

Python初級數據分析員證書

(六) 數據分析及可視化專案

13. 數據分析專案

Demo 10 - Amazon Bestseller



Review

- Statistics
- Hypothesis testing
- Algebra
- Linear regression
- Propositional logic
- Python
- R
- SQL
- Pandas, NumPy, SciPy
- Data Visualization, Matplotlib, Seaborn, Plotly
- Dashboard Visualization, Business Intelligence
- Storytelling



13. 數據分析專案 Data Analysis Project – Demo 10

Chapter Summary

- Scenario
- Wordcloud
- Missingno
- Data Import
- Ratings
- Reviews
- Author

Scenario

By analysing a database of Amazon Top 50 Bestselling Books 2009-2019, you are asked to get inspired for books trading for a local online book shop.



WordCloud

WordCloud are useful for visualizing common words in a text or data set.

```
1 pip install wordcloud
```

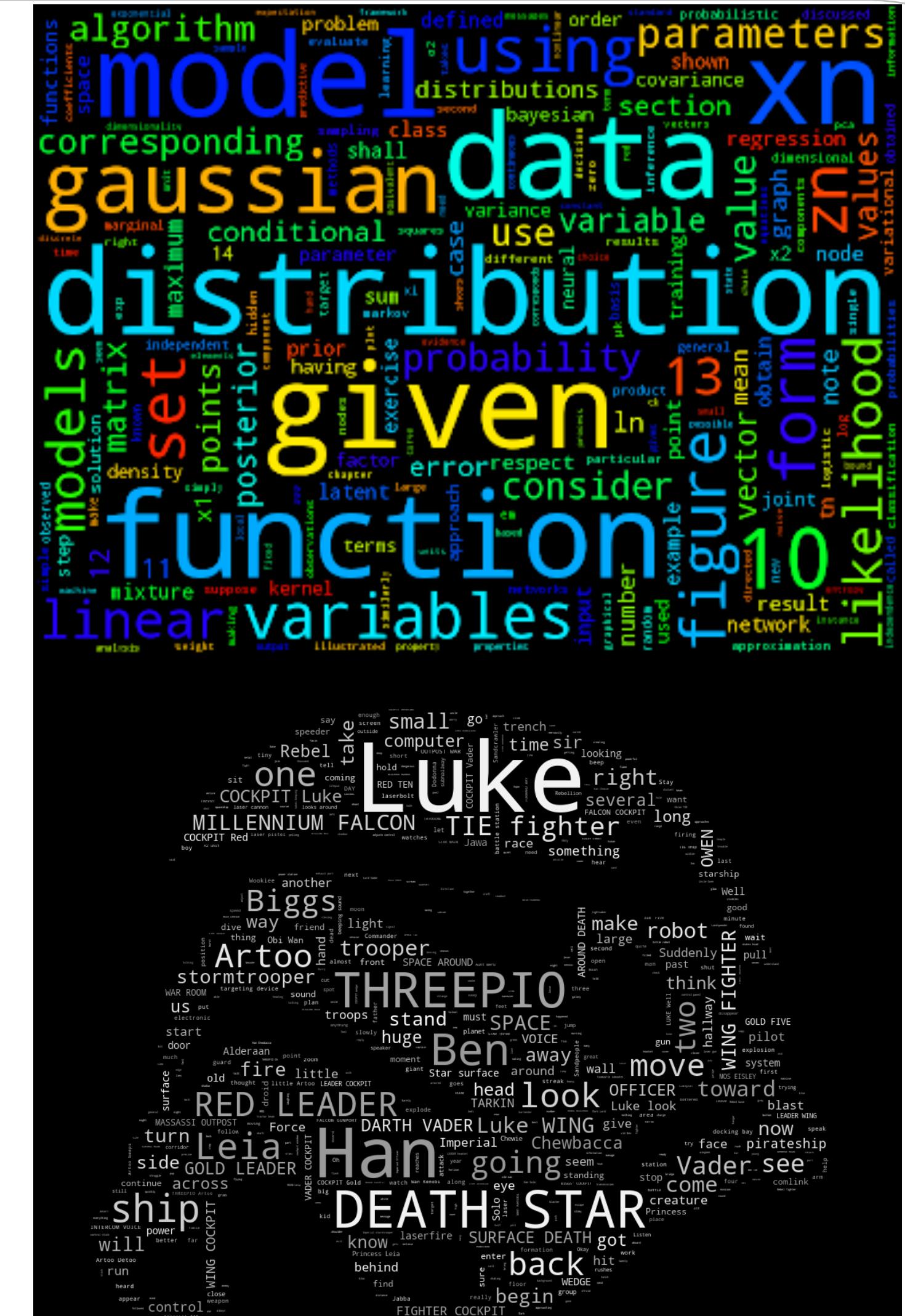
Collecting wordcloud

