DMWR Final Project: An Exploration of F1 Data

Avi Kodali 2022-05-05

Part 1: Ergast Data

Prix

stop time.

Circuit

pandemic. Therefore, I thought it would be fun to work with Formula 1 data, and see what I could find through analyzing it. The first part of my final project is using the historical F1 data from https://ergast.com/mrd/. This is an API that has access to F1 data from its inception in 1950 until now. This includes things like race results, pit stop times, driver data, season data. Pitstop Data

One of my biggest passions is Formula 1 racing. In fact, I have watched each and every race religiously since the start of the Covid-19

I decided to first look at pit stop data from this API, which has data specifically from 2011 until now. My first step was to attempt to use the API to get the pit stop data. However, this proved to be very difficult/inefficient because each query to the API for the pit stop data only allows access to the pit stop data of a specific year and round number of an F1 season. This would mean that I would have to have to combine the pit stop dataframes of roughly 200 races at some point, as well as combine that data with other data from Ergast such as constructor names, race locations, etc. Therefore, I decided to explore the site more to see if I had any other options, and I discovered that I

could download the data from the Ergast database from https://ergast.com/mrd/db/. This made obtaining the dataframe I wanted more obtainable. So using the CSV database tables from Ergast, I cleaned and combined the different tables to gain my final pit stop dataframe,

which part of can be seen below. raceName raceDate circuitName circuitAltitude dateTime year round driverName driverCode constructorName stop lap duration milliseconds seconds Toro Rosso Australian 2011-Albert Park 10 2011-03- 2011 1 Jaime ALG 1 1 26.898 26898 26.898 27 Grand 03-27 **Grand Prix** Alguersuari 17:05:23 Circuit Prix MSC Australian 2011-Albert Park 10 2011-03- 2011 1 Michael 1 1 25.021 25021 25.021 Mercedes 03-27 **Grand Prix** Grand Schumacher Prix Circuit 17:05:52 Australian 2011-Albert Park 10 2011-03- 2011 1 Mark WEB Red Bull 1 11 23.426 23426 23.426 Grand 03-27 **Grand Prix** 27 Webber 17:20:48 Prix Circuit Australian 2011-Albert Park 10 2011-03- 2011 1 ALO 1 12 23.251 23251 23.251 Fernando Ferrari Grand 03-27 **Grand Prix** 27 Alonso Prix Circuit 17:22:34 MAS Australian 2011-Albert Park 10 2011-03- 2011 1 Felipe Ferrari 1 13 23.842 23842 23.842 03-27 **Grand Prix** 27 Massa Grand 17:24:10 Prix Circuit Australian 2011-Albert Park 10 2011-03- 2011 1 Rubens BAR Williams 1 13 23.643 23643 23.643 03-27 27 Grand **Grand Prix** Barrichello

write.csv(pitstops, file = "pitstops_full.csv") Now with this cleaned data, I can do some analysis to see what I can find. First I want to see which constructor have the fastest average pit

Now I can save this clean pitstop dataframe into a csv file so that anyone can use it to start their own analysis.

17:24:29

constructorName avgPitStopTime numberOfRaces 24.23619 24.44421 85 Lotus Lotus F1 32.46375 285

HRT 32.67777 150 33.92462 241 Caterham 233 Marussia 34.39684 Toro Rosso 44.48698 683 Sauber 44.54653 621 50.84577 609 Force India Renault 63.26183 403 From this table, it seems that Virgin is the constructor with the fastest average pit stop time. However, a big part of this is probably because the Virgin only raced for two years (2010, 2011). I think I should also point out that there are two Lotus teams because they are technically two different teams (Lotus turned into Lotus F1 in the 2012 F1 season). Another interesting insight from this table is that most of these fastest average pit stop times are from teams that do not currently exist right now, i.e. teams that raced in the early 2010s. I think this might be because average pit stop times overall have increased from 2011 to now. However, I think it is important to note that the

summary statistics for this table include outliers that could have been caused by pit stop problems, retirements, etc. Therefore, I decided to filter out these outliers and see how the results change. constructorName avgPitStopTime numberOfRaces Mercedes 23.56088 845 Red Bull 23.57966 853 Ferrari 23.81331 839

