



BIG DATA AND BUSINESS INTELLIGENCE MODULE

CIS4008-N-BF1-2022

Formula One World Championship Analysis

Module Leader : Dr Annalisa Occhipinti

Lab Tutor : Ms. Nissy Mathews

Submitted By : Ajay Kuruvilla Johnson (W9521662)

Submission Date: 11/01/2023

Contents

A. Executive Summary:	3
B. Body	6
Introduction:	6
• Data Source :.....	6
• Business Questions/Problems:	8
C. Finding based on analysis and evaluation.....	9
D. Conclusions with Business Questions and Answers	27
Recommendations	28
SECTION 2 : APPENDIX	30
ICA – Appendix: BI Design	31
A. Data Pre-Processing or Data Cleansing:	31
B. Data Modeling via Star Schema:	40
C. DAX and M Language	50
D. Dashboard.....	54
E. Self- Assessment	58

A. Executive Summary:

Formula One with its rich history is becoming a more global sport. The report focuses on how accidents and finishes vary over the years from (1950 – 2022).

We attempt to find certain characteristics and KPIs for a driver to be successful and using AI we try to forecast driver wins for the next 3 years.

We also try to find out what it takes to become a world champion and how different teams have been dominant in different eras of F1.

Findings and Conclusions

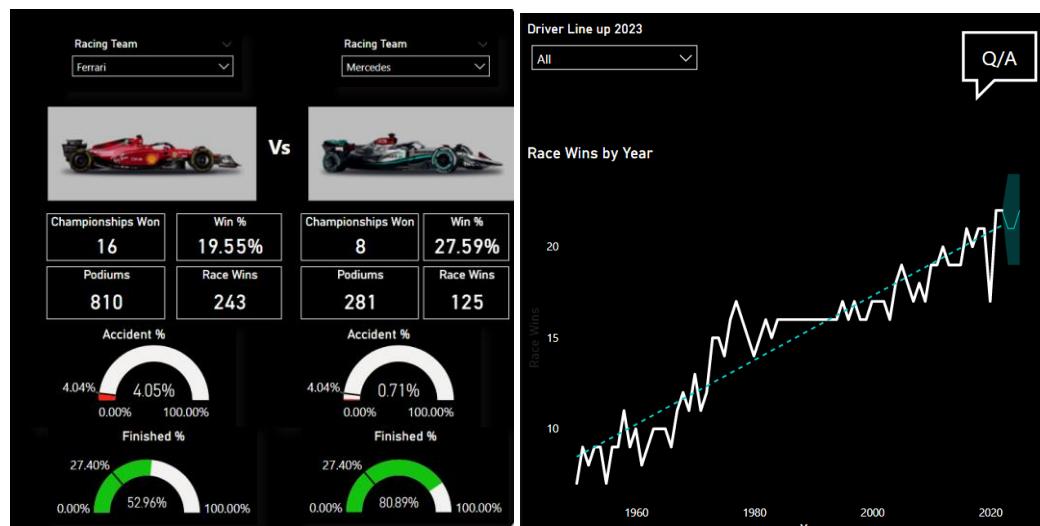
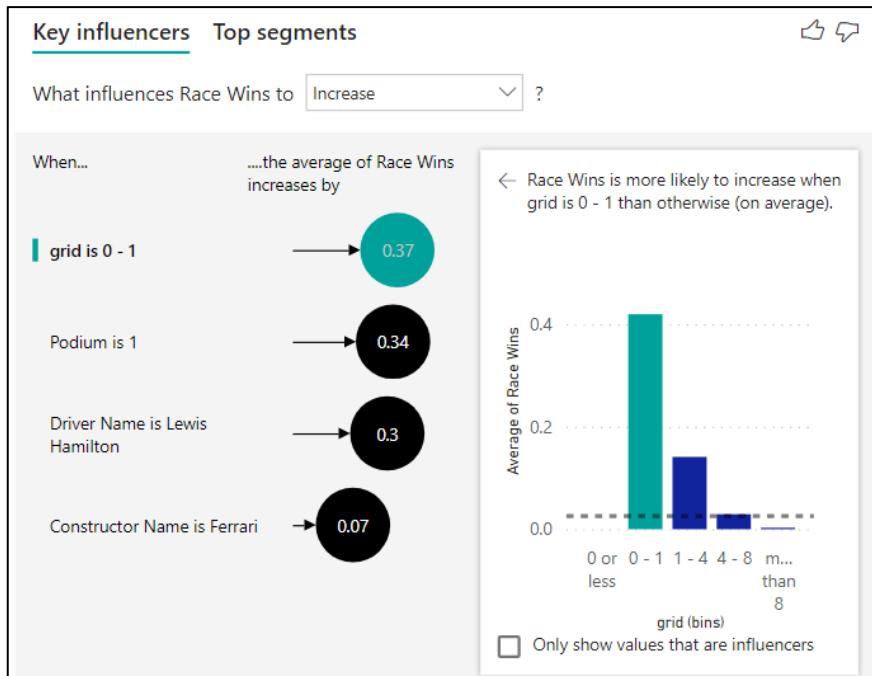
1. The most successful driver is Lewis Hamilton with 103 wins to his name and is a **seven time** world championship winner, tieing the record with Micheal Schumacher.
2. The most successful team is Ferrari because they have been in F1 since the beginning. Mercedes and RedBull are dominant in the Hybrid era.
3. The accident percentage has drastically reduced from 6.42% in 1950s it has reduced to 1.82% in 2022.
4. The finished percentage has drastically increased from 17.2% in 1950 to 63.41% in 2022.
5. In the past 10 years, Monaco has been the track with the highest accident percentage with 7.14% . On the other hand, Spa is the race track which has a high relative finish rate of 67.39%

The graphs below gives a quick understanding of the type of analysis done in this report.

Due to time constraints I have only chosen the top 20 driver to analyse in detail.
And the top 5 constructors

Apart from this we try to use AI to forecast driver wins for the year 2023

- Some Graphs from the Main Dashboard



Recommendations:

1. Teams with vast resources like Ferrari and McLaren are not able to perform at the top level. The number of (DNFs) – Did Not finish is higher due to car issues compared to Mercedes and Red Bull. The cars built should be reliable and hence top teams need to invest in the future to get their winning strategy right.
2. Winning Races is paramount to getting the prized Drivers and Constructors Championships and hence the main factors to look into that is to qualify in top 3 grid position in qualifying to give a racer maximum chances of winning on race day.
3. The Average age of a championship winning driver is 31.85 years. This means that the driver has been nurtured within a team and has gotten comfortable in order for him to give his best. Teams and senior management play a vital role in therefore recruiting bright talent and nurturing them so that they can become future world champions. This strategy has worked out well for red bull with Max Verstappen.

B.Body

Introduction:

Formula One is one of the biggest global sports and has a rich history originating in the 1950's. Winning races consistently leads to achieving the prestigious Driver's and Constructor's Championship titles. Hence racing teams are always on the lookout to understand and analyse factors which result in race wins. Using this dashboard, we can dive deep into some metrics which makes drivers win and also reasons for not winning as well.

Formula One is considered to be one of the dangerous sports with drivers at times having fatal consequences. We will explore the trend on accident rates and how this could help drivers be more efficient and consistent in their performance.

The steps of data pre-processing, data modelling and data visualisation with some emphasis on AI and analytics would have been showcased in this dashboard.

- **Data Source :**

The dataset is taken from Kaggle from the link : [Formula 1 World Championship \(1950 - 2022\) | Kaggle](#).

The Championship Data was taken from Wikipedia : [List of Formula One World Drivers' Champions - Wikipedia](#) and [List of Formula One World Constructors' Champions - Wikipedia](#)

The dataset from kaggle has **14** tables namely:

1. **Circuits:** circuitId,circuitRef,name,location,country,lat,lng,alt,url.
2. **Constructor_results:**
constructorResultsId,raceId,constructorId,points,status.
3. **Constructor_standings:**
constructorStandingsId,raceId,constructorId,points,position,positionText,wins.
4. **Constructors:** constructorId,constructorRef,name,nationality,url.
5. **Driver_standings:**
driverStandingsId,raceId,driverId,points,position,positionText,wins.
6. **Drivers:**
driverId,driverRef,number,code,forename,surname,dob,nationality,url.
7. **Lap Times :** raceId,driverId,lap,position,time,milliseconds.
8. **Pit stops :** raceId,driverId,stop,lap,time,duration,milliseconds.
9. **Qualifying :**
qualifyId,raceId,driverId,constructorId,number,position,q1,q2,q3
10. **Races :**
raceId,year,round,circuitId,name,date,time,url,fp1_date,fp1_time,fp2_date,fp2_time,fp3_date,fp3_time,quali_date,quali_time,sprint_date,sprint_time

11. Results :

resultId,raceId,driverId,constructorId,number,grid,position,positionText,positionOrder,points,laps,time,milliseconds,fastestLap,rank,fastestLapTime,fastestLapSpeed,statusId

12. Seasons: year,url

13. Sprint results:

resultId,raceId,driverId,constructorId,number,grid,position,positionText,positionOrder,points,laps,time,milliseconds,fastestLap, fastestLapTime , statusid.

14 status: statusid, status

The extra Championships data is taken from Wikipedia and has 2 tables:

1. Driver Championships
2. Constructor Championships

• The reason for choosing this dataset:

I am an ardent fan of Formula One and I know the sport quite well since I watch the races regularly. The possibility of finding some exciting insights in the data was quite motivating. Hence, I was excited on the opportunity to do an in depth analysis in PowerBI with this dataset.

• Will this dataset help in developing specific business skills?

I believe by doing this Big Data analysis, I got to understand how to analyse a huge dataset and how to break it down into different sections. The importance of pre-processing and data modelling became much more apparent to me.

The skills I have gained such as looking for KPI's and Performance metrics is of high importance and can be applied to my further career aspiration as a business analyst or data analyst.

- **Business Questions/Problems:**

This is a very important part of my project. The business questions I have taken are given below:

1. Have race accidents decreased over the years (1950 -2022)?
2. How has the finished % of drivers changed over the years (1950 -2022)?
3. Does Qualifying P1 (Position 1) have any impact on winning on race day?
4. Which racetrack has the highest risk of accidents on average and average finishes?
5. Can we forecast the number of race wins for drivers in 2023?
6. Which driver has been the most successful and is there any defining reason why he has been so successful?
7. Among the constructors, which teams are performing the best in the recent turbo era (2014 -2021)?
8. The Formula One sport was initially a European sport, how has that changed in recent years?
9. What is the average age for winning a driver's championship?
10. Which driver and team have the best winning percentage?

These questions are answered in the key findings section

- **Key User Groups :**

This dashboard is mainly intended for senior management and team principles who can use the historical data of drivers and teams and run useful analysis on every performance metric.

- **Specific Features and it's need:**

The main features I will be looking at is the safety of cars and how accidents have reduced or increased over time and what are the reasons for the same. This has massive applications as F1 has seen some fatal accidents in the past and making the sport more safe and fun is the objective.

- **Most Important KPI**

The main KPI would be **race wins** and **win%**. As a racing driver, the sport is extremely competitive and the goal is to secure a championship. That's where the money and fame lies. To get the championship one has to win and get as many points as possible in a given season.

C.Finding based on analysis and evaluation

This is the most important section of the project and I have explained with screenshots of the various visuals used in my dashboard.

- HOME Page:

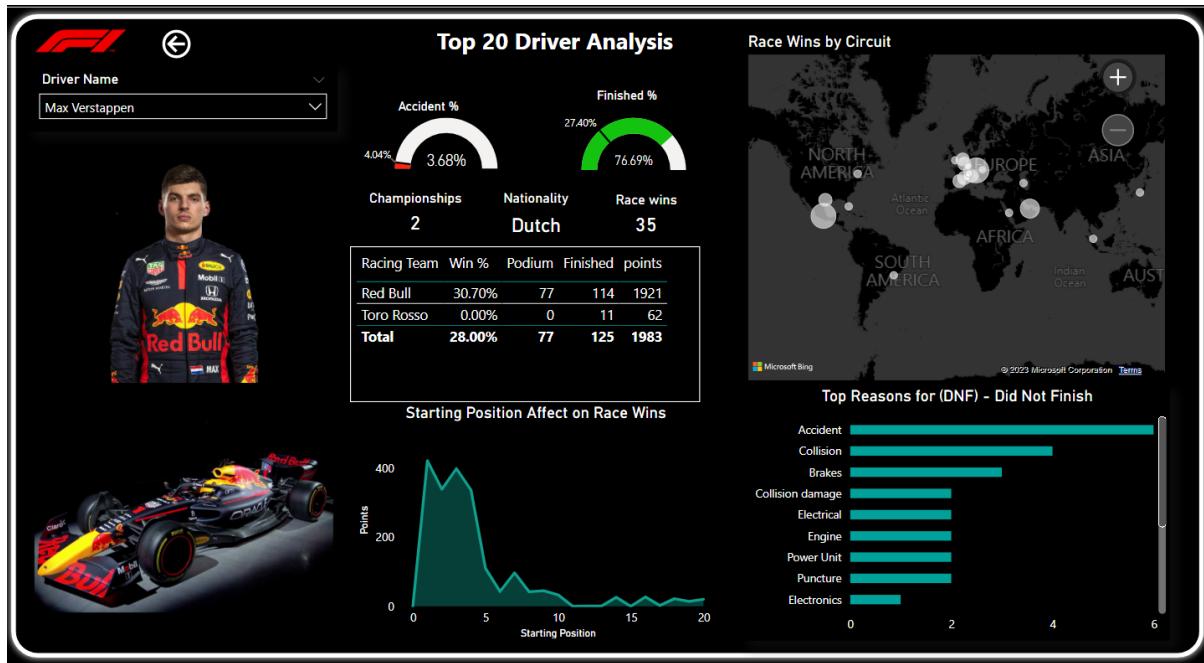


- This is a home page which use different buttons for navigation.
- A background image of F1 legend ,’Ayrton Senna’ is kept as an homage to his legacy.
- Each card is a navigation button with a background image given to make it visually pleasing. The same is applied for the other buttons as well.



Navigation Button

- DRIVER Analysis:



Overview



- Used here is the 'Image by CloudScope' visual.
- Using this we can make a dynamic image change according to the selection in the slicer.
- As the slicer selects another driver the image and their car will change dynamically.
- I had to get the image URL which I have put in a separate table called 'Driver and Car Images' and then made a 1:1 data modelling relationship to get this visual working.
- Only the top 20 Drivers were selected as it would have been a herculean task to get the images for all 848 drivers!

- Gauge Visual:



This is nice visual to show the threshold value, the minimum and maximum value.

My main business question and finding is based on the accident and finished % of drivers, constructors and the circuits.

DAX measures have been used to create this visual:

1 Accidents =												
2												
	IF([Results>Status] = "Accident", 1, 0)											
accId	grid	position	points	Status	Race Wins	Accidents						
29	4	0	0	Engine	0	0						
1 Finished = if([Results>Status] = "Finished", 1 , 0)												
accId	grid	position	points	Status	Race Wins	Accidents	Finished					
29	4	0	0	Engine	0	0	0					

Key Findings:

1. Able to calculate Accident % and finished % of Drivers.

- Table visual:

Racing Team	Win %	Podium	Finished	points
Red Bull	30.70%	77	114	1921
Toro Rosso	0.00%	0	11	62
Total	28.00%	77	125	1983

This table is chosen as it easy to show multiple multiple columns In an elegant without me having to add a lot of cards to the page.

