analysis between the total number of cards (1070) and the total number of assists (648) and the total number of goals scored (940) as well. Team management and coaches should ensure adequate number of trainings are done particularly by the team players and goalkeepers who are in the bottom half of the league.



Figure 1: card showing total assists



Figure 2: card showing total goals

Also, it could be generally asserted and assumed that players could be said to be in their prime at age twenty-seven (27) and age twenty-eight (28) because players ranging between these ages accounted for the most goals scored as well. Presumably, it could be assumed that the reason

for English players being the most in the English premier league is because it is a league that is being played in their home country although they are very high performing, hardworking and talented professionals who contribute immensely to the performance of their respective teams hence, they do have an upper hand in terms of getting easily scouted than other players who play in different countries.

Data source description and BI Requirements.

Introduction

Established and started in 1992, the English premier league is known to be the topmost domestic football league in the world, having the most cost-effective and high-income broadcasting right deals globally. The league is played annually by 20 teams who face each other twice within a calendar season, that is one home and one away fixture while having the top four on the league table qualify and compete to play in the UEFA champions league. The 2020/21 premier league season was a special, unpredictable and challenging season due to a global pandemic that the world faced, that is the Covid-19 pandemic. Despite all the challenges and hurdles, the teams and players and the premier league as a franchise provided fans and supporters with a fun, entertaining and action-packed season.

This BI project is tailored towards analysing the English premier league 2020/21 season and evaluating how the season played out and analysing club performance as well as player performance and involvement and contributions to a challenging but competitive and fun season.

Data Source and Description

The dataset used for this project is the English premier league 2020/21 which was obtained from Kaggle, a credible and widely accepted website by data analyst and professionals globally: <https://www.kaggle.com/datasets/rajatrc1705/english-premier-league202021>

The dataset contains one table with eighteen (18) columns and five hundred and thirty-three (533) rows which contains information on crucial stats about the English premier league 2020/21 edition. The dataset contains information on all the players that played in the English premier league and basic and standard stats about these players like the goals they scored, the assists they made, their expected goals and expected assists as well, passes attempted, pass accuracy and the teams these players play for as well as their nationality.

Column name	Description

Name	Name of the player
Club	Club the player plays for
Position	The position a player plays regularly
Nationality	The country the player comes from.
Age	The age of the player
Matches	The number of matches played by each player.
Starts	The number of times the player was named in the starting eleven (11).
Mins	The number of minutes played by the player.
Goals	The number of goals scored by the player
Assists	The number of times the player has assisted other players to score.
Passes attempted	The number of passes attempted by the player
Perc Passes Completed	The number of times that the player accurately passed to his teammate
Penalty goals	The number of penalties scored by the player
Penalty attempted	The number of penalties played by the player
xG	Expected number of goals from the player in a match
xA	Expected number of assists from the player in a match
Yellow cards	The number of yellow cards a player got for indiscipline, technical fouls, or other minor fouls
Red cards	The number of red cards a player got for accumulating two yellow cards or one major foul

Table: English premier league 2020/21 columns and descriptions.

Screenshots of the dataset are presented below:

The screenshot shows a Microsoft Excel spreadsheet titled "EPL_20_21.xlsx". The table has 38 rows and approximately 20 columns. The columns include: A1 (row header), Name, Club, Nationality, Position, Age, Matches, Starts, Mins, Goals, Assists, Passes, At Per cent, Pass Penalty, G Penalty, xG, xA, xB, xC, xD, xE, xF, xG, xH, xI, xJ, xK, xL, xM, xN, xO, xP, xQ, xR, xS, xT, xU, xV, xW, xX, xY, xZ, AA, AB, AC. The data consists of player statistics from the 2020-21 season, such as Mason Mount's 21 goals and 10 assists, or Timo Werner's 29 goals and 10 assists.

Figure 3: English premier league dataset

Rationale for choosing this dataset

This dataset was chosen for several reason but most importantly it was chosen for the following reasons:

- Challenging of skillsets and knowledge acquired: This dataset was selected because of how it challenged my knowledge acquired over the period of my studies and how it also put my acquired skillsets to the test. The dataset pushed me to limits because it tested my data cleaning and pre-processing skills as well my data analytics skills and most importantly how to analyse big data problems.
- Demonstration of data analytics skills: Another reason why this dataset was chosen was because the size of the dataset that is the number of columns and rows as well as the number of missing values aided in the demonstration of data analytics skills acquired so far. Also, it helped in developing my analytical skills using DAX and M language as well that is, using DAX to create custom columns and new measures.
- Credibility of the source of the dataset: The source of the dataset that is Kaggle, is a globally accepted website, that is acknowledged by millions and

professionals worldwide and in active use in academia and industry as well. So, the credibility of the source of the dataset played a major in the selection of this dataset.

- Interest and passion for subject topic and the accuracy of the dataset due to how recent it is: The interest and passion that I have for football also played a huge role in the selection of this dataset because I am a huge football fan and always take the opportunity to look further and analyse anything that has to do with football, so this played a role. And, how recent the dataset is also played a role. The dataset contains stats about the English premier league 2020/21 season. This means that it displays current information, and it is painting a picture of modern football.
- Easily understandable dataset: For any dataset to be analysed, the understanding of what the dataset entails is very essential to the analysis and addressing the problems to come up with solutions. The dataset gave sufficient information needed for the analysis and also had an extensive number of rows and columns which gave detailed insights and information as to what the whole dataset is about, and this contributed to making the dataset easier to analyse.

BI Requirements\Questions

The project aims to examine strategically the English premier league 2020/21 season while trying to unveil the relationship between the different characteristics of the football matches played and the teams' performances as well as their players contributions to the seasonal performance of the teams. The analysis is primarily focused on answering the following questions:

- At what age are players said to be in their prime based on the number of matches they play and their contribution in terms of goals and assists?
- Which teams had the best players in the 2020/21 season per their goals and assists?
- Which nationality had the most players playing in the English premier league 2020/21 season?
- Which positions contributed the most to their respective teams in terms of goals and assists?
- What metrics are used in determining who the best players were for the English premier league 2020/21 season?

Finding based on analysis and evaluation

A variety of different visualization were created, and they were all tailored towards answering the business questions asked earlier on in the report and these visualizations are presented below.

- At what age are players said to be in their prime based on the number of matches they play and their contribution in terms of goals and assists?

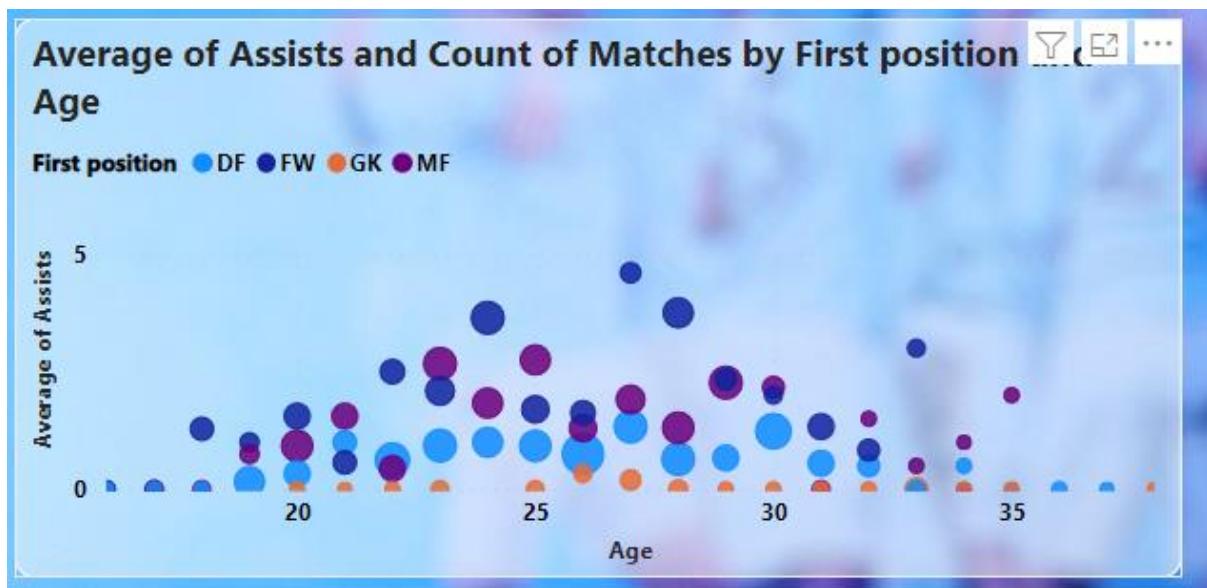


Figure 4: scatter plot from positions dashboard

A scatter chart showing at what age players average the most assist and get the most playing time with matches per the positions they play.

The assists, matches and age column of the EPL 20_21_22 table and positions column of the position column was used.

The scatter plot was used to represent this visual because it best represents two quantitative measures for different categories through the positioning of a single point and to determine if there is a linear relationship between the values.

- Which teams had the best players in the 2020/21 season per their goals and assists?

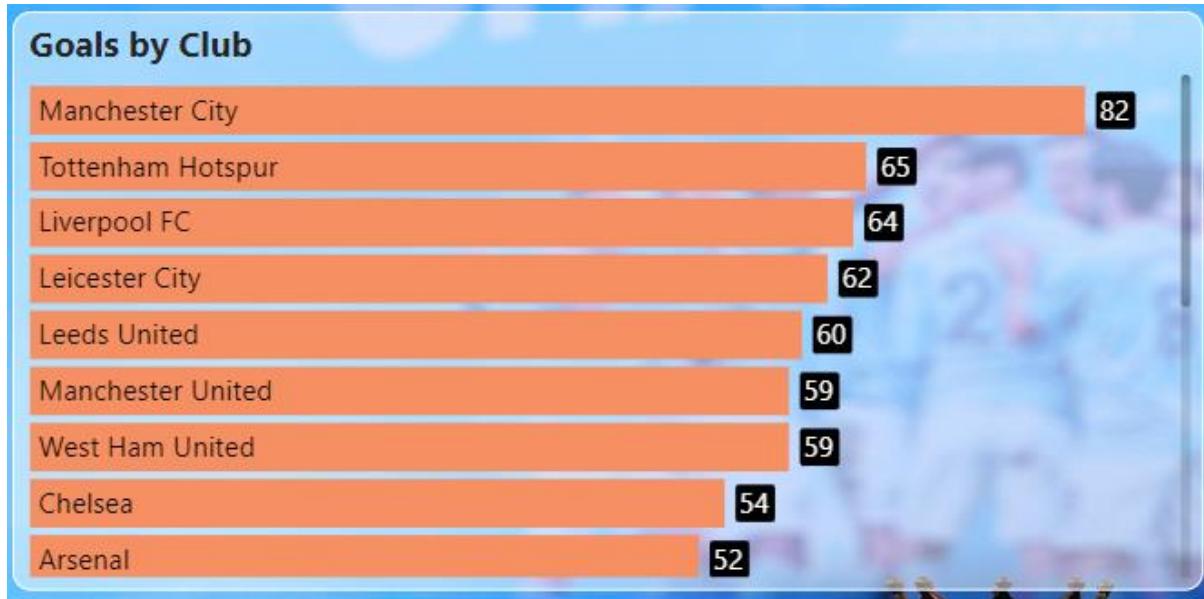


Figure 5: Bar chart for club performance analysis dashboard

A horizontal bar chart showing the number of goals that were scored by various clubs, this chart communicates the fact that Manchester city had the best players for the 2020/21 season judging from the number of goals they scored followed by Tottenham and Liverpool.

The goals column from the EPL_20_21_22 table was used as well as the club's column from the club table.

The bar chart was used because the chart is visualizing a data type that represents quantitative values against categories.

- Which nationality had the most players playing in the English premier league 2020/21 season?

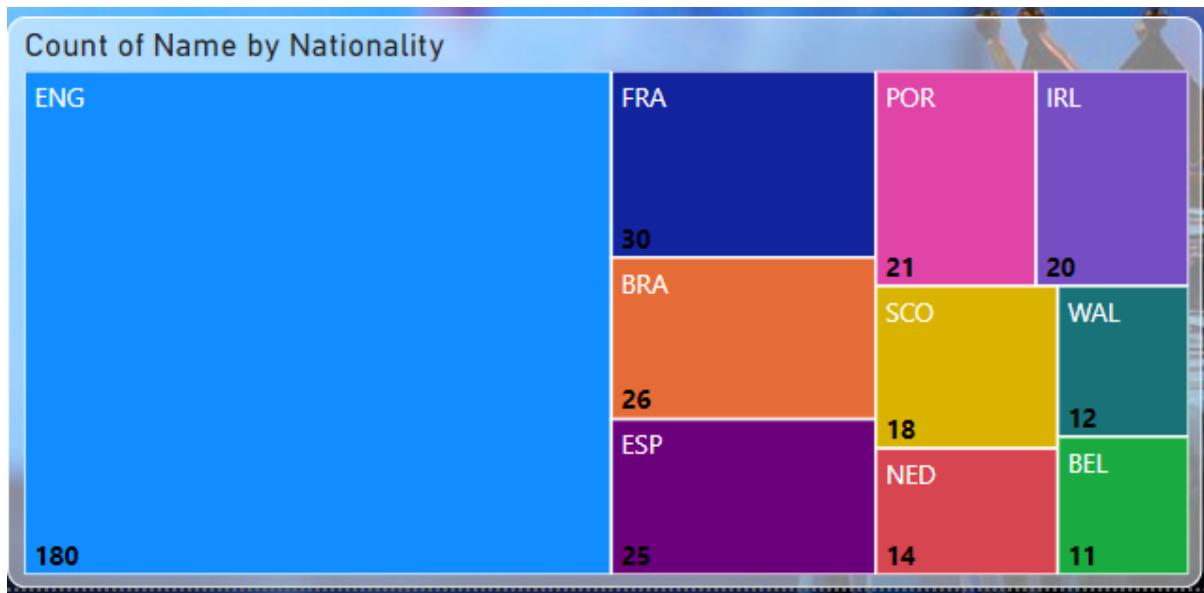


Figure 6: Treemap from Nationality analysis dashboard

A treemap representing the top ten (10) countries with the most players in the league for the 2020/21 season, and it can be inferred from the visualization that England has the most players in the league with one hundred and eighty (180) players.

The name column from the players table together with the nationality column from the nationality table were used in creating this chart.

The treemap was used in visualizing because the dataset represents quantitative values by displaying one constituent part of a whole. Colour has been used to indicate different categorical groupings.

- Which positions contributed the most to their respective teams in terms of goals and assists?



Figure 7: Stacked bar chart from the position analysis dashboard



Figure 8:drilldown radial bar chart from positions analysis

Two charts that is the stacked bar chart and the drilldown radial bar chart were used in displaying this information which is communicating the positions that contributed the most to the performance of their respective teams through goals and assists.

The goals column from the EPL_20_21_22 table and the position column from the position table were used for the stacked bar chart. And the assists column from the EPL_20_21_22 table and the position column from the position table were used in the creation of this visualization.

The bar chart was used because the chart is visualizing a data type that represents quantitative values against categories.

Key Findings

- Manchester city had the highest assists at 55, while Sheffield united had the lowest assist at 133. And out of the overall assist's percentage, Manchester city accumulated 8.49% of the general assists, and the assists range between all 20 clubs was 13-55.
- English premier league 2020/21 season winners Manchester city had the most goals with eighty-two (8)2 goals which shows that goals are very crucial in winning the league at the end of the season.
- Sheffield united scored the least goals nineteen (19), which shows they won few games which eventually led to them getting relegated.
- One thousand and seventy goals (1070) yellow cards were issued to players throughout the season which indicated that it was a very aggressive season.
- Players aged 27 and below could be said to be the playmakers of the league because they had the most assist at the end of the season.
- Forwards had the highest number of goals accounting for 59.26% of the total goals scored followed by midfielders.
- England had the most players with one hundred and eighteen (180) players playing in the league and with the most players the English players could be said to be the most aggressive players because they accumulated the most cards three hundred and eighty-eight (388) for the season.