Downloading wordcloud-1.9.1.1-cp39-cp39-mac

```
Requirement already satisfied: numpy>=1.6.1 from wordcloud) (1.23.1)
```

Requirement already satisfied: pillow in /Users/dcloud/.local/lib/python3.7/site-packages (9.2.0)

Requirement already satisfied: matplotlib in wordcloud) (3.5.2)



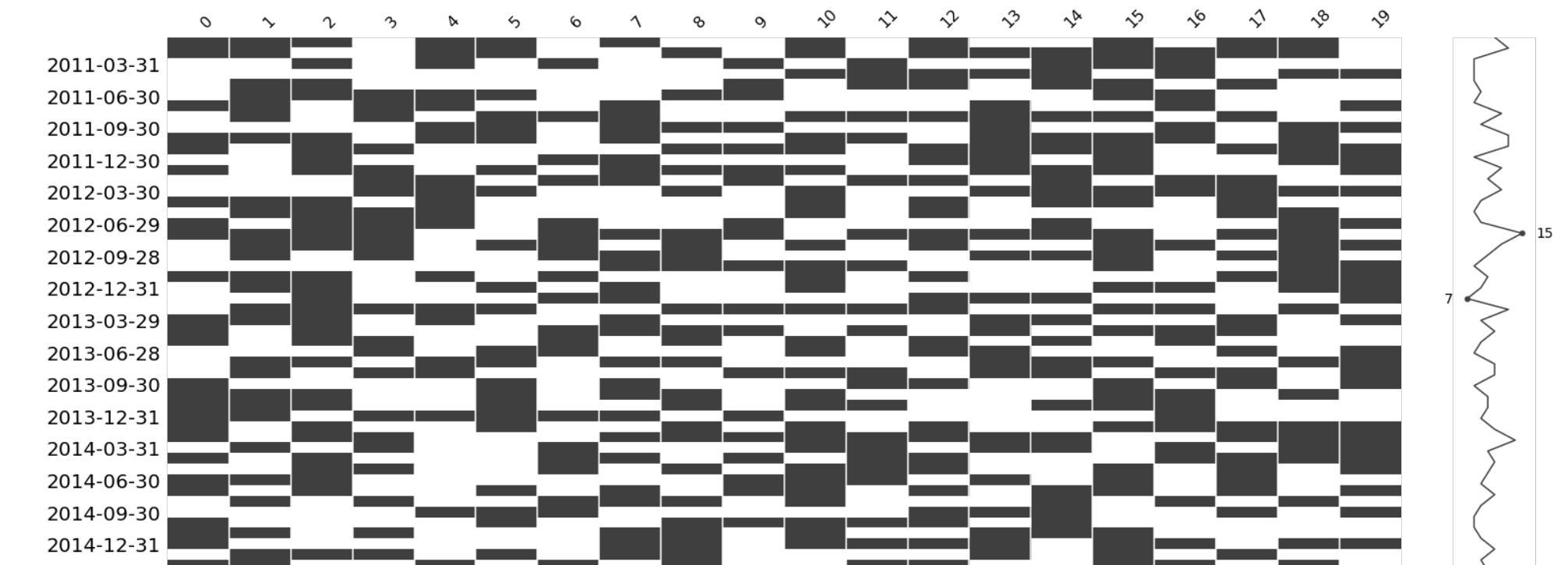
Missingno

Missingno is an excellent and simple to use Python library that provides a series of visualisations to understand the presence and distribution of missing data within a pandas DataFrame.

```
1 pip install missingno
```

```
Collecting missingno
  Downloading missingno-0.5.2-py3-none-any.whl
Requirement already satisfied: numpy in /Users/ingno) (1.23.1)
Requirement already satisfied: matplotlib in /Umissingno) (3.5.2)
Requirement already satisfied: scipy in /Users/ingno) (1.10.1)

