

Introduction

This project is an analysis of my personal Spotify streaming history from 6/4/2020 - 11/3/2021. In this project, I examine my streaming data via song, artist, and date/time dimensions to uncover insights and patterns about my listening behaviors. Below are some of the questions that I will address:

- Who are my most listened to artists?
- What are my top songs?
- How much content do I listen to per day?
- How long are the songs that I listen to?
- How does my listening behavior change depending on the time of day?
- How does my listening behavior change depending on the day of week? Is this the same for songs vs. podcasts?
- During what months do I listen to content the most?
- Which do I listen to more, songs or podcasts?

In addition to answering these questions, I will also discuss potential factors that could explain why I got the results that I did.

Below is a data dictionary that defines the data found in each column:

Column Name	Definition
endTime	Date and time of when the stream ended in UTC format (Coordinated Universal Time zone).
artistName	Name of "creator" for each stream (e.g. the artist name if a music track).
trackName	Name of items listened to or watched (e.g. title of music track or name of video).
UniqueID	Concatenation of artistName and trackName to create a unique ID for each song/artist combination.
msPlayed	How many milliseconds the track was listened to.
minutesPlayed	How many minutes the track was listened to.
Language	The language of the track.
Date	The year, date, and month in which the track was listened to (YYYY-MM-DD).
Year	The year in which the track was listened to (YYYY).
Month	The month in which the track was listened to (MM).
Week	The week of the year in which the track was listened to (1-53).
Day of Month	The day of the month in which the track was listened to (1-31).
Day of Week (Number)	The day of the week in which the track was listened to (0-6).
Day of Week (Name)	The day of the week in which the track was listened to (Monday-Sunday).
Hour of Day	The hour of the day in which the track was listened to (0-23).
Year + Week	The year and week in which the track was listened to (YYYY-WW).
Year + Month	The year and month in which the track was listened to (YYYY-MM).
Hour of Day - Adjusted	The hour of the day in which the track was listened to (0-23), adjusted to align with a given timezone.
Day of Week (Number) - Adjusted	The day of the week in which the track was listened to (0-6), adjusted to align with a given timezone.
Day of Week (Name) - Adjusted	The day of the week in which the track was listened to (Monday-Sunday), adjusted to align with a given timezone.

Transforming the Data

Importing the necessary packages and transforming the CSV file into usable data:

1

2020-06-05 13:32:00

Pardon My Take

Dana White, Booger McFarland And Sour Grapes D...

46634

0.78

...

7

Friday

2020-23

2020-06

2

2020-06-06 02:46:00

Miles Jaye

Let's Start Over

Miles Jaye: Let's Start Over

1677

0.03

...

6

Saturday

2020-23

2020-06

rows x 17 columns

#fixing an issue with the data where having "s" in the artist's or track's name triggered the regex

import warnings

warnings.filterwarnings('ignore')

df['artistName'] = df['artistName'].str.replace('\s','')\n

df['trackName'] = df['trackName'].str.replace('\s','')\n

df['UniqueID'] = df['artistName'] + " " + df['trackName']

Artist Analysis

low that our data has been cleaned and organized, we will analyze the data surrounding the artists to whom I listen.

Who are my top artists?

artist_minutes = pd.DataFrame(df.groupby('artistName')['minutesPlayed'].sum()\\\n .sort_values(ascending = False).head(10))\nround(artist_minutes,2)

minutesPlayed

artistName	minutesPlayed
Pardon My Take	10,768.23
Mac Miller	2,802.1
Tyler, The Creator	1,617.97
Olivia Rodrigo	996.14
Billie Eilish	900.27
Louis The Child	891.23
FKJ	884.77
BROCKHAMPTON	798.15
Kanye West	748.62
TroyBoi	706.58

plt.rcParams["figure.figsize"]=10,6\nsns.barplot(data = artist_minutes.index, x = artist_minutes.minutesPlayed,\\\n color = 'lime', estimator = np.sum)\nsns.set(rc={'axes.facecolor':'black', 'figure.facecolor':'white'})\nplt.ylabel('Artist Name', weight='bold')\nplt.xlabel('Minutes Played', weight='bold')\nplt.rcParams["axes.labelsize"] = 15\nplt.title('Top 10 Artists by Minutes Played', weight='bold').set_fontsize('20')\nplt.show()

Top 10 Artists by Minutes Played

Artist Name

Pardon My Take

Mac Miller

Tyler, The Creator

Olivia Rodrigo

Billie Eilish

Louis The Child

BROCKHAMPTON

Kanye West

TroyBoi

0

2500

5000

7500

10000

12500

15000

17500

20000

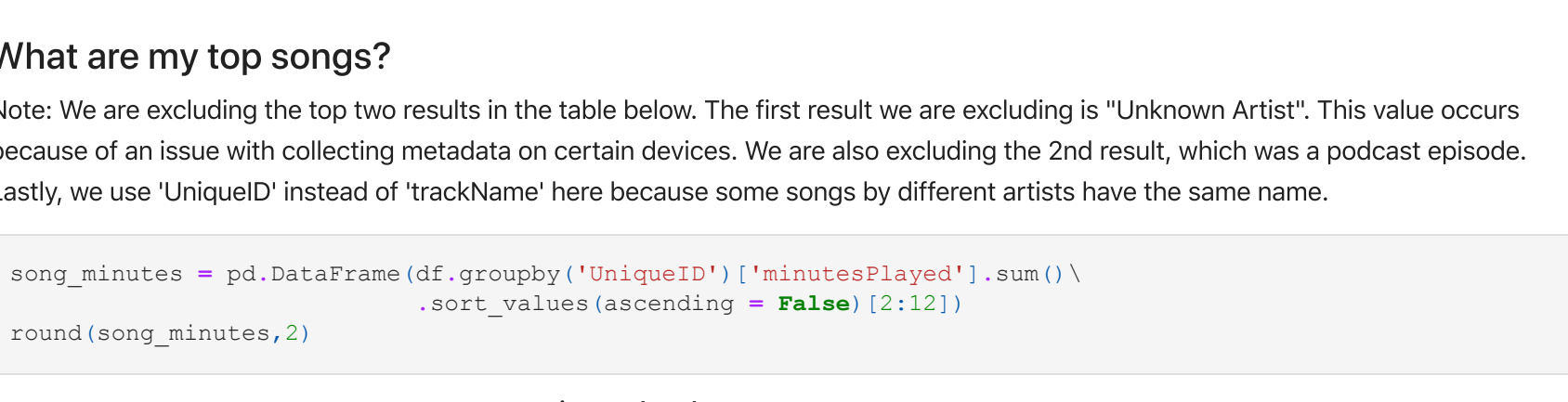
Minutes Played

as we can see in the above table and bar chart, my top artist by far is the podcast "Pardon My Take," followed by artists Mac Miller and Tyler, The Creator. Seeing Pardon My Take at the top is expected because the average length of a podcast is much longer than the average length of a song and I listen to the podcast for hours every week.

Now that our data has been cleaned and organized, we will analyze the data surrounding the artists to whom I listen.

Who are my top artists?