- English players accounted for the 34.80% of the total number of cards issued during the English premier league 2020/21 season.
- There are players from fifty-seven (57) countries who played in the English premier league 2020/21 season.

Recommendations.

- Data on dates could have been provided to allocate a more detailed analysis report. With data on game weeks and matchdays, the analysis of this project could have been drilled down to the number of wins and losses per the games played and a cumulative column could have been created for top half and bottom half of the league table with details on individual match days.
- Teams should focus more on being attacking if they want to win the league or remain in the league because football has evolved, and it has been transformed into an attacking sport rather than a defensive one.
- Team management and coaches should ensure adequate number of trainings are done particularly by the team players and goalkeepers who are in the bottom half of the league.
- Sheffield United and Fulham should endeavour to develop and focus more on attacking prospects to strengthen their goal scoring ratio.
- Metrics used in calculating the existing xA and xG do not favour outfield player that is goalkeepers and defenders hence it should be looked at and revised to accommodate all players of the league.

Conclusion

This project has really helped me in developing and gaining extra knowledge not only on powerBi, but it has also broadened my scope on data analytics, and it has helped me gain much insight and skills in data analytics and pre-processing. Because I did not have prior knowledge and much understanding before this project, but this project has helped me develop and practice many skills that I never knew before.

I encountered multiple problems and challenges right from creating of missing values and nulls values, to data cleaning/pre-processing, and then to creation of new columns and new measures using DAX expressions and M language. During the development and exploration of visualizations, I realized that there were blanks existing in some of my dimension tables when I used a slicer to visualize the dataset, but these blanks were not present when I looked at the dataset in the power query editor and the data view mode on power BI. I removed blanks from rows, removed duplicates and removed errors as well but the blanks were still present. I then went the extra mile to research and watch videos online on how to get rid of these blanks from my dataset, but all of these efforts did not have any effect on the blanks. This is the main and major challenge that I was faced with during the development of my project, and I was unable to rectify this problem. But with the help of my module leader and my lab tutor and some individual research through the internet and YouTube videos, I was able to go through and surpass all the challenges and I look forward to applying and the knowledge gained, and the skills acquired in the coming semesters and during my professional experience as an analyst as I know these skills acquired will go a long way to help me in the industry.

SCHOOL OF COMPUTING, ENGINEERING AND DIGITAL TECHNOLOGIES

MSC. DATA SCIENCE WITH ADVANCED PRACTICE

BAAH NKANSAH KWAKU

C2082617



BIG DATA AND BUSINESS INTELLIGENCE

CIS 4008-N

ENGLISH PREMIER LEAGUE 2020/21 ANALYSIS

APPENDIX: BUSINESS INTELLIGENCE DESIGN

SECTION 2

JANUARY 10, 2023

Data Pre-Processing or Data Cleansing

Loading the data

The first step in the data pre-processing and cleaning stage of this analysis is to load the dataset that was used for the project, in this case the English premier league dataset in Microsoft PowerBI. This can be done primarily in two ways which are, selecting the file type and format from “Get Data” menu and then going ahead to import it or to create a blank query and have M – Language imports the dataset. For the purpose of this project, the dataset was imported by selecting “Get Data”.

The Microsoft powerBi tool launched and opened.

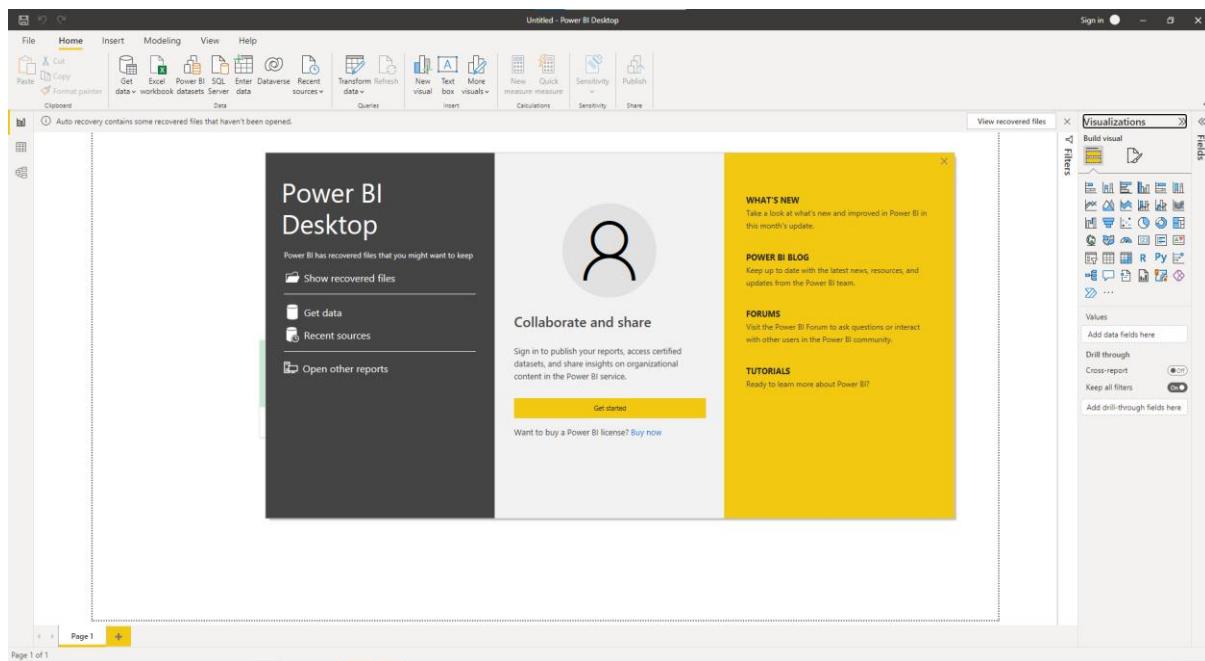


Figure 9: launching PowerBi tool

The data is loaded by clicking on “Get data” from the Home Tab of Microsoft powerBi. After clicking on “Get data”, a dropdown containing different formats and

data sources appear. The dataset for this analysis is in the comma separated values (csv. files) format, hence select “Text/csv” from the dropdown that appears.

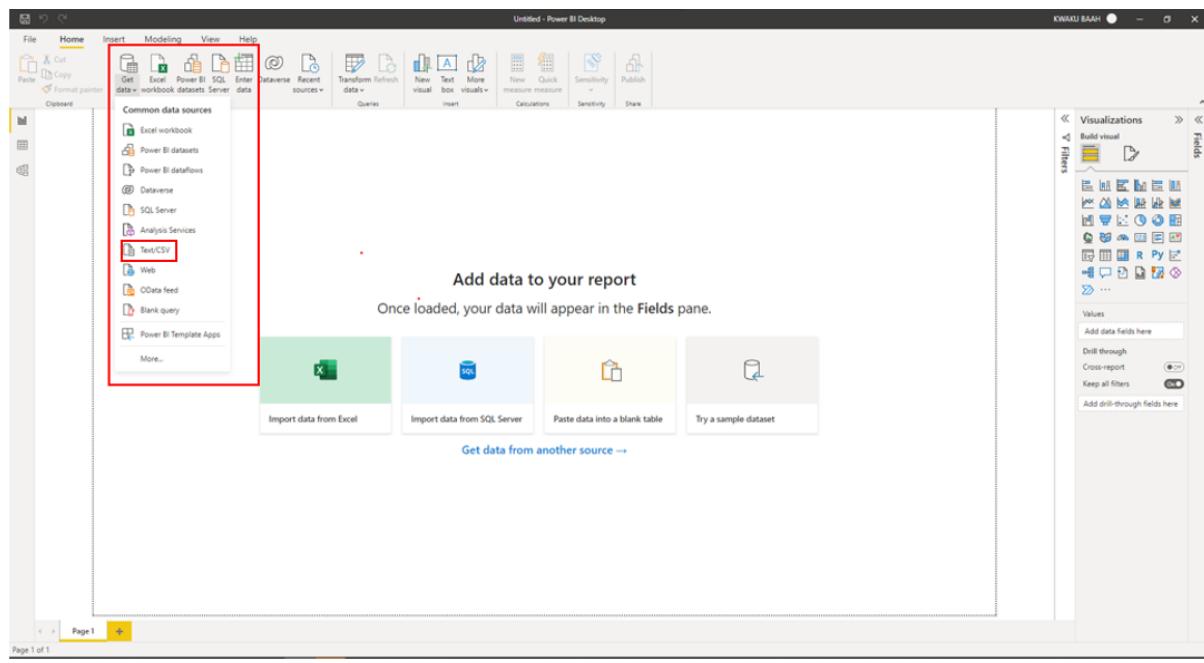


Figure 10: loading the dataset.

This opens a dialogue box that displays the English premier league data with the option to load, transform data or cancel.

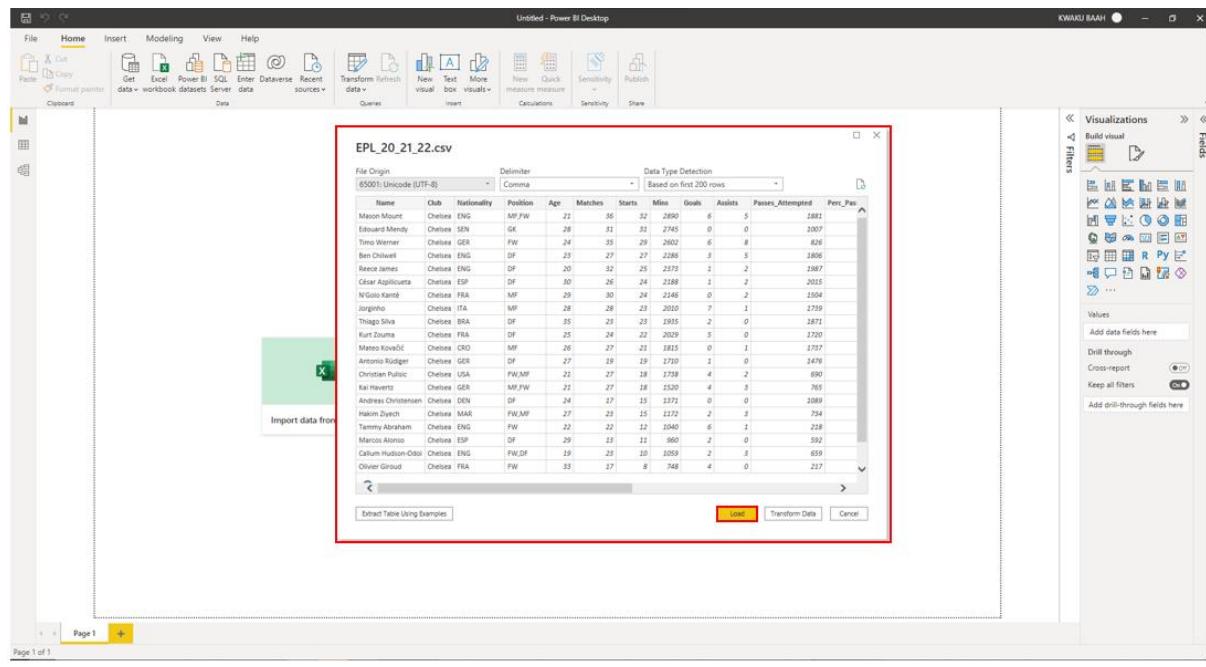


Figure 11: Get data dialogue box.

The process of loading the data was completed after the load option was selected in the previous dialogue box.

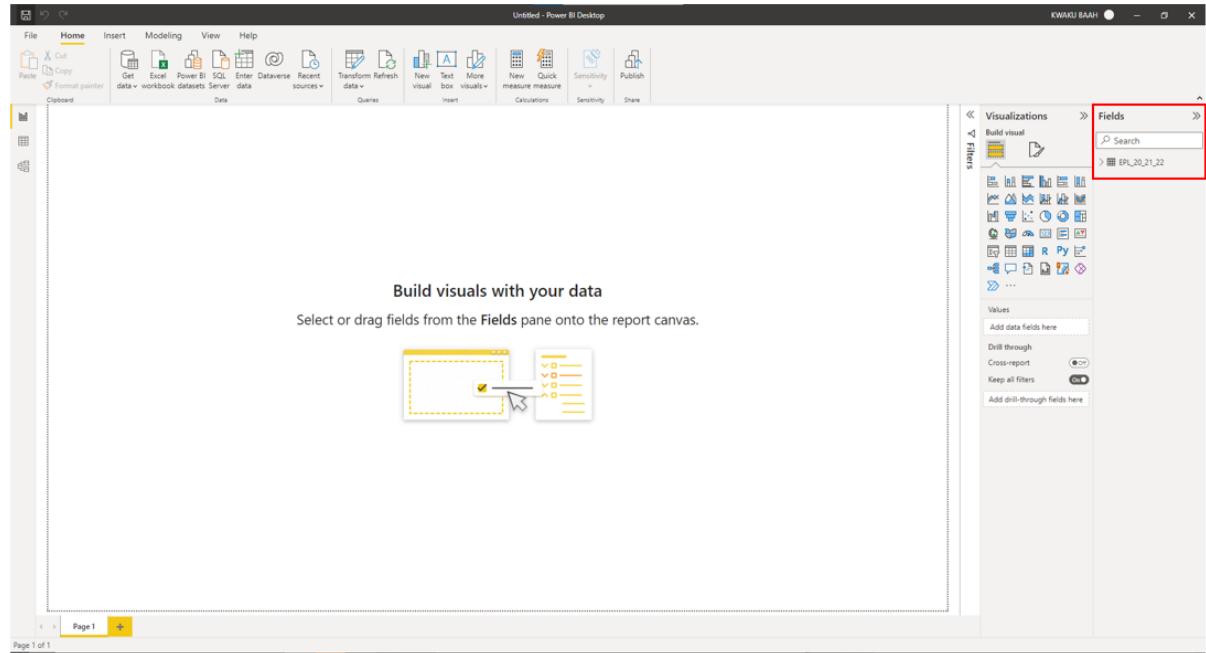


Figure 12: Dataset loaded

Data cleaning and Pre-processing

The next step is the cleaning and pre-processing of the dataset that was just imported and loaded into the Microsoft powerBi tool. The cleaning and pre-processing of the dataset is done using the power query editor. To get the power query editor, click on “Transform data” in the home Tab. Then select “Transform data” from the dropdown that displays

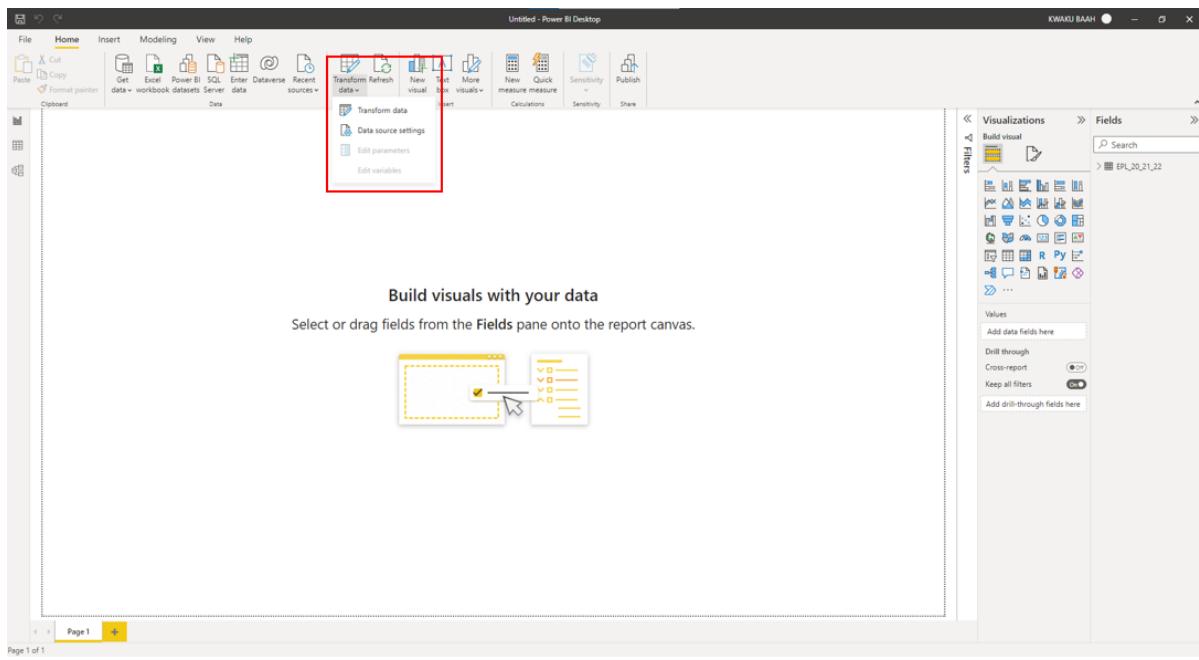


Figure 13: Transform data for power query editor

This opens the power query editor dialogue box.