McLaren 24.05513 828 593 Force India 24.05836 Virgin 24.23619 77 Lotus F1 24.23822 281 390 24.29716 Renault 24.44421 85 Lotus Toro Rosso 24.49106 671 This table is interesting because the top three constructors with the fastest average pit stop times (Mercedes, Red Bull, Ferrari) are the top three constructors in the past decade, i.e they are the teams that mainly fought for the F1 World Championship. Furthermore, a lot of the constructors in this table are still older teams that don't race right now. With the insights from these tables, I can try plotting the data. First, I want to see how average pit stop time has changed over time by constructor. Average Pit Stop Time from 2011 to 2022

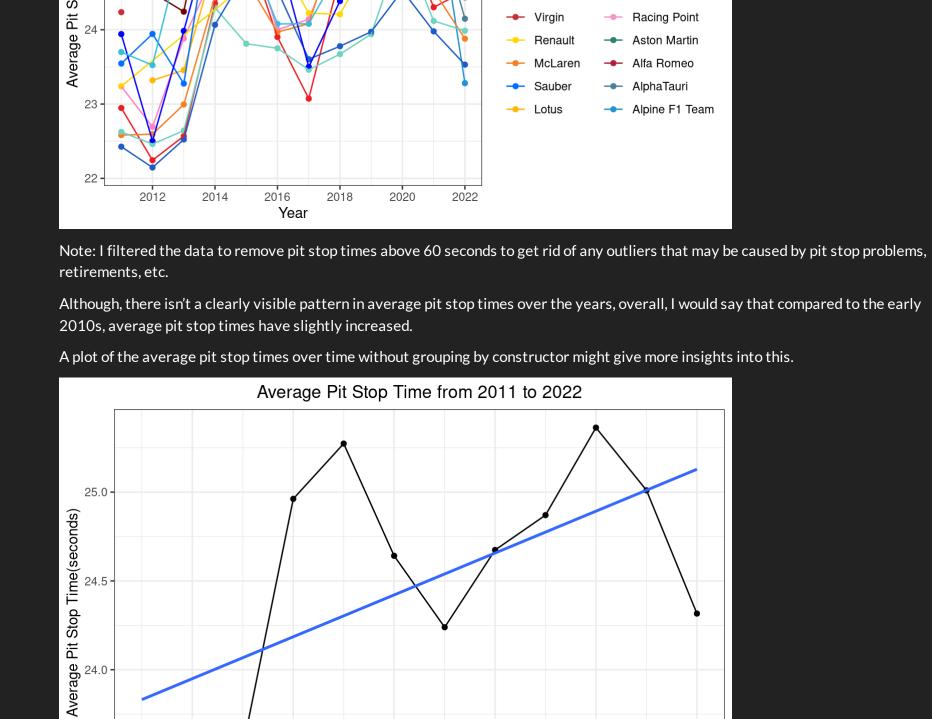
-- Lotus F1 - Marussia Manor Marussia --- Haas F1 Team

Constructor

Toro Rosso → HRT

Force India

Racing Point Aston Martin



2016

Year

Distribution of Pit Stop Times By Constructor

2014

Now I want to look at the distribution of the pit stop times by constructor.