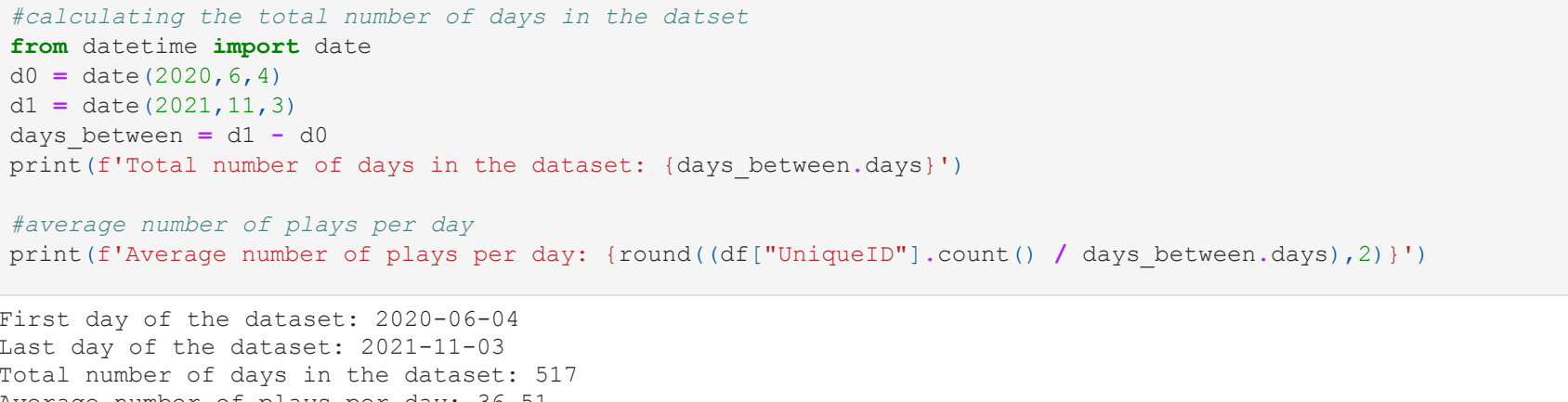
In (6):	<pre>artist_minutes = pd.DataFrame(df.groupby('artistName')['minutesPlayed'].sum()) round(artist_minutes,2)</pre>																						
Out (6):	<table><tr><th>artistName</th><th>minutesPlayed</th></tr><tr><td>Pardon My Take</td><td>19768.23</td></tr><tr><td>Mac Miller</td><td>2,802.1</td></tr><tr><td>Tyler, The Creator</td><td>1,873.97</td></tr><tr><td>Olivia Rodrigo</td><td>996.14</td></tr><tr><td>Billie Eilish</td><td>900.27</td></tr><tr><td>Louis The Child</td><td>891.23</td></tr><tr><td>FKJ</td><td>884.77</td></tr><tr><td>BROCKHAMPTON</td><td>798.15</td></tr><tr><td>Kanye West</td><td>748.62</td></tr><tr><td>TroyBoi</td><td>706.58</td></tr></table>	artistName	minutesPlayed	Pardon My Take	19768.23	Mac Miller	2,802.1	Tyler, The Creator	1,873.97	Olivia Rodrigo	996.14	Billie Eilish	900.27	Louis The Child	891.23	FKJ	884.77	BROCKHAMPTON	798.15	Kanye West	748.62	TroyBoi	706.58
artistName	minutesPlayed																						
Pardon My Take	19768.23																						
Mac Miller	2,802.1																						
Tyler, The Creator	1,873.97																						
Olivia Rodrigo	996.14																						
Billie Eilish	900.27																						
Louis The Child	891.23																						
FKJ	884.77																						
BROCKHAMPTON	798.15																						
Kanye West	748.62																						
TroyBoi	706.58																						
In (44):	<pre>plt.rcParams["figure.figsize"]=10,6 sns.barplot(data = song_minutes, y = artist_minutes.index, x = artist_minutes.minutesPlayed, \ color = 'lime', estimator = np.sum) sns.set(rc={'axes.facecolor':'black', 'figure.facecolor':'white'}) plt.ylabel('Artist Name', weight='bold') plt.xlabel('Minutes Played', weight='bold') plt.rcParams["axes.labelsize"] = 15 plt.title('Top 10 Artists by Minutes Played', weight='bold').set_fontsize('20') plt.show()</pre>																						



As we can see in the above table and bar chart, my top artist by far is the podcast "Pardon My Take," followed by artists Mac Miller and Tyler, The Creator. Seeing Pardon My Take at the top is expected because the average length of a podcast is much longer than the average length of a song and I listen to the podcast for hours every week.

Who are my top artists (excluding podcasts)?

In (8):	<pre>#this is done to exclude my top result, which is a podcast artist_minutes_songs = pd.DataFrame(df.groupby('artistName')['minutesPlayed'].sum()) round(artist_minutes_songs,2)</pre>																						
Out (8):	<table><tr><th>artistName</th><th>minutesPlayed</th></tr><tr><td>Mac Miller</td><td>2,802.1</td></tr><tr><td>Tyler, The Creator</td><td>1,873.97</td></tr><tr><td>Olivia Rodrigo</td><td>996.14</td></tr><tr><td>Billie Eilish</td><td>900.27</td></tr><tr><td>Louis The Child</td><td>891.23</td></tr><tr><td>FKJ</td><td>884.77</td></tr><tr><td>BROCKHAMPTON</td><td>798.15</td></tr><tr><td>Kanye West</td><td>748.62</td></tr><tr><td>TroyBoi</td><td>706.58</td></tr><tr><td>Tom Misch</td><td>667.52</td></tr></table>	artistName	minutesPlayed	Mac Miller	2,802.1	Tyler, The Creator	1,873.97	Olivia Rodrigo	996.14	Billie Eilish	900.27	Louis The Child	891.23	FKJ	884.77	BROCKHAMPTON	798.15	Kanye West	748.62	TroyBoi	706.58	Tom Misch	667.52
artistName	minutesPlayed																						
Mac Miller	2,802.1																						
Tyler, The Creator	1,873.97																						
Olivia Rodrigo	996.14																						
Billie Eilish	900.27																						
Louis The Child	891.23																						
FKJ	884.77																						
BROCKHAMPTON	798.15																						
Kanye West	748.62																						
TroyBoi	706.58																						
Tom Misch	667.52																						
In (9):	<pre>plt.rcParams["figure.figsize"]=10,6 sns.barplot(data = artist_minutes_songs, y = artist_minutes_songs.index, \ color = 'lime', estimator = np.sum) sns.set(rc={'axes.facecolor':'black', 'figure.facecolor':'white'}) plt.ylabel('Artist Name', weight='bold') plt.xlabel('Minutes Played', weight='bold') plt.rcParams["axes.labelsize"] = 15 plt.title('Top 10 Artists by Minutes Played (Excluding Podcasts)', \ weight='bold').set_fontsize('20') plt.show()</pre>																						



We now have a clearer picture of my top artists for songs (ranked by minutes played) after excluding the Pardon My Take podcast from the data. Mac Miller and Tyler, The Creator have the most minutes listened to for over 4.5 hours during the timeframe of this dataset!

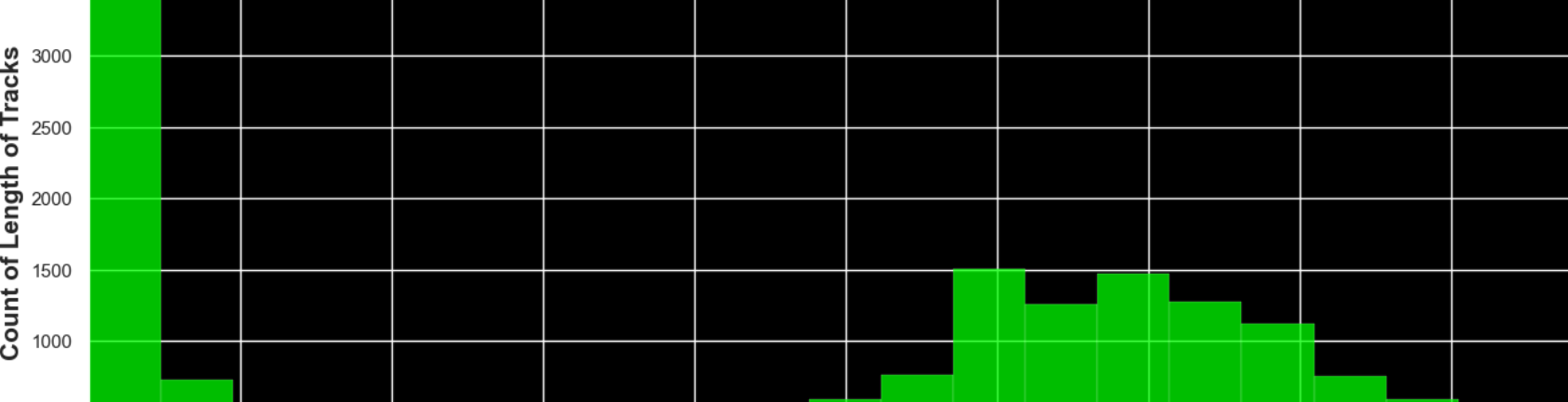
While "Dang" is one of my favorite songs, I am also designing a video game level using this song, which would account for some of its plays.