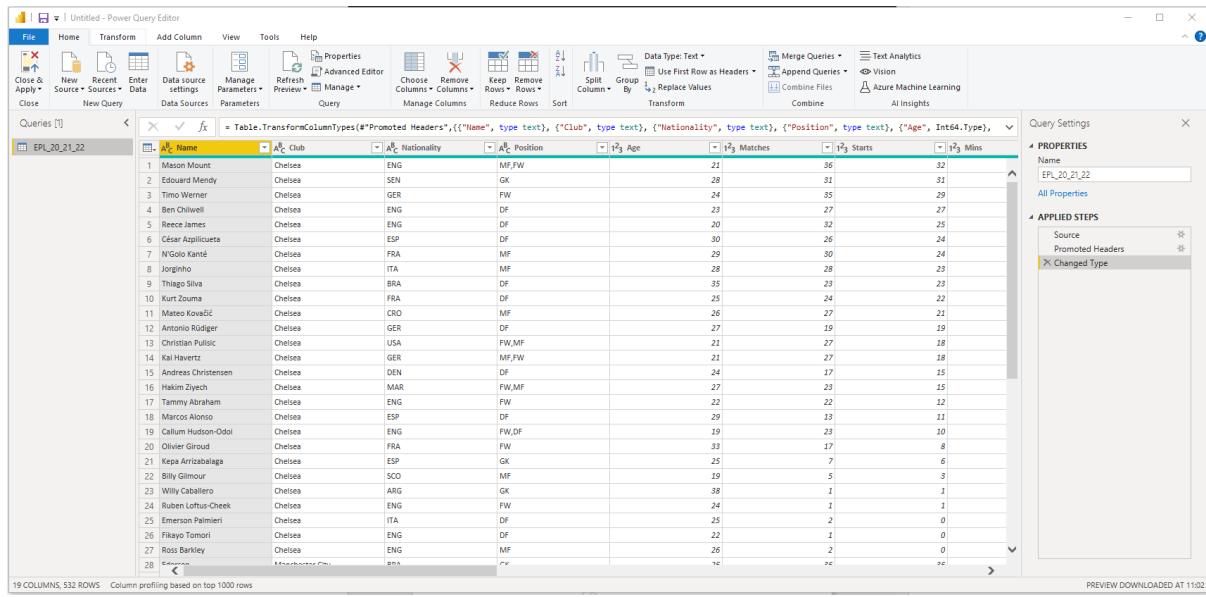


Figure 14: Power query editor dialogue box

Removing Errors/Nulls/Columns

The player ID column was contained errors as initially detected by powerBi. Eventually, the root cause of the error was revealed, and it was because the player ID column contained qualitative values when it should have contained quantitative values ideally and powerBi was unable to rectify these errors automatically.

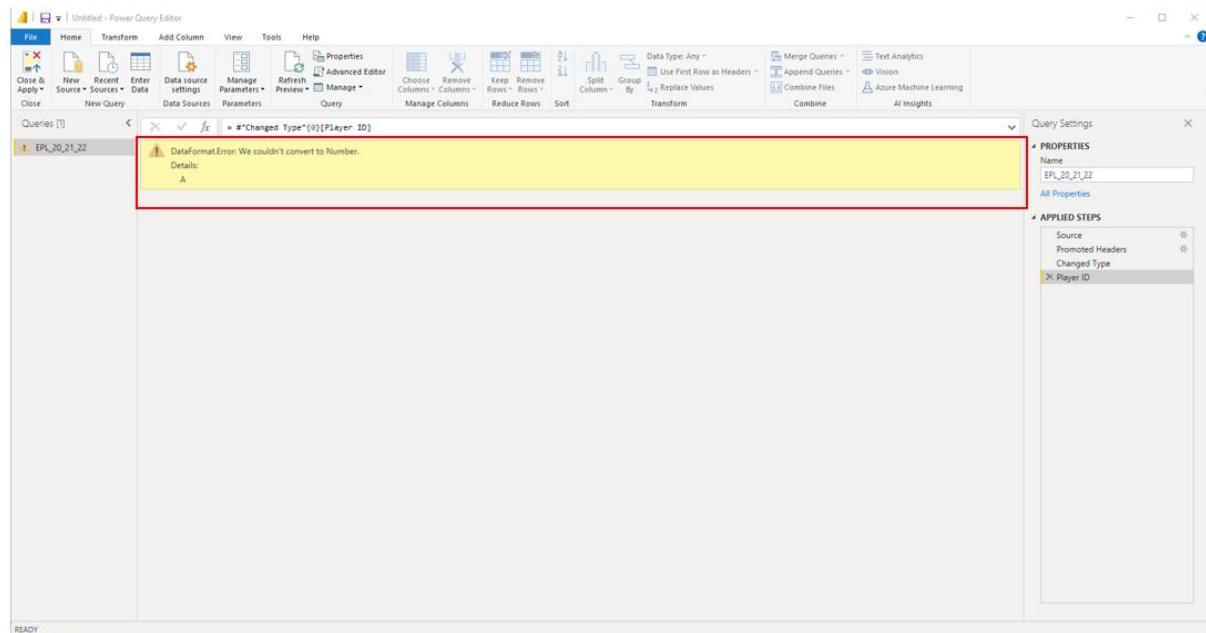


Figure 15: Errors detected by PowerBi.

The column was then removed from the dataset as the absence of this column from the dataset was not going to have a huge impact.

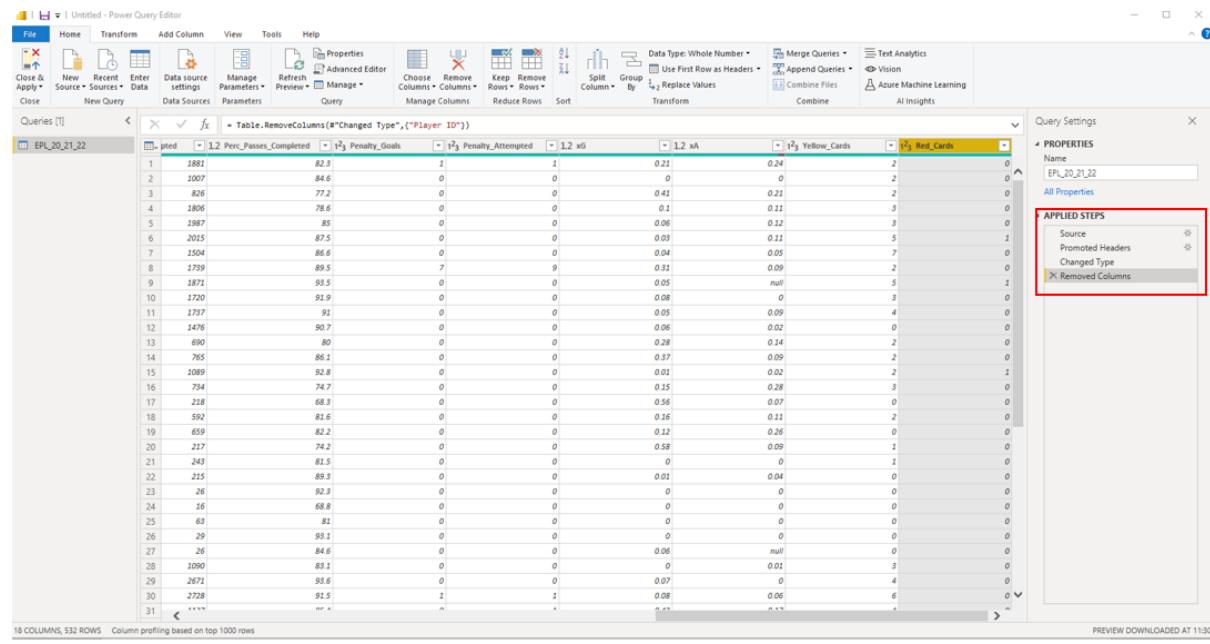


Figure 16: Removing columns

Empty rows as well as rows containing null values and blanks were also removed because the percentages that these null values and empty rows covered in the dataset was not marginal and they were not really essential for the analysis of this project and their presence in the dataset was not going to contribute marginally. Null values in the xA column were removed. This was done by right clicking on the arrow next to the xA column and unticking the “null” values from the dropdown that displayed.