Key Findings:

1. Able to see the Win % of Drivers with different teams.

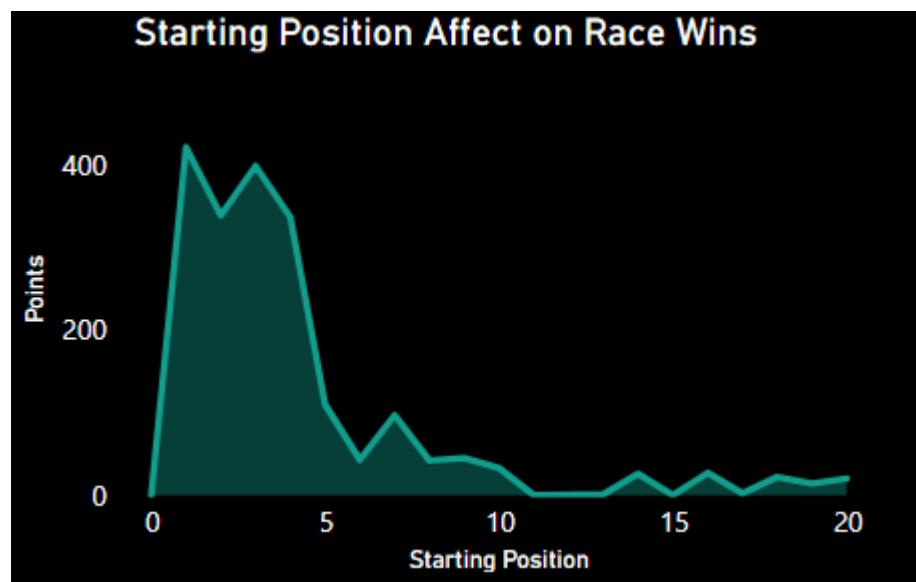
The Win% , Podium and Finished Columns use DAX Measures and Columns.

```
1 Win % = [sum(Results[Race Wins])/sum(Results[Finished])]
```

```
1 Podium = IF  
2 | OR(if(Results[position] = 1,1,  
3 | IF(Results[position] = 2,1,  
4 | IF(Results[position]=3,1,0))),0),1,0)
```

grid	position	points	Status	Race Wins	Accidents	Finished	Podium
20	A	n	Engine	0	0	0	0

- Area Chart visual:

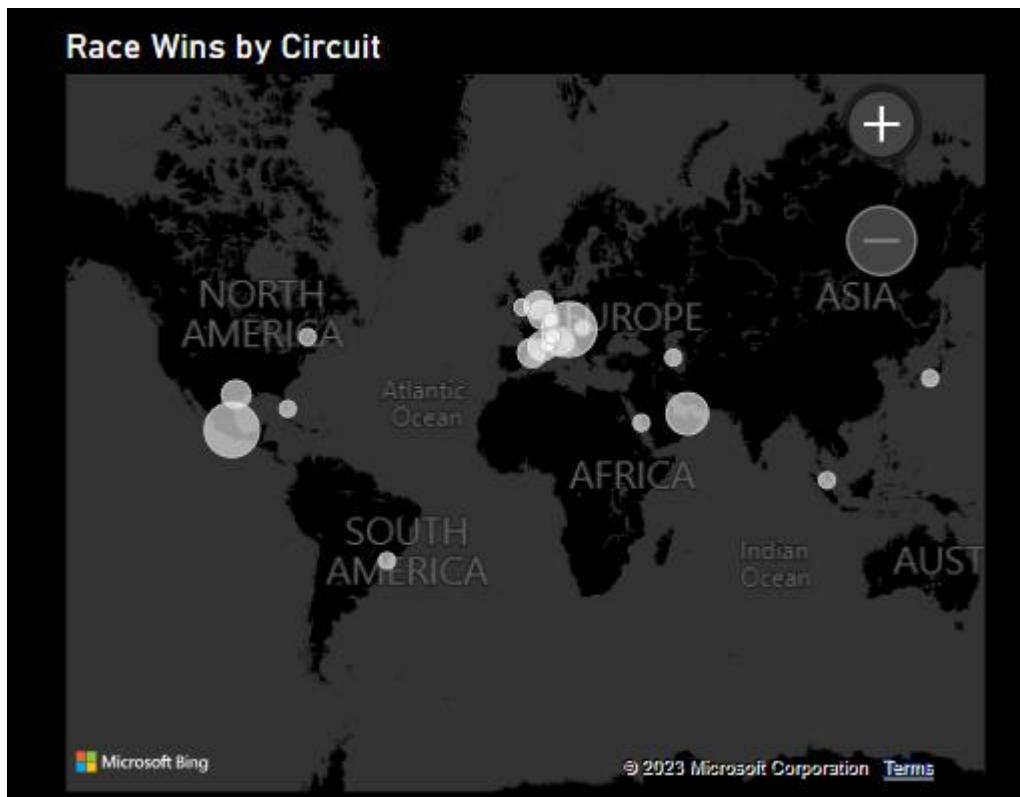


Having the y-axis as the points column from results and in the x-axis we have the 'grid' column from results table. This chart is chosen as it is easy to interpret and visually pleasing to look at.

Key Findings:

1. We see from the visual that the starting position of 1 , 2 and 3 have the most chances of winning maximum points and this significantly reduces if the starting position is 5th and above.
2. The chart will be different for different racers. Some racers prefer to lead while others perform well during a chase or race battle.

- Map Chart:



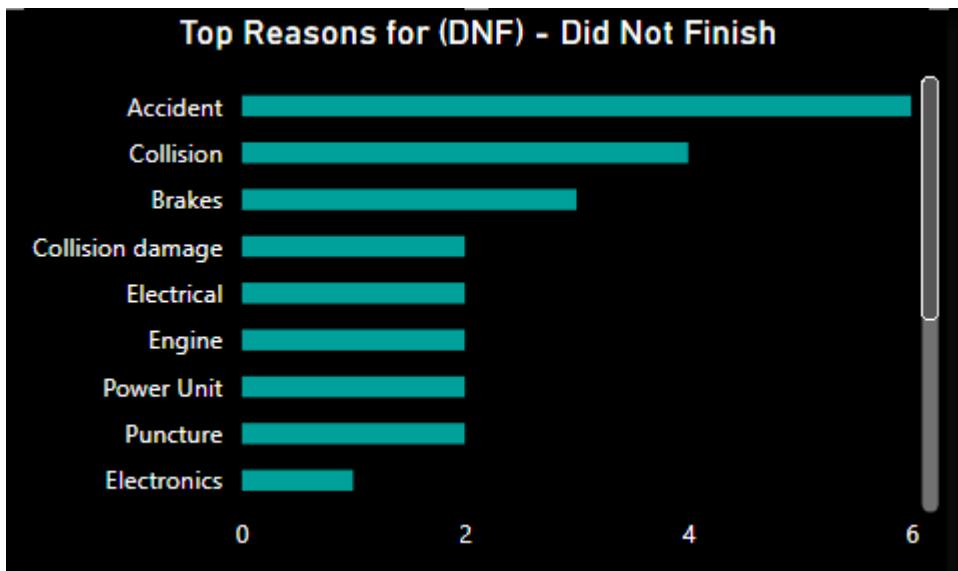
This chart is taken specifically to showcase clearly which **circuits** are a driver's strong point. The size of the **bubble** denotes the number of race wins by a driver. The map is interactive and moves accordingly to the driver selected.

A gray scale was chosen with white bubbles was chosen to give a contrast to the grayish background.

Key Findings:

1. Most racers are dominant in their home tracks but the drivers who are successful in historically proven tracks such as Monza and Silverstone are generally classified as elite racers.

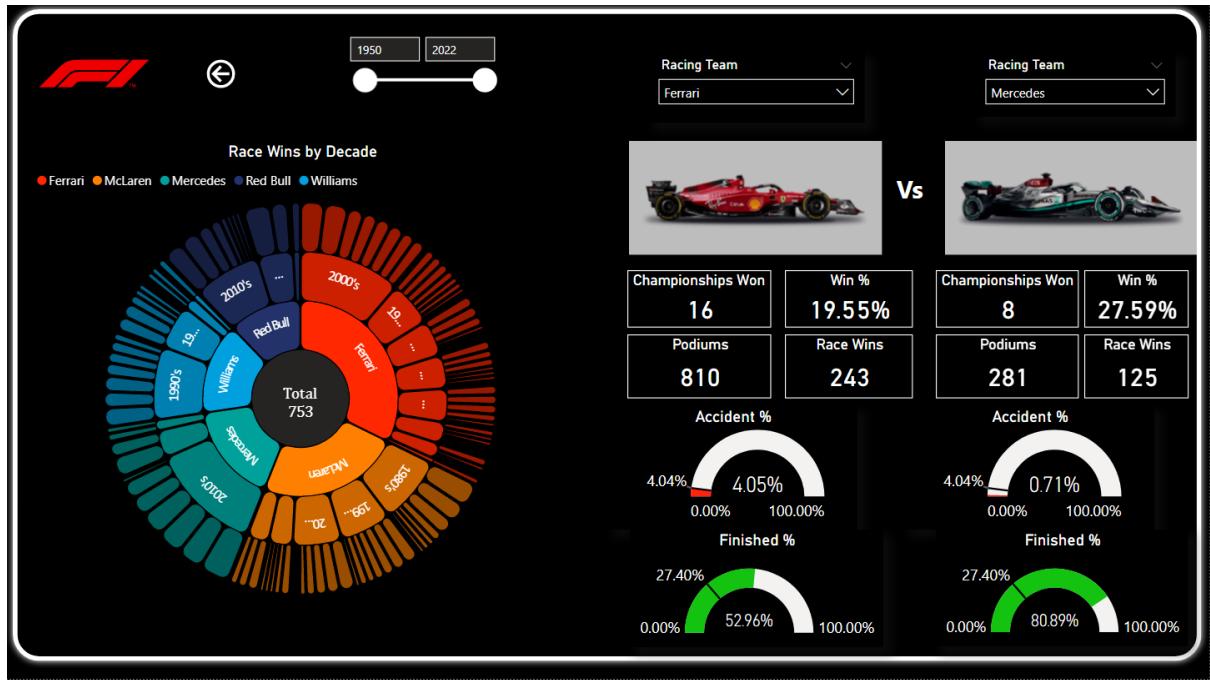
- **Horizontal Bar Chart:** This is a standard bar graph. It is chosen specifically to understand which are main reasons for Driver (DNF) – Did not finish.



Key Findings:

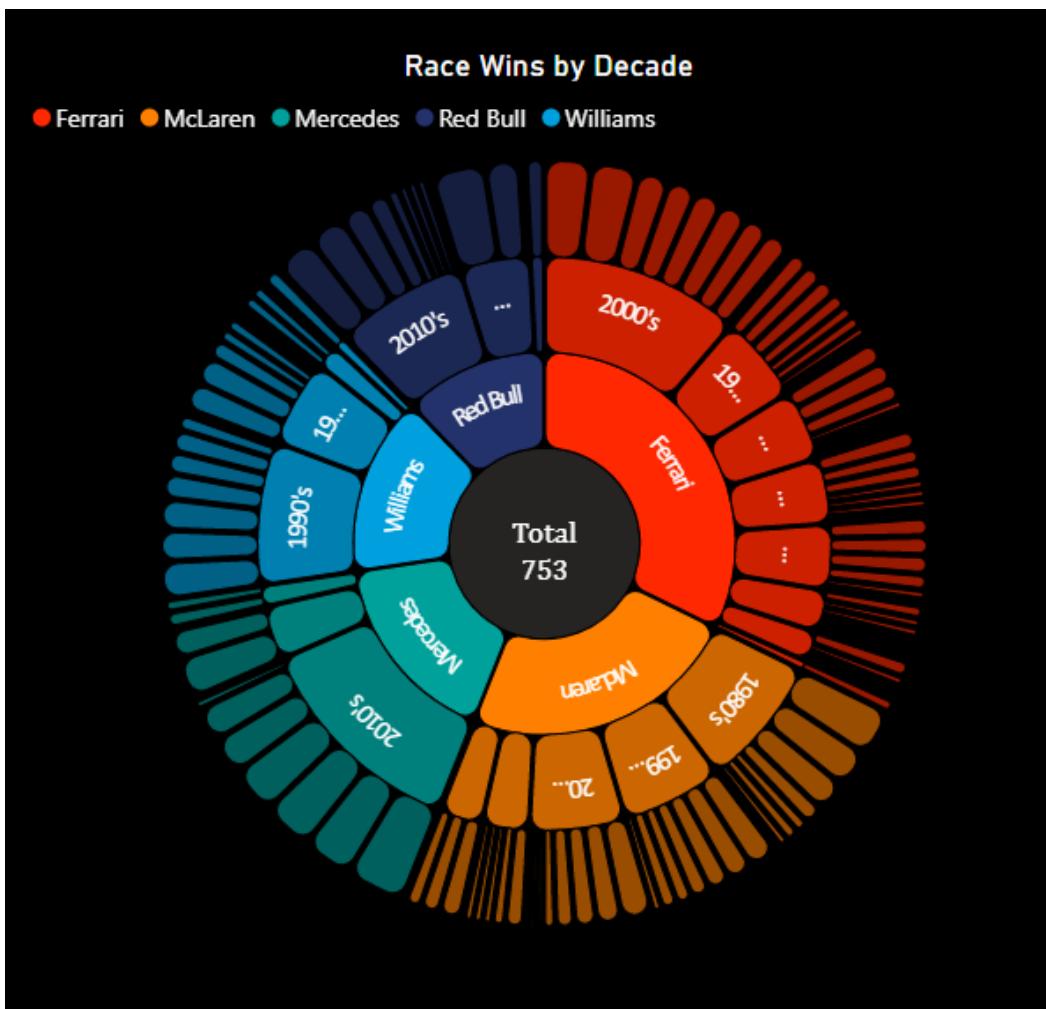
1. This is an important visual as we can see have an insight into their driving style and also how efficient their cars are.
2. When the reason is ‘Engine’ or ‘Electrical’ failure at the top, we can understand the car is not fully reliable.
3. On Driver traits, ‘Lewis Hamilton’ has very few (DNF’s) and are mostly due to collision by another racer. Hence it can be said he drives much more safely than for example, ‘Max Verstappen’ who is more of an aggressive driver by comparison.

- CONSTRUCTOR Analysis (Page 2):



Overview

- The Zoomable Sunburst chart: (Not Covered In Lesson)



This is a zoomable sunburst chart and this is the external chart which I have chosen. I have chosen this particular graph as I wanted to display the dominance of certain teams over the 72 years of F1 history.

It has two levels and it shows the race wins by Constructors over the decades.

Key Findings:

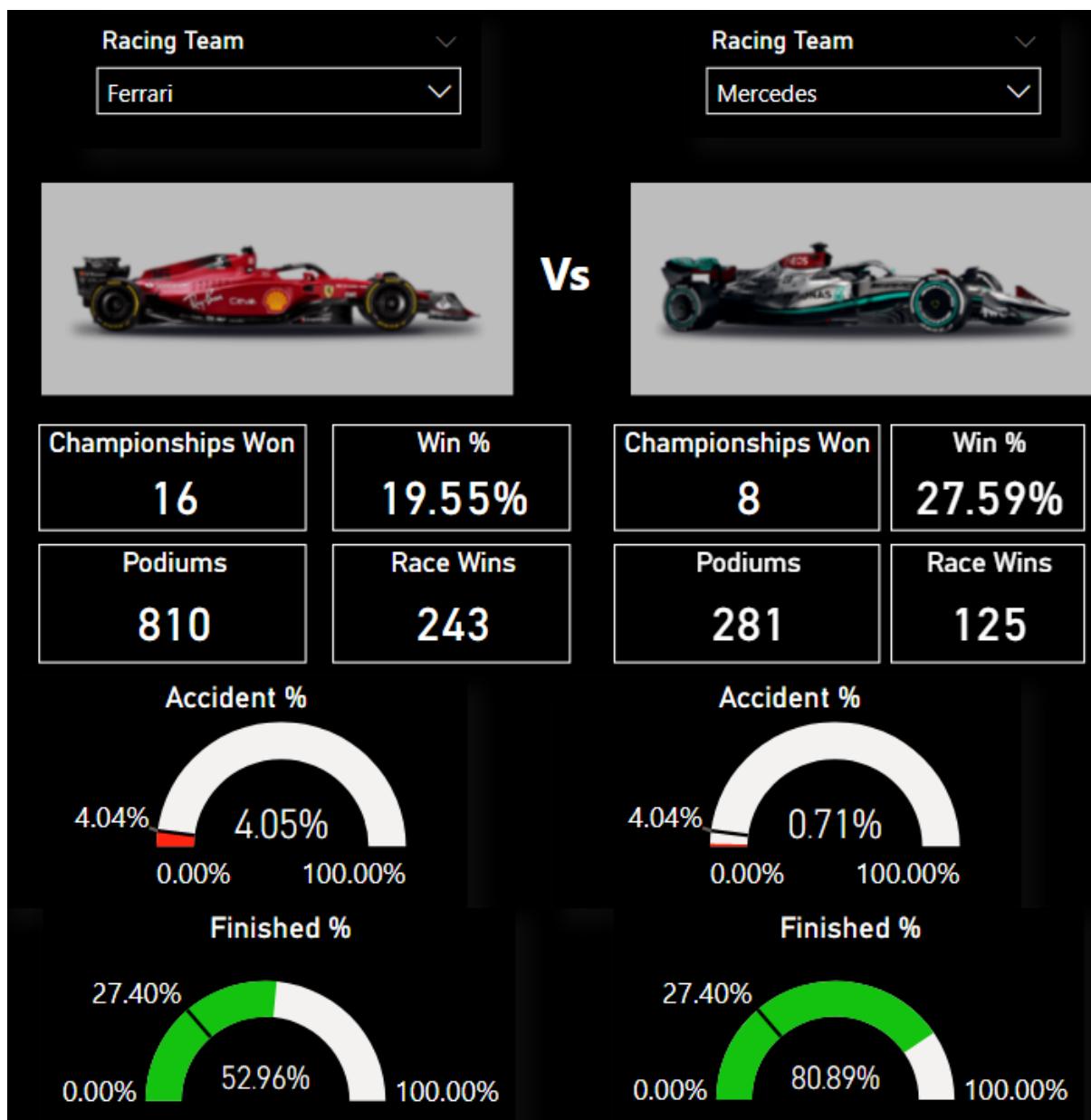
1. From this visual we can see that in the Hybrid Era (2014-2021) Mercedes have been completely dominant with Lewis Hamilton breaking records.
2. The years between (2010 – 2013) belonged to Red Bull.
3. And the years between (2000 – 2004) and some would argue till 2008 was dominated by 'Ferrari'.

Hence this insight can be displayed nicely using the sunburst chart.

For this a DAX Column 'Decade' has been used to create this visual.

```
1 Decade = IF(Results[Year] < 1960, "1950's",
2             IF(AND(Results[Year] < 1970, Results[Year] >= 1960), "1960's",
3                 IF(AND(Results[Year] < 1980, Results[Year] >= 1970), "1970's",
4                     IF(AND(Results[Year] < 1990, Results[Year] >= 1980), "1980's",
5                         IF(AND(Results[Year] < 2000, Results[Year] >= 1990), "1990's",
6                             IF(AND(Results[Year] < 2010, Results[Year] >= 2000), "2000's",
7                                 IF(AND(Results[Year] < 2020, Results[Year] >= 2010), "2010's",
8                                     IF(AND(Results[Year] < 2030, Results[Year] >= 2020), "2020's", "Unknown"))))))))
```

- Comparison using Cards, Gauge axis and and Image visual:



This is used to compare the stats and figures of different teams. (Only the Top 5 teams are used).