Time Analysis

Now that we have analyzed my top artists and songs, we will analyze the data surrounding the time and dates during which I listen.

How many total minutes of music and podcasts did I listen to?

In (12):	<pre>#Total minutes played (milliseconds divided by 60,000 to convert to minutes) round(df['msPlayed'].sum())/60000)</pre>
Out (12):	63552
In (11):	<pre>plt.rcParams["figure.figsize"]=10,6 sns.barplot(data = song_minutes, y = song_minutes.index, x = song_minutes.minutesPlayed, color = 'lime', \ estimator = np.sum) sns.set(rc={'axes.facecolor':'black', 'figure.facecolor':'white'}) plt.ylabel('Song Name', weight='bold') plt.xlabel('Minutes Played', weight='bold') plt.rcParams["axes.labelsize"] = 15 plt.title('Top 10 Songs by Minutes Played', weight='bold').set_fontsize('20') plt.show()</pre>



As we can see in the above table and bar chart, my top song is Mac Miller's "Dang!", followed by Olivia Rodrigo's "drivers license" and Ariana Grande's "Stuck with U." I listened to the song "Dang!" for over 4.5 hours during the timeframe of this dataset!

While "Dang" is one of my favorite songs, I am also designing a video game level using this song, which would account for some of its plays.

How many minutes of music and podcasts did I listen to per day, on average?

In (13):	<pre>#Finding the first and last dates, respectively, in the dataset print(f'First day of the dataset: {df["Date"][0].to_string(index=False)}') #2020-06-04 print(f'Last day of the dataset: {df["Date"][-1].to_string(index=False)}') #2021-11-03 #Calculating the total number of days in the dataset from datetime import date d0 = date(2020,6,4) d1 = date(2021,11,3) days_between = d1 - d0 print(f'Total number of days in the dataset: {days_between.days}') #Average number of plays per day print(f'Average number of plays per day: {round((df["UniqueID"].count() / days_between.days),2)}')</pre>
Out (13):	First day of the dataset: 2020-06-04 Last day of the dataset: 2021-11-03 Total number of days in the dataset: 517 Average number of plays per day: 36.51

As we can see from the above output, I listen to about 36.5 tracks per day.

How many minutes of music and podcasts did I listen to per day, on average?

In (14):	<pre># Average minutes listened to per day print(f'Average minutes listened to per day: {round(df['msPlayed'].sum()/(60000/517))}') Average minutes listened to per day: 123</pre>
Out (14):	123

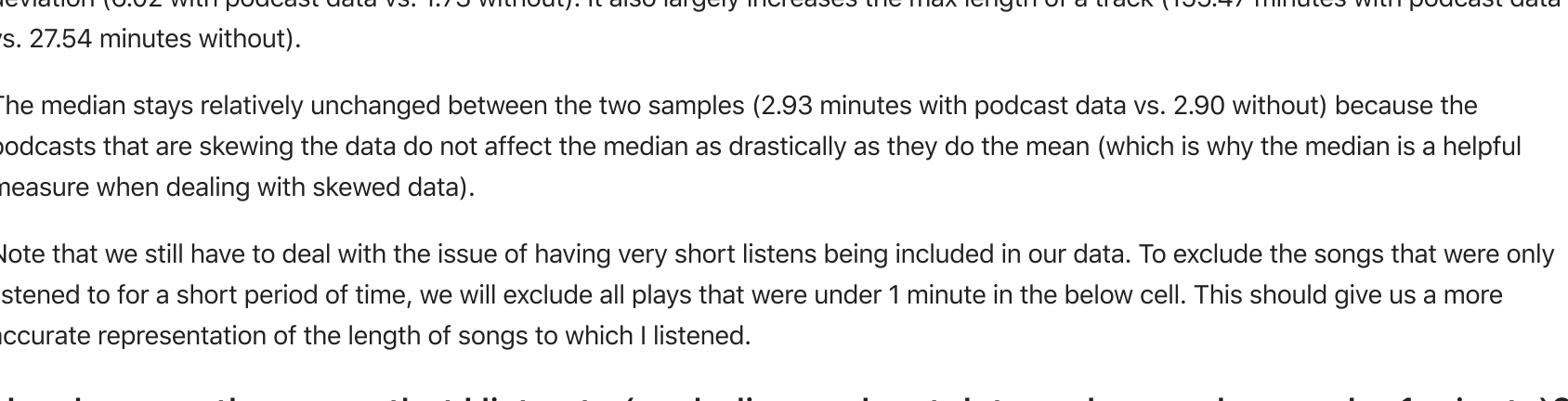
As we can see from the above output, I listen to about 123 minutes (just over 2 hours) per day.

How long is the average piece of media that I listen to?

In (15):	<pre># All data including podcasts pd.DataFrame(round(df['minutesPlayed'].describe(),2))</pre>																		
Out (15):	<table><tr><th></th><th>minutesPlayed</th></tr><tr><td>count</td><td>18,878.0</td></tr><tr><td>mean</td><td>3.37</td></tr><tr><td>std</td><td>6.02</td></tr><tr><td>min</td><td>0.0</td></tr><tr><td>25%</td><td>0.59</td></tr><tr><td>50%</td><td>2.93</td></tr><tr><td>75%</td><td>3.76</td></tr><tr><td>max</td><td>165.47</td></tr></table>		minutesPlayed	count	18,878.0	mean	3.37	std	6.02	min	0.0	25%	0.59	50%	2.93	75%	3.76	max	165.47
	minutesPlayed																		
count	18,878.0																		
mean	3.37																		
std	6.02																		
min	0.0																		
25%	0.59																		
50%	2.93																		
75%	3.76																		
max	165.47																		

In (16): df['minutesPlayed'].median()

Out (16): 2.93

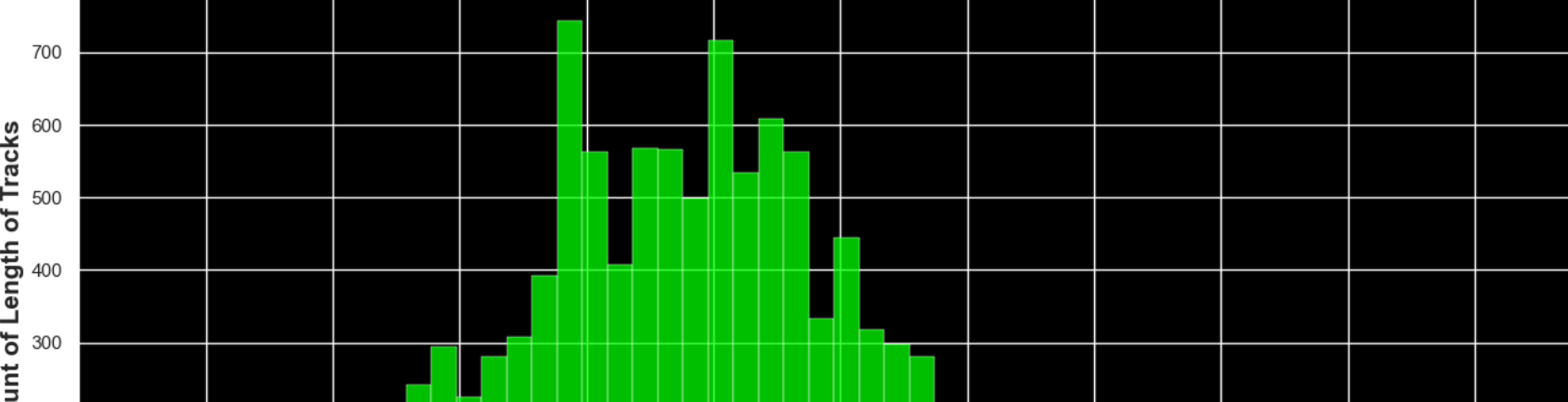


As we can see from table and graph above, the mean track length (songs and podcasts) that I listen to is 3.37 minutes, and the median is 2.93 minutes.