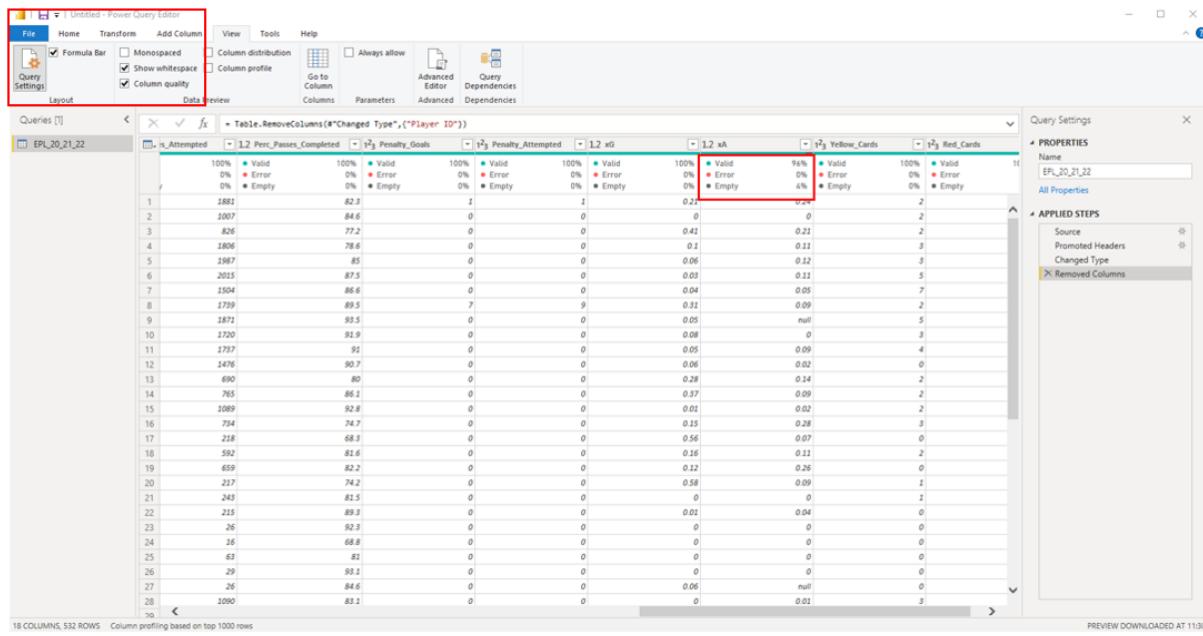


Figure 17: Removing blanks

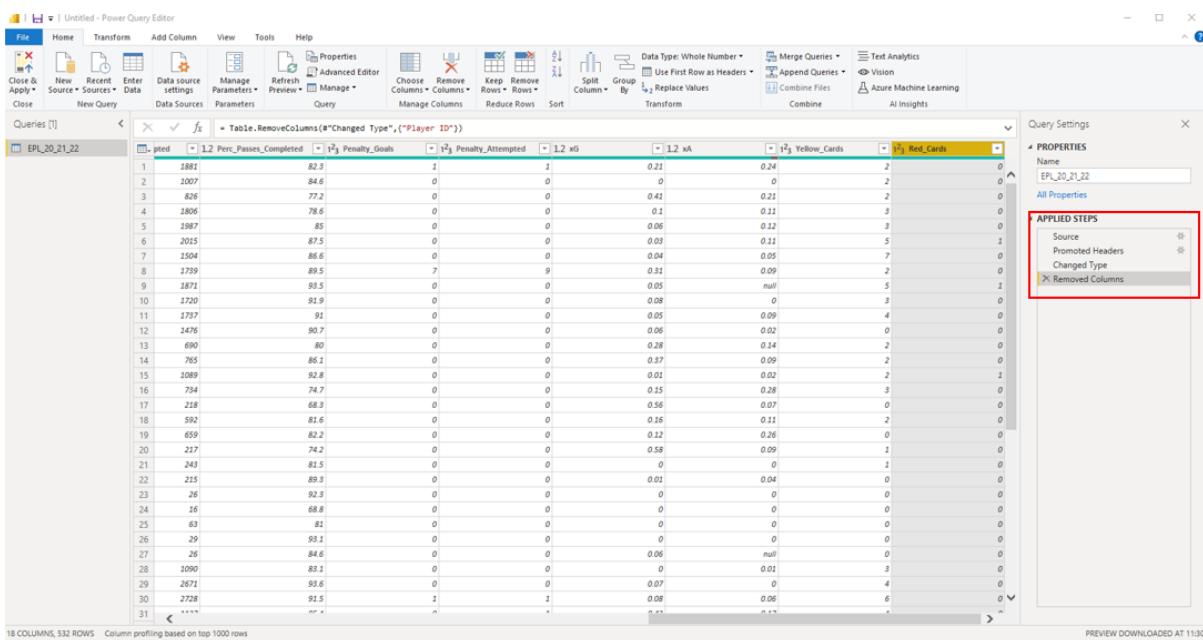


Figure 18: Errors removed

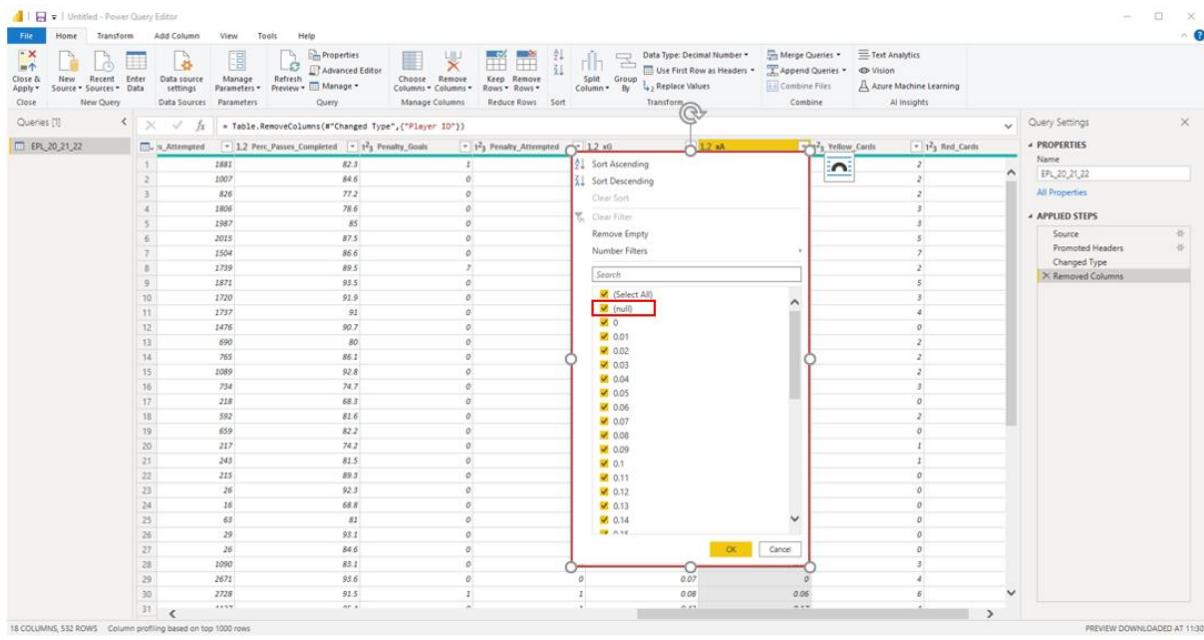


Figure 19: Unticking null values

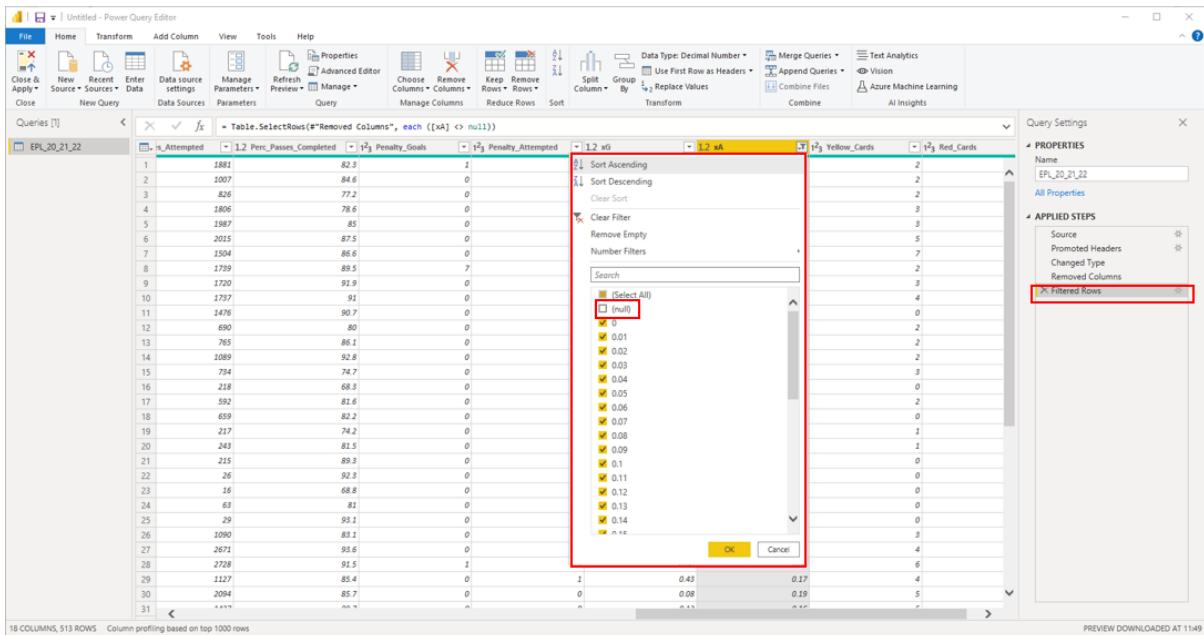


Figure 20: null values removed.

Creating and renaming columns

The index column was created and renamed as ID. This was done by clicking on “Add column” from the home tab and selecting Index column from the Add column tab and choosing “from 1” from the pop-up box that appears. No other columns were deleted as all the columns in the dataset were needed for the analysis

The screenshot shows the Microsoft Power Query Editor interface. The ribbon at the top has tabs for File, Home, Transform, Add Column, View, Tools, and Help. The 'Add Column' tab is highlighted with a red box. The main area displays a table with 18 columns and 513 rows, representing player statistics. The columns are labeled: AP₁ Name, AP₂ Club, AP₃ Nationality, AP₄ Position, AP₅ Age, AP₆ Matches, AP₇ Starts, and AP₈ Miss. The data includes players like Mason Mount, Edouard Mendy, Timo Werner, etc. On the right side, there are sections for 'Query Settings', 'PROPERTIES' (Name: EPL_20_21_22), and 'APPLIED STEPS'. The 'APPLIED STEPS' section shows a step named 'Filtered Rows'. At the bottom left, it says '18 COLUMNS, 513 ROWS' and 'Column profiling based on top 1000 rows'. At the bottom right, it says 'PREVIEW DOWNLOADED AT 11:49'.

Figure 21: Adding index column

The screenshot shows the Microsoft Power Query Editor interface. The top menu bar includes File, Home, Transform, Add Column, View, Tools, and Help. The 'Transform' tab is selected, with the 'Add Column' dropdown open, showing options like 'Conditional Column', 'Index Column', and 'Duplicate Column'. The main area displays a table with 513 rows and 19 columns, representing EPL 20_21_22 data. The columns include 'Penalty_Goals', 'Penalty_Attempted', 'Yellow_Cards', 'Red_Cards', and an 'Index' column. The bottom right corner of the editor shows 'PREVIEW DOWNLOADED AT 12:10'. On the right side, there's a 'Query Settings' pane with sections for 'PROPERTIES' (Name: EPL_20_21_22) and 'APPLIED STEPS'. The 'APPLIED STEPS' section is expanded, showing a list of steps: 'Source', 'Promoted Headers', 'Changed Type', 'Removed Columns', 'Filtered Rows', and 'Added Index', with the 'Added Index' step highlighted by a red box.

Figure 22: Index column added.

BI Data Modelling via Star Schema - Facts and Dimensions.

After the data cleaning and pre-processing had been done, the next step was to create dimensions tables and establish a relationship between the dimension tables and the fact table.

Creating Dimension Tables.

The club dimension table was the first-dimension table that was created. In order to create this, we reopened the power query editor, and duplicate the EPL 2020/21 table. After the table had been duplicated, the duplicated table was renamed “club”, only the club column was to be used for this table hence all the other columns were removed. This was done by selecting the club column and right clicking and choosing to remove other columns.

Queries [2] fx = Table.TransformColumnTypes(#"Promoted Headers", {{"Name", type text}, {"Club", type text}, {"Nationality", type text}, {"Position", type text}, {"Age", Int64.Type}, {"Matches", Int64.Type}, {"Starts", Int64.Type}, {"Minutes", Int64.Type}, {"Goals", Int64.Type}})

Club

Name	Club	Nationality	Position	Age	Matches	Starts	Minutes	Goals
Mason Mount	Chelsea	ENG	MF/FW	21	36	32	2890	
2 Edward Mendy	Chelsea	SEN	GK	28	31	31	2745	
3 Timo Werner	Chelsea	GER	FW	24	35	29	2602	
4 Ben Chilwell	Chelsea	ENG	DF	23	27	27	2386	
5 Reece James	Chelsea	ENG	DF	20	32	29	2373	
6 César Aspitarte	Chelsea	ESP	DF	30	26	24	2188	
7 Ngolo Kanté	Chelsea	FRA	MF	29	30	24	2146	
8 Jorginho	Chelsea	ITA	MF	28	28	23	2010	
9 Thiago Silva	Chelsea	BRA	DF	35	23	23	1845	
10 Kurt Zouma	Chelsea	FRA	DF	25	24	22	2029	
11 Mateo Kovacic	Chelsea	CRO	MF	28	27	21	1815	
12 Antonio Rudiger	Chelsea	GER	DF	27	19	19	1710	
13 Christian Pulisic	Chelsea	USA	FW/MF	21	27	18	1738	
14 Kai Havertz	Chelsea	GER	MF/FW	21	27	18	1520	
15 Andreas Christensen	Chelsea	DEN	DF	24	17	15	1371	
16 Hakim Ziyech	Chelsea	MAR	FW/MF	27	23	15	1272	
17 Tammy Abraham	Chelsea	ENG	FW	22	22	12	1040	
18 Marcos Alonso	Chelsea	ESP	DF	29	13	11	960	
19 Callum Hudson-Odoi	Chelsea	ENG	FW/DF	19	23	10	1059	
20 Oliver Giroud	Chelsea	FRA	FW	33	17	8	748	
21 Kepa Arrizabalaga	Chelsea	ESP	GK	25	7	6	550	
22 Billy Gilmour	Chelsea	ENG	MF	19	5	3	261	
23 Willy Caballero	Chelsea	ARG	GK	38	1	1	90	
24 Ruben Loftus-Cheek	Chelsea	ENG	FW	24	1	1	60	
25 Emerson Palmieri	Chelsea	ITA	DF	25	2	0	90	
26 Fábio Tomori	Chelsea	ENG	DF	22	1	0	45	
27 Ross Barkley	Chelsea	ENG	MF	26	2	0	42	
28 Ederson	Manchester City	BRA	GK	28	36	36	3240	
29 Ruben Dias	Manchester City	POR	DF	23	32	32	2843	
30 Rodri	Manchester City	ESP	MF	24	34	31	2748	
31 Raheem Sterling	Manchester City	ENG	FW	25	31	28	2556	
32 João Pedro	Manchester City	POR	DF	26	28	27	2299	
33 Bernardo Silva	Manchester City	POR	MF/FW	25	26	24	2065	
34 İkay Göndogan	Manchester City	GER	MF	29	28	23	2029	
35 Kevin De Bruyne	Manchester City	BEL	MF	29	25	23	1997	
36 Riyad Mahrez	Manchester City	ALG	FW	29	27	23	1849	
37 Gabriel Jesus	Manchester City	BRA	FW	23	29	22	2063	
38 Kyle Walker	Manchester City	ENG	DF	30	24	22	1946	

19 COLUMNS, 532 ROWS Column profiling based on top 1000 rows PREVIEW DOWNLOADED AT 10:07

Figure 23: EPL 2020/21_22 table duplicated and renamed

Queries [2] fx = Table.TransformColumnTypes(#"Promoted Headers", {{"Name", type text}, {"Club", type text}, {"Nationality", type text}, {"Position", type text}, {"Age", Int64.Type}, {"Matches", Int64.Type}, {"Starts", Int64.Type}, {"Minutes", Int64.Type}, {"Goals", Int64.Type}})

Club

Name	Club	Nationality	Position	Age	Matches	Starts	Minutes	Goals
Mason Mount	Chelsea	ENG	MF/FW	21	36	32	2890	
2 Edward Mendy	Chelsea	SEN	GK	28	31	31	2745	
3 Timo Werner	Chelsea	GER	FW	24	35	29	2602	
4 Ben Chilwell	Chelsea	ENG	DF	23	27	27	2386	
5 Reece James	Chelsea	ENG	DF	20	32	25	2373	
6 César Aspitarte	Chelsea	ESP	DF	30	26	24	2188	
7 Ngolo Kanté	Chelsea	FRA	MF	29	30	24	2146	
8 Jorginho	Chelsea	ITA	MF	28	28	23	2010	
9 Thiago Silva	Chelsea	BRA	DF	35	23	23	1845	
10 Kurt Zouma	Chelsea	FRA	DF	25	24	22	2029	
11 Mateo Kovacic	Chelsea	CRO	MF	28	27	21	1815	
12 Antonio Rudiger	Chelsea	GER	DF	27	19	19	1710	
13 Christian Pulisic	Chelsea	USA	FW/MF	21	27	18	1738	
14 Kai Havertz	Chelsea	GER	MF/FW	21	27	18	1520	
15 Andreas Christensen	Chelsea	DEN	DF	24	17	15	1371	
16 Hakim Ziyech	Chelsea	MAR	FW/MF	27	23	15	1272	
17 Tammy Abraham	Chelsea	ENG	FW	22	22	12	1040	
18 Marcos Alonso	Chelsea	ESP	DF	29	13	11	960	
19 Callum Hudson-Odoi	Chelsea	ENG	FW/DF	19	23	10	1059	
20 Oliver Giroud	Chelsea	FRA	FW	33	17	8	748	
21 Kepa Arrizabalaga	Chelsea	ESP	GK	25	7	6	550	
22 Billy Gilmour	Chelsea	ENG	MF	19	5	3	261	
23 Willy Caballero	Chelsea	ARG	GK	38	1	1	90	
24 Ruben Loftus-Cheek	Chelsea	ENG	FW	24	1	1	60	
25 Emerson Palmieri	Chelsea	ITA	DF	25	2	0	90	
26 Fábio Tomori	Chelsea	ENG	DF	22	1	0	45	
27 Ross Barkley	Chelsea	ENG	MF	26	2	0	42	
28 Ederson	Manchester City	BRA	GK	28	36	36	3240	
29 Ruben Dias	Manchester City	POR	DF	23	32	32	2843	
30 Rodri	Manchester City	ESP	MF	24	34	31	2748	
31 Raheem Sterling	Manchester City	ENG	FW	25	31	28	2556	
32 João Pedro	Manchester City	POR	DF	26	28	27	2299	
33 Bernardo Silva	Manchester City	POR	MF/FW	25	26	24	2065	
34 İkay Göndogan	Manchester City	GER	MF	29	28	23	2029	
35 Kevin De Bruyne	Manchester City	BEL	MF	29	25	23	1997	
36 Riyad Mahrez	Manchester City	ALG	FW	29	27	23	1849	
37 Gabriel Jesus	Manchester City	BRA	FW	23	29	22	2063	
38 Kyle Walker	Manchester City	ENG	DF	30	24	22	1946	

19 COLUMNS, 532 ROWS Column profiling based on top 1000 rows PREVIEW DOWNLOADED AT 10:07

Figure 24: columns removed

The duplicates in the club column were removed by right clicking on club column and selecting “remove duplicates”.

The screenshot shows the Power Query Editor interface with the 'Club' table selected. A context menu is open over the table, with the 'Remove Duplicates' option highlighted. The 'APPLIED STEPS' pane on the right shows the step 'Removed Other Columns' has been removed.

Figure 25: Duplicates removed from the club table

A custom column was created for club ID using the index column. This was done by clicking “Add column” and selecting Index column and choosing “from 1”. This new column was renamed Club ID.

The screenshot shows the Power Query Editor interface with the 'Club' table selected. A red box highlights the newly added 'Index' column, which contains numerical values from 1 to 20 corresponding to each club name. The 'APPLIED STEPS' pane on the right shows the step 'Added Index' has been applied.

Figure 26: Index column created for club ID

The screenshot shows the Power Query Editor interface with the following details:

- File**, **Home**, **Transform**, **Add Column**, **View**, **Tools**, **Help** menu items.
- Toolbars for **Columns**, **Formulas**, **Invoke Custom Functions**, **General**.
- Buttons for **From Text**, **From Number**, **From Date & Time**.
- Icons for **Trigonometry**, **Statistics**, **Standard Scientific**, **Rounding**, **Date**, **Time**, **Duration**, **Information**, **Text**, **Analytics**, **Vision**, **Azure Machine Learning**, **AI Insights**.
- Queries [3]** pane: EPL_20_21_22, Club, Nationality.
- Club** table preview:

	Club	Club ID
1	Chelsea	1
2	Manchester City	2
3	Manchester United	3
4	Liverpool FC	4
5	Leicester City	5
6	Watford United	6
7	Tottenham Hotspur	7
8	Arsenal	8
9	Leeds United	9
10	Everton	10
11	Arton Villa	11
12	Newcastle United	12
13	Wolverhampton Wanderers	13
14	Crystal Palace	14
15	Southampton	15
16	Brighton	16
17	Burnley	17
18	Fulham	18
19	West Bromwich Albion	19
20	Sheffield United	20
- Query Settings** pane: Name set to Club.
- APPLIED STEPS** pane: Shows steps like Promoted Headers, Changed Type, Removed Other Columns, Removed Duplicates, Added Index, and Renamed Columns.
- Bottom status bar: 2 COLUMNS, 20 ROWS, PREVIEW DOWNLOADED AT 1034.