23.5

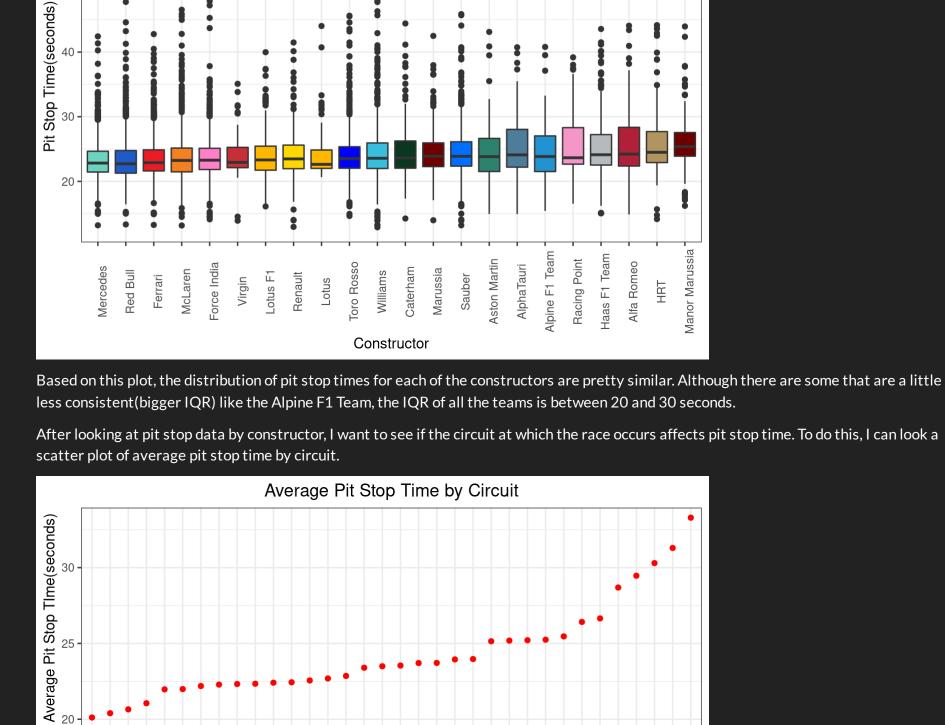
temperature, humidity, etc.).

Now I want to see how circuit altitude affects average pit stop time.

Based off of this plot, it is definitely clear that compared to the early 2010s, average F1 pit stop times have increased over time.

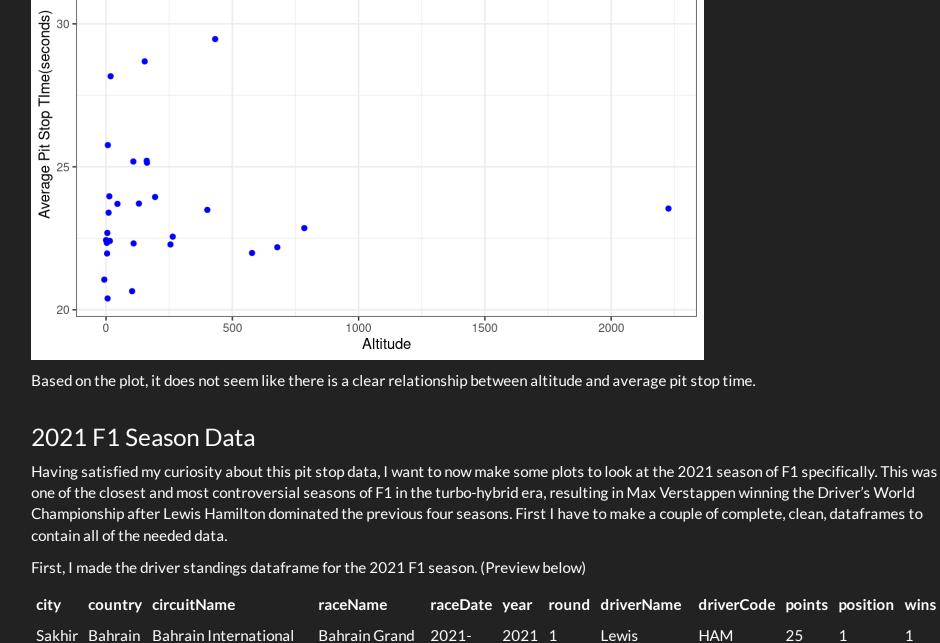
2020

2022



Circuit There definitely is some difference for pit stop times between different circuits, although I do have to point out that Circui Park Zandvoort has the fastest average pit stop time because it is a very new circuit that only had one race in 2021. This is a pretty interesting insight. These differences could be due to the length of the pit lane at each circuit or perhaps the overall environment at each circuit(like