There are two important callouts to be made here. The first is that this data includes podcasts, which have longer averages, which vary between songs. Second, there are a large number of listens that are very close to 0 minutes, which may be due to Spotify keeping records of very brief plays (for example, if I listen to a song for 5 seconds and decide to skip it). Because of these callouts, the above data is not indicative of the true length of songs to which I listen, so we need to change the data to more accurately reflect the data surrounding the songs to which I listen.

How long are the songs that I listen to (excluding podcast data)?

In (18):	<pre># Only song data (excluding podcasts) songs_only = pd.DataFrame(df.loc[df['artistName'] != 'Pardon My Take','minutesPlayed']) round(songs_only.describe(),2)</pre>																		
Out (18):	<table><tr><th></th><th>minutesPlayed</th></tr><tr><td>count</td><td>17,852.0</td></tr><tr><td>mean</td><td>2.45</td></tr><tr><td>std</td><td>1.73</td></tr><tr><td>min</td><td>0.0</td></tr><tr><td>25%</td><td>0.55</td></tr><tr><td>50%</td><td>2.9</td></tr><tr><td>75%</td><td>3.68</td></tr><tr><td>max</td><td>27.58</td></tr></table>		minutesPlayed	count	17,852.0	mean	2.45	std	1.73	min	0.0	25%	0.55	50%	2.9	75%	3.68	max	27.58
	minutesPlayed																		
count	17,852.0																		
mean	2.45																		
std	1.73																		
min	0.0																		
25%	0.55																		
50%	2.9																		
75%	3.68																		
max	27.58																		
In (19):	songs_only.median()																		
Out (19):	minutesPlayed 2.9 dtype: float64																		



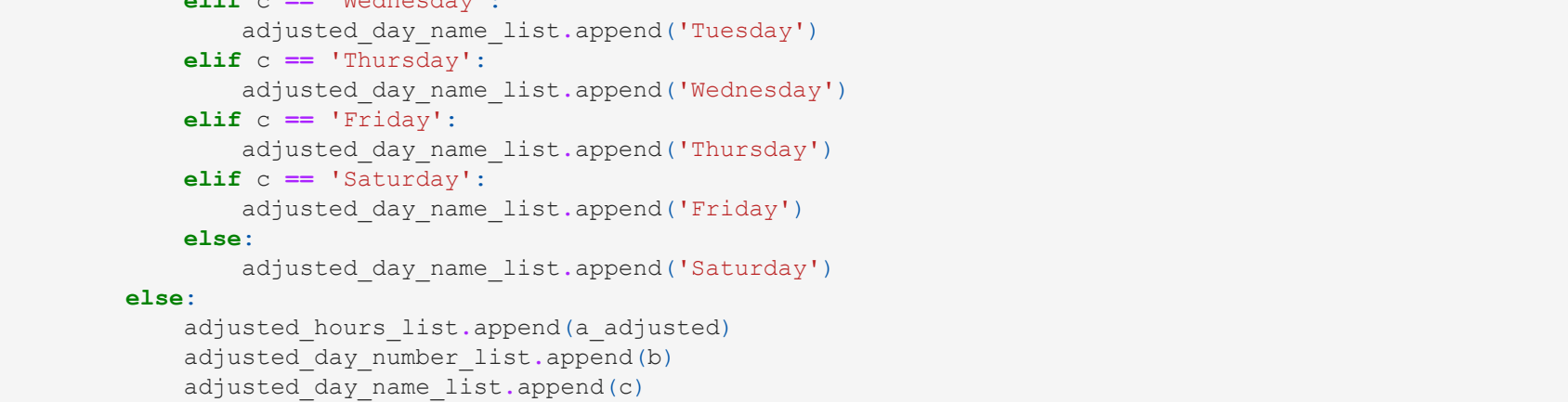
As we can see from the above table and graph, the mean track length (songs only) that I listen to is 2.45 minutes and the median is 2.90 minutes. As expected, this shows that the podcasts skew the mean more heavily towards the right and vastly increase the standard deviation (6.02 with podcast data vs. 1.73 without). It also largely increases the max length of a track (155.47 minutes with podcast data vs. 27.54 minutes without).

The median stays relatively unchanged between the two samples (2.93 minutes with podcast data vs. 2.90 without) because the podcasts that are skewing the data do not affect the median as drastically as they do the mean (which is why the median is a helpful measure when dealing with skewed data).

Noted that we still have to deal with the issue of having very short listens being under 1 minute. To exclude the songs that were only listened for a short period of time, we will exclude all plays that were under 1 minute in the below cell. This should give us a more accurate representation of the length of songs to which I listened.

How long are the songs that I listen to (excluding podcast data and song plays under 1 minute)?

In (21):	<pre>songs_over_1_min = pd.DataFrame(df.loc[(df['artistName'] != 'Pardon My Take') & (df['minutesPlayed'] >= 1),\ round(songs_over_1_min.describe(),2)])</pre>																		
Out (21):	<table><tr><th></th><th>minutesPlayed</th></tr><tr><td>count</td><td>12,705.0</td></tr><tr><td>mean</td><td>3.36</td></tr><tr><td>std</td><td>1.15</td></tr><tr><td>min</td><td>1.0</td></tr><tr><td>25%</td><td>2.76</td></tr><tr><td>50%</td><td>3.32</td></tr><tr><td>75%</td><td>3.92</td></tr><tr><td>max</td><td>27.58</td></tr></table>		minutesPlayed	count	12,705.0	mean	3.36	std	1.15	min	1.0	25%	2.76	50%	3.32	75%	3.92	max	27.58
	minutesPlayed																		
count	12,705.0																		
mean	3.36																		
std	1.15																		
min	1.0																		
25%	2.76																		
50%	3.32																		
75%	3.92																		
max	27.58																		
In (22):	songs_over_1_min.median()																		
Out (22):	minutesPlayed 3.36 dtype: float64																		



As we can see from the above table and graph, the mean track length (songs only) that I listen to is 2.45 minutes and the median is 2.90 minutes. As expected, this shows that the podcasts skew the mean more heavily towards the right and vastly increase the standard deviation (6.02 with podcast data vs. 1.73 without). It also largely increases the max length of a track (155.47 minutes with podcast data vs. 27.54 minutes without).

The median stays relatively unchanged between the two samples (2.93 minutes with podcast data vs. 2.90 without) because the podcasts that are skewing the data do not affect the median as drastically as they do the mean (which is why the median is a helpful measure when dealing with skewed data).

Noted that we still have to deal with the issue of having very short listens being under 1 minute. To exclude the songs that were only listened for a short period of time, we will exclude all plays that were under 1 minute in the below cell. This should give us a more accurate representation of the length of songs to which I listened.

How long are the songs that I listen to (excluding podcast data and song plays under 1 minute)?

In (21):	<pre>songs_over_1_min = pd.DataFrame(df.loc[(df['artistName'] != 'Pardon My Take') & (df['minutesPlayed'] >= 1),\ round(songs_over_1_min.describe(),2)])</pre>																		
Out (21):	<table><tr><th></th><th>minutesPlayed</th></tr><tr><td>count</td><td>12,705.0</td></tr><tr><td>mean</td><td>3.36</td></tr><tr><td>std</td><td>1.15</td></tr><tr><td>min</td><td>1.0</td></tr><tr><td>25%</td><td>2.76</td></tr><tr><td>50%</td><td>3.32</td></tr><tr><td>75%</td><td>3.92</td></tr><tr><td>max</td><td>27.58</td></tr></table>		minutesPlayed	count	12,705.0	mean	3.36	std	1.15	min	1.0	25%	2.76	50%	3.32	75%	3.92	max	27.58
	minutesPlayed																		
count	12,705.0																		
mean	3.36																		
std	1.15																		
min	1.0																		
25%	2.76																		
50%	3.32																		
75%	3.92																		
max	27.58																		
In (22):	songs_over_1_min.median()																		
Out (22):	minutesPlayed 3.36 dtype: float64																		



As we can see from the above table and graph, the mean track length (songs only) that I listen to is 2.45 minutes and the median is 2.90 minutes. As expected, this shows that the podcasts skew the mean more heavily towards the right and vastly increase the standard deviation (6.02 with podcast data vs. 1.73 without). It also largely increases the max length of a track (155.47 minutes with podcast data vs. 27.54 minutes without).