So this is designed to see which teams dominated in which year. All of this can be controlled by a year slicer to show data of particular years.

This visual using multiple visuals to portray a story. This is particularly useful to understand which teams are more safe and more accident prone. Which team has a better winning %.

Key Findings:

1. We see that Mercedes in recent year have a very high Win % and their finished Races % is much higher than any other which means their Accident % is also extremely low.
2. Red Bull and Mercedes in the past 10 years have been more successful than all the other teams.
3. Ferrari has fallen short in recent years and their last drivers championship was in 2008.

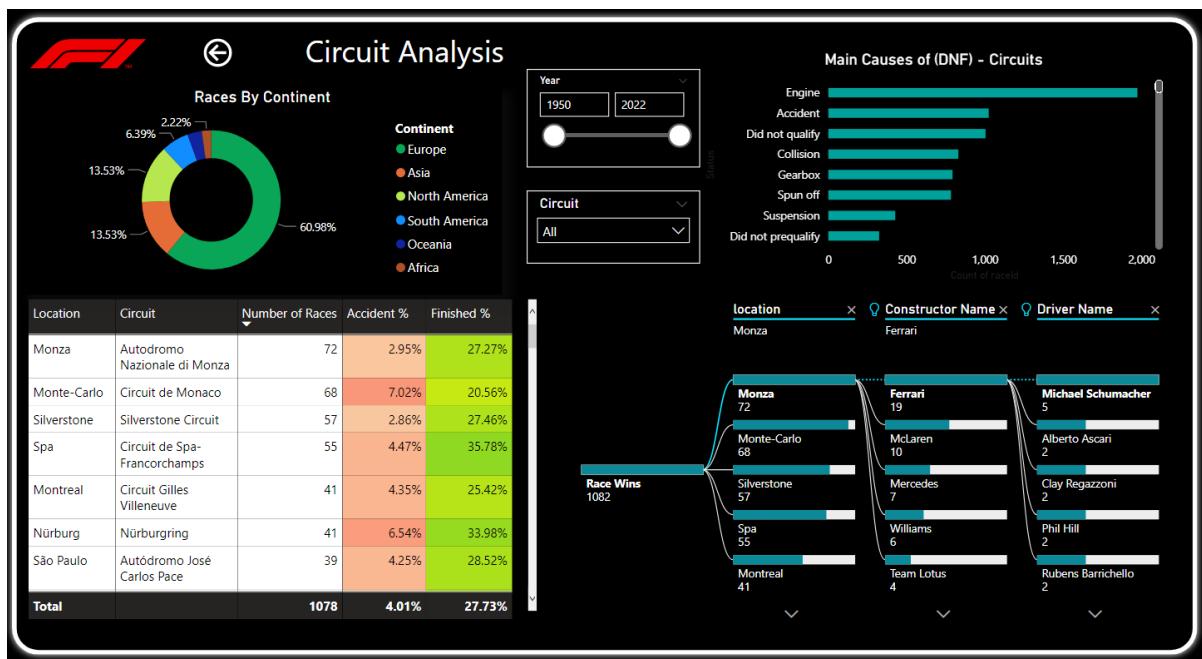
Images of the cars and the URL was used to put in a table called ‘Driver and Car Images’ and a 1:1 relationship is done with Drivers Data and Constructors data.

- Slicer



Standard Slicer with the seasons (Race Year) are given so as to drilldown on a particular season or maybe an certain time period or decade.

- CIRCUIT ANALYSIS: (Page 3)



Overview

- Table Visual with Conditional Formatting:

This visual is best to display multiple columns of data side by side for several rows.

The conditional formatting is done so that we can sort according to their descending or ascending values.

Location	Circuit	Number of Races	Accident %	Finished %
Casablanca	Ain Diab	1	20.00%	20.00%
Pescara	Pescara Circuit	1	12.50%	25.00%
Heusden-Zolder	Zolder	10	12.41%	15.96%
California	Long Beach	8	12.27%	15.91%
Barcelona	Montjuïc	4	11.90%	13.10%
Indianapolis	Indianapolis Motor Speedway	19	11.20%	31.66%
Detroit	Detroit Street Circuit	7	10.47%	15.18%
Styria	Zeltweg	1	10.00%	10.00%
Total		1078	4.01%	27.73%

Since accident % and finished % were quite important to me, this visual made it super easy and makes it easy to understand for the viewer as well to understand the circuits with high accident rates and finished rates.

Key Findings:

1. If we take data from the past 10 years, we see that the most dangerous track is monaco's circuit with a Accident % of 7.26%
2. The track with most average finishes is Spa with 66.14%.

- **Decomposition tree visual:**



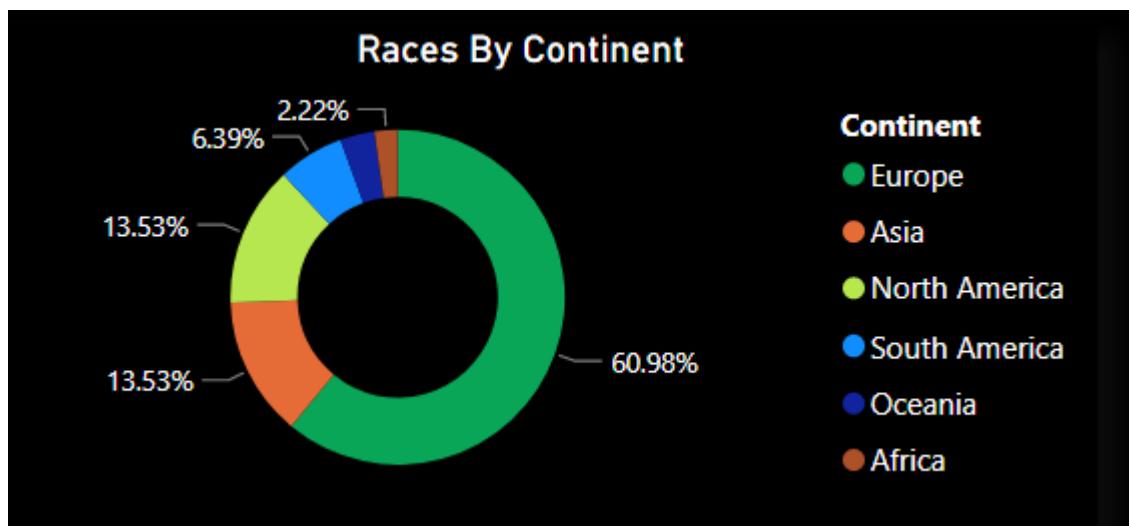
This visual has been taken as it displays relative value for the Constructors and Drivers who have been most successful in a certain circuit. The 'bulb' next to the Constructor and driver Name heading signifies that using AI 'high value' or 'Low value', power BI automatically detects the different input and gives an output, which, in this case is of high value.

This visual is also helpful to see how other drivers or constructors compare to the top value.

Key Findings:

1. Silverstone which is the track in Great Britain has most wins by Lewis Hamilton. And the UK is the home base for Mercedes.
2. Monza has Michael Schumacher in race wins which is the home track of Ferrari.

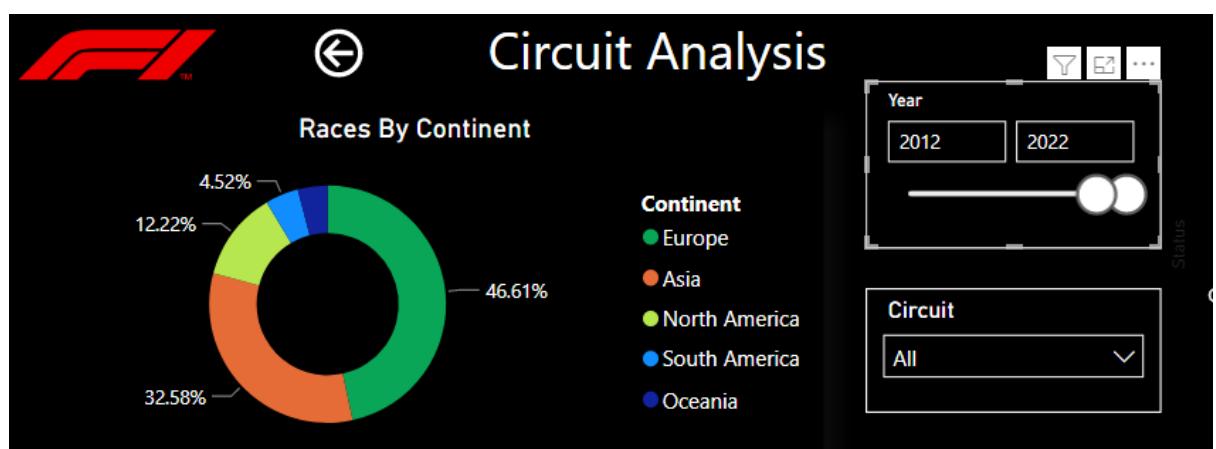
- **Pie chart:**



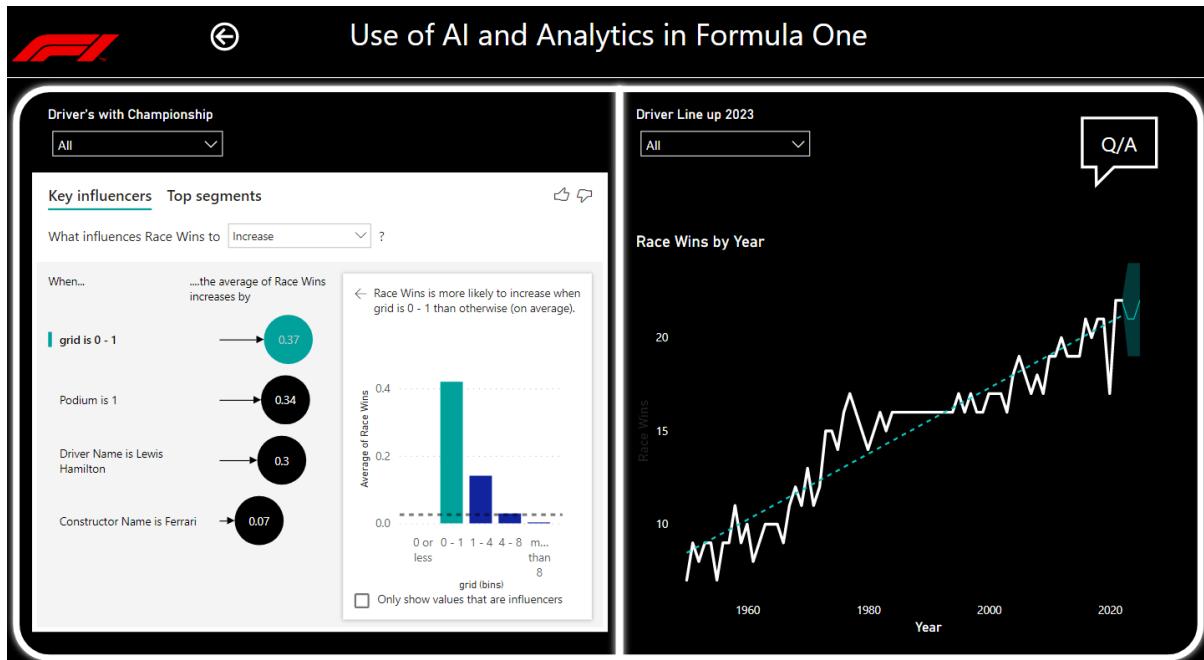
This is a standard pie chart visualised that gives an idea of which continent hosts more of the F1 races. This chart helped me answer the question of how in recent years, the sport is become much more global.

Key Findings:

1. In past 10 years there has been more races happening in the Asian continent and some in North America making F1 truly a global sport.



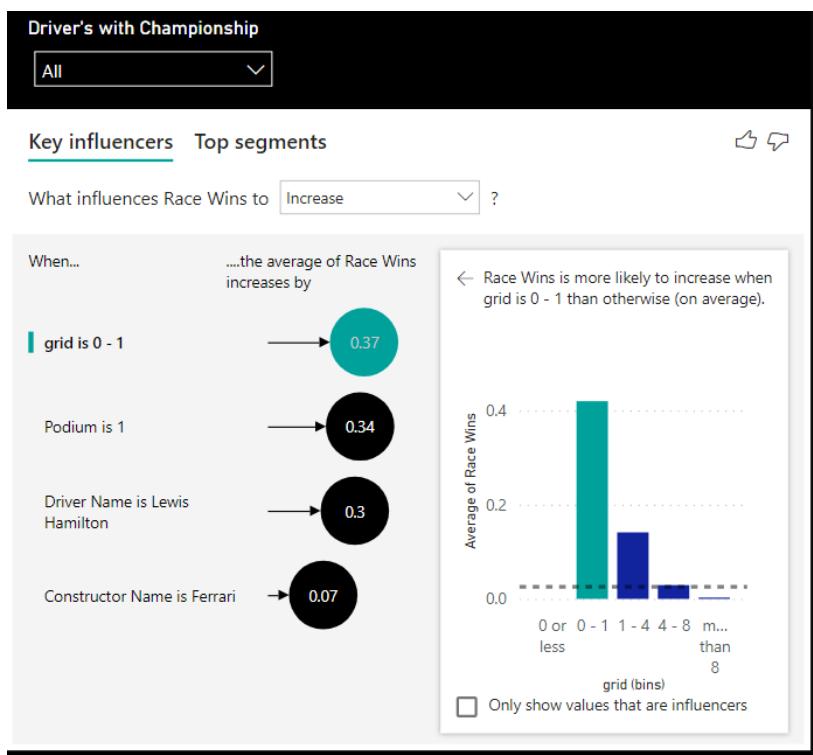
- Use of AI and Analytics in Formula One (Page 4)



Overview

I have kept this page mainly only for charts which have AI being used.

- **KPI Influencers Visual:**

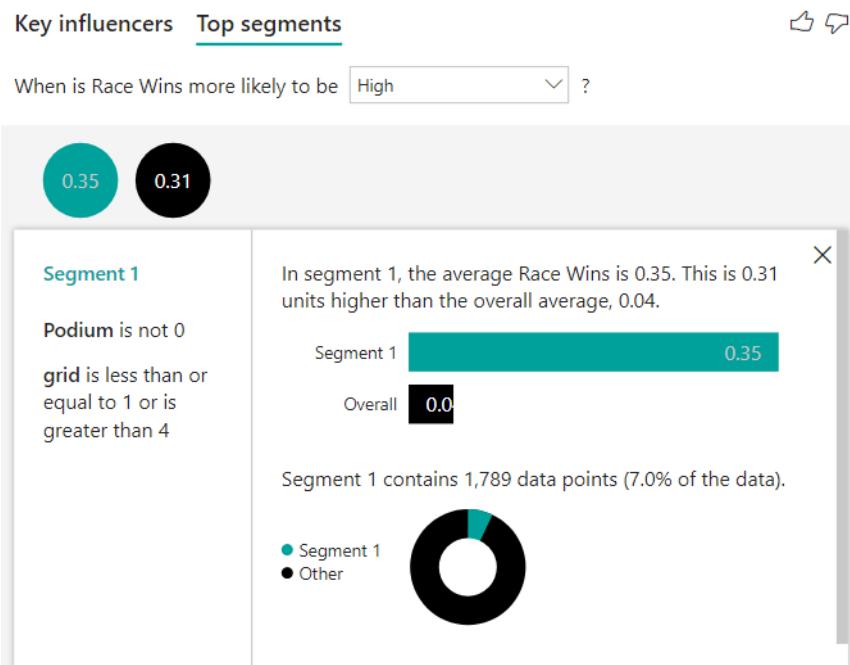


This visual helps us understand what are the factors that can help increase Race wins.