Figure 27: Index column renamed as club ID

The second-dimension table created was the nationality table. This was also done by duplicating the EPL 2020/21 table and renaming it “Nationality”. All columns were removed except for the nationality column because it was the only column needed for this table. This was done by selecting the nationality column and right clicking on it and selecting remove other columns. The duplicates in the nationality column were also removed, this was done by right clicking on the nationality column and selecting remove duplicates.

The screenshot shows the Power Query Editor interface with the 'Club' table selected. A context menu is open over the first row of the table, with the 'Duplicate' option highlighted. The 'APPLIED STEPS' pane on the right shows a single step named 'Promoted Headers' with the 'Changed Type' checkbox checked.

Figure 28: Duplicating the EPL table to create nationality table

The screenshot shows the Power Query Editor interface with the 'Club' table selected. A context menu is open over the first row of the table, with the 'Remove' option highlighted. The 'APPLIED STEPS' pane on the right shows a single step named 'Promoted Headers' with the 'Changed Type' checkbox checked.

Figure 29:removing other columns

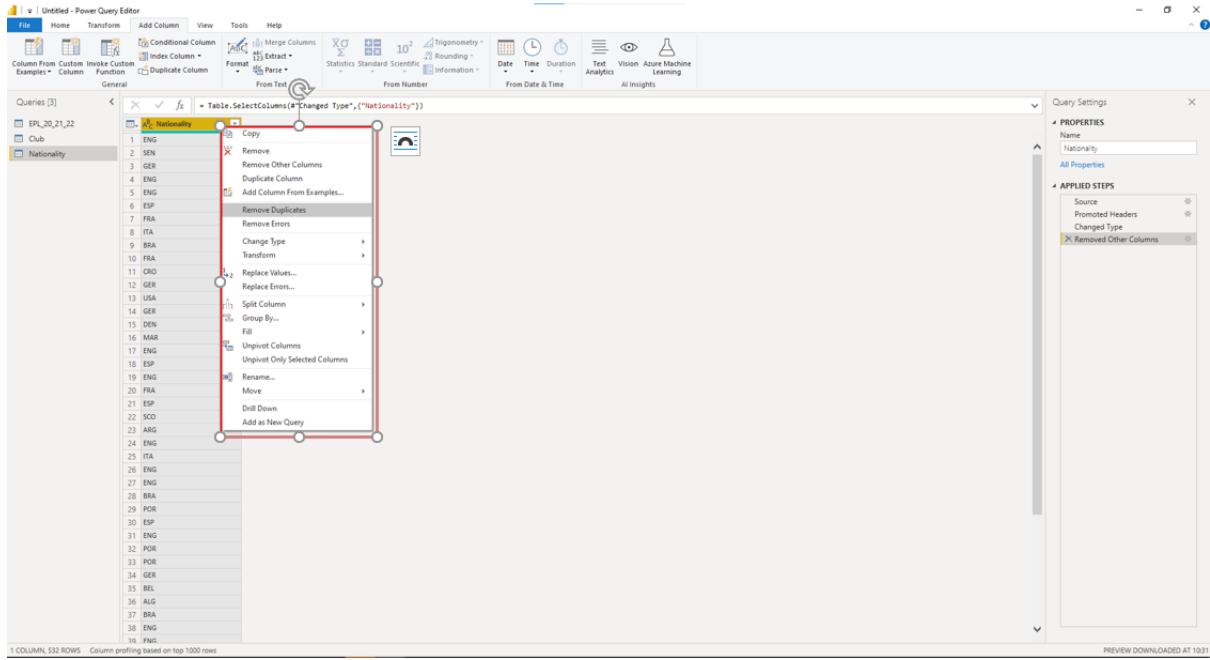


Figure 30: Removing duplicates

An ID column was also added to this table by clicking on “Add column” on the home tab and selecting “index” and choosing “from 1” in the pop-up box that appears. This column was renamed “Nationality Id”.

Figure 31: Creating index for nationality ID

The screenshot shows the Microsoft Power Query Editor interface. In the center, there is a table named 'Nationality' with two columns: 'Index' and 'Nationality ID'. The 'Index' column lists country abbreviations like ENG, SEN, GER, ESP, FRA, ITA, BRA, CRO, USA, DEN, MAR, SCO, ARG, POR, BEL, ALG, UKR, NED, SWE, UGD, SRB, WAL, CIV, NGA, EGY, TUR, CMR, GUI, SUR, JPN, IRL, GRE, NNR, GHA, AUT, AAM, ARA, CZE, and PKR. The 'Nationality ID' column contains numerical values ranging from 1 to 39. A red box highlights the 'Nationality ID' column. On the right side, the 'Query Settings' pane shows the 'Name' field set to 'Nationality'. The 'APPLIED STEPS' pane lists the steps taken: 'Source', 'Promoted Headers', 'Changed Type', 'Removed Other Columns', 'Removed Duplicates', 'Added Index', and 'Renamed Columns'. The status bar at the bottom indicates '2 COLUMNS, 59 ROWS' and 'Column profiling based on top 1000 rows'.

Figure 32: renaming index as club ID

This process was repeated and used to create two more dimension tables, Players, and positions tables.

Adding index columns to fact table.

The index (Id) columns that were created for all the dimension tables were added to the fact table in order to be able to create relationships. The dimension tables had to be merged with the fact table in order to make it possible. This was done by clicking on “Merge Queries” on the home tab and creating a connection between them on their shared columns. for the first-dimension table, that is the club dimension table, the shared column was club, hence the club column was selected on both the dimension table and the fact table using the “left outer” as the join kind. After this had been done, the fact table and club dimension table are merged, by clicking on “ok”. The illustration is detailed below

The screenshot shows the Power Query Editor interface with a query named 'EPL_20_21_22'. The 'Merge Queries' step is highlighted with a red box in the ribbon bar. The main pane displays a table with 532 rows and 19 columns, representing player statistics. The columns include Name, Club, Nationality, Position, Age, Matches, Starts, and Minutes. The 'APPLIED STEPS' pane on the right shows the 'Changed Type' step applied to the query.

Figure 33: Merging queries for the fact table and dimension table.

Merge

X

Select a table and matching columns to create a merged table.

EPL_20_21_22



Name	Club	Nationality	Position	Age	Matches	Starts	Mins	Goals	Assists	Passe
Mason Mount	Chelsea	ENG	MF,FW	21	36	32	2890	6	5	
Edouard Mendy	Chelsea	SEN	GK	28	31	31	2745	0	0	
Timo Werner	Chelsea	GER	FW	24	35	29	2602	6	8	
Ben Chilwell	Chelsea	ENG	DF	23	27	27	2286	3	5	

Club



Club	club ID
Chelsea	1
Manchester City	2
Manchester United	3
Liverpool FC	4
Leicester City	5

Join Kind

Left Outer (all from first, matching from second)

Use fuzzy matching to perform the merge

↳ Fuzzy matching options

✓ The selection matches 532 of 532 rows from the first table.

OK

Cancel

Figure 34: Club column merged for the fact table and the club dimension table.

Now, the column is expanded, and the club column is unticked leaving only the club Id column and this is done by right clicking on the arrow next to the new column. The column is renamed as “club Id”

The screenshot shows the Power Query Editor interface with the 'Club' column expanded. A red box highlights the 'Club' checkbox in the 'Select All Columns' section of the 'Applied Steps' pane. The 'Properties' pane shows the query name as 'EPL_20_21_22'. The 'Applied Steps' pane lists 'Source', 'Promoted Headers', 'Changed Type', and 'Merged Queries'. The main table view shows columns like 'Jals', 'Penalty_Attempted', 'Yellow_Cards', 'Red_Cards', 'Player ID', and 'Club'. The 'Club' column contains values such as 1, 0, 0, 0, 0, etc.

Figure 35: Club ID column expanded

The screenshot shows the Power Query Editor interface with the 'Club' column renamed. A red box highlights the 'Renamed Columns' section in the 'Applied Steps' pane. The 'Properties' pane shows the query name as 'EPL_20_21_22'. The 'Applied Steps' pane lists 'Source', 'Promoted Headers', 'Changed Type', 'Merged Queries', and 'Expanded Club.1'. The main table view shows the renamed 'Club' column, which now contains values like '0 A', '0 A', '0 A', '0 A', etc. The table also includes columns for 'Jals', 'Penalty_Attempted', 'Yellow_Cards', 'Red_Cards', 'Player ID', and the renamed 'Club' column.

Figure 36: Club ID column renamed.

This process is repeated for all the other three (3) dimension tables.

The next step to do after repeating for all the other dimension tables is, to remove all the columns that are in the dimension tables from the fact table with the exception of the IDs and this is done by pressing and holding the control button and selecting all these columns and right clicking one of the columns and selecting “remove columns” from the dropdown box that appears.

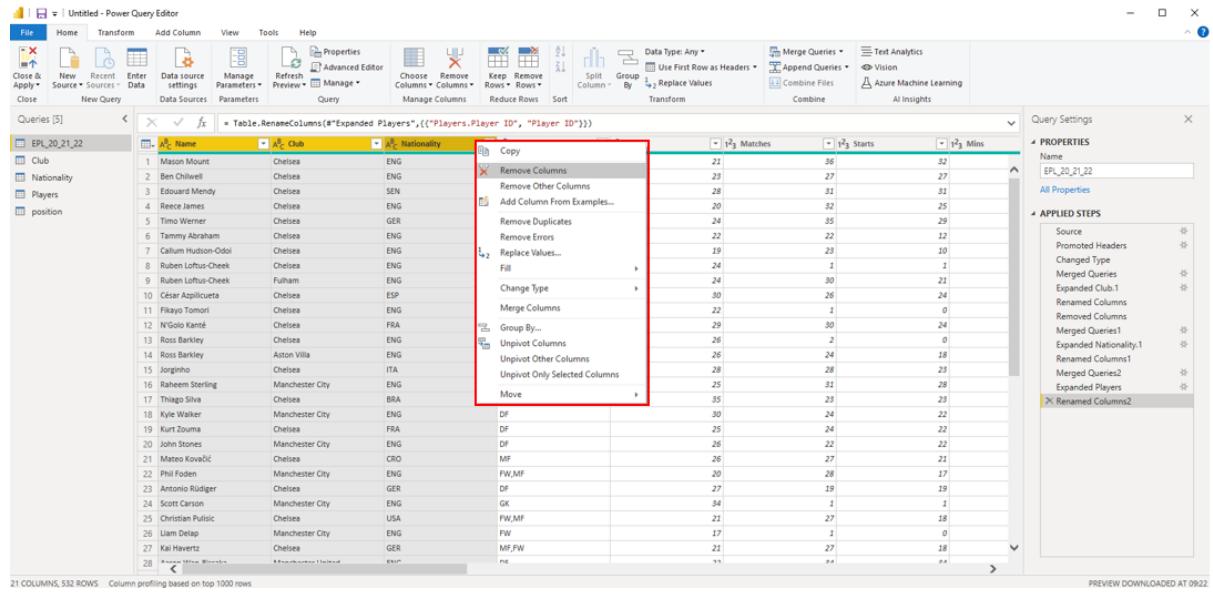


Figure 37: Removing columns that are in the dimension table that is also in the fact table.

Duplicating and splitting columns.

For the positions dimension table, the position column was duplicated by right clicking on the position column and selecting duplicate column from the dropdown box that displays. The column was then split into two using the delimiter, this was done by selecting the “split using delimiter” on the home tab. After the dialogue box opens, select split column using comma, then click “ok”. The first position and second position columns were then created.

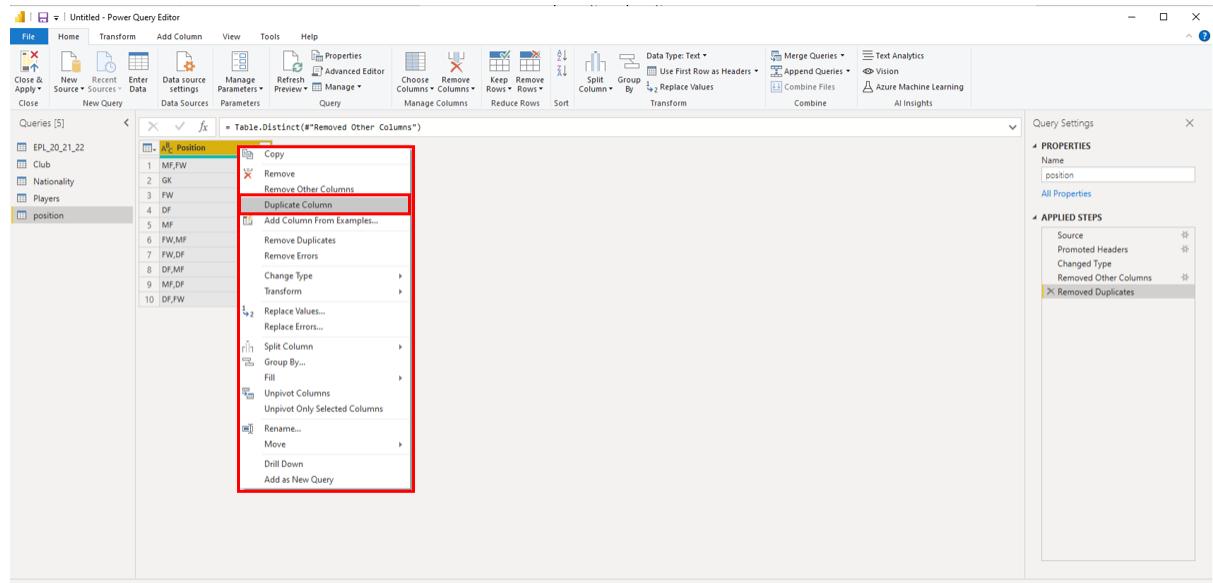


Figure 38:Duplicating positions column

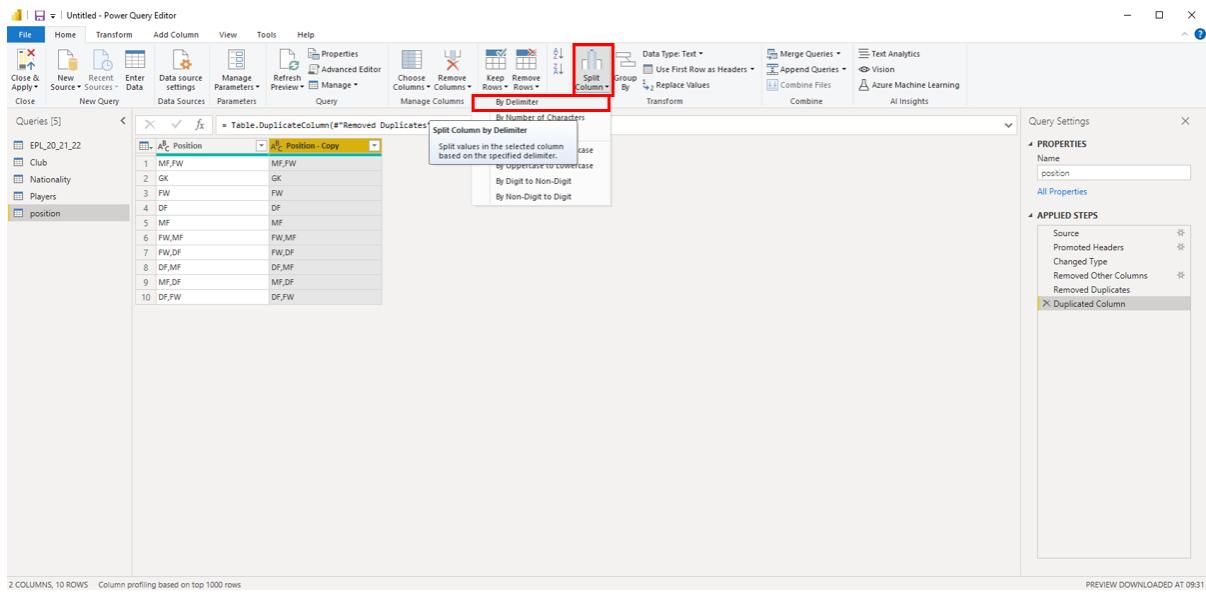


Figure 39: Selecting “split column by” on the Home Tab.