The median stays relatively unchanged between the two samples (2.93 minutes with podcast data vs. 2.90 without) because the podcasts that are skewing the data do not affect the median as drastically as they do the mean (which is why the median is a helpful measure when dealing with skewed data).

Noted that we still have to deal with the issue of having very short listens being under 1 minute. To exclude the songs that were only listened for a short period of time, we will exclude all plays that were under 1 minute in the below cell. This should give us a more accurate representation of the length of songs to which I listened.

How long are the songs that I listen to (excluding podcast data and song plays under 1 minute)?

In (21):	<pre>songs_over_1_min = pd.DataFrame(df.loc[(df['artistName'] != 'Pardon My Take') & (df['minutesPlayed'] >= 1),\ round(songs_over_1_min.describe(),2)])</pre>																		
Out (21):	<table><tr><th></th><th>minutesPlayed</th></tr><tr><td>count</td><td>12,705.0</td></tr><tr><td>mean</td><td>3.36</td></tr><tr><td>std</td><td>1.15</td></tr><tr><td>min</td><td>1.0</td></tr><tr><td>25%</td><td>2.76</td></tr><tr><td>50%</td><td>3.32</td></tr><tr><td>75%</td><td>3.92</td></tr><tr><td>max</td><td>27.58</td></tr></table>		minutesPlayed	count	12,705.0	mean	3.36	std	1.15	min	1.0	25%	2.76	50%	3.32	75%	3.92	max	27.58
	minutesPlayed																		
count	12,705.0																		
mean	3.36																		
std	1.15																		
min	1.0																		
25%	2.76																		
50%	3.32																		
75%	3.92																		
max	27.58																		
In (22):	songs_over_1_min.median()																		
Out (22):	minutesPlayed 3.36 dtype: float64																		



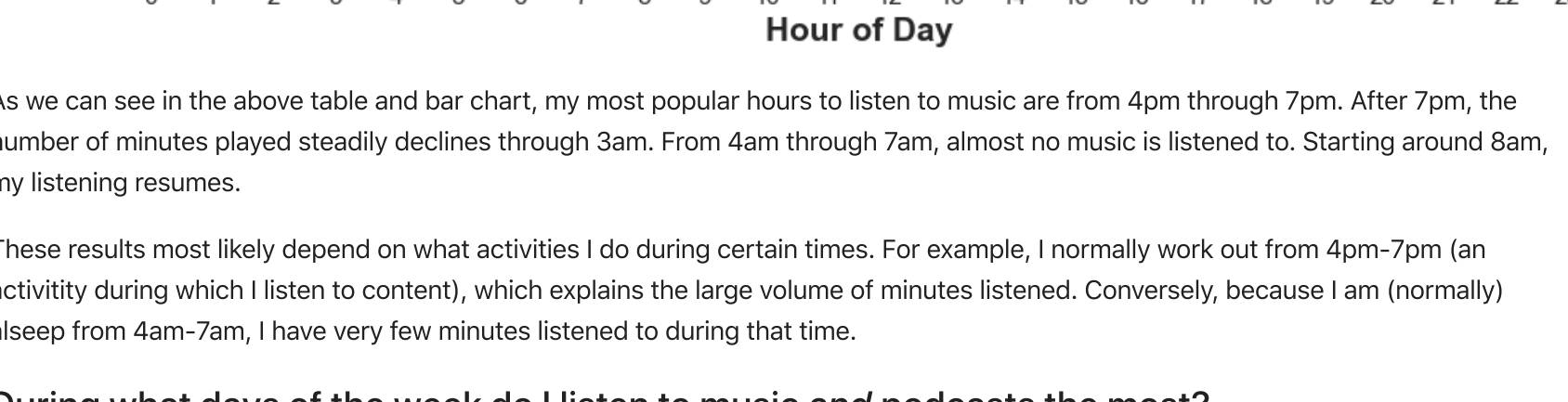
As we can see from the above table and graph, the mean track length (songs only) that I listen to is 2.45 minutes and the median is 2.90 minutes. As expected, this shows that the podcasts skew the mean more heavily towards the right and vastly increase the standard deviation (6.02 with podcast data vs. 1.73 without). It also largely increases the max length of a track (155.47 minutes with podcast data vs. 27.54 minutes without).

The median stays relatively unchanged between the two samples (2.93 minutes with podcast data vs. 2.90 without) because the podcasts that are skewing the data do not affect the median as drastically as they do the mean (which is why the median is a helpful measure when dealing with skewed data).

Noted that we still have to deal with the issue of having very short listens being under 1 minute. To exclude the songs that were only listened for a short period of time, we will exclude all plays that were under 1 minute in the below cell. This should give us a more accurate representation of the length of songs to which I listened.

How long are the songs that I listen to (excluding podcast data and song plays under 1 minute)?

In (21):	<pre>songs_over_1_min = pd.DataFrame(df.loc[(df['artistName'] != 'Pardon My Take') & (df['minutesPlayed'] >= 1),\ round(songs_over_1_min.describe(),2)])</pre>																		
Out (21):	<table><tr><th></th><th>minutesPlayed</th></tr><tr><td>count</td><td>12,705.0</td></tr><tr><td>mean</td><td>3.36</td></tr><tr><td>std</td><td>1.15</td></tr><tr><td>min</td><td>1.0</td></tr><tr><td>25%</td><td>2.76</td></tr><tr><td>50%</td><td>3.32</td></tr><tr><td>75%</td><td>3.92</td></tr><tr><td>max</td><td>27.58</td></tr></table>		minutesPlayed	count	12,705.0	mean	3.36	std	1.15	min	1.0	25%	2.76	50%	3.32	75%	3.92	max	27.58
	minutesPlayed																		
count	12,705.0																		
mean	3.36																		
std	1.15																		
min	1.0																		
25%	2.76																		
50%	3.32																		
75%	3.92																		
max	27.58																		
In (22):	songs_over_1_min.median()																		
Out (22):	minutesPlayed 3.36 dtype: float64																		



As we can see from the above table and graph, the mean track length (songs only) that I listen to is 2.45 minutes and the median is 2.90 minutes. As expected, this shows that the podcasts skew the mean more heavily towards the right and vastly increase the standard deviation (6.02 with podcast data vs. 1.73 without). It also largely increases the max length of a track (155.47 minutes with podcast data vs. 27.54 minutes without).

The median stays relatively unchanged between the two samples (2.93 minutes with podcast data vs. 2.90 without) because the podcasts that are skewing the data do not affect the median as drastically as they do the mean (which is why the median is a helpful measure when dealing with skewed data).

Noted that we still have to deal with the issue of having very short listens being under 1 minute. To exclude the songs that were only listened for a short period of time, we will exclude all plays that were under 1 minute in the below cell. This should give us a more accurate representation of the length of songs to which I listened.

How long are the songs that I listen to (excluding podcast data and song plays under 1 minute)?

In (21):	<pre>songs_over_1_min = pd.DataFrame(df.loc[(df['artistName'] != 'Pardon My Take') & (df['minutesPlayed'] >= 1),\ round(songs_over_1_min.describe(),2)])</pre>																		
Out (21):	<table><tr><th></th><th>minutesPlayed</th></tr><tr><td>count</td><td>12,705.0</td></tr><tr><td>mean</td><td>3.36</td></tr><tr><td>std</td><td>1.15</td></tr><tr><td>min</td><td>1.0</td></tr><tr><td>25%</td><td>2.76</td></tr><tr><td>50%</td><td>3.32</td></tr><tr><td>75%</td><td>3.92</td></tr><tr><td>max</td><td>27.58</td></tr></table>		minutesPlayed	count	12,705.0	mean	3.36	std	1.15	min	1.0	25%	2.76	50%	3.32	75%	3.92	max	27.58
	minutesPlayed																		
count	12,705.0																		
mean	3.36																		
std	1.15																		
min	1.0																		
25%	2.76																		
50%	3.32																		
75%	3.92																		
max	27.58																		
In (22):	songs_over_1_min.median()																		
Out (22):	minutesPlayed 3.36 dtype: float64																		

As we can see from the above table and graph, the mean track length (songs only) that I listen to is 2.45 minutes and the median is 2.90 minutes. As expected, this shows that the podcasts skew the mean more heavily towards the right and vastly increase the standard deviation (6.02 with podcast data vs. 1.73 without). It also largely increases the max length of a track (155.47 minutes with podcast data vs. 27.54 minutes without).