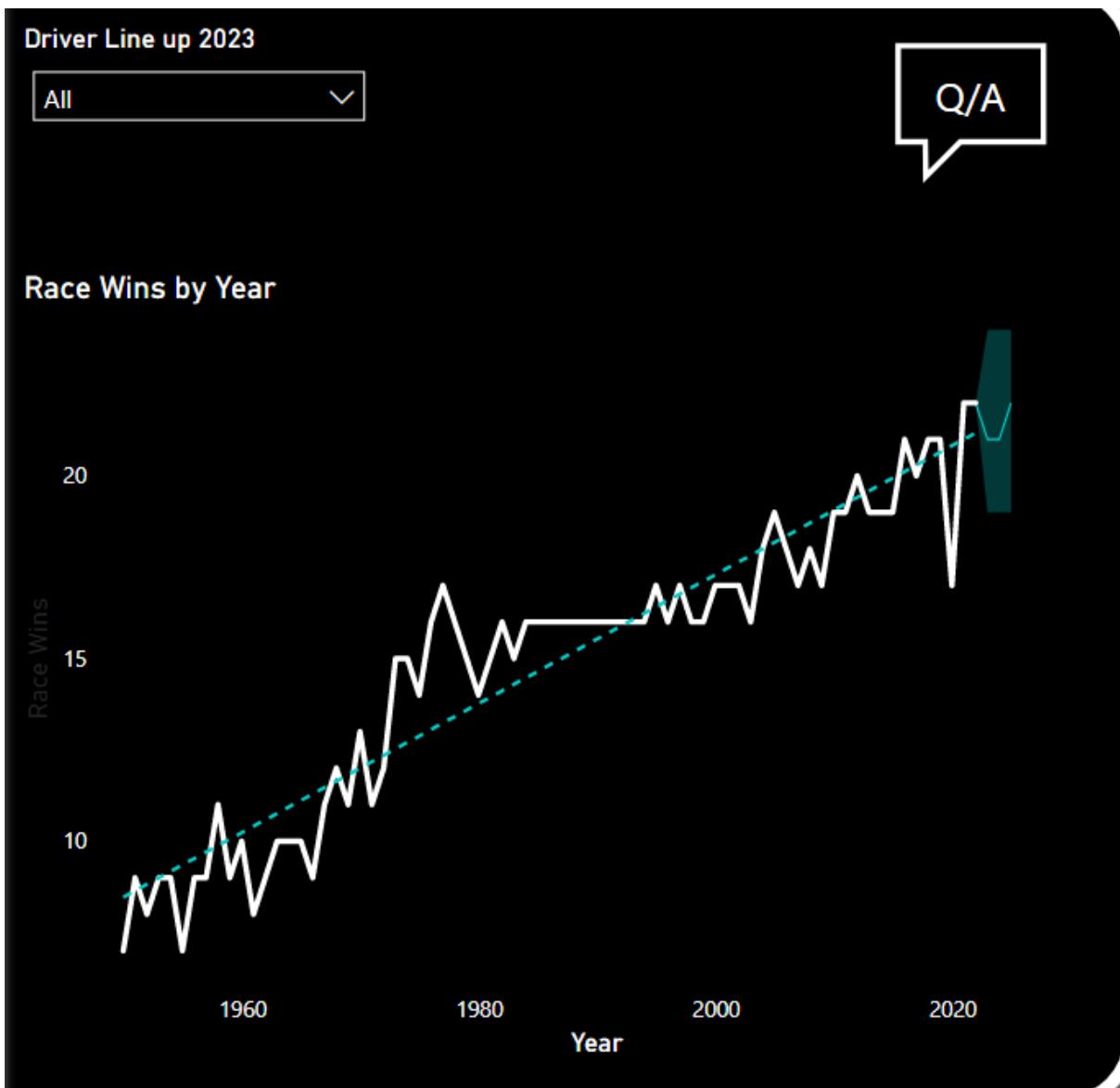
Key Findings:

1. If the starting grid position is 1, the chances of winning are much higher. This is especially true for circuits like Monaco.
2. The chances of winning a championship is by accumulating most amount of points. And the way to do that is to get is by either winning the race or having a podium (1,2,3 finish)

This visual also gives us Top Segments which are very helpful for analysis.



- **Line Chart with Forecasting:**



The area at the end of the line chart is the forecasted part. I have forecasted race wins for 3 more from 2022. It is fun to have a look into the future, although this is not very accurate, it gives a finding based on the previous historical data.

Key Findings:

1. Lewis Hamilton is forecasted to win only 2 races in 2023.

- The Q/A Visual:



I have kept the Q/A as a button to conserve space.

This is a really cool visual. I have taught the AI some Q/A and it has come up with some suggested graphs.

Manage terms			
Manage the terms and definitions you've taught Q&A.			
Term ↑	Definition ⓘ	Modified	Actions ⓘ
best	name is best if total race wins is greater than 25	04/01/2023	ⓘ
country	nationality	04/01/2023	ⓘ
driver	name	04/01/2023	ⓘ
finished percentage	finished %	04/01/2023	ⓘ
good	name is good if total race wins is high	04/01/2023	ⓘ
nation	country	04/01/2023	ⓘ
Race Name	<i>Race Name</i>	04/01/2023	
Team name	<i>Constructor</i>	04/01/2023	
win percentage	win %	04/01/2023	ⓘ
Wins	race win	04/01/2023	ⓘ

Teaching some terms to AI

Show the trend of finished% over the years X Add

Showing results for Finished % sorted by result year

Reorder your suggested questions

- Show circuit location by accident %
- Show the trend of finished% over the years
- Show the trend of accident % over the years
- What is the average age of driver's champions?
- Driver total championships
- Circuit Location by races map

Save

Preview your result

Show the trend of finished% over the years

60%

40%

20%

1950 2000

Year

Suggested and trained questions

More in less

<input type="checkbox"/> Championship Data	<input checked="" type="checkbox"/>
<input type="checkbox"/> Circuits	<input checked="" type="checkbox"/>
<input type="checkbox"/> Constructors Data	<input checked="" type="checkbox"/>
<input type="checkbox"/> Continent	<input checked="" type="checkbox"/>
<input type="checkbox"/> Driver And Car Images	<input type="checkbox"/>
<input type="checkbox"/> Driver Data	<input checked="" type="checkbox"/>
<input type="checkbox"/> Race Car Images	<input type="checkbox"/>
<input type="checkbox"/> Races	<input checked="" type="checkbox"/>
<input type="checkbox"/> Results	<input checked="" type="checkbox"/>

Choosing the right columns to be displayed so that the graphs are displayed accordingly

Ask a question about your data X

Try one of these to get started

Show circuit location by accident %

Show the trend of finished% over the years

Show the trend of accident % over the years

What is the average age of driver's champions?

Driver total championships

Circuit Location by races map

top constructors by race wins

Who are top drivers by race wins?

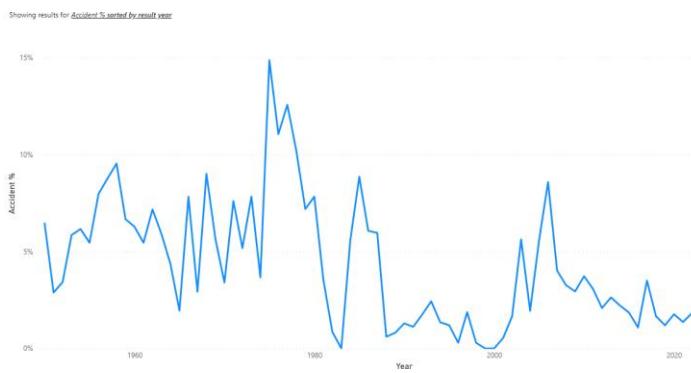
[Show fewer suggestions](#)

Suggested and trained questions

D. Conclusions with Business Questions and Answers

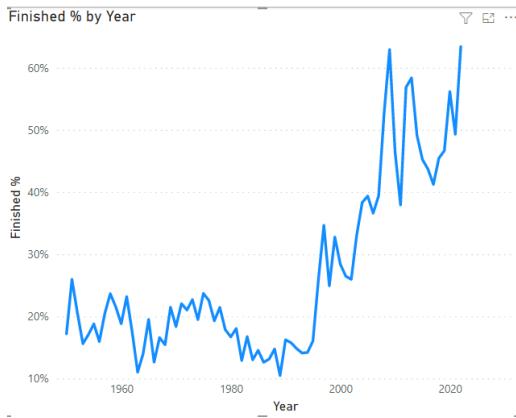
1. Have race accidents decreased over the years (1950 -2022)?

Answer : Yes they have, I was able to find this from the Q/A graph. From 6.42% in 1950s it has reduced to 1.82% in 2022.



2. How has the finished % of drivers changed over the years (1950 -2022)?

From 17.2% in 1950 to 63.41% in 2022.



3. Does Qualifying P1 (Position 1) have any impact on winning on race day?

Yes qualifying P1 puts a significant advantage for the driver to lead the race.

4. Which racetrack has the highest risk of accidents on average and average finishes?

In the past 10 years:

Highest Accident % - Monaco

Highest Finish % - Spa

5. Can we forecast the number of race wins for drivers in 2023?

Yes, Lewis Hamilton – 2 wins

Max Verstappen – 14 wins

6. Which driver has been the most successful and is there any defining reason why he has been so successful?

Lewis Hamilton as per the KPI Influencer Chart. His accident % over his career is 0.32% and finished % is 87.10%. Main reasons of DNF are Collision and Gearbox issues which shows how good his driving calibre is and only if a car issues comes he does not finish.

7. Among the constructors, which teams are performing the best in the recent turbo era (2014 -2021)?

Mercedes with 8 Consecutive Constructors' Championships.

8. The Formula One sport was initially a European sport, how has that changed in recent years?

In the past 10 years, Asia hosts 32.8% of the races compared to only 13.53% in the past.

9. What is the average age for winning a driver's championship?

31.85 years

10. Which driver and team have the best winning percentage?

Mercedes with 27.59% and Hamilton with 38.15%

Recommendations

1. Ferrari and McLaren, two teams with vast resources, are unable to compete at the highest level. Compared to Mercedes and Red Bull, there are more drivers who (DNF) - Did Not Finish - the race because of car problems. The cars constructed should be dependable, thus top teams must invest in the future to develop a winning strategy.
2. The major consideration is to qualify in the top 3 grid positions in qualifying to provide a racer the best chance of victory on race day. Winning races is essential to obtaining the coveted Drivers and Constructors Championships.
3. A championship-winning driver is 31.85 years old on average. This indicates that the driver has grown up in a team environment and is now at ease enough to give it his all. Teams and senior management are crucial in finding and developing talented individuals who have the potential to become future world champions. Red Bull has had success with this tactic while using Max Verstappen.
4. The circuit at Monaco has the highest average risk of having accidents or collisions at 7.14% in the past 10 years. This is due to street style circuit.

Although this has been the circuit layout for more than 50 years a change in the track should be considered as there is always some sort of collision bound to happen due to its narrow roads. Safety of the drivers is the main priority and other newer circuits such as the Red Bull ring, Bahrain International Circuit and the Abu Dhabi Yas Marina circuit have accident percentages namely 0%, 0.8% and 0.88% respectively. Newer safety principles are considered while building these race tracks.

5. It would be nice to see more races happening in North America and Asia and other parts of the world. Asia is slowly climbing up in the percentages of races hosted and the same would be nice to see other parts of the world.
6. Drivers need to be mindful of rash and aggressive racing because as per the data, the chances of not finishing the race is quite high and you need to be in the race to get points for the Drivers Championships. Taking the example of Lewis Hamilton: He has one of the lowest accident ratios and has the highest finished percentages. He is the most successful driver in F1 after Michael Schumacher.

SECTION 2 : APPENDIX

ICA – Appendix: BI Design

A. Data Pre-Processing or Data Cleansing:

Data Pre-Processing and Data Cleansing is an integral part of any analysis and visualisation. From the Formula One dataset, I have taken only selected tables so that it's relevant to my business questions. So once the data has been uploaded to Power BI, the steps to rectify such as removing /N values, null values, renaming and removing columns, changing data types, rectifying errors, merging tables etc.

Demonstrated below are all the pre-processing steps with the 'result' column taken as an example

The steps are given below:

- **Replacing '/N' values:**

Before Pic:

1	resultid	raceid	driver ref	constructor ref	number	grid	position	positionText	positionOrder	points	laps	time	
2	1	18	hamilton	mclaren	22	1	1	1	1	10	58	34:50.6	
3	2	18	heidfeld	bmw_sauber	3	5	2	2	2	8	58	5:47.8	
4	3	18	rostberg	williams	7	7	3	3	3	6	58	8:16.3	
5	4	18	alonso	renault	5	11	4	4	4	5	58	11:18.1	
6	5	18	kovalainen	mclaren	23	3	5	5	5	4	58	15:04.4	
7	6	18	nakajima	williams	8	13	6	6	6	3	57	/N	
8	7	18	bourdais	toro_rosso	14	17	7	7	7	2	55	/N	
9	8	18	raikkonen	ferrari	1	15	8	8	8	1	53	/N	
10	9	18	kubica	bmw_sauber	4	2	W	R	9	0	47	/N	
11	10	18	glock	toyota	12	18	W	R	10	0	43	/N	
12	11	18	sato	super_aguri	18	19	W	R	11	0	32	/N	
13	12	18	piquet_jr	renault	6	20	W	R	12	0	30	/N	
14	13	18	massa	ferrari	2	4	W	R	13	0	29	/N	
15	14	18	couthard	red_bull	9	8	W	R	14	0	25	/N	
16	15	18	trulli	toyota	11	6	W	R	15	0	19	/N	
17	16	18	sutil	force_india	20	22	W	R	16	0	8	/N	
18	17	18	webber	red_bull	10	14	W	R	17	0	0	/N	
19	18	18	button	honda	16	12	W	R	18	0	0	/N	
20	19	18	davidson	super_aguri	19	21	W	R	19	0	0	/N	
21	20	18	vettel	toro_rosso	15	9	W	R	20	0	0	/N	
22	21	18	fisichella	force_india	21	16	W	R	21	0	0	/N	
23	22	18	barichello	honda	17	10	W	R	22	0	58	/N	
24	23	19	raikkonen	ferrari	1	2	1	1	10	56	31:18.6		
25	24	19	kubica	bmw_sauber	4	4	2	2	2	8	56	19:57	
26	25	19	kovaleinen	mclaren	23	8	3	3	3	6	56	38:45	
27	26	19	trulli	toyota	11	3	4	4	4	5	56	45:832	
28	27	19	hamilton	mclaren	22	9	5	5	5	4	56	46:548	
29	28	19	heidfeld	bmw_sauber	3	5	6	6	6	5	56	49:532	
30	29	19	webber	red_bull	10	6	7	7	7	2	56	+1:08.130	
31	30	19	massa	renault	2	3	8	8	8	1	56	+1:10.041	
32	31	19	couthard	red_bull	9	12	9	9	9	0	56	+1:18.220	
33	32	19	button	honda	15	11	10	10	10	0	56	+1:26.214	
34	33	19	piquet_jr	renault	6	13	11	11	11	0	56	+1:32.202	
35	34	19	fisichella	force_india	21	17	12	12	12	0	55	/N	
36	35	19	barichello	honda	17	14	13	13	13	0	55	/N	
37	36	19	rostberg	williams	7	16	14	14	14	0	55	/N	
38	37	19	davidson	super_aguri	19	21	15	15	15	0	55	/N	

Fig 1: Shows an image of the 'results' table before replacing '/N' values

Replacing '/N' values in column 'fastestLap' and 'rank' with '0' so that the null values will not having any meaning while doing aggregate functions.

Fig 2: An image showing how to '/N' values are replaced with 0

The same process has been done for other columns in results,

- **Renaming columns:**

Fig 3: The 'Constructor Chassis' column has been renamed to 'Constructor Champions'

- **Changing the data type:**

Changing the data type from text to the appropriate format of whole number

The screenshot shows a Power Query Editor window with a table of data. The 'rank' column is highlighted with a red oval, indicating it is selected for modification. A context menu is open over the column, with the 'Whole Number' option highlighted. Other options visible in the menu include Decimal Number, Fixed decimal number, Percentage, Date/Time, Date, Time, Date/Time/Timezone, Duration, Text, True/False, and Binary.

positionOrder	t ² ₃ points	t ² ₃ laps	A _C time	t ² ₃ milliseconds	A _C fastestLap	A _C rank	A _C fastestLapTime	t ² ₃ fastest
1	10	58 34.50.6		5690616 39		12	Decimal Number	5
2	8	58 5.478		5696094 41		12	Fixed decimal number	5
3	6	58 8.163		5698779 41		12	Whole Number	5
4	5	58 17.181		5707797 58		12	Percentage	5
5	4	58 18.014		5708630 43		12	Date/Time	5
6	3	57		null 50		12	Date	5
7	2	55		null 22		12	Time	5
8	1	53		null 20		12	Date/Time/Timezone	5
9	0	47		null 15		12	Duration	5
10	0	43		null 23		12	Text	5
11	0	32		null 24		12	True/False	5
12	0	30		null 20		12	Binary	5
13	0	29		null 23		12	Using Locale...	5
14	0	25		null 21		12		5
15	0	19		null 18		17		01:32.0
16	0	8		null 8				
17	0	0		null 0				
18	0	0		null n				

Fig 4 : A image showing how to change the data type

- **Rectifying Errors:**

✓ Results
▶ 25,840 rows loaded. 6 errors.

Fig 4: An image showing the errors while loading the dataset

The screenshot shows a Power Query Editor window with a table of data. The 'grid' column is highlighted with a red box, indicating it is selected for modification. A context menu is open over the column, with the 'Replace Errors' option highlighted. A tooltip for 'Replace Errors' says 'Replace all errors in the currently selected columns with the specified value.' The 'grid' column contains several 'Error' values, which are being replaced by '0'.

Row Number	t ² ₃ resultId	t ² ₃ raceld	A _C driver ref	A _C constructor ref	t ² ₃ number	t ² ₃ grid
1	17716	17716	732 marsh	brm	Error	
2	17717	17717	732 davis	porsche	Error	
3	17718	17718	732 abate	lotus-climax	Error	
4	17719	17719	732 burgess	cooper-climax	Error	
5	17940	17940	741 bordeu	lotus-climax	Error	
6	20320	20322	728 hailwood	lola	Error	

Fig 5: An image showing error rectification in the results column

Replace Errors with '0' as the values causing the error '\N' and since it's a number data type it shows up an error. Hence replacing it with 0 so that the errors are removed.

The screenshot shows a Power BI interface with a table and a modal dialog. The table has columns: A^B_C driver ref, A^B_C constructor ref, 1²₃ number, 1²₃ grid, and 1²₃ position. The 'number' column contains several 'Error' entries. A modal dialog titled 'Replace Errors' is open, prompting the user to enter a value to replace errors. The 'Value' field contains '0'. There are 'OK' and 'Cancel' buttons at the bottom of the dialog.

Fig 6: Replacing Error with 0 as the number column shows

After Error Correction:

The screenshot shows the same table after error correction. The 'number' column now contains all '0' values instead of 'Error'.