Average Pit Stop Time vs Altitude



Circuit Prix 03-28 Hamilton Bahrain Grand 2021-Sakhir Bahrain Bahrain International 2021 1 Max VER Circuit Prix 03-28 Verstappen Sakhir Bahrain Bahrain International Bahrain Grand 2021-2021 1 Valtteri BOT

Prix

Bahrain Grand 2021-

Bahrain Grand 2021-

Circuit

Circuit

Sakhir Bahrain Bahrain International Circuit Bahrain Grand

it. It looks a lot like the graphs the official F1 Instagram account posted.

scale_color_manual(values = drivers_2021_colors) + scale_linetype_manual(values = drivers_2021_lines) +

theme(axis.text.x.bottom = element_text(angle = 90), plot.title = element_text(hjust = 0.5),

axis.title.x = element_blank(),

ggplot(aes(x = fct_reorder(raceName, raceDate), y = points,

Now I can get started on working on a couple of plots.

driver_standings_2021 %>%

geom_line() +

Sakhir Bahrain Bahrain International

Sakhir Bahrain Bahrain International

Circuit Prix 03-28 Sakhir Bahrain Bahrain International Bahrain Grand 2021-LEC 2021 1 Charles Circuit Prix 03-28 Leclerc Next, I made the constructor standings dataframe for the 2021 F1 season. (Preview below) city country circuitName raceName raceDate year round constructorName points position wins Sakhir Bahrain Bahrain International Circuit Bahrain Grand 2021 1 2021-Mercedes 41 Prix 03-28 Sakhir Bahrain Bahrain International Circuit Bahrain Grand 2021-2021 1 Red Bull 03-28 Sakhir Bahrain Bahrain International Circuit Bahrain Grand 2021-2021 1 McLaren 03-28 Sakhir Bahrain Bahrain International Circuit Bahrain Grand 2021-2021 1 0 Ferrari 03-28 Prix Sakhir Bahrain Bahrain International Circuit Bahrain Grand 2021-2021 1 AlphaTauri

03-28

2021-

03-28

I have included the code for the F1 World Driver's championship plot because it took me a lot of time to make, and personally I am proud of

color = driverCode, linetype = driverCode, group = driverCode)) +

2021 1

Aston Martin

03-28

03-28

2021 1

2021 1

Bottas

Lando Norris NOR

Sergio Pérez PER

25

18

0

0

legend.text = element_text(size = 8)) + labs(title = "2021 F1 Driver's World Championship", 2021 F1 Driver's World Championship 400 -Driver — HAM --- MAZ --- BOT VET 300 -- STR ···· KUB -- SAI — GAS --· oco --· RIC - MSC

This plot definitely show how the driver's championship was pretty close, as the top two drivers (Hamilton and Verstappen) exchanged places in terms of points total. And that was not the only close battle, as there are several driver who were close to each other throughout