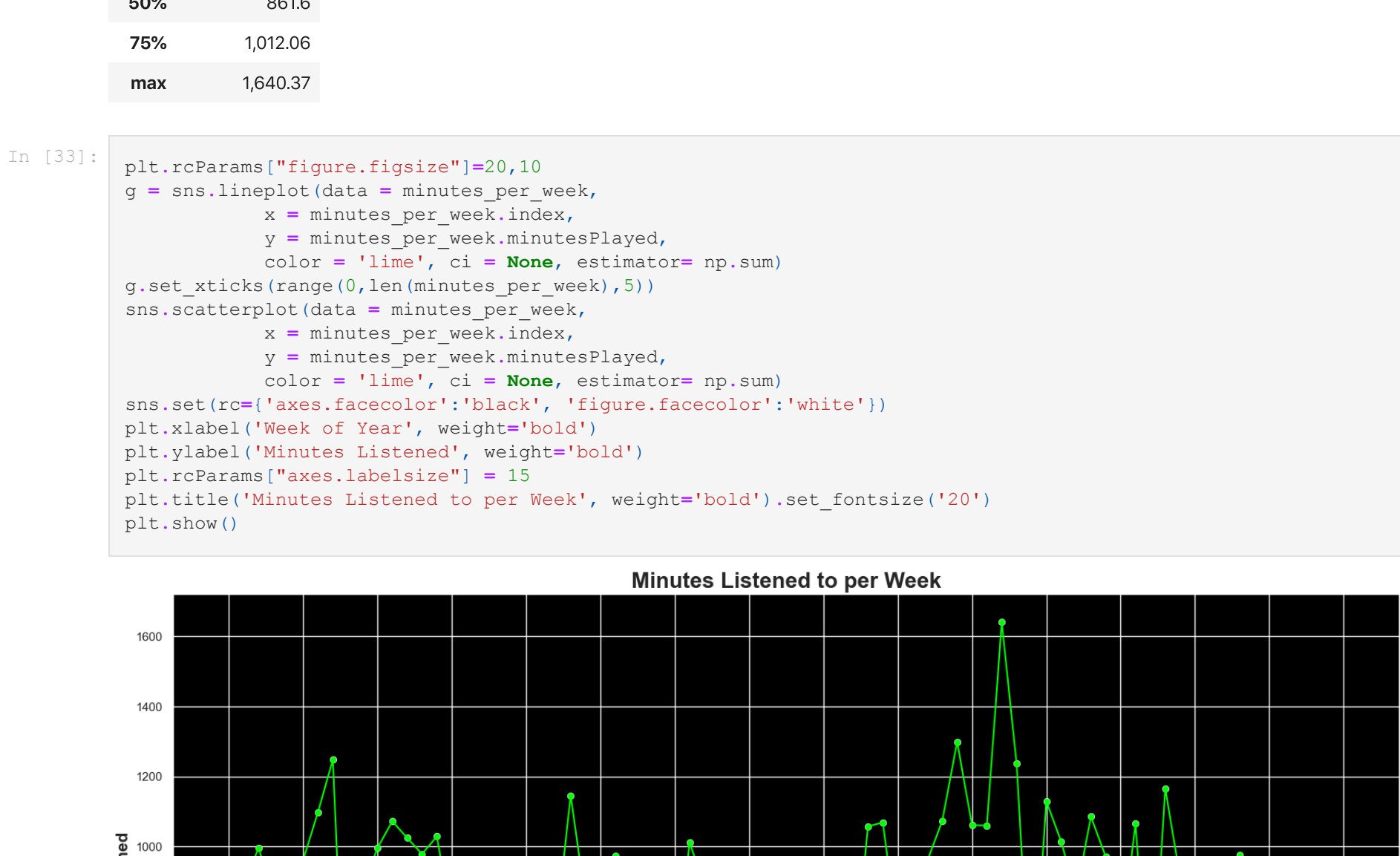
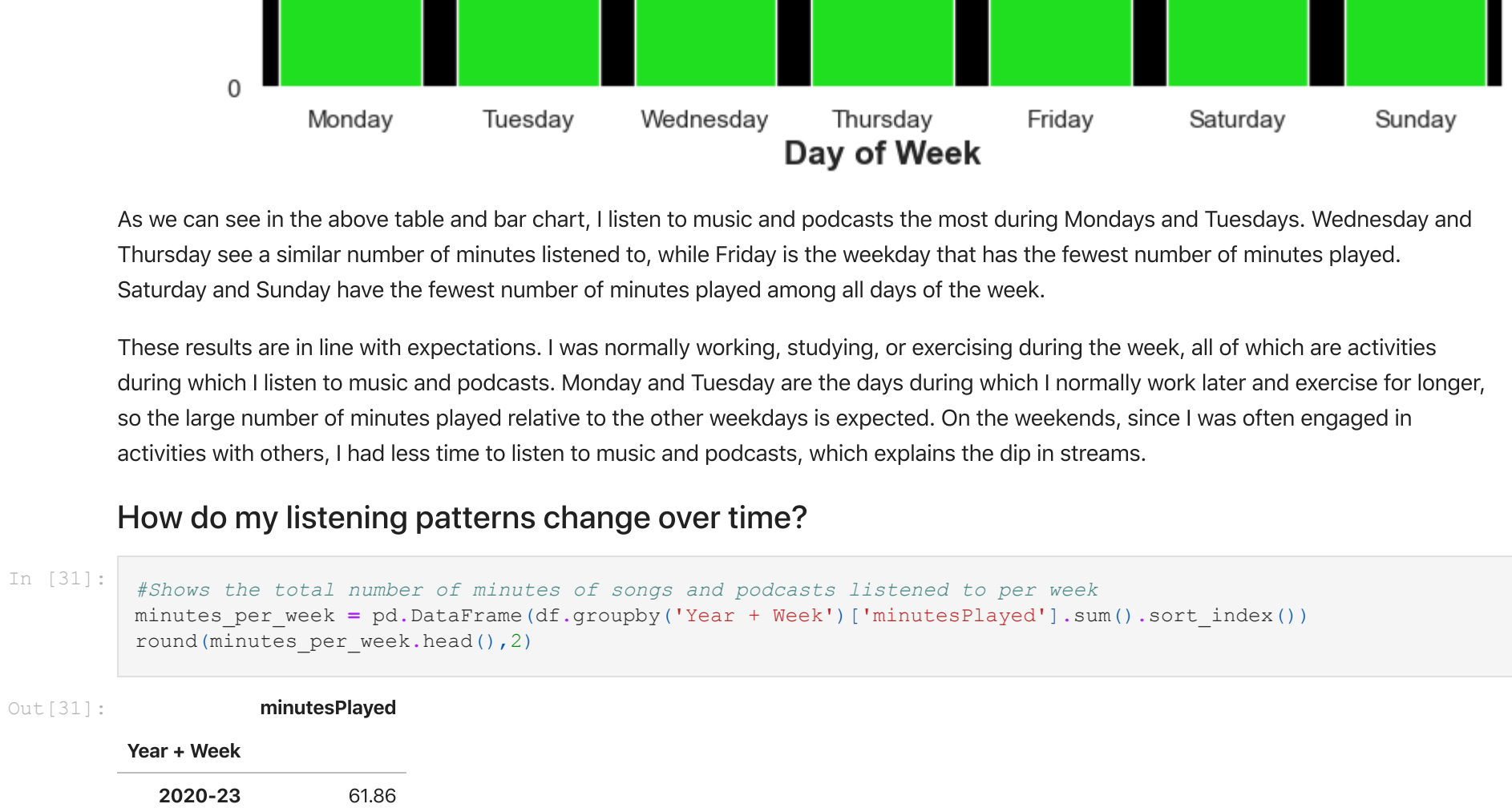
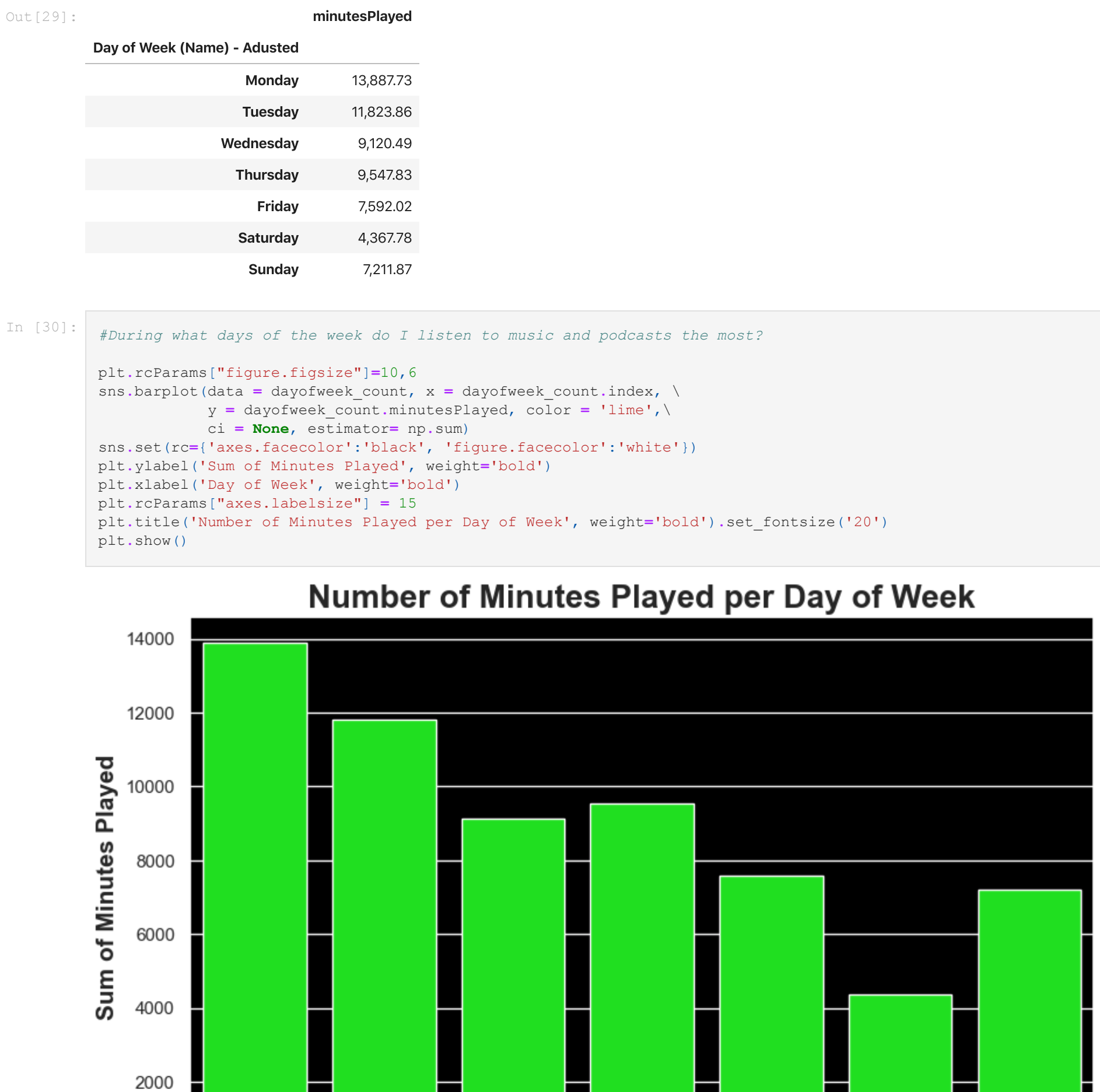
The median stays relatively unchanged between the two samples (2.93 minutes with podcast data vs. 2.90 without) because the podcasts that are skewing the data do not affect the median as drastically as they do the mean (which is why the median is a helpful measure when dealing with skewed data).