1.2 Row Number	1 ² ₃ resultId	1 ² ₃ raceId	A ^B _C driver ref	A ^B _C constructor ref	1 ² ₃ number	1 ² ₃ grid
17716	17716	732	marsh	brm	0	
17717	17717	732	davis	porsche	0	
17718	17718	732	abate	lotus-climax	0	
17719	17719	732	burgess	cooper-climax	0	
17940	17940	741	bordeu	lotus-climax	0	
20320	20322	728	hailwood	lola	0	

Fig 7: Image after error correction

- **Removing columns:**

The screenshot shows the Power Query Editor interface. A context menu is open over a cell in the 'time' column of a table. The menu options include 'Copy', 'Remove', 'Remove Other Columns', 'Duplicate Column', 'Add Column From Examples...', 'Remove Duplicates', 'Remove Errors', 'Change Type', and 'Transform'. The 'Remove' option is highlighted with a red box.

Fig 8: Image showing how to remove columns

• Merged Queries :

The screenshot shows the Power Query Editor interface with a red box highlighting the 'Merge Queries' option in the ribbon's 'Home' tab. A context menu is open over a table, with 'Merge Queries' also highlighted. The table contains several columns labeled with Greek letters like τ^2_3 , τ^3_3 , etc.

Fig 9: An Image showing the steps of merging

Merge

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

Results

resultId	raceId	driver ref	constructor ref	number	grid	position	positionText	positionOrder
1	18	hamilton	mclaren	22	1	1	1	
2	18	heidfeld	bmw_sauber	3	5	2	2	
3	18	rosberg	williams	7	7	3	3	
4	18	alonso	renault	5	11	4	4	
5	18			22	2	5	5	

Driver Data

Driver Id	driverRef	number	code	forename	surname	dob	nationality
1	hamilton	44	HAM	Lewis	Hamilton	07/01/1985	British
2	heidfeld	1N	HEI	Nick	Heidfeld	10/05/1977	German
3	rosberg	6	ROS	Nico	Rosberg	27/06/1985	German
4	alonso	14	ALO	Fernando	Alonso	29/07/1981	Spanish
5	kovalainen	1N	KOV	Heikki	Kovalainen	19/10/1981	Finnish

Join Kind

Left Outer (all from first, matching from second)

Use fuzzy matching to perform the merge

▷ Fuzzy matching options

Fig 10: Columns with the same name needs to be chosen to complete merging

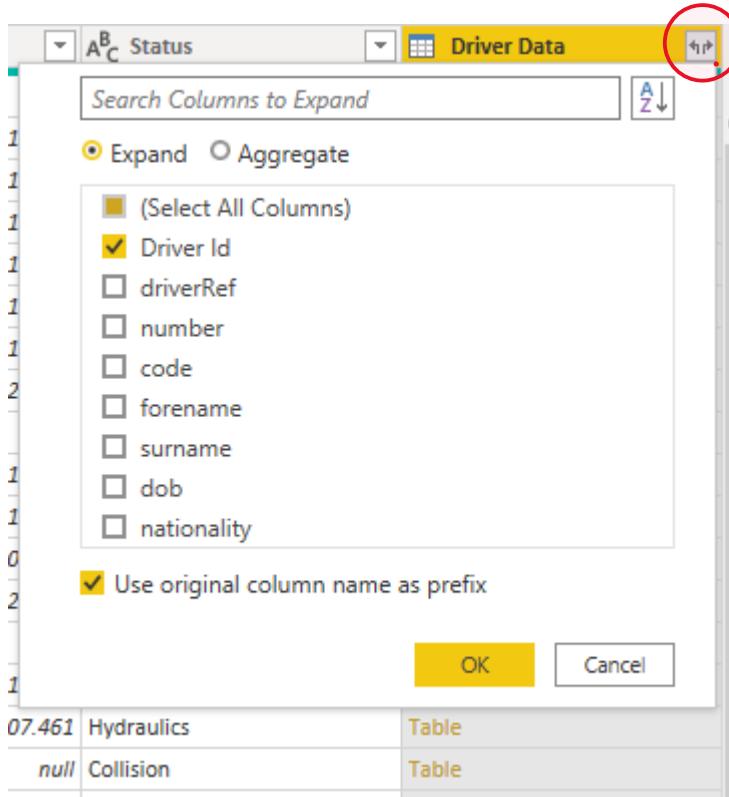


Fig 11 : Driver Id is the column we need and is expanded to get the required values

Line	1.2 fastestLapSpeed	A C Status	1 2 3 Driver Data.Drivers.Driver Id
	218.3	Finished	1
	209.033	Finished	1
	203.969	+1 Lap	1
	204.323	Finished	1
	222.085	Finished	1
	153.152	Finished	1
	202.559	Collision	1
	205.022	Finished	1
	199.398	Finished	1
	216.552	Finished	1

Fig 12 : Image after Merged Query is completed

• Results Table

The screenshot shows a Power BI interface with a results table. The table has columns: raceid, Driver ID, grid, position, points, Status, and Constructor ID. The 'Applied Steps' pane on the right lists several steps, including Replaced Value1, Replaced Value2, and Replaced Value3.

Fig 13: Results Table after pre-processing with the 'applied steps' image:

In a similar way, all the steps of pre-processing have been done for the other columns in the dataset.

• Drivers Data Table:

Before:

The screenshot shows a Power BI interface with a driver data table. The table includes columns such as understandingId, record, points, position, positionText, name, driverRef, number, code, and more. The 'Applied Steps' pane on the right shows various steps like Promoted Headers, Changed Type, and Removed Columns.

After:

The screenshot shows a Power BI interface with a driver data table after processing. The table includes columns such as Driver ID, driverRef, surname, name, dob, and nationality. The 'Applied Steps' pane on the right shows steps like Reordered Columns1, Removed Columns1, and Removed Columns2.

• Constructors Data Table:

Before:

constructorResultsId	Race Name	Year	points	status	constructorRef	name	nationality
1	Australian Grand Prix	2008	14	N	mclaren	McLaren	British
2	1				bmw_sauber	BMW Sauber	German
3	2		8	N			
4	3		9	N	williams	Williams	British
5	4		5	N	renault	Renault	French
6	5		2	V	toro_rosso	Toro Rosso	Italian
7	6		1	V	ferrari	Ferrari	Italian
8	7		0	V	toyota	Toyota	Japanese
9	8		0	V	super_aguri	Super Aguri	Japanese
10	9		0	V	red_bull	Red Bull	Austrian
11	10		0	V	force_india	Force India	Indian
12	11		0	V	honda	Honda	Japanese
13	12		10	V	ferrari	Ferrari	Italian
14	13		11	V	bmw_sauber	BMW Sauber	German
15	14		10	V	mclaren	McLaren	British

After:

ConstructorID	constructorRef	constructorName	nationality
1	mclaren	McLaren	British
2	bmw_sauber	BMW Sauber	German
3	williams	Williams	British
4	renault	Renault	French
5	toro_rosso	Toro Rosso	Italian
6	ferrari	Ferrari	Italian
7	toyota	Toyota	Japanese
8	super_aguri	Super Aguri	Japanese
9	red_bull	Red Bull	Austrian
10	force_india	Force India	Indian
11	honda	Honda	Japanese
12	spyker	Spyker	Dutch
13	mf1	MF1	Russian
14	spyker_mf1	Spyker MF1	Dutch
15	sauber	Sauber	Swiss
16	rar	R&R	Austrian

• Races Table:

Before:

record	year	round	circuited	name	date	time	url
1	2009	1	1	Australian Grand Prix	29/03/2009	06:00:00	http://en.wikipedia.org/wiki/2009_Australian_Grand_Prix
2	2009	2	2	Malaysian Grand Prix	05/04/2009	09:00:00	http://en.wikipedia.org/wiki/2009_Malaysian_Grand_Prix
3	2009	3	17	Chinese Grand Prix	19/04/2009	07:00:00	http://en.wikipedia.org/wiki/2009_Chinese_Grand_Prix
4	2009	4	3	Bahrain Grand Prix	26/04/2009	12:00:00	http://en.wikipedia.org/wiki/2009_Bahrain_Grand_Prix
5	2009	5	4	Spanish Grand Prix	10/05/2009	12:00:00	http://en.wikipedia.org/wiki/2009_Spanish_Grand_Prix
6	2009	6	6	Monaco Grand Prix	24/05/2009	12:00:00	http://en.wikipedia.org/wiki/2009_Monaco_Grand_Prix
7	2009	7	5	Turkish Grand Prix	07/06/2009	12:00:00	http://en.wikipedia.org/wiki/2009_Turkish_Grand_Prix
8	2009	8	9	British Grand Prix	21/06/2009	12:00:00	http://en.wikipedia.org/wiki/2009_British_Grand_Prix
9	2009	9	20	German Grand Prix	12/07/2009	12:00:00	http://en.wikipedia.org/wiki/2009_German_Grand_Prix
10	2009	10	11	Hungarian Grand Prix	26/07/2009	12:00:00	http://en.wikipedia.org/wiki/2009_Hungarian_Grand_Prix
11	2009	11	12	European Grand Prix	23/08/2009	12:00:00	http://en.wikipedia.org/wiki/2009_European_Grand_Prix
12	2009	12	13	Belgian Grand Prix	30/08/2009	12:00:00	http://en.wikipedia.org/wiki/2009_Belgian_Grand_Prix
13	2009	13	14	Italian Grand Prix	13/09/2009	12:00:00	http://en.wikipedia.org/wiki/2009_Italian_Grand_Prix
14	2009	14	15	Singapore Grand Prix	27/09/2009	12:00:00	http://en.wikipedia.org/wiki/2009_Singapore_Grand_Prix
15	2009	15					
16	2009	16					

After:

record	year	round	circuited	raceName	date
1	1	2009	1	1. Australian Grand Prix	29/03/2009
2	2	2009	2	2. Malaysian Grand Prix	05/04/2009
3	3	2009	3	3. Chinese Grand Prix	19/04/2009
4	4	2009	4	3. Bahrain Grand Prix	26/04/2009
5	5	2009	5	4. Spanish Grand Prix	10/05/2009
6	6	2009	6	6. Monaco Grand Prix	24/05/2009
7	7	2009	7	5. Turkish Grand Prix	26/07/2009
8	8	2009	8	8. German Grand Prix	24/09/2009
9	9	2009	9	20. German Grand Prix	22/07/2009
10	10	2009	10	21. Hungarian Grand Prix	26/07/2009
11	11	2009	11	22. European Grand Prix	21/08/2009
12	12	2009	12	15. Belgian Grand Prix	08/09/2009
13	13	2009	13	14. Italian Grand Prix	12/09/2009
14	14	2009	14	25. Singapore Grand Prix	27/09/2009
15	15	2009	15	22. Japanese Grand Prix	04/10/2009
16	16	2009	16	18. Brazilian Grand Prix	18/10/2009

• Championship Data

Before:

Year	Driver	Age	Driver Team	Constructor's Champions
1990	Giuseppe Farina	44	Alfa Romeo	
1991	Juan Fangio	40	Alfa Romeo	
1992	Alberto Ascari	34	Ferrari	
1993	Alberto Ascari	35	Ferrari	
1994	Juan Fangio	43	Maserati	
1995	Juan Fangio	43	Mercedes	
1996	Juan Fangio	44	Mercedes	
1997	Juan Fangio	45	Ferrari	
1998	Juan Fangio	46	Maserati	
1999	Mike Hawthorn	29	Ferrari	Vanwall
2000	Jack Brabham	33	Cooper	Cooper
2001	Jack Brabham	34	Cooper	Cooper
2002	Phil Hill	34	Ferrari	Ferrari
2003	Graham Hill	33	BRM	BRM
2004	Jim Clark	27	Lotus	Lotus
2005	John Surtees	30	Ferrari	Ferrari

After:

Championship ID	Year	Driver	Age	Driver Team	Constructor's Champions	Constructor ID
1	2	1974 Emerson Fittipaldi	27	McLaren	McLaren	1
2	1	1984 Niki Lauda	28	McLaren	McLaren	1
3	5	1985 Alan Prost	30	McLaren	McLaren	1
4	4	1988 Ayrton Senna	28	McLaren	McLaren	1
5	5	1989 Ayrton Senna	34	McLaren	McLaren	1
6	6	1990 Ayrton Senna	30	McLaren	McLaren	1
7	7	1992 Ayrton Senna	32	McLaren	McLaren	2
8	8	1998 Mika Häkkinen	30	McLaren	McLaren	1
9	9	1980 Alan Jones	34	Williams	Williams	3
10	10	1982 Nelson Piquet	29	Brabham	Williams	3
11	11	1986 Alan Prost	32	McLaren	Williams	3
12	12	1987 Nelson Piquet	35	Williams	Williams	3
13	13	1992 Nigel Mansell	39	Williams	Williams	3
14	14	1993 Alain Prost	38	Williams	Williams	3
15	15	1994 Michael Schumacher	25	Benetton	Williams	3

• Circuit Data

Before :

Column1	Column2	Column3	Column4	Column5	Column6	Column7	Column8	Column9	Column10	Column11	Column12	Column13	Column14	Column15
1	circuits	circuitid	name	location	country	lat	lon	alt	url					
2		albert_park	Albert Park Grand Prix Circuit	Melbourne	Australia	-37.8407	144.968	10	http://en.wikipedia.org/wiki/Albert_Park_Grand_Prix_Circuit					
3		sepang	Sepang International Circuit	Kuala Lumpur	Malaysia	2.70683	101.718	18	http://en.wikipedia.org/wiki/Sepang_International_Circuit					
4		bahrain	Bahrain International Circuit	Sakhir	Bahrain	26.0325	50.51305	7	http://en.wikipedia.org/wiki/Bahrain_International_Circuit					
5		catalunya	Circuit de Barcelona-Catalunya	Montmeló	Spain	41.37	2.26111	109	http://en.wikipedia.org/wiki/Circuit_de_Barcelona-Catalunya					
6		istanbul	Istanbul Park	Istanbul	Turkey	40.9517	29.405	130	http://en.wikipedia.org/wiki/Istanbul_Park					
7		monaco	Circuit de Monaco	Monte-Carlo	Monaco	43.7347	7.42056	7	http://en.wikipedia.org/wiki/Circuit_de_Monaco					
8		villeneuve	Circuit Gilles Villeneuve	Montreal	Canada	45.5	-75.5228	13	http://en.wikipedia.org/wiki/Circuit_Gilles_Villeneuve					
9		magny_cours	Circuit de Nevers Magny-Cours	Magny-Cours	France	46.8644	3.16361	228	http://en.wikipedia.org/wiki/Circuit_de_Nevers_Magny-Cours					
10		silverstone	Silverstone Circuit	Silverstone	UK	52.0786	-1.01894	153	http://en.wikipedia.org/wiki/Silverstone_Circuit					
11		hockenheimring	Hockenheimring	Hockenheim	Germany	49.3278	8.56583	103	http://en.wikipedia.org/wiki/Hockenheimring					
12		hungaroring	Hungaroring	Budapest	Hungary	47.5789	19.2466	264	http://en.wikipedia.org/wiki/Hungaroring					
13		valencia	Valencia Street Circuit	Valencia	Spain	39.4589	-0.331667	4	http://en.wikipedia.org/wiki/Valencia_Street_Circuit					
14		spa	Circuit de Spa-Francorchamps	Spa	Belgium	50.4372	5.97139	401	http://en.wikipedia.org/wiki/Circuit_de_Spa-Francorchamps					
15		monza	Autodromo Nazionale di Monza	Monza	Italy	45.6156	9.28111	162	http://en.wikipedia.org/wiki/Autodromo_Nazionale_di_Monza					
16		marina_bay	Marina Bay Street Circuit	Marina Bay	Singapore	1.2914	103.864	18	http://en.wikipedia.org/wiki/Marina_Bay_Street_Circuit					

After:

circuitid	circuitRef	name	location	lat	lon	alt	continent	country
1	3	bahrain	Bahrain International Circuit	26.025	50.5106	7	Asia	Bahrain
2	16	fuj	Fuji Speedway	35.3717	138.927	583	Asia	Japan
3	22	suzuka	Suzuka Circuit	34.4831	136.542	45	Asia	Japan
4	28	okayama	Okayama International Circuit	34.915	134.222	266	Asia	Japan
5	37	shanghai	Shanghai International Circuit	31.2339	121.455	2	Asia	China
6	20	ysam	Yas Marina	25.9904	29.9767	1450	Africa	South Africa
7	56	george	Prince George Circuit	-33.0498	27.8735	15	Africa	South Africa
8	35	wongam	Korean International Circuit	34.7333	126.417	0	Asia	Korea
9	64	ain_dhab	Ain Dhab	33.796	-7.6875	19	Africa	Morocco
10	68	buldh	Buldh International Circuit	28.5487	77.5531	19	Asia	India
11	2	sepang	Sepang International Circuit	3.7083	102.739	18	Asia	Malaysia
12	78	losail	Losail International Circuit	21.6319	39.1044	15	Asia	Qatar
13	71	sotchi	Sochi Autodrom	43.4057	39.9578	2	Asia	Russia
14	77	jeddah	Jeddah Corniche Circuit	21.6319	39.1044	15	Asia	Saudi Arabia
15	25	marina_bay	Marina Bay Street Circuit	1.2914	103.864	18	Asia	Singapore

B. Data Modeling via Star Schema:

Among the 16 tables provided from Kaggle source, I have chosen only selected tables.