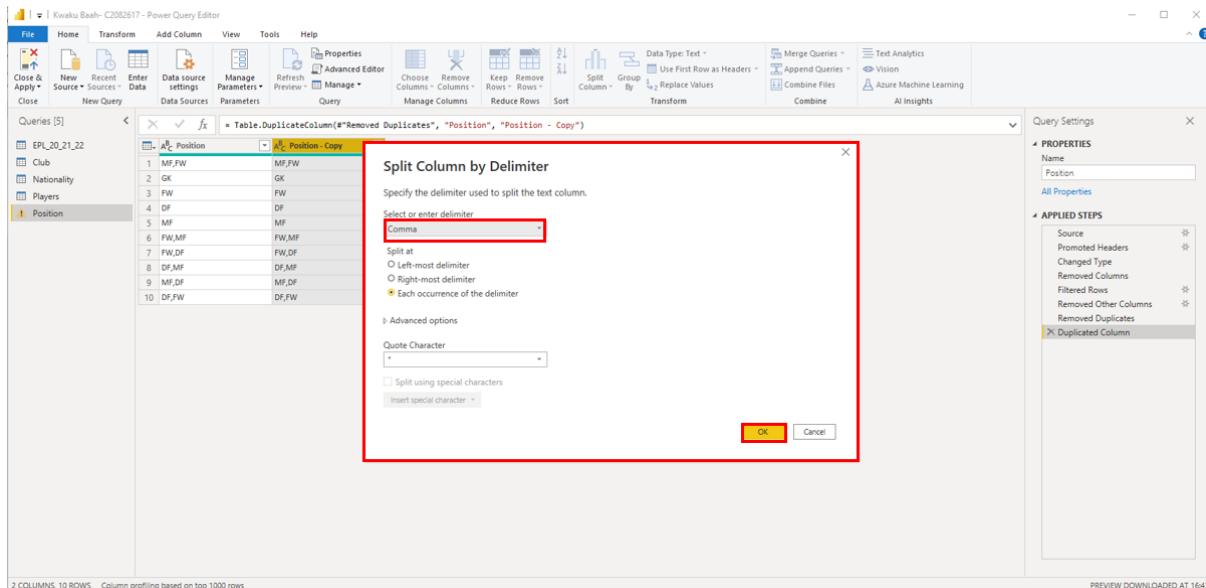


Figure 40: Positions column split using “Split column using delimiter”

For the secondary position column, the null values were replaced with N/A. this was done by right clicking on secondary position column and selecting replace values from the dropdown box that appears. Now, the null values were replaced with N/A in the “replace values” dialogue box that appeared.

The screenshot shows the Power Query Editor interface with the 'Transform' tab selected. A context menu is open over the 'Position' column, specifically over the 'second position'. The menu path 'Replace Values...' is highlighted with a red box. Other options visible in the menu include Copy, Remove, Remove Other Columns, Duplicate Column, Add Column From Examples, Remove Duplicates, Remove Errors, Change Type, Transform, Replace Values..., Replace Errors..., Split Column, Group By..., Fill, Unpivot Columns, Unpivot Other Columns, Unpivot Only Selected Columns, Rename..., Move, Drill Down, and Add as New Query.

Figure 41:Selecting replace values from the dropdown.

The screenshot shows the Power Query Editor interface with the 'Transform' tab selected. A 'Replace Values' dialog box is open over the 'Position' column. The 'Value To Find' field contains 'null' and the 'Replace With' field contains 'NA'. The 'OK' button is highlighted with a red box. The dialog also includes an 'Advanced options' section and a 'Cancel' button.

Figure 42:Replacing null values with N/A.

Creating relationships.

In order to properly analyse and visualise this project and meet the requirements of this project, relationships had to be created between the fact table and the dimension tables. This was done by selecting “manage relationships” from the home tab of the model view of the powerBi and “new” is selected from the dialogue box that appears.

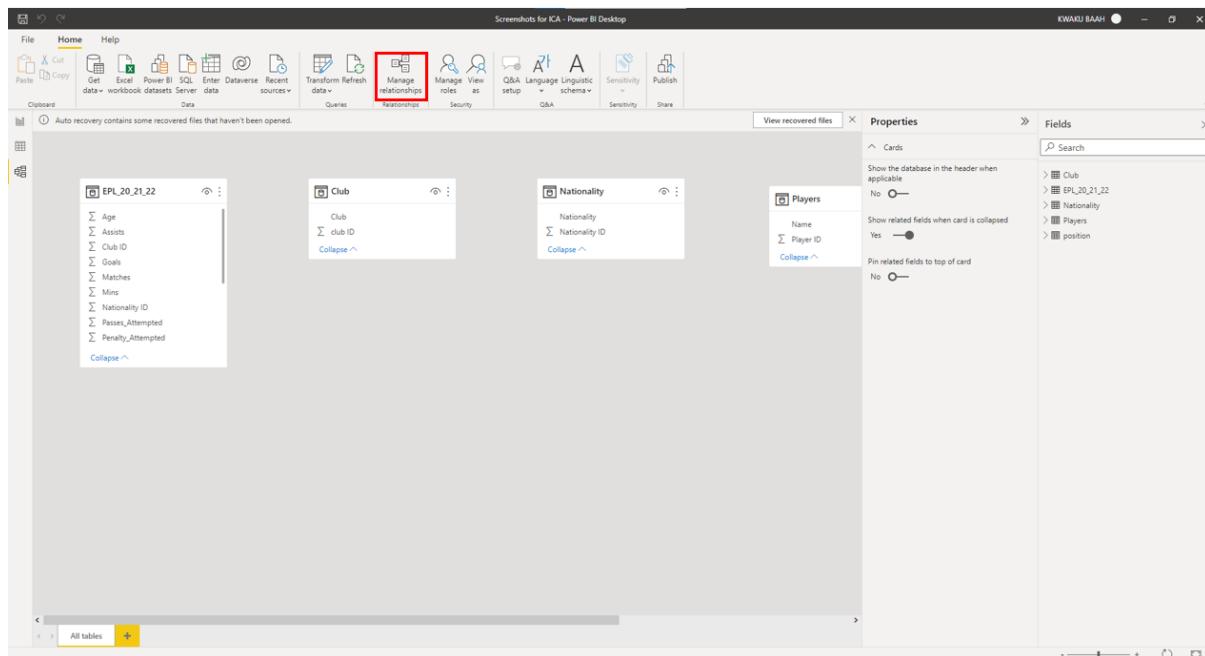


Figure 43:Selecting to manage relationships.

Manage relationships

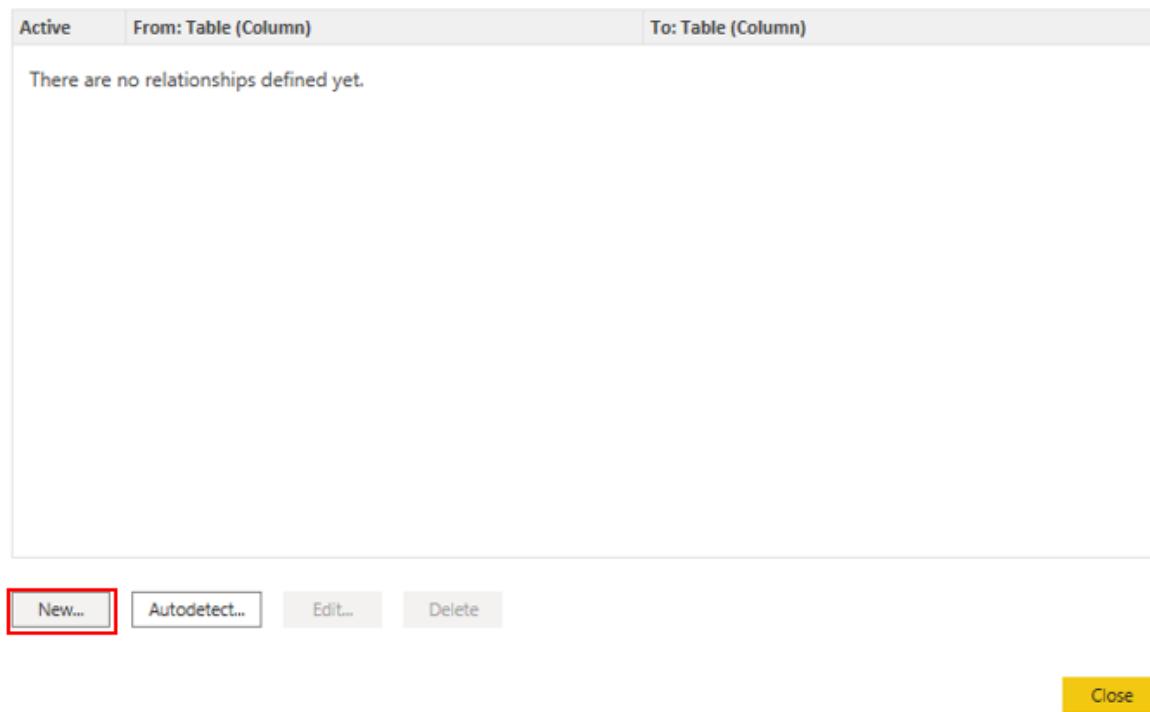


Figure 44:Select “new” from the manage relationships dialogue box.

The first relationship is between the “EPL 2020/21” that is the fact table, and the “club” dimension table is connected by the “club id” column. A relationship with cardinality “many to one” and a “single” cross filter direction is created between the fact table and dimension table.

Create relationship

X

Select tables and columns that are related.

EPL_20_21_22

ses_Completed	Penalty_Goals	Penalty_Attempted	xG	xA	Yellow_Cards	Red_Cards	Club ID	Name
84.6	0	0	0	0	0	2	0	1
68.8	0	0	0	0	0	0	0	1
93.1	0	0	0	0	0	0	0	1

Club

Club	club ID
Chelsea	1
Manchester City	2
Manchester United	3

Cardinality Cross filter direction

Many to one (*;1) Single

Make this relationship active Apply security filter in both directions

Assume referential integrity

OK Cancel

Figure 45:Relationship created between the fact table and club dimension table.

Again, this process is repeated for all the three-dimension tables as well.

A star schema model, which is where one fact table is connected to many dimension tables is the type of data that was developed after all the relationships were created.

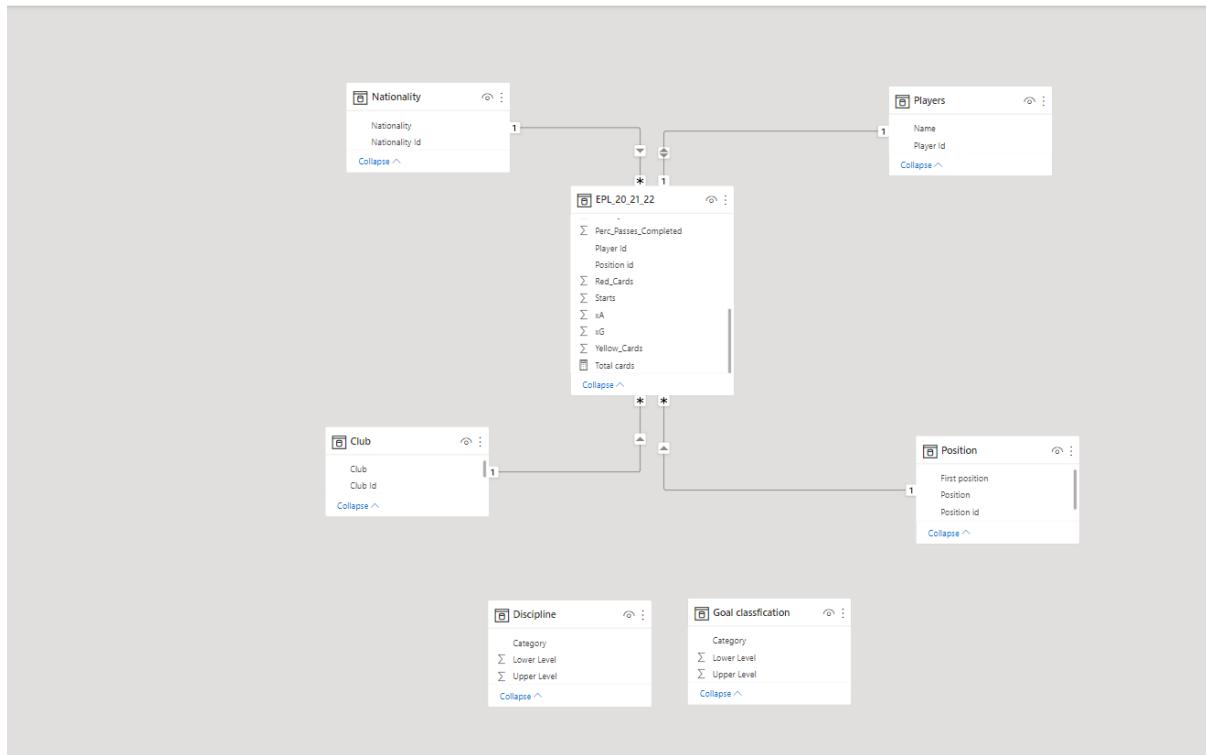


Figure 46:Star Schema model created for relationships.

DAX and M Language

The creation of DAX was the last stage in the data pre-processing stage, and DAX was used to generate two (2) measures and a column that is, Total club goals column, average goals per club measure and Total cards measure. This column and measures were created in order to go a step in simplifying the analysis and making the analysis easier to understand as well using their metrics. M languages were also created using the blank query and these were done so that the goal scorers could be categorised in four types in order to classify and be able to tell who the very best are, and it was also used in generating a table that told a story about how aggressive players were for this particular season judging by the number of cards they accumulated over the season.

The first DAX that was created was the Total club goals column and this was created using the DAX expression below.

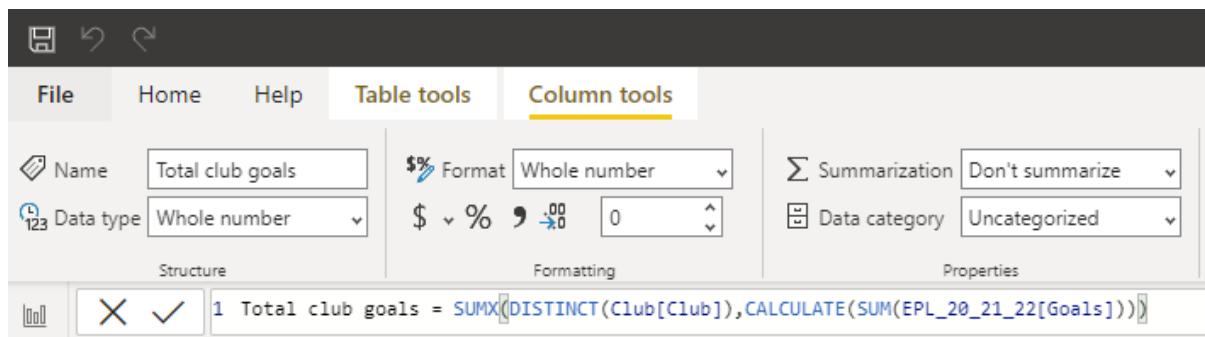


Figure 47:DAX expression for total club goals.