Noted that we still have to deal with the issue of having very short listens being under 1 minute. To exclude the songs that were only listened for a short period of time, we will exclude all plays that were under 1 minute in the below cell. This should give us a more accurate representation of the length of songs to which I listened.

How long are the songs that I listen to (excluding podcast data and song plays under 1 minute)?

In (21):	<pre>songs_over_1_min = pd.DataFrame(df.loc[(df['artistName'] != 'Pardon My Take') & (df['minutesPlayed'] >= 1),\ round(songs_over_1_min.describe(),2)])</pre>																		
Out (21):	<table><tr><th></th><th>minutesPlayed</th></tr><tr><td>count</td><td>12,705.0</td></tr><tr><td>mean</td><td>3.36</td></tr><tr><td>std</td><td>1.15</td></tr><tr><td>min</td><td>1.0</td></tr><tr><td>25%</td><td>2.76</td></tr><tr><td>50%</td><td>3.32</td></tr><tr><td>75%</td><td>3.92</td></tr><tr><td>max</td><td>27.58</td></tr></table>		minutesPlayed	count	12,705.0	mean	3.36	std	1.15	min	1.0	25%	2.76	50%	3.32	75%	3.92	max	27.58
	minutesPlayed																		
count	12,705.0																		
mean	3.36																		
std	1.15																		
min	1.0																		
25%	2.76																		
50%	3.32																		
75%	3.92																		
max	27.58																		
In (22):	songs_over_1_min.median()																		
Out (22):	minutesPlayed 3.36 dtype: float64																		

As we can see from the above table and graph, the mean track length (songs only) that I listen to is 2.45 minutes and the median is 2.90 minutes. As expected, this shows that the podcasts skew the mean more heavily towards the right and vastly increase the standard deviation (6.02 with podcast data vs. 1.73 without). It also

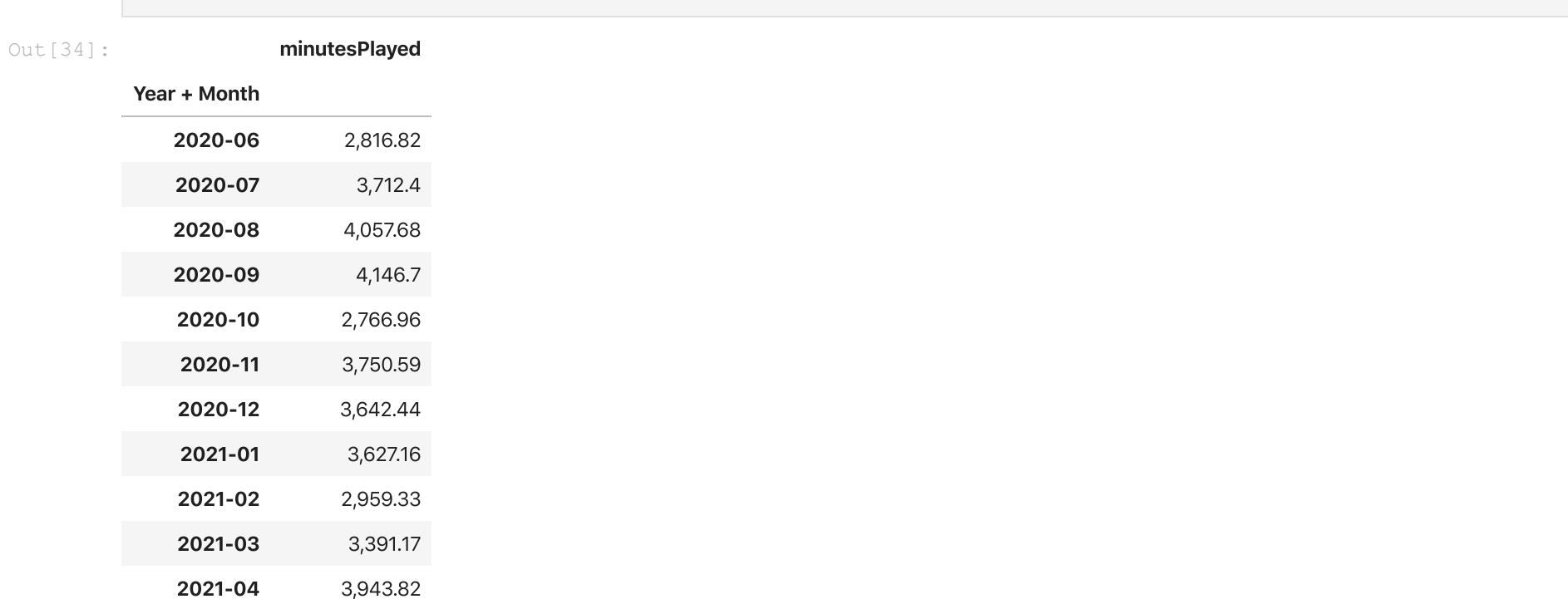


As we can see in the above table and line chart, the number of minutes listened to by week can vary drastically from week to week but overall appears to be generally stable.