The screenshot displays a data modeling interface with seven tables arranged in two rows:

- Top Row:**
 - circuits**: Columns include alt, circuitid, circuitRef, country, lat, lng, location, name, and url.
 - constructor_results (Updated)**: Columns include constructorRef, constructorResultsid, name, nationality, points, Race Name, status, and Year.
 - qualifying**: Columns include constructorid, driverid, number, position, q1, q2, q3, qualifyid, and raceld.
 - results (Updated)**: Columns include constructor ref, driver ref, fastestLap, fastestLapSpeed, fastestLapTime, grid, laps, milliseconds, and number.
- Bottom Row:**
 - driver_standings (Updated)**: Columns include code, dob, driverRef, driverStandingsid, forename, nationality, number, points, position, positionText, raceld, and surname.
 - races**: Columns include circuitid, date, fp1_date, fp1_time, fp2_date, fp2_time, fp3_date, fp3_time, name, qual_date, qual_time, raceld, round, sprint_date, sprint_time, time, url, and year.
 - pit_stops**: Columns include driverid, lap, milliseconds, Minutes, raceld, Seconds, and stop.

Fig 14: Before Data Modeling Pic:

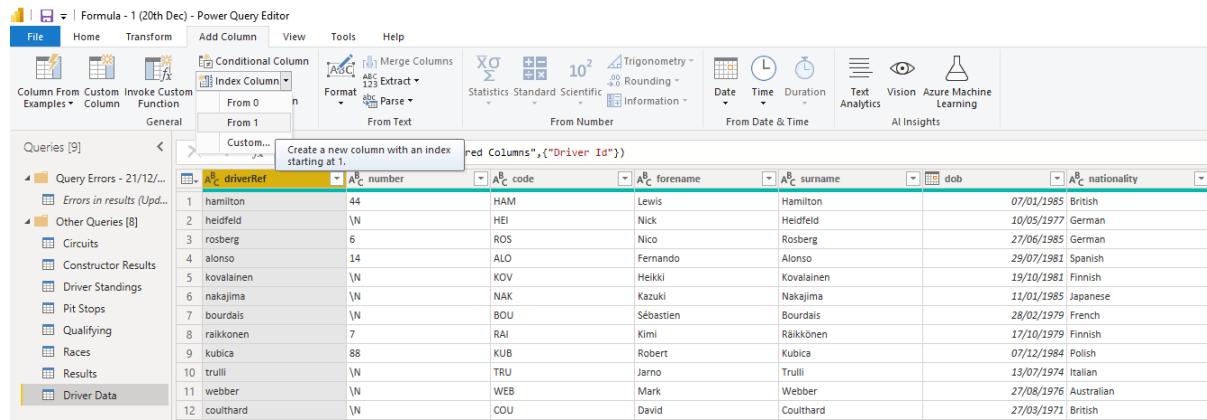
Since the dataset already has some structure, I have re-created some tables such as drivers data and constructors data to display the normalisation technique. I have then created an index column for each of these two tables and then did a merge query to connect these new columns. I have also discarded the ‘pitstops’ and ‘qualifying’ tables as I am not using them for my visualisation.

The steps for normalisation and creating new relationships is given below:

1. Normalisation

Creating a new table called ‘Drivers Data’

1. Duplicating the ‘Driver standings’ data and removing unwanted columns

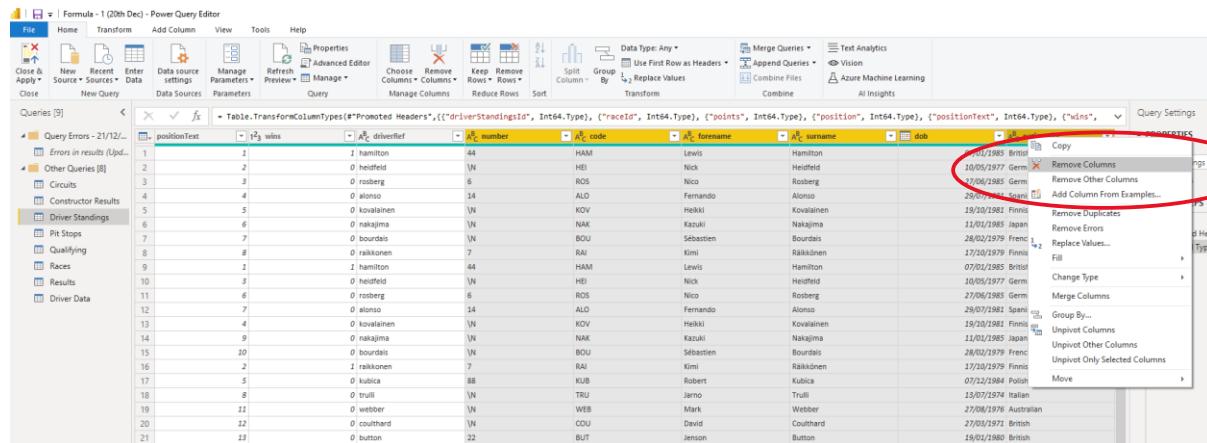


The screenshot shows the Power Query Editor interface with the 'Driver Data' query selected. The table contains 12 rows of driver information, including name, number, code, forename, surname, date of birth, and nationality. The columns are labeled A_C number, A_C code, A_C forename, A_C surname, dob, and A_C nationality.

	A _C number	A _C code	A _C forename	A _C surname	dob	A _C nationality
1	hamilton	44	HAM	Lewis	Hamilton	07/01/1985 British
2	heidfeld	\N	HEI	Nick	Heidfeld	10/05/1977 German
3	rosberg	6	ROS	Nico	Rosberg	27/06/1985 German
4	alonso	14	ALO	Fernando	Alonso	29/07/1981 Spanish
5	kovalainen	\N	KOV	Heikki	Kovalainen	19/10/1981 Finnish
6	nakajima	\N	NAK	Kazuki	Nakajima	11/01/1985 Japanese
7	bourdais	\N	BOU	Sébastien	Bourdais	28/02/1979 French
8	raikkonen	7	RAI	Kimi	Räikkönen	17/10/1979 Finnish
9	kubica	88	KUB	Robert	Kubica	07/12/1984 Polish
10	trulli	\N	TRU	Jarno	Trulli	13/07/1974 Italian
11	webber	\N	WEB	Mark	Webber	27/08/1976 Australian
12	coulthard	\N	COU	David	Coulthard	27/03/1971 British

Fig 15: Before Pic of ‘Driver Standings’ table:

2. Removing unnecessary columns in Driver Standings



The screenshot shows the Power Query Editor with a transformation step applied to the 'Driver Data' query. The step uses the formula: Table.TransformColumnTypes(Table.PromoteHeaders(“driverStandings”, {{"raceId”, Int64.Type}, {"resultId”, Int64.Type}, {"points”, Int64.Type}, {"position”, Int64.Type}, {"positionText”, Int64.Type}, {"wins”, Int64.Type}})). The 'wins' column is highlighted in yellow, indicating it is selected for removal. A context menu is open over the 'wins' column, with the 'Remove Columns' option circled in red.

	A _C number	A _C code	A _C forename	A _C surname	dob	A _C nationality
1	1	hamilton	44	Hamilton	07/01/1985 British	
2	2	heidfeld	\N	Heidfeld	10/05/1977 German	
3	3	rosberg	6	Rosberg	27/06/1985 German	
4	4	alonso	14	Fernando	29/07/1981 Spanish	
5	5	kovvalainen	\N	Heikki	19/10/1981 Finnish	
6	6	nakajima	\N	Kazuki	11/01/1985 Japanese	
7	7	bourdais	\N	Sébastien	28/02/1979 French	
8	8	raikkonen	7	Räikkönen	17/10/1979 Finnish	
9	1	hamilton	44	Lewis	07/01/1985 British	
10	3	heidfeld	\N	Heidfeld	10/05/1977 German	
11	6	rosberg	6	Rosberg	27/06/1985 German	
12	7	alonso	14	Fernando	29/07/1981 Spanish	
13	4	kovvalainen	\N	Heikki	19/10/1981 Finnish	
14	9	nakajima	\N	Kazuki	11/01/1985 Japanese	
15	10	bourdais	\N	Sébastien	28/02/1979 French	
16	2	raikkonen	7	Räikkönen	17/10/1979 Finnish	
17	5	kubica	88	Robert	07/12/1984 Polish	
18	8	trulli	\N	Trulli	13/07/1974 Italian	
19	11	webber	\N	Mark	27/08/1976 Australian	
20	12	coulthard	\N	David	27/03/1971 British	
21	13	button	22	Jenson	19/01/1980 British	

Fig 16: Removing unnecessary columns in ‘Driver Standing’s table

3. Creating Index column driver Id:

We are selecting the index column from 1

The screenshot shows the Power BI desktop interface with the 'Create a new column with an index' dialog open. The 'From 1' dropdown is circled in red. The table preview shows a new column 'Index' highlighted with a red box.

	DriverRef	number	code	forename	surname	dob	nationality	Index
1	hamilton	44	HAM	Lewis	Hamilton	07/01/1985	British	1
2	heidfeld	\N	HEI	Nick	Heidfeld	10/05/1977	German	2
3	rosberg	6	ROS	Nico	Rosberg	27/06/1985	German	3
4	alonso	14	ALO	Fernando	Alonso	29/07/1981	Spanish	4
5	kovalainen	\N	KOV	Heikki	Kovalainen	19/10/1982	Finnish	5
6	nakajima	\N	NAK	Kazuki	Nakajima	11/01/1985	Japanese	6
7	bourdais	\N	BOU	Sébastien	Bourdais	28/02/1979	French	7
8	raikkonen	7	RAI	Kimi	Räikkönen	17/10/1979	Finnish	8
9	kubica	88	KUB	Robert	Kubica	07/12/1984	Polish	9
10	trulli	\N	TRU	Jarno	Trulli	18/07/1974	Italian	10
11	webber	\N	WEB	Mark	Webber	27/06/1986	Australian	11
12	coulthard	\N	COU	David	Coulthard	27/02/1971	British	12
13	button	22	BUT	Jenson	Button	19/07/1980	British	13
14	piquet_ir	\N	PQI	Nelson	Piquet Jr.	25/07/1985	Brazilian	14

Fig 17:Image showing the main options to choose for creating an index column

After pic:

The screenshot shows the Power BI desktop interface with the table after creating the 'Driver Id' column. The 'Driver Id' column is highlighted with a red box.

	Driver Id	DriverRef	number	code	forename	surname	dob	nationality
1	1	hamilton	44	HAM	Lewis	Hamilton	07/01/1985	British
2	2	heidfeld	\N	HEI	Nick	Heidfeld	10/05/1977	German
3	3	rosberg	6	ROS	Nico	Rosberg	27/06/1985	German
4	4	alonso	14	ALO	Fernando	Alonso	29/07/1981	Spanish
5	5	kovvalainen	\N	KOV	Heikki	Kovalainen	19/10/1982	Finnish
6	6	nakajima	\N	NAK	Kazuki	Nakajima	11/01/1985	Japanese
7	7	bourdais	\N	BOU	Sébastien	Bourdais	28/02/1979	French
8	8	raikkonen	7	RAI	Kimi	Räikkönen	17/10/1979	Finnish

Fig 18:And after pic of the new 'Drivers ID' column

Creating a new Constructors Data Table:

The same steps shown in creating the 'Drivers ID' table has been used to create the Constructors Data table.

The constructors data table has been created from the **constructors results** table shown below.

	constructorResultsId	Race Name	Year	points	status	constructorRef	name	nationality
1	1	Australian Grand Prix	2008	14 \N	mclaren	McLaren	British	
2	2	Australian Grand Prix	2008	8 \N	bmw_sauber	BMW Sauber	German	
3	3	Australian Grand Prix	2008	9 \N	williams	Williams	British	
4	4	Australian Grand Prix	2008	5 \N	renault	Renault	French	
5	5	Australian Grand Prix	2008	2 \N	toro_rosso	Toro Rosso	Italian	
6	6	Australian Grand Prix	2008	1 \N	ferrari	Ferrari	Italian	
7	7	Australian Grand Prix	2008	0 \N	toyota	Toyota	Japanese	
8	8	Australian Grand Prix	2008	0 \N	super_aguri	Super Aguri	Japanese	
9	9	Australian Grand Prix	2008	0 \N	red_bull	Red Bull	Austrian	
10	10	Australian Grand Prix	2008	0 \N	force_india	Force India	Indian	

Fig 19:Before pic of the Constructor Results table

	Constructor Id	constructorRef	name	nationality
1	1	mclaren	McLaren	British
2	2	bmw_sauber	BMW Sauber	German
3	3	williams	Williams	British
4	4	renault	Renault	French
5	5	toro_rosso	Toro Rosso	Italian
6	6	ferrari	Ferrari	Italian
7	7	toyota	Toyota	Japanese
8	8	super_aguri	Super Aguri	Japanese
9	9	red_bull	Red Bull	Austrian
10	10	force_india	Force India	Indian

Fig 20::After pic of the Constructors Data table

2. Creating New relationships

- Showing data relationship between Drivers Data and Results table

Step 1: Select 'Transform Data'

The screenshot shows the Power BI ribbon with several icons and sections. The 'Transform data' icon, which is a grid with a pencil, is highlighted with a red circle. Below the ribbon, there are three preview panes: 'Images' (empty), 'Continent' (listing Continent and Country), and 'Circuits' (listing alt, circuitId, circuitRef, Continent, Country, lat, lng, location, name). The 'Continent' pane is currently selected.

Fig 21:First step is to select Transform data to create a new relationship

Step 2: Click on 'New'

The screenshot shows the 'Manage relationships' dialog box. At the top, there are tabs for 'Active', 'From: Table (Column)', and 'To: Table (Column)'. The main area displays the message 'There are no relationships defined yet.' At the bottom, there are four buttons: 'New...', 'Autodetect...', 'Edit...', and 'Delete'. The 'New...' button is circled in red. A yellow 'Close' button is located at the bottom right.

Step 3:

Select the key column with which you would connect the two tables. Creating relationship with Drivers Data and Results table.

We choose the cardinality as One to Many Or Many to One and the cross filter direction as 'Single'.

×

Edit relationship

Select tables and columns that are related.

Results											
rid	position	points	Status	Race Wins	Accidents	Finished	Podium	Decade	Driver ID	Co	Ca
4	0	0	Engine	0	0	0	0	2000's	8		
4	0	0	Engine	0	0	0	0	2000's	19		
2	0	0	Engine	0	0	0	0	2000's	32		

Driver Data						
driverRef	forename	surname	dob	nationality	Driver Id	Driver Name
andretti	Michael	Andretti	05 October 1962	American	114	Michael Andretti
cheever	Eddie	Cheever	10 January 1958	American	149	Eddie Cheever
sullivan	Danny	Sullivan	09 March 1950	American	185	Danny Sullivan

Cardinality **Cross filter direction**

Make this relationship active Apply security filter in both directions
 Assume referential integrity

OK **Cancel**

Fig 22: 'Driver ID' is the key which connects these two tables

We choose the Cardinality as 'Many to One' or 'One to Many' because we want data to connect from a dimension table such as 'Driver Data' which has all unique values (Single) to the table having ('Many').

We choose the cross filter to be 'Single' because I have two or more tables that also have lookup tables and the data modelling is a Snowflake Schema.

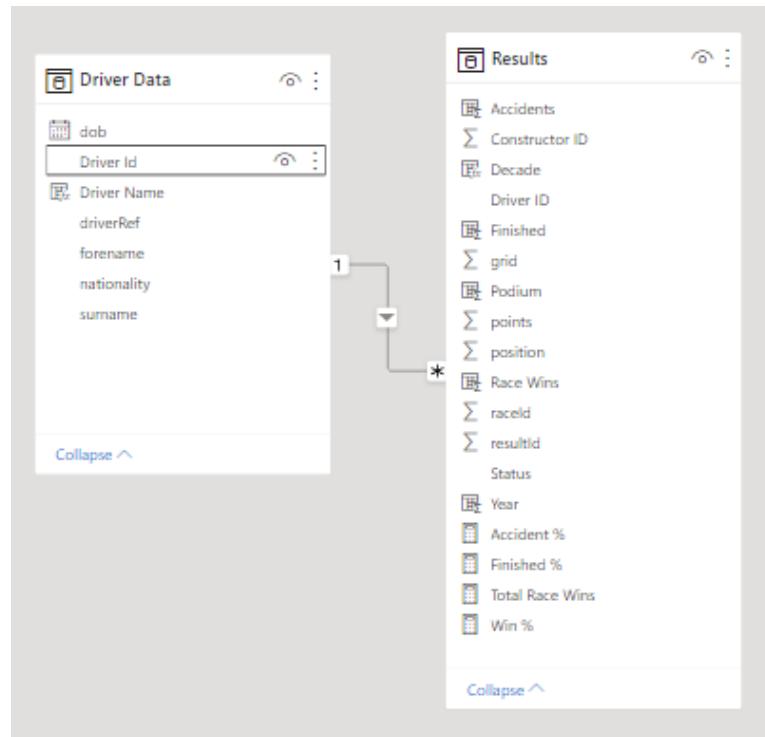


Fig 23: An after pic showing the relationship between 'Driver Data' and 'Results' table

- **Creating the relationship between Constructors Data and Results table:**

Edit relationship

Select tables and columns that are related.