```



Import needed library

```

1 import numpy as np
2 import pandas as pd
3
4 # Basic Visualization tools
5 import matplotlib
6 import matplotlib.pyplot as plt
7 plt.style.use('ggplot')
8 plt.rcParams['figure.dpi'] = 300
9 import seaborn as sns
10 sns.set_palette('husl')
11
12 plt.rc('figure', figsize=(17,13))
13 import plotly.express as px
14 import plotly.graph_objs as go
15 import plotly.offline as pyo
16 from plotly.subplots import make_subplots
17
18 # Special Visualization
19 import wordcloud, missingno
20 from wordcloud import WordCloud # wordcloud
21 import missingno as msno # check missing value
22 #import networkx as nx

```

```

23
24 # Plotly visualization
25 import plotly.offline as py
26 from plotly.offline import iplot, init_notebook_mode
27 import plotly.graph_objs as go
28 # Required to use plotly offline in jupyter notebook
29 py.init_notebook_mode(connected = True)
30
31 #importing plotly and cufflinks in offline mode
32 import cufflinks as cf
33 import plotly.offline
34 cf.go_offline()
35 cf.set_config_file(offline=False, world_readable=True)
36
37 # Display markdown formatted output like bold, italic bold etc
38 from IPython.display import Markdown
39 def bold(string):
40     display(Markdown(string))

```

Data Import

```
1 data = pd.read_csv('bestsellers.csv')
2 data.head()
```

	Name	Author	User Rating	Reviews	Price	Year	Genre
0	10-Day Green Smoothie Cleanse	JJ Smith	4.7	17350	8	2016	Non Fiction
1	11/22/63: A Novel	Stephen King	4.6	2052	22	2011	Fiction
2	12 Rules for Life: An Antidote to Chaos	Jordan B. Peterson	4.7	18979	15	2018	Non Fiction
3	1984 (Signet Classics)	George Orwell	4.7	21424	6	2017	Fiction
4	5,000 Awesome Facts (About Everything!) (Natio...	National Geographic Kids	4.8	7665	12	2019	Non Fiction

```
1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 550 entries, 0 to 549
Data columns (total 7 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   Name         550 non-null    object  
 1   Author       550 non-null    object  
 2   User Rating  550 non-null    float64
 3   Reviews      550 non-null    int64  
 4   Price        550 non-null    int64  
 5   Year         550 non-null    int64  
 6   Genre        550 non-null    object  
dtypes: float64(1), int64(3), object(3)
memory usage: 30.2+ KB
```

DF.describe(include='all')

Apart from using `data.describe()` for numerical data structure viewing, we may use `data.describe(include=['O'])` for string data viewing, and `data.describe(include='all')` for all