There are sometimes large drops in minutes played from one week to the next, such as from 2020-26 (the 26th week of 2020) to 2020-27 (the 27th week of 2020). With a standard deviation of approximately 267 minutes per week, there is certainly some volatility within the data.

There are various factors that could have influenced why certain weeks are higher than others, such as if I exercised more or drove more in certain weeks (both activities during which I heavily listen to music/podcasts).

During what months do I most often listen to music *and* podcasts?



It appears that the monthly data is distributed relatively uniformly but with certain waves that occur throughout the months.

We can see from this data that the month in which I listened the most content was May 2021, while the month in which I listened to the least amount of content was October 2020 (we disregard November 2021 since it does not have a full month's worth of data). There is a 79% increase in minutes listened from October 2020 to May 2021!

There appear to be three different waves over the months with respect to minutes listened. The first peak starts in June 2020 and ends in October 2020, the second peak begins in November 2020 and ends in February 2021, and the third peak begins in March 2021 and ends in October 2021. There are various factors that could have influenced why certain months peaked, such as if I exercised more or drove more in certain months (both activities during which I heavily listen to music/podcasts) or if I was introduced to new music in certain months (which would cause me to listen more than normal).

During what days of the week do I listen to *podcasts* the most?

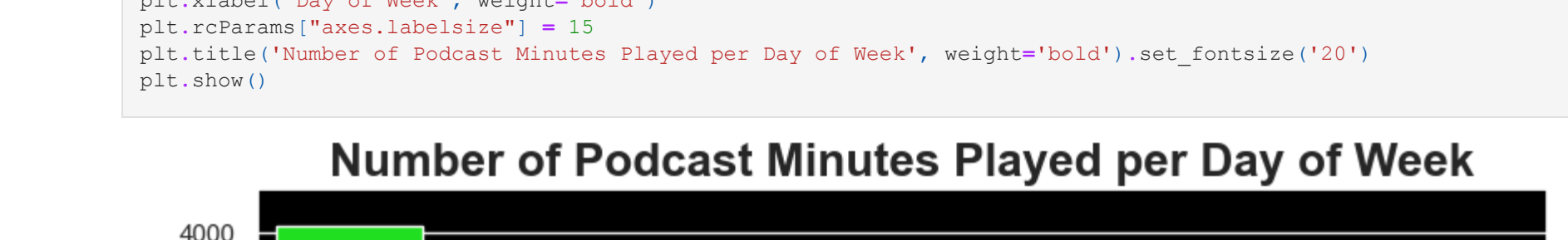
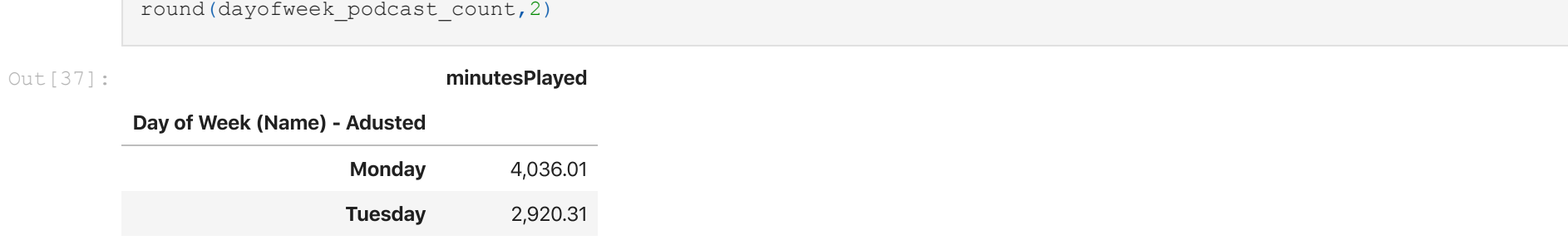


We can see in the above table and bar chart that I most often listen to podcasts on Monday, Tuesday, Wednesday, Thursday, and Friday have a similar number of minutes listened to per day, while Saturday and Sunday once again have the fewest number of minutes played.

As mentioned in the day of week chart for songs and podcast data, I was normally working, studying, or exercising during the week, all of which are activities during which I listen to podcasts. Monday is the day during which I normally work the latest and exercise for the longest, so the large number of minutes played relative to the other weekdays is expected. One other reason why Monday might have such a relatively large volume of minutes played besides the reasons mentioned above is that the length of episodes released on Mondays are typically longer than those released on Wednesdays and Fridays.

It is important to note that new podcasts for Pardon My Take are typically released every Monday, Wednesday, and Friday morning.

How many minutes per month did I spend listening to songs vs. podcasts?



As we can see from the above lineplots and table, I listen to music more often than podcasts per month. The difference in minutes listened between songs and podcast varies greatly from month to month—October 2021 only saw a difference of 85 minutes, while June 2021 saw a large difference of 3,017 minutes. The average difference between minutes listening to songs vs. minutes listening to podcasts is 1,419 minutes.

The large average gap between songs and podcast minutes is in line with expectations. Podcasts have a ceiling of how many minutes I can listen to, because once I finish the most recent episode, I have no new content to listen to. I will never run out of songs to listen to, on the other hand. The other major factor is simply content preference—I am more often in the mood to listen to music than I am to listen to podcasts.

Takeaways

We have navigated through nearly 1.5 years worth of my streaming history. Here are some takeaways that were uncovered during the process:

- Pardon My Take is my most listened to artist overall, while Mac Miller is my top artist for songs. This is because I listened to Pardon My Take for hours every week and because Mac Miller is my favorite musician.
- My top song was Dang! (feat. Anderson .Paak) by Mac Miller. While this is one of my favorite songs, I am also designing a video game level using this song, which would account for some of its plays.
- I listened to 63,352 minutes of content over 517 days for an average number of 123 minutes listened to per day. The daily number of minutes listened to depends on factors such as if I have been exposed to new music/podcast episodes and what activities I'm doing on that day (e.g., exercising, driving, studying, etc.).
- When excluding podcast plays and listens under 1 minute, the average length of my songs are 3.36 minutes long.
- I listen to music the most from 4pm through 7pm, while from 4am through 7am I listen to almost no music. This depends largely on what activities I do during certain times. For example, I normally work out from 4pm-7pm (an activity during which I listen to content), which explains the large volume of minutes listened. Conversely, because I am (normally) asleep from 4am-7am, I have very few minutes listened to during that time.
- Monday is the day of the week during which I listen to the most content (songs and podcasts). I listen to the least amount of content on the weekends (especially Saturday). Similar to the explanation behind my hour-of-day data, my day-of-week results largely depend on what activities I do during which days. During the week I engage in activities where I more often listen to content (e.g., exercising, driving, and studying), while on the weekends my activities often involve others (with whom I am less likely to listen to content). For the podcast data specifically, the podcast release dates of Monday, Wednesday, and Friday also impact the days on which I listen to episodes.
- The month during which I listened to the most music was May 2021, and the month during which I listened to the least music was October 2020. While likely less affected by weekly activities such as exercise and working, the minutes played per month could change depending on when I am introduced to new music (which would cause me to listen to more music than normal).
- I listen to songs more often than podcasts (an average difference of 1,419 minutes per month). This is most likely because there are only a finite number of new podcast episodes to which I can listen, but there is an unlimited number of songs to which I can listen.

Thank you for reading this analysis. I hope you enjoyed it!