The second DAX that was created was the Average goals per club measure and this was accomplished using the DAX expression below,

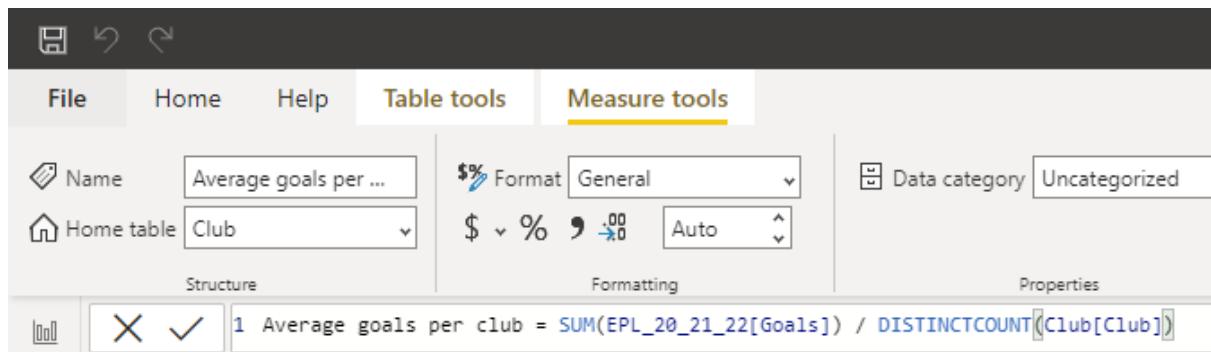


Figure 48:DAX expression for Average club goals.

Finally, total cards measure was the third and last DAX that was created and the expression below was the DAX expression used in creating it.

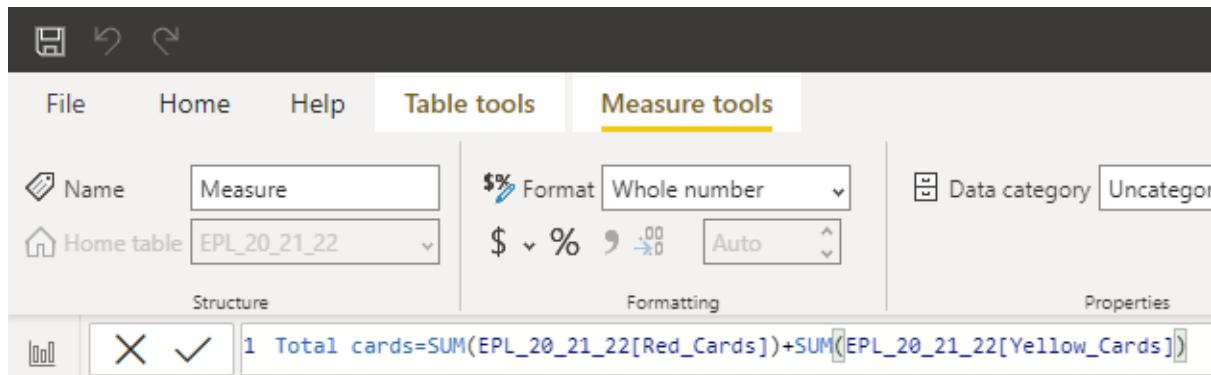


Figure 49: DAX expression for total cards.

Also, an M language was created for goal scorer classification by right clicking on “Get data” and select blank query editor from the dropdown. And the M-language below was used in creating the query,

The screenshot shows the Microsoft Power Query Editor interface. The main area displays the M language code for the 'Goal classification' query:

```

let GoalClassification = Table.FromRecords({
    [Category = "Top Scorer", Lower Level = 21, Upper Level = 25],
    [Category = "Mid Scorer", Lower Level = 12, Upper Level = 20],
    [Category = "Low Scorer", Lower Level = 6, Upper Level = 11],
    [Category = "Basic Scorer", Lower Level = 0, Upper Level = 5]
})
in Goal_Classification

```

The preview pane shows a sample of the data:

Player Name	Position	Goal Classification
Alexander	DF	27
Sead Kolasinac	DF	28
Shkodran Mustafi	DF	29
Hector Bellerín	DF	25
Roberto Soriano	MF	1
Andrea Belotti	MF	1
Emre Can	MF	0
Sead Kolasinac	MF	0
Shkodran Mustafi	MF	47
Roberto Soriano	MF	0
Andrea Belotti	MF	0
Emre Can	MF	16
Sead Kolasinac	GK	0
Emre Can	GK	0
Sead Kolasinac	GK	12

The properties pane on the right shows the query's name is 'Goal classification'.

Figure 50:M language for goal classification.

Dashboard

A variety of different visualizations were explored and used in the creation of the dashboards. These dashboards and visualizations were tailored towards answering the business questions that analysis aims to answer hence the particular choice of every chart to visualize particular data types.

Dashboard 1: Homepage

The dashboard being displayed below which is the Homepage was created using a combination of icons, textbox, and blank page navigators. Blank page navigators were used in creating the pages that is, Club Performance Analysis, Player Performance Analysis, Position Analysis and AI Insights. These pages also serve as buttons for page navigators which help in navigating through the pages they are meant to display. A textbox was used for the module name, student name, student ID and course code as well. The home button is a downloaded image which also serves as a page navigator that takes the report back to the homepage. The premier league and Teesside university logos are also downloaded images. The title of the page was also created using the textbox. The image of the first dashboard is displayed below in fig. 51.

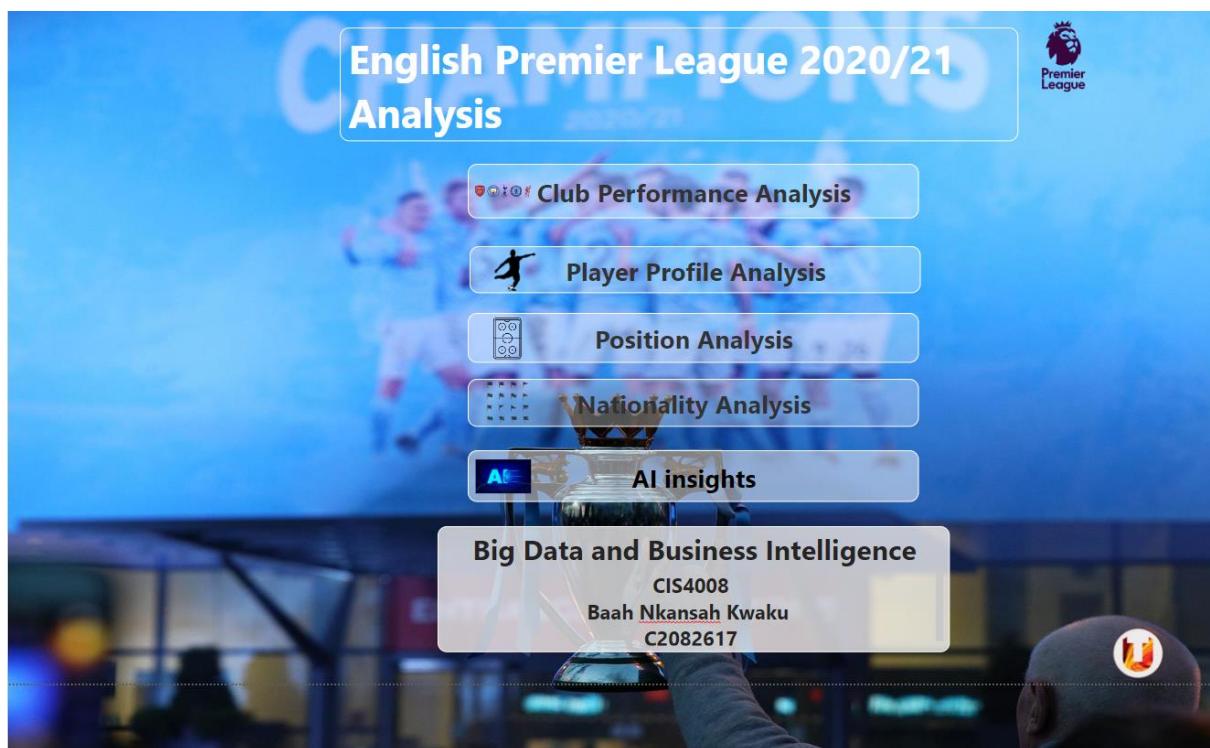


Figure 51:Homepage dashboard

Dashboard 2: Club Performance Analysis.

The below dashboard also displays a combination of different visualizations together with five (5) cards and one (1) slicer. The first chart is a horizontal bar chart which talks about the goals that were scored by the respective teams in the English Premier League during the 2020/21 season. The second chart is an Area chart which discusses the number of assists that each premier league club had for the English Premier League 2020/21 season. The third visualization is an animated chart that looks at the sum of the total club goals scored by each club for this particular season. The Drilldown radial bar chart was used to visualize the data type. A DAX expression was created for a new measure that was used in the creation of this chart. DAX expression created was,

Total club goals

$$= \text{SUMX}(\text{DISTINCT}(\text{CLUB}[\text{CLUB}], \text{CALCULATE}(\text{SUM}(\text{EPL 202122}[GOALS])))$$

The five cards created were used to give an idea of the number of clubs, the number of matches played, the number of yellow cards, the average goals and the average goals per club, while slicer was created for a dropdown list of the clubs that played the premier league season and this slicer works on the three (3) charts and five (5) cards as well.

The back, home and next button are also page navigators that goes back to the previous page, the homepage, and the next page respectively.

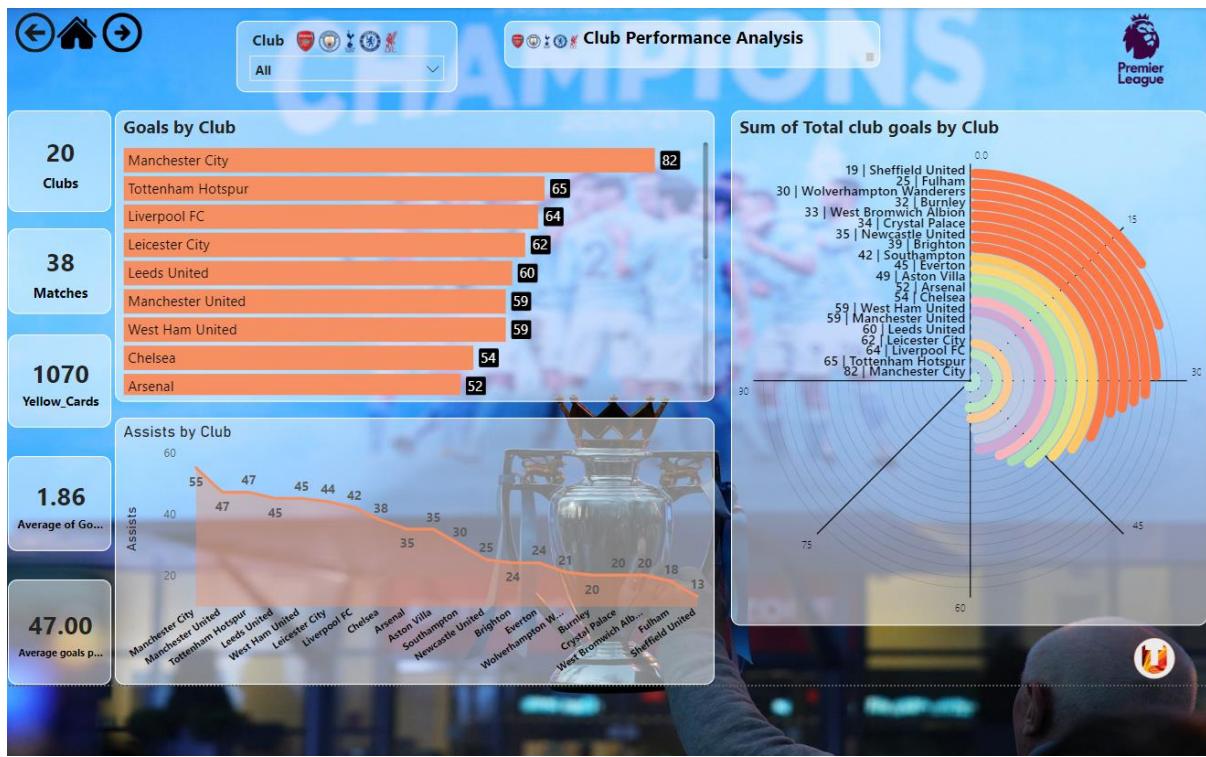


Figure 52:Club Performance Analysis Dashboard

Dashboard 3: Player Profile Analysis.

This is a special and different dashboard that was created with twelve (12) cards and three (3) slicers, that contains details and information about every player. The three (3) slicers which are, players, discipline and goal classification all work on navigating and giving an insight into the cards. The 12 cards contain data on Age of the players, the number of assists they had as well as the number of goals they scored, the passes completed, the number of red and yellow cards, their expected goals, the number of starts and the number of matches they played, and the penalty attempts they had and the number of penalty goals they scored.

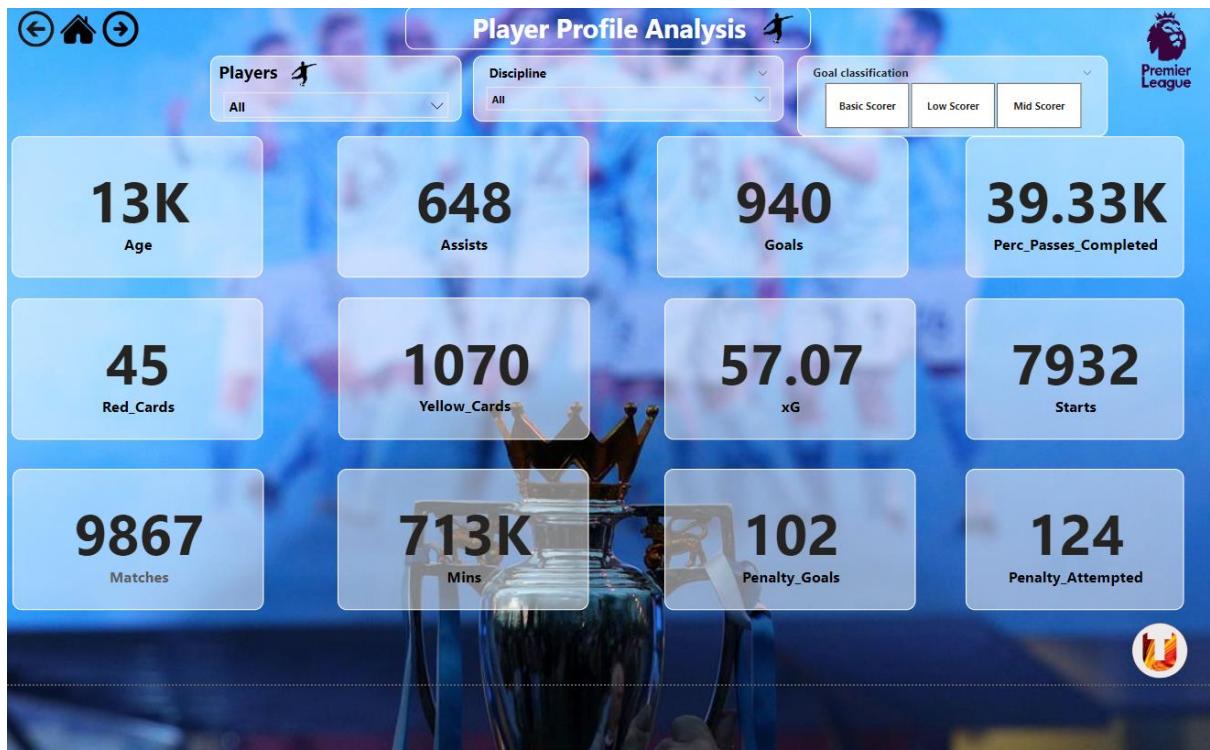


Figure 53: Player Profile Analysis dashboard.