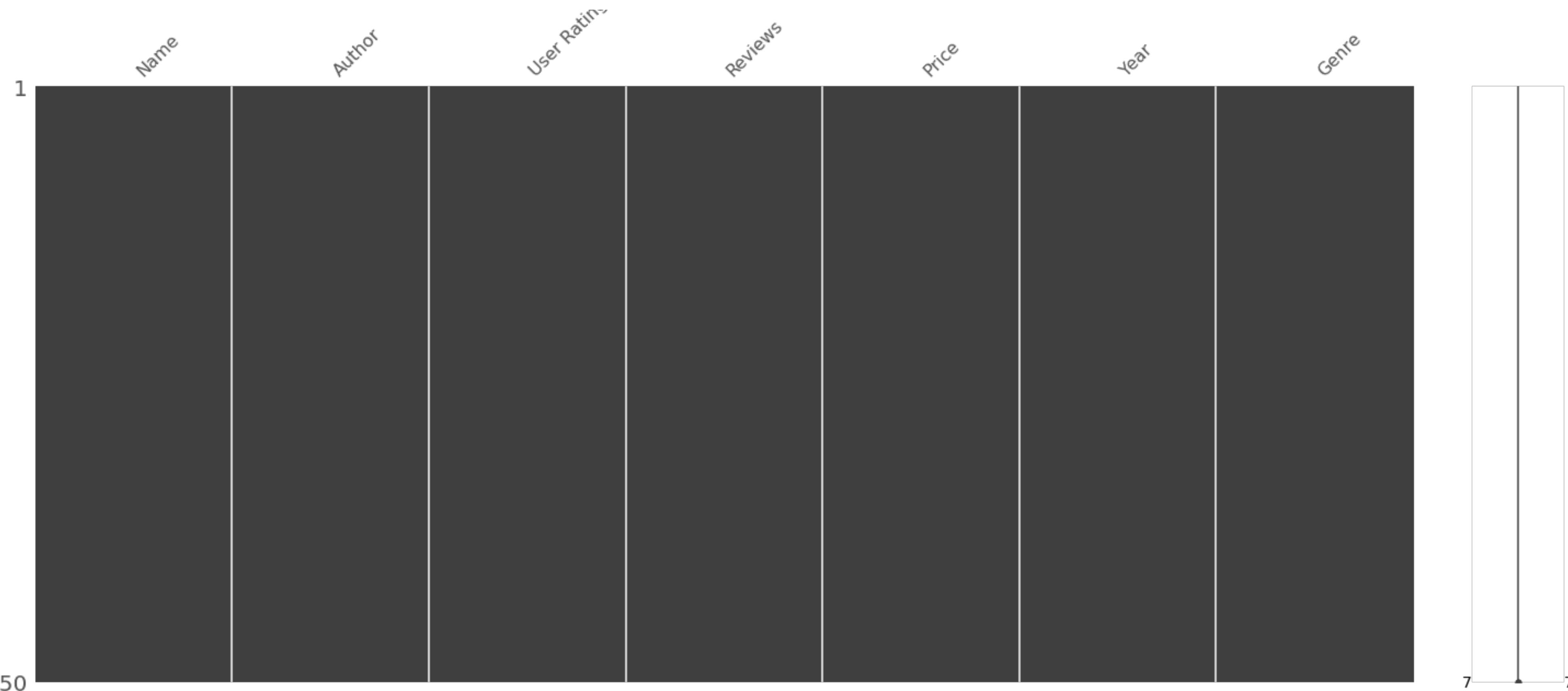
```
1 data.describe(include='all')
```

	Name	Author	User Rating	Reviews	Price	Year	Genre
count	550	550	550.000000	550.000000	550.000000	550.000000	550
unique	351	248	NaN	NaN	NaN	NaN	2
top	Publication Manual of the American Psychologic...	Jeff Kinney	NaN	NaN	NaN	NaN	Non Fiction
freq	10	12	NaN	NaN	NaN	NaN	310
mean	NaN	NaN	4.618364	11953.281818	13.100000	2014.000000	NaN
std	NaN	NaN	0.226980	11731.132017	10.842262	3.165156	NaN
min	NaN	NaN	3.300000	37.000000	0.000000	2009.000000	NaN
25%	NaN	NaN	4.500000	4058.000000	7.000000	2011.000000	NaN
50%	NaN	NaN	4.700000	8580.000000	11.000000	2014.000000	NaN
75%	NaN	NaN	4.800000	17253.250000	16.000000	2017.000000	NaN
max	NaN	NaN	4.900000	87841.000000	105.000000	2019.000000	NaN