Results

points	Status	Race Wins	Accidents	Finished	Podium	Decade	Driver ID	Constructor ID	Ye
0	Engine	0	0	0	0	2000's	8	6	
0	Engine	0	0	0	0	2000's	19	6	
0	Engine	0	0	0	0	2000's	32	6	

Constructors Data

constructorRef	Constructor Name	nationality	Constructor Id
mclaren	McLaren	British	1
williams	Williams	British	3
bar	BAR	British	16

Cardinality: Many to one (*:1) **Cross filter direction**: Single

Make this relationship active Apply security filter in both directions
 Assume referential integrity

OK **Cancel**

Fig 23: 'Constructor ID' is the key which connects these two tables

'Constructor ID' is the key which connects these two tables

Create relationship

Select tables and columns that are related.

Races

raceId	year	round	circuitId	Race Name	date
13	2009	13	14	Italian Grand Prix	13 September 2009
31	2008	14	14	Italian Grand Prix	14 September 2008
48	2007	13	14	Italian Grand Prix	09 September 2007

Results

resultId	raceId	grid	position	points	Status	Race Wins	Accidents	Finished	Year	Podium
246	29	4	0	0	Engine	0	0	0	2008	
13	18	4	0	0	Engine	0	0	0	2008	
1114	69	2	0	0	Engine	0	0	0	2006	

Cardinality: One to many (1:*)

Cross filter direction: Single

Make this relationship active

Assume referential integrity

Apply security filters in both directions

OK Cancel

'raceid' is the key which connects these two tables

- Creating Relationship between custom columns created for 'Race Car Images' to 'Constructors Data'

Create relationship

Select tables and columns that are related.

Race Car Images

Constructor	Team Link
McLaren	https://www.formula1.com/content/dam/fom-website...
Williams	https://www.formula1.com/content/dam/fom-website...
Ferrari	https://www.formula1.com/content/dam/fom-website...

Constructors Data

constructorRef	Constructor Name	nationality	Constructor Id
mclaren	McLaren	British	1
williams	Williams	British	3
bar	BAR	British	16

Cardinality: One to one (1:1)

Cross filter direction: Both

Make this relationship active

Assume referential integrity

OK Cancel

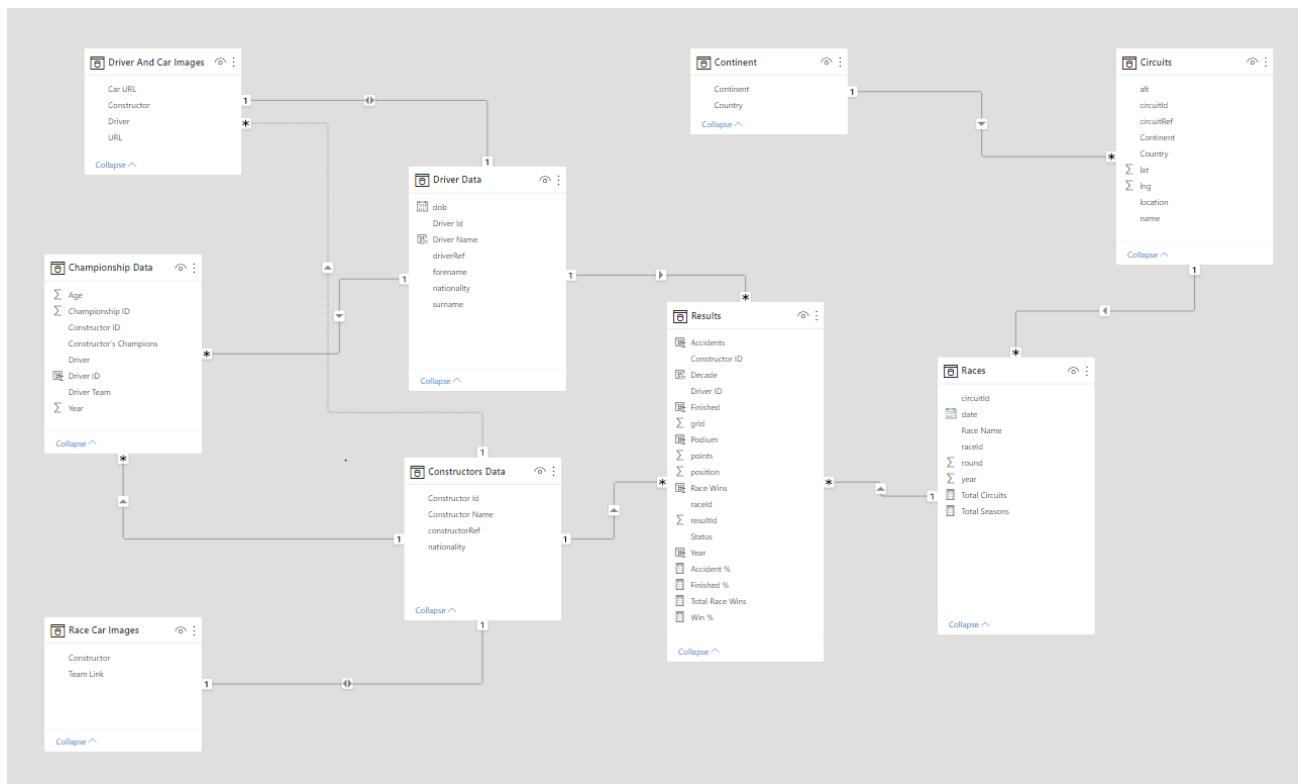
In this case we have a One to One (1:1) relationship because the Race Car Images is a unique column with a purpose only for the images and the Cross filter direction is 'Both' so that the information can flow both ways.

Hence using the same method, all the relationships have been created

All the relationships and their keys:

1. Driver Data to Results Data **Key : Driver ID:**
Cardinality : One to Many
Cross Filter Direction : Single
2. Constructors Data to Results Data **Key : Constructor ID:**
Cardinality : One to Many
Cross Filter Direction : Single
3. Races to Results Data **Key : raceid**
Cardinality : One to Many
Cross Filter Direction : Single
4. Races to Results Data **Key : raceid**
Cardinality : One to Many
Cross Filter Direction : Single
5. Circuits to Races Data **Key : circuitid**
Cardinality : One to Many
Cross Filter Direction : Single
6. Championship Data to Driver Data **Key : Driver ID**
Cardinality : Many to One
Cross Filter Direction : Single
7. Championship Data to Constructors Data **Key : Constructor ID**
Cardinality : Many to One
Cross Filter Direction : Single
8. Race Car Images to Constructor Data **Key : Driver Name**
Cardinality : One to One
Cross Filter Direction : Both
9. Driver and Car Images to Constructor Data **Key : Constructor Name**
Cardinality : Many to One
Cross Filter Direction : Single
10. Driver and Car Images to Driver Data **Key : Driver Name**
Cardinality : One to One
Cross Filter Direction : Both

Final Screenshot after creating the relationships:



C. DAX and M Language

DAX:

DAX and M language has been used to create custom columns and measures so as to get the right visualisation.

- **Decade:** A column named '**Decade**' has been created which is used in the zoomable sunburst chart so to understand which teams are dominant and also to keep the chart clean instead of having multiple years scattered throughout.

The screenshot shows the Power BI Data Editor interface. At the top, there is a code editor window containing DAX code for creating a 'Decade' column based on the year. Below the code editor is a preview table with columns: raceId, grid, position, points, Status, Race Wins, Accidents, Finished, Year, Podium, Decade, Driver ID, Constructor ID. Two rows of data are shown: one for raceId 29 and one for raceId 18. The 'Decade' column is highlighted in yellow in both the code and the preview table.

```

1 Decade = IF(Results[Year] < 1960, "1950's",
2             IF(AND(Results[Year] < 1970, Results[Year] >= 1960), "1960's",
3                 IF(AND(Results[Year] < 1980, Results[Year] >= 1970), "1970's",
4                     IF(AND(Results[Year] < 1990, Results[Year] >= 1980), "1980's",
5                         IF(AND(Results[Year] < 2000, Results[Year] >= 1990), "1990's",
6                             IF(AND(Results[Year] < 2010, Results[Year] >= 2000), "2000's",
7                                 IF(AND(Results[Year] < 2020, Results[Year] >= 2010), "2010's",
8                                     IF(AND(Results[Year] < 2030, Results[Year] >= 2020), "2020's", "Unknown")))))))))
9

```

raceId	grid	position	points	Status	Race Wins	Accidents	Finished	Year	Podium	Decade	Driver ID	Constructor ID
29	4	0	0	Engine	0	0	0	2008	0	2000's	8	6
18	4	0	0	Engine	0	0	0	2008	0	2000's	19	6

- **Podium:** A column named '**Podium**' has been created to understand how many podiums drivers achieve. This is especially useful in the 'KPI Influencers' visual.

The screenshot shows the Power BI Data Editor interface. At the top, there is a code editor window containing DAX code for creating a 'Podium' column based on the position. Below the code editor is a preview table with columns: raceId, grid, position, points, Status, Race Wins, Accidents, Finished, Year, Podium, Decade, Driver ID, Constructor ID. Two rows of data are shown: one for raceId 29 and one for raceId 18. The 'Podium' column is highlighted in yellow in both the code and the preview table.

```

1 Podium = IF
2             OR(if(Results[position] = 1,1,
3                 IF(Results[position] = 2,1,
4                     IF(Results[position]=3,1,0))),0,1,0)

```

raceId	grid	position	points	Status	Race Wins	Accidents	Finished	Year	Podium	Decade	Driver ID	Constructor ID
29	4	0	0	Engine	0	0	0	2008	0	2000's	8	6
18	4	0	0	Engine	0	0	0	2008	0	2000's	19	6

- **Race Wins:** A column named '**Race wins**' is created to understand the total wins by drivers and constructors. This column has been extremely useful in creating many visual for the driver and constructor analysis.

The screenshot shows the Power BI Data Editor interface. At the top, there is a code editor window containing DAX code for creating a 'Race Wins' column based on the position. Below the code editor is a preview table with columns: raceId, grid, position, points, Status, Race Wins, Accidents, Finished, Year, Podium, Decade, Driver ID, Constructor ID. Two rows of data are shown: one for raceId 29 and one for raceId 18. The 'Race Wins' column is highlighted in yellow in both the code and the preview table.

```

1 Race Wins =
2             IF(Results[position] = 1, 1 ,0)

```

raceId	grid	position	points	Status	Race Wins	Accidents	Finished	Year	Podium	Decade	Driver ID	Constructor ID
29	4	0	0	Engine	0	0	0	2008	0	2000's	8	6
18	4	0	0	Engine	0	0	0	2008	0	2000's	19	6

- **Accidents:** A column named '**Accidents**' is created to understand how the average accident percentage over the years in formula one history.

1 Accidents =							
2 IF(Results[Status] = "Accident", 1, 0)							
aceld	grid	position	points	Status	Race Wins	Accidents	F
29	4	0	0	Engine	0	0	
..

- Driver name: A column named '**Driver Name**' is created to join driver forename and surname together.

```
1 Driver Name = 'Driver Data'[forename] & " " & 'Driver Data'[surname]
```

forename	surname	dob	nationality	Driver Id	Driver Name
Michael	Andretti	05 October 1962	American	114	Michael Andretti
Eddie	Cheever	10 January 1958	American	149	Eddie Cheever
Danny	Sullivan	09 March 1950	American	185	Danny Sullivan

- **Driver ID**

A column named '**Driver ID**' is created by using the '**LOOKUPVALUE**' function.

1 Driver ID =							
2 LOOKUPVALUE('Driver Data'[Driver Id],'Driver Data'[Driver Name],'Championship Data'[Driver])							
Driver	Age	Driver Team	Constructor's Champions	Constructor ID	Championship ID	Driver ID	
Fittipaldi	27	McLaren	McLaren	1	1	214	
Lauda	35	McLaren	McLaren	1	2	174	
..

Measures:

- Win %

```
1 Win % = [sum(Results[Race Wins])/sum(Results[Finished])]
```



- Total Race Wins

```
Total Race Wins = SUM([Results[Race Wins]])
```



- Total Circuits

```
Total Circuits = DISTINCTCOUNT([Circuits[circuitId]])
```



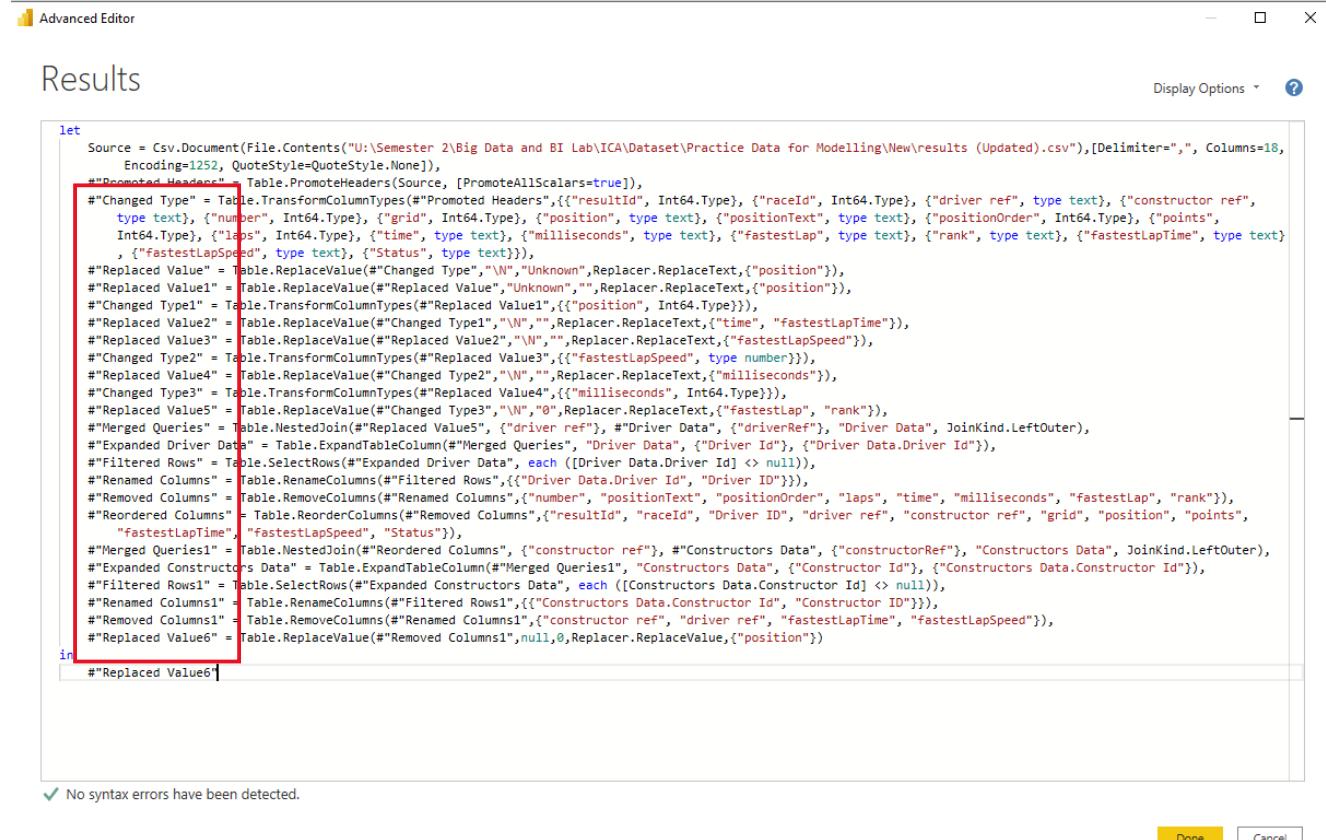
- Total Seasons

```
Total Seasons = DISTINCTCOUNT([Races[year]])
```