> Constructor — Mercedes Red Bull

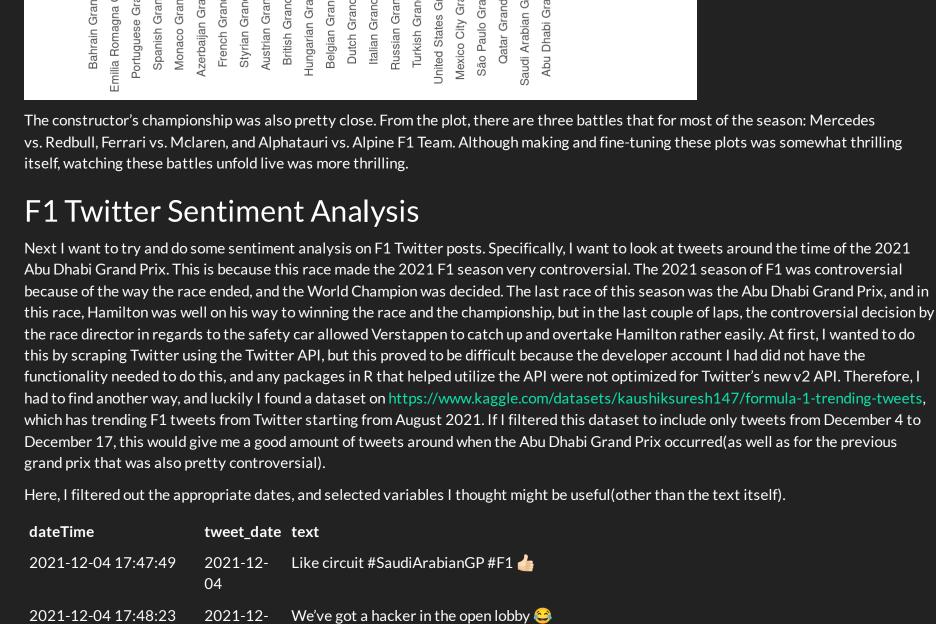
Williams

Haas F1 Team

— Aston Martin — Alfa Romeo — AlphaTauri — Alpine F1 Team

the season. Let's see if we can find the same thing from the constructor's championship.

2021 F1 Constructor's World Championship



04 start. Eager beavers. #F1 2021-12-2021-12-04 17:48:30 This could be interesting #Q3 #F1 04 2021-12-04 17:48:31 2021-12- #Q3 start #F1 https://t.co/SMZ98ba3Ln 04

2021-12-04 17:48:24 2021-12- Dangerous track. Many collisions avoided only for lucky #SaudiArabiaGP #F1 #SaudiArabianGP

2021-12- Leclerc and the two Mercedes are waiting at the end of the pitlane for the green light for Q3 to

26089

4271

2411

2040

3155

-2864

570

951

1093

1510

04

04

#F1 #SaudiArabianGP

https://t.co/x2qkapKBc2

2021-12-04 17:48:28

2021-12-05

2021-12-06

2021-12-07

2021-12-08

2021-12-09

#Quali

Now I have to make a data	rame that follows the tidy text form	nat.	
tweet_date		word	
2021-12-04		like	
2021-12-04		circuit	
2021-12-04		saudiarabiangp	
2021-12-04		f1	
2021-12-04		we've	
2021-12-04		got	
With this tidy data, I can tr	y and perform sentiment analysis.		
tweet_date	negative	positive	sentiment
2021-12-04	2912	3242	330

2037	4472	2435	
3826	9303	5477	
54224	79801	25577	
8758	12476	3718	
4121	5854	1733	
2953	4327	1374	
4719	6906	2187	
1921	2910	989	
, , , , , , , , , , , , , , , , , , , ,			
Sentiment of Trending	F1 Tweets		
Sentiment of Trending	F1 Tweets		
	54224 8758 4121 2953 4719 1921	3826 9303 54224 79801 8758 12476 4121 5854 2953 4327 4719 6906 1921 2910	3826930354775422479801255778758124763718412158541733295343271374471969062187

28953

3701

1460

947

1645

The plot is interesting. First, it looks like December 5 is the only outwardly negative day, which is interesting because it was the day of the Saudi Arabian Grand Prix. This negative score might be because that race was filled with a lot of crashes, and highlighting it all was contact between Verstappen and Hamilton, which could have been dirty driving. This may be why the day was negative. Another thing that is interesting is that December 12 is super positive. This might be because of the overall excitment for the last, championship deciding race. Furthermore, it does not seem like the few days after this controversial race was negative, like I thought it would be. A factor in this might be that I only used trending tweets, so those tweets might not be representative of the true sentiments on Twitter. Or maybe I made an error in the sentiment analysis process. Or maybe it simply really wasn't as controversial as I thought. One last thing I want to try is to make a wordcloud to visualize the words from the tweets.

Date

holy prettydangerous



used a lot, but in the context of that tweet, Verstappen winning might have been negative. Also, December 12 might be so positive because there was an influx of tweets proclaiming Verstappen as race winner and World Champion, so that might be why it was so positive. Furthermore, although there are a lot of positive words, there are also plenty of negative words, especially curse words. Since the occurrence of those positive words is higher, that might be why we do not have a lot of days that are clearly negative. With this, my final project covering the exploration of F1 data is done. I found a lot of insights through my analysis of various sets of F1 data, and it was pretty enjoyable. I hope you found what I did interesting as well.