Check missing data with missingno

```
1 msno.matrix(data)
```

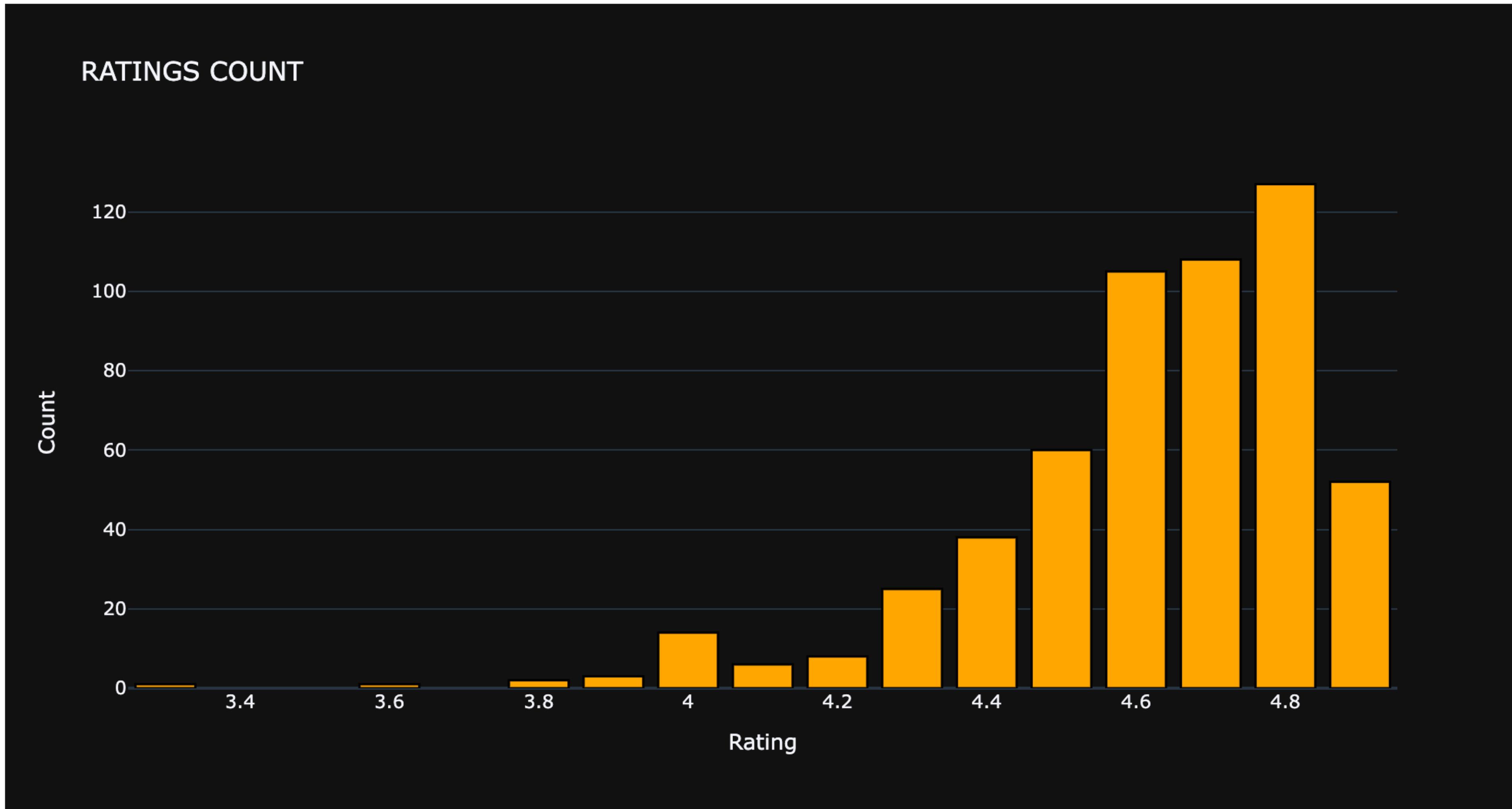
<AxesSubplot:>



Bestseller rating counts

```
1 temp_df = data['User Rating'].value_counts().reset_index()
2
3 # create trace1
4 trace1 = go.Bar(x = temp_df['index'], y = temp_df['User Rating'],
5                  marker = dict(color = 'rgb(255,165,0)',
6                                line=dict(color='rgb(0,0,0)',width=1.5)))
7 layout = go.Layout(template= "plotly_dark",title = 'RATINGS COUNT',
8                      xaxis = dict(title = 'Rating'), yaxis = dict(title = 'Count'))
9 fig = go.Figure(data = [trace1], layout = layout)
10 fig.show()
11 bold("**MOST OF THE RATINGS ARE IN THE RANGE OF 4.6 TO 4.8**")
```