M Language

- Results Table:



```
let
    Source = Csv.Document(File.Contents("U:\Semester 2\Big Data and BI Lab\ICA\Dataset\Practice Data for Modelling\New\results (Updated).csv"),[Delimiter=",", Columns=18, Encoding=1252, QuoteStyle=QuoteStyle.None]),
    #"Promoted Headers" = Table.PromoteHeaders(Source, [PromoteAllScalars=true]),
    #"Changed Type" = Table.TransformColumnTypes(#"Promoted Headers",{{"resultId", Int64.Type}, {"raceId", Int64.Type}, {"driver ref", type text}, {"constructor ref", type text}, {"number", Int64.Type}, {"grid", Int64.Type}, {"position", type text}, {"positionText", type text}, {"positionOrder", Int64.Type}, {"points", Int64.Type}, {"laps", Int64.Type}, {"time", type text}, {"milliseconds", type text}, {"fastestLap", type text}, {"rank", type text}, {"fastestLapTime", type text}, {"fastestLapSpeed", type text}, {"Status", type text}}),
    #"Replaced Value" = Table.ReplaceValue(#"Changed Type","\\N","Unknown",Replacer.ReplaceText,{"position"}),
    #"Replaced Value1" = Table.ReplaceValue(#"Replaced Value","Unknown","",Replacer.ReplaceText,{"position"}),
    #"Changed Type1" = Table.TransformColumnTypes(#"Replaced Value1",{{"position", Int64.Type}}),
    #"Replaced Value2" = Table.ReplaceValue(#"Changed Type1","\\N","",Replacer.ReplaceText,{"time", "fastestLapTime"}),
    #"Replaced Value3" = Table.ReplaceValue(#"Replaced Value2","\\N","",Replacer.ReplaceText,{"fastestLapSpeed"}),
    #"Changed Type2" = Table.TransformColumnTypes(#"Replaced Value3",{{"fastestLapSpeed", type number}}),
    #"Replaced Value4" = Table.ReplaceValue(#"Changed Type2","\\N","",Replacer.ReplaceText,{"milliseconds"}),
    #"Changed Type3" = Table.TransformColumnTypes(#"Replaced Value4",{{"milliseconds", Int64.Type}}),
    #"Replaced Value5" = Table.ReplaceValue(#"Changed Type3","\\N","0",Replacer.ReplaceText,{"fastestLap", "rank"}),
    #"Merged Queries" = Table.NestedJoin(#"Replaced Value5", {"driver ref"}, #Driver Data, {"driverRef"}, "Driver Data", JoinKind.LeftOuter),
    #"Expanded Driver Data" = Table.ExpandTableColumn(#"Merged Queries", "Driver Data", {"Driver Id"}, {"Driver Data.Driver Id"}),
    #"Filtered Rows" = Table.SelectRows(#"Expanded Driver Data", each ([Driver Data.Driver Id] <> null)),
    #"Renamed Columns" = Table.RenameColumns(#"Filtered Rows",{{"Driver Data.Driver Id", "Driver ID"}}, {"fastestLapTime", "fastestLapSpeed", "Status"}),
    #"Removed Columns" = Table.RemoveColumns(#"Renamed Columns",{"number", "positionText", "positionOrder", "laps", "time", "milliseconds", "fastestLap", "rank"}),
    #"Reordered Columns" = Table.ReorderColumns(#"Removed Columns", {"resultId", "raceId", "Driver ID", "driver ref", "constructor ref", "grid", "position", "points", "fastestLapTime", "fastestLapSpeed", "Status"}),
    #"Merged Queries1" = Table.NestedJoin(#"Reordered Columns", {"constructor ref"}, #Constructors Data, {"constructorRef"}, "Constructors Data", JoinKind.LeftOuter),
    #"Expanded Constructors Data" = Table.ExpandTableColumn(#"Merged Queries1", "Constructors Data", {"Constructor Id"}, {"Constructors Data.Constructor Id"}),
    #"Filtered Rows1" = Table.SelectRows(#"Expanded Constructors Data", each ([Constructors Data.Constructor Id] <> null)),
    #"Renamed Columns1" = Table.RenameColumns(#"Filtered Rows1", {"Constructors Data.Constructor Id", "Constructor ID"}), {"Renamed Columns1", {"constructor ref", "driver ref", "fastestLapTime", "fastestLapSpeed"}},
    #"Removed Columns1" = Table.RemoveColumns(#"Renamed Columns1", {"constructor ref", "driver ref", "fastestLapTime", "fastestLapSpeed"}),
    #"Replaced Value6" = Table.ReplaceValue(#"Removed Columns1",null,0,Replacer.ReplaceValue,{"position"})
in
#"Replaced Value6"
```

✓ No syntax errors have been detected.

Done Cancel

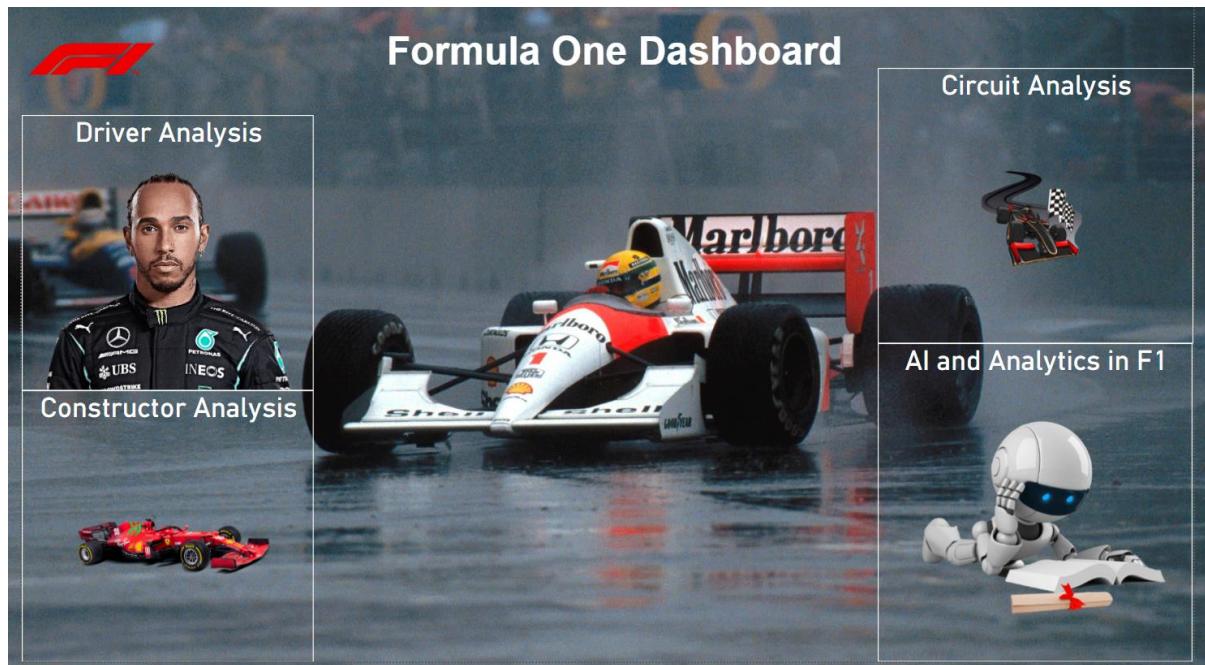
The highlighted region shows the pre-processing steps which are done in m-language this includes Replace Values, Fixing errors and Merged Queries etc.

In the same way we can see how m-language is used in other tables as well.

D. Dashboard

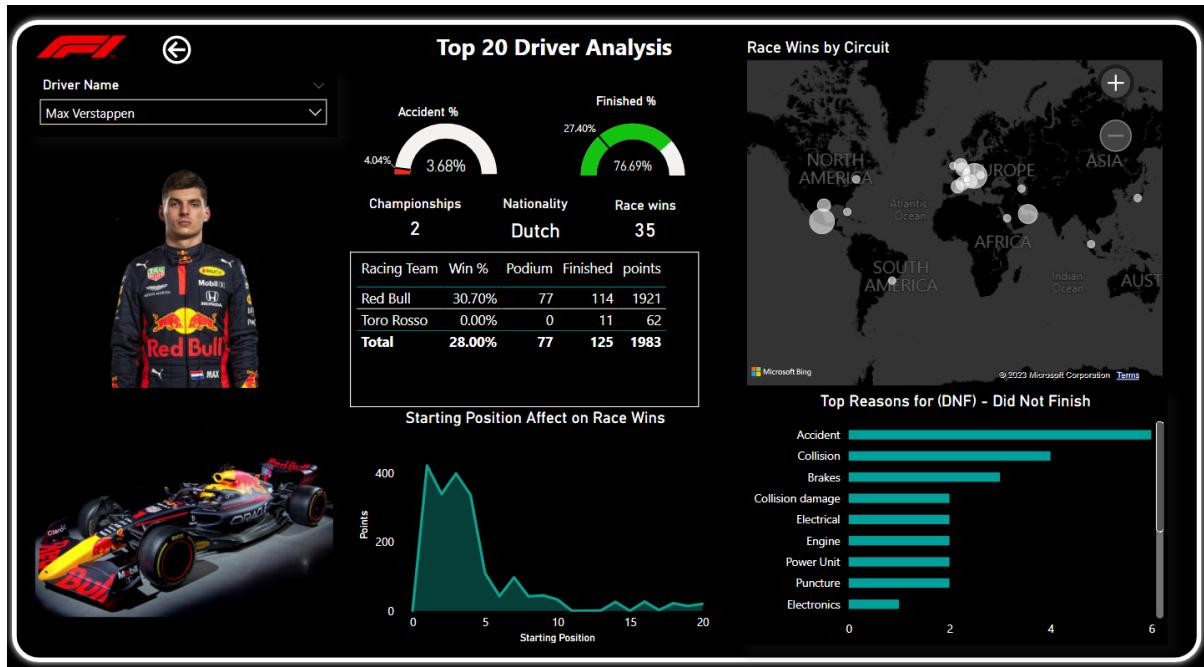
- Power BI dashboard (full collection of visuals) and how the content of the Power BI pages is organised.

Home Page :



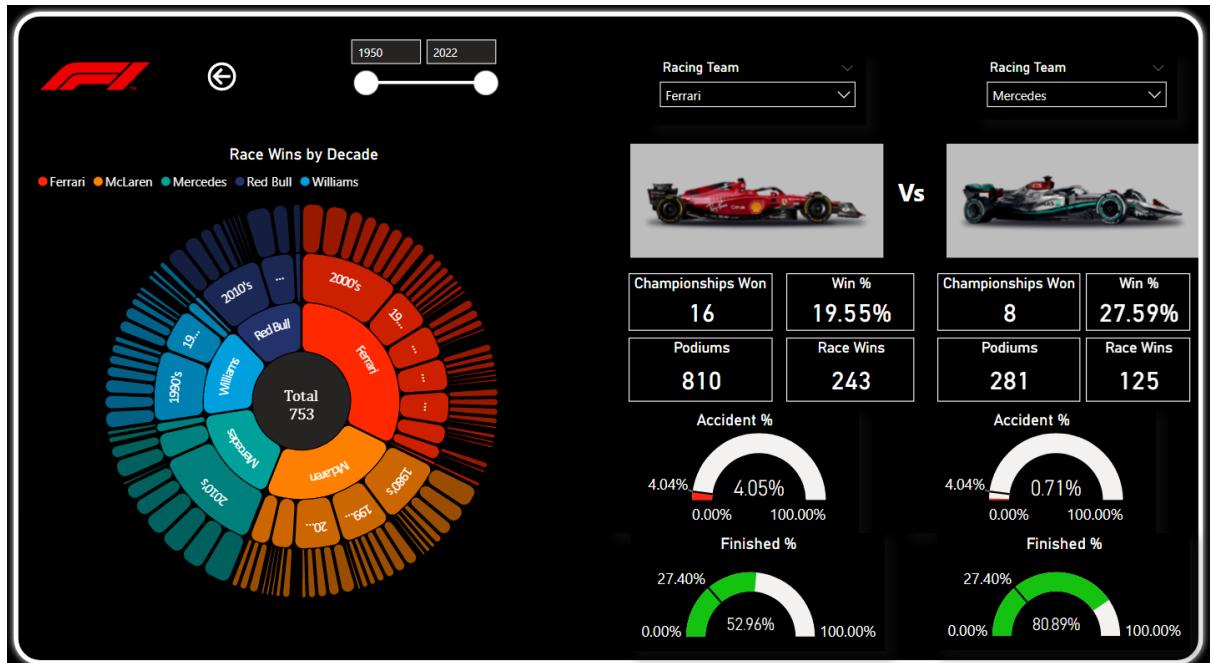
- This is the home page of the dashboard. All the 4 pages of the dashboard is linked to the home page.
- I have taken first the Driver Analysis to be my first page because the sport mainly revolves around the drivers.
- The second page is Constructor Analysis. The teams of each of these drivers is analysed here.
- The third page is circuit analysis of all the different race tracks in F1 and analysis based on track location is done here.
- AI and Analytics in F1 is the last page. It uses AI related charts to forecast and see what factors influences certain performance metrics like race wins.

Driver Analysis:



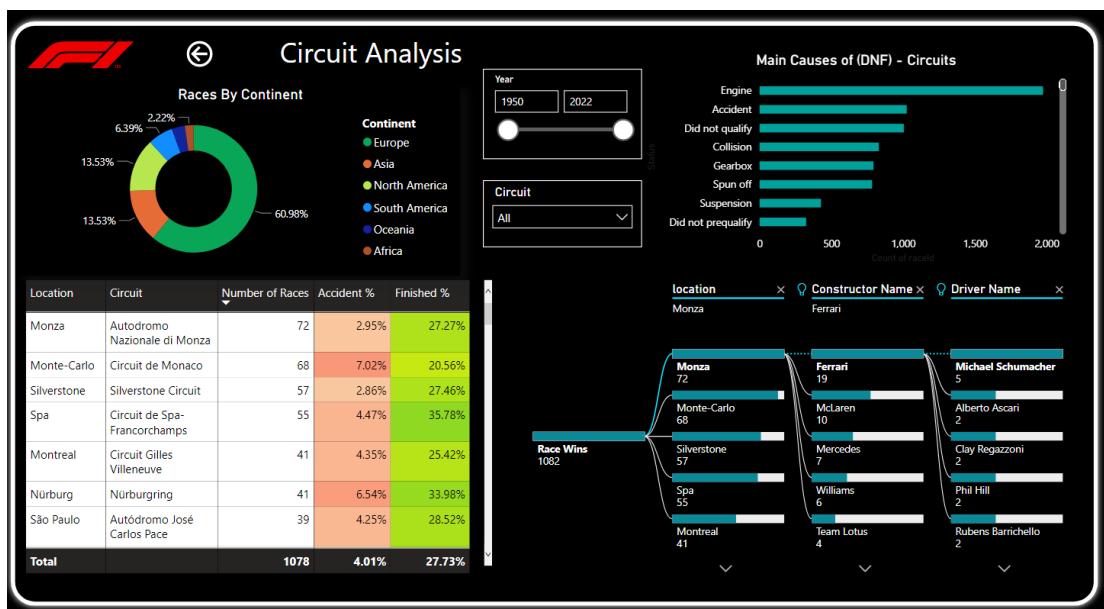
- This is the first page of my dashboard.
- It gives an in-depth analysis of the top 20 Drivers sorted by the race wins.
- The design is done with the consideration of normal human eye movement going from the left side of the screen to the right.
- Hence the driver image and their respective cars are to the left with the slicer.
- Then the gauge axis is kept in the middle which shows accident % and finished % along with the cards and table showing individual driver stats. An area chart is kept beneath it as well to show the influence of starting points on race wins.
- To the right we have the map chart which shows the circuits each driver has won a race.
- The graph to the bottom right is a horizontal bar chart showing top reasons for (DNF) – Did Not Finish.

Constructor Analysis:



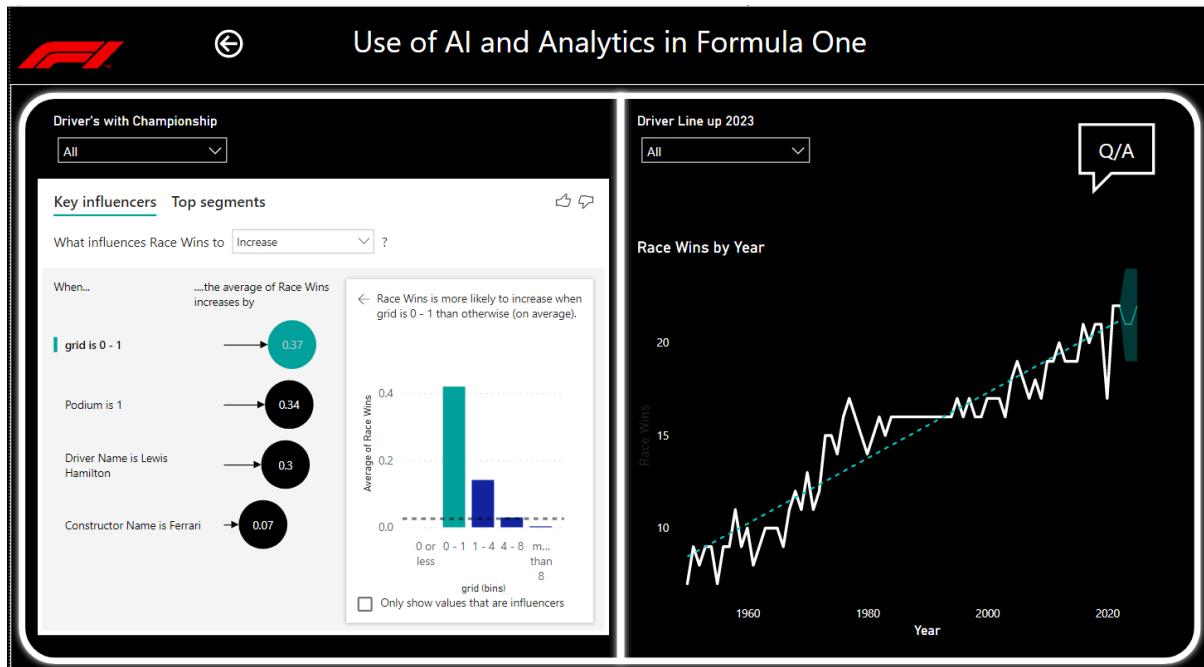
- This is the page analysing the Constructors or the team analysis.
- The animated sunburst is kept to the left as it is the most important chart in this page.
- To the right I have made a comparison between any two teams. (Only top 5 teams are selected).
- The cards and gauge axis move with respect to the input given.

Circuit analysis:



- The Circuit analysis is the 3rd page of the dashboard. The table and the decomposition tree chart are the most import visuals in this graph.
- The bar chart on the top right shows different causes of DNF are linked to the year and circuit.
- The pie chart shows the circuits by continent level.

Use of AI and Analytics in F1:



- Graphs in which AI plays a major role is kept in this page.
- It has been ordered with no specific organising technique in mind.

E. Self-Assessment

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	95
Data Pre-processing and Data Modelling	Many pre-processing steps have been applied. The data model is well-structured	100
Dax and M language	Both DAX and M Language have been extensively used in the report	100
Dashboard Design	The dashboard contains a variety of charts, including advanced ones not covered in the module.	95
Average		Add below the average of the four cells above: 97.5