Dashboard 4: Positions Analysis.

This dashboard is made up of a variety of different charts and visualizations. The first chart is a scatter chart that is analysing the average number of assists as against their position and the number of matches played per their ages. This was done to get insight into the age at which players could be said to be in their prime per the position they are playing in their respective teams. The second chart is a stacked bar chart looking at the number of goals scored by various main positions of the players. A card was also created to get a fair idea of the average of goals in a match. The third chart is represented by a drilldown radial bar chart. This visualization takes a look at which position had the greatest number of assists and this was done so as to get insight into which position contributes the most to the performance of the team. An Aster plot was used for the last chart which is visualizing the count of red cards for the respective positions for the 2020/21 season of the premier league. This was also done to gain knowledge about the most aggressive positions for this particular season. A slicer was also created for the various positions and this slicer reacts on the four (4) charts in this dashboard.

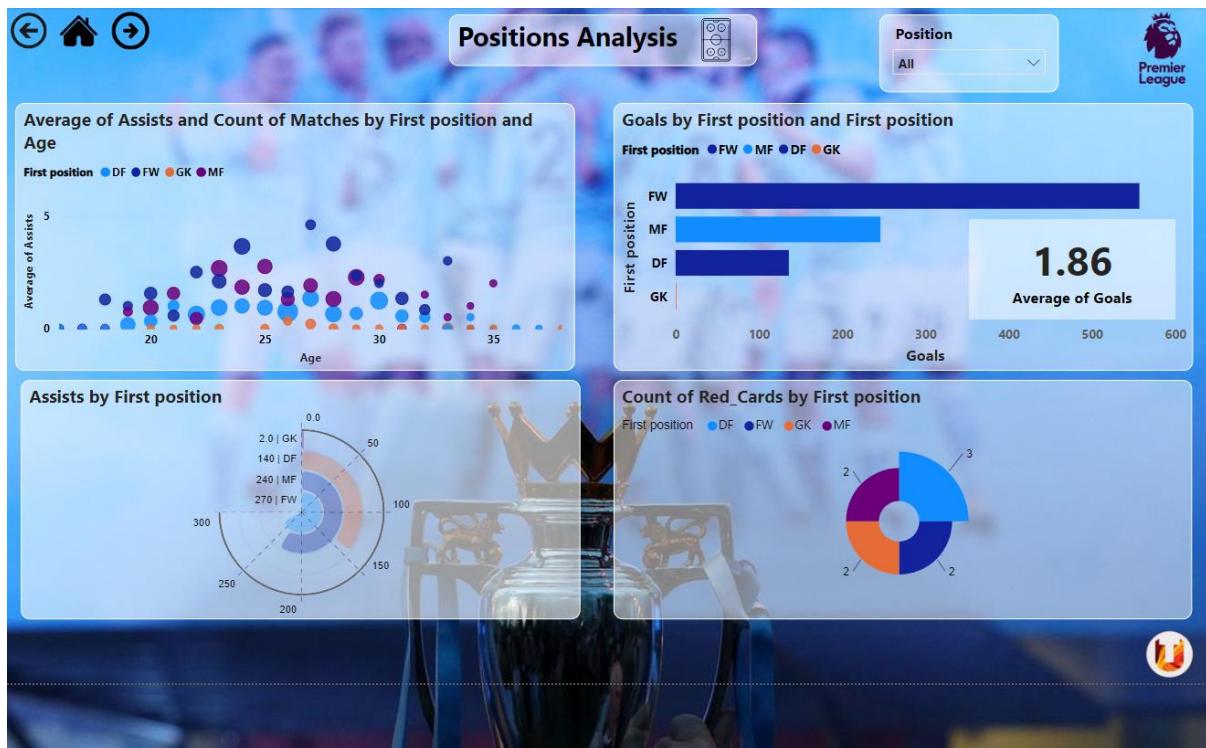


Figure 54:Positions Dashboard

Dashboard 5: Nationality Analysis

The nationality analysis dashboard contains four (4) charts, two (2) cards and two (2). The first chart is a map, and it was used because it is a chart that is representing and visualizing spatial data. The chart looks at the number of goals that were scored by various players through their nationality. The second visualization is an animated bar chart race which is visualizing the total number of cards and the matches played against the nationality of the players, this was done to figure out the nationality with the most aggressive players in the English Premier League 2020/21 season. The third chart is a treemap looking at the top 10 countries with the most players in the premier

league 2020/21 season. The last chart is a stacked bar chart analysing the total number of cards against the nationality of the players.

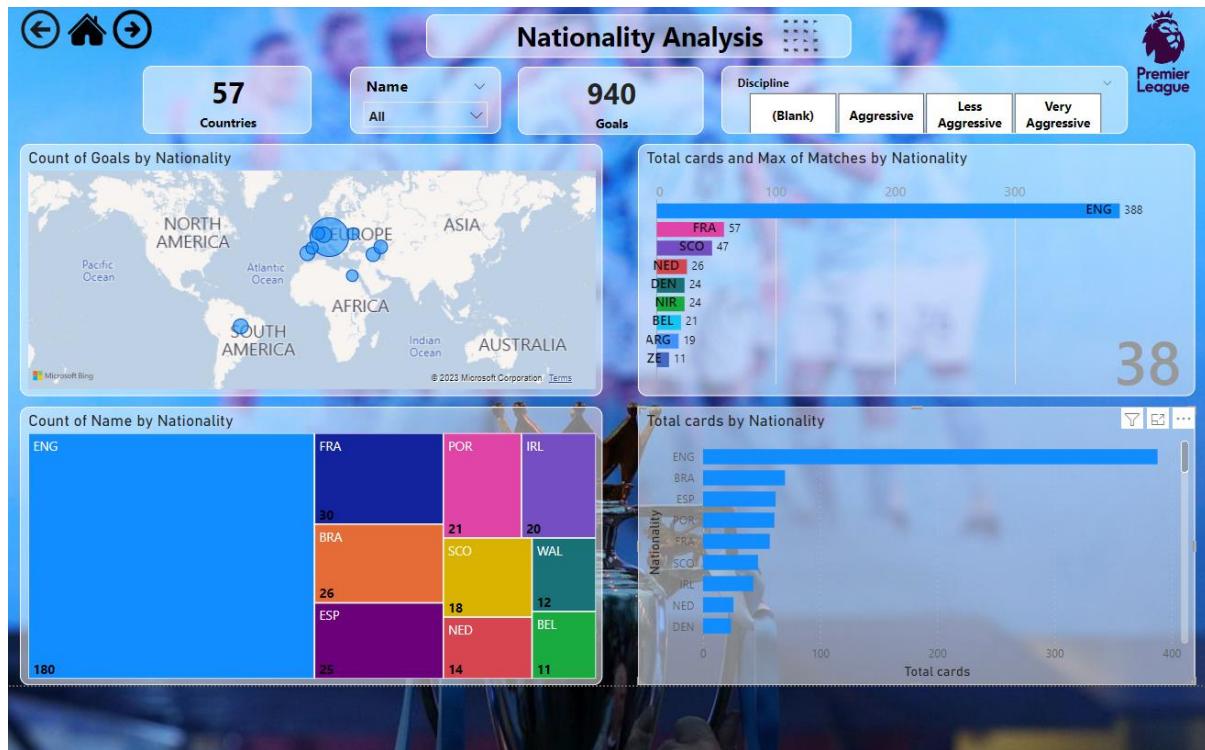


Figure 55:Nationality dashboard

Dashboard 6: AI Insights

This dashboard is a uniquely designed dashboard for insights into artificial intelligence. There were four (4) visualizations that were created for the AI Insights dashboard. The first chart is a decomposition tree that was used to analyse goals based on parameters like positions, club, and Nationality. The second chart is the key influencers chart, this chart was also used to analyse position as well while using goals, assists, matches, mins, penalties scored and age as parameters of explanation and also to analyse what influences decrease or increase in these parameters and also hidden patterns in the data. The third chart was a bar chart but focusing on distribution changes for goals and assists by various position. By right clicking on the bar chart and selecting “Analyze” from the pop-up box, there are charts that appear that help

identify outliers in goals and assists by various positions. The last chart is the Q&A chart, this visualization has certain in-built suggested questions generated by powerBi that comes with answers in the form of charts as well. Also, the “Teach Q&A” function was used in asking questions and determining how the questions are to be answered and submitting them.

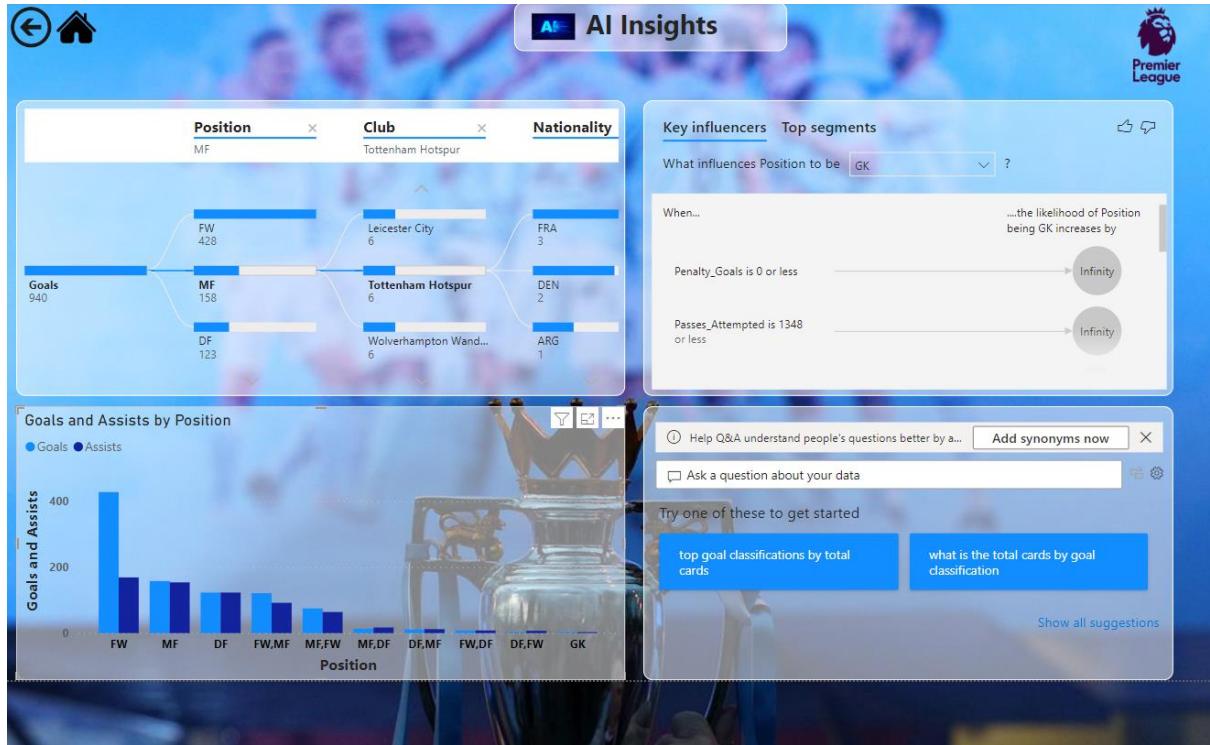


Figure 56:AI Insights dashboard

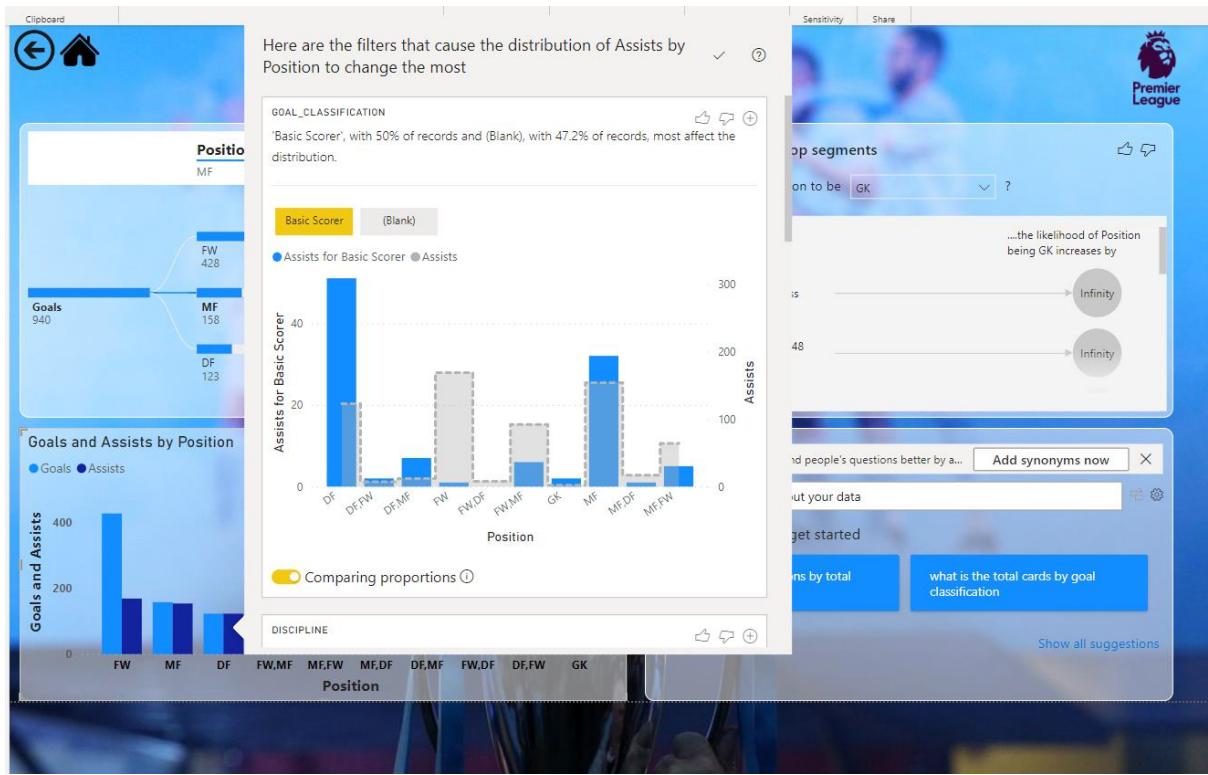


Figure 57: Further step for AI on bar chart.

REFERENCES

- (2021) *Premier League Season Review 2020/21*. Available at:
<https://www.premierleague.com/season-review> (Accessed: January 10, 2023).
- EY (no date). Available at: https://assets.ey.com/content/dam/ey-sites/ey-com/pt_br/topics/ey-economic-advisory-/ey-premier-league-economic-and-social-impact-january-2019.pdf (Accessed: January 10, 2023).
- Mwangi., B. (2021) *Exploring the English premier league 2020–21 season*, Medium.
Heartbeat. Available at: <https://heartbeat.comet.ml/exploring-the-english-premier-league-2020-21-season-ad10c6c83f1c> (Accessed: January 10, 2023).

Use the table below to **self-assess** your work. This will help reflect on your work.

You must keep this table in your report.

Report Section	Description	Grade your work from 0 to 100
Report Structure	The report is well-written, and it contains all the relevant sections	87
Data Pre-processing and Data Modelling	Many pre-processing steps have been applied. The data model is well-structured	60
Dax and M language	Both DAX and M Language have been extensively used in the report	75
Dashboard Design	The dashboard contains a variety of charts, including advanced ones not covered in the module.	82
Average		Add below the average of the four cells above: 76