Bestseller rating counts

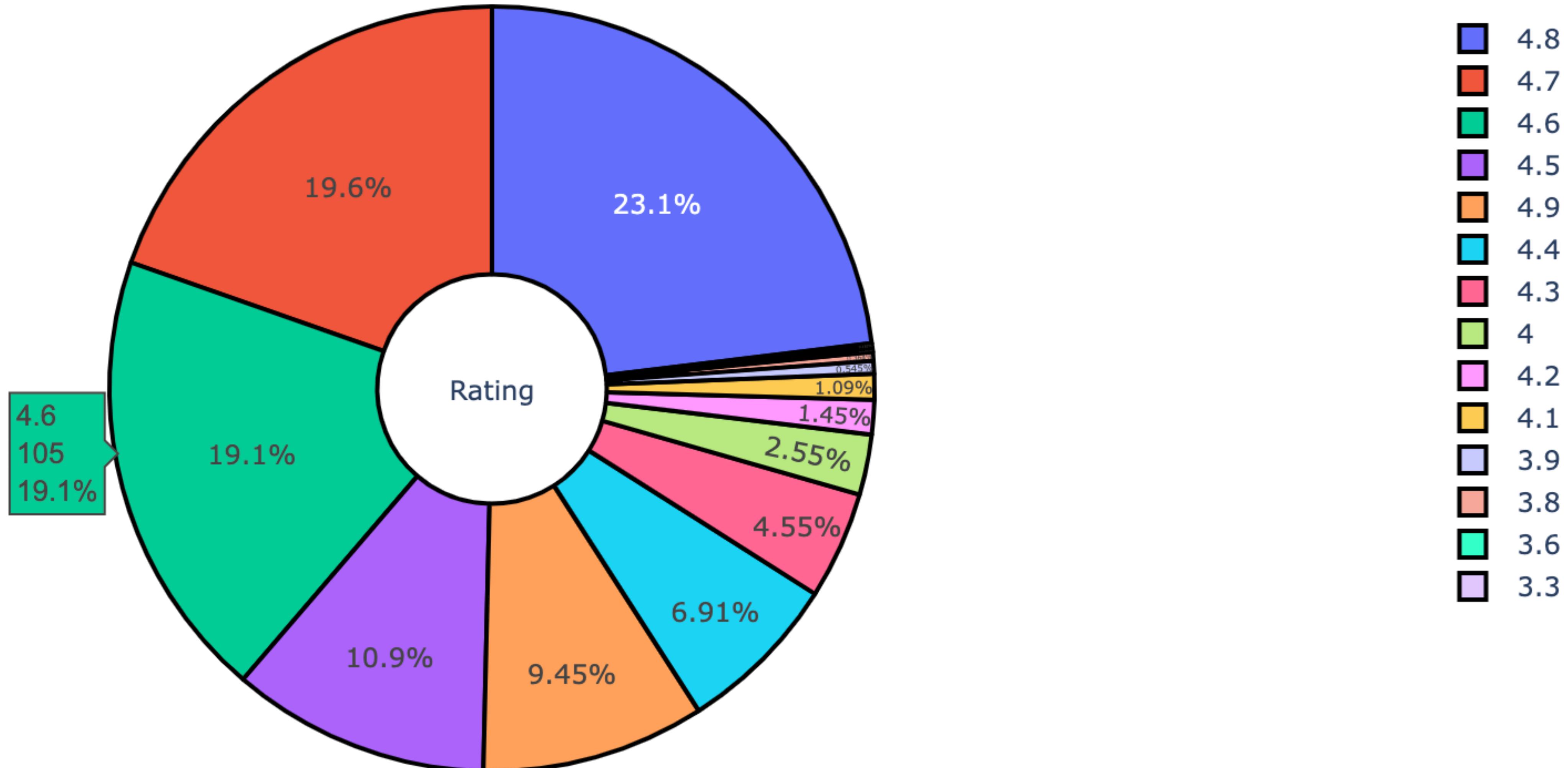


MOST OF THE RATINGS ARE IN THE RANGE OF 4.6 TO 4.8

Bestseller rating counts

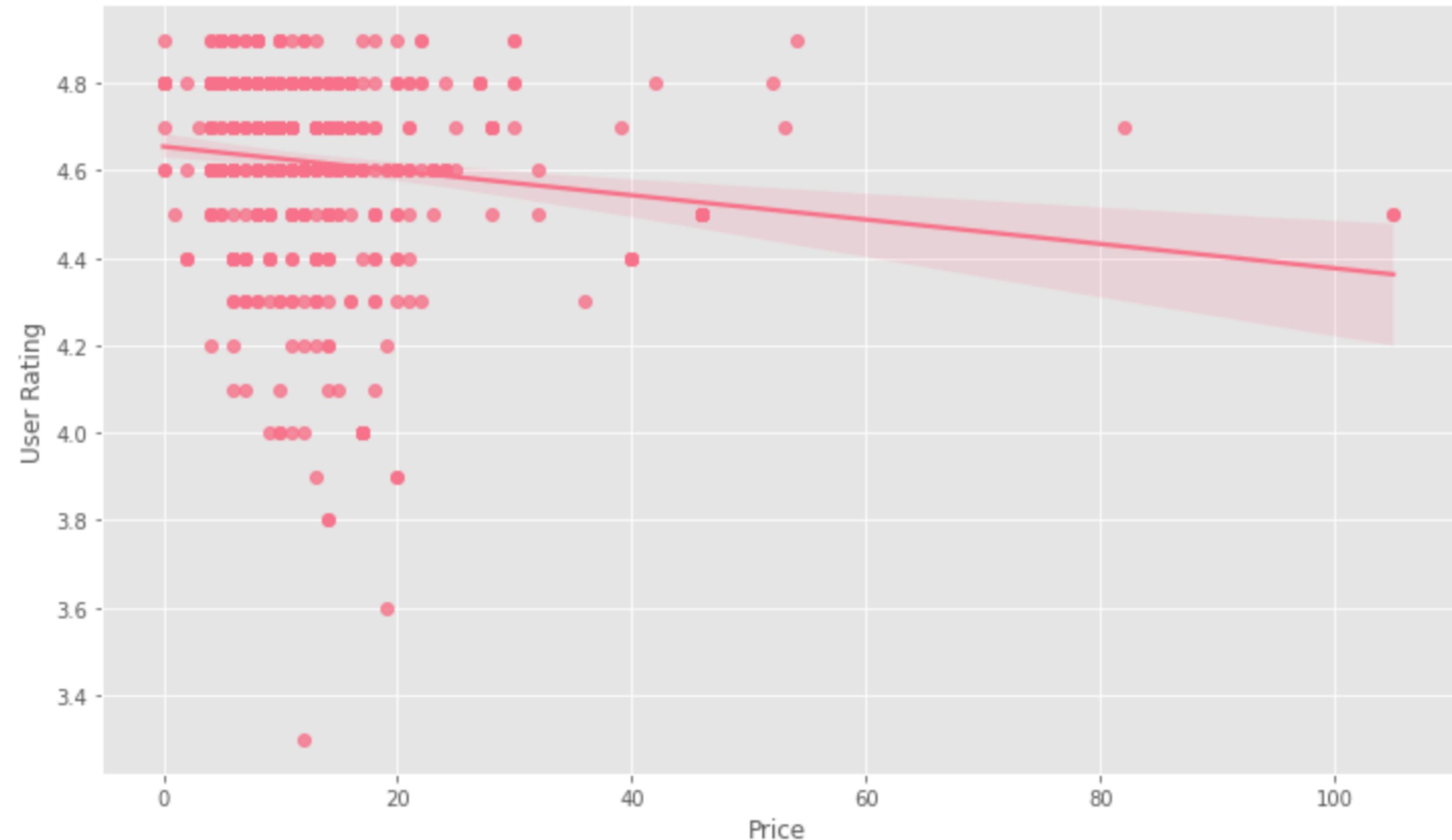
```
1 def pie_plot(cnt_srs, title):
2     labels=cnt_srs.index
3     values=cnt_srs.values
4     trace = go.Pie(labels=labels,
5                      values=values,
6                      title=title,
7                      hoverinfo=['percent+value'],
8                      textinfo='percent',
9                      textposition='inside',
10                     hole=0.3,
11                     showlegend=True,
12                     marker=dict(colors=plt.cm.viridis_r(np.linspace(0, 1, 14)),
13                                 line=dict(color='#000000', width=2) )
14                 )
15     return trace
16 py.iplot([pie_plot(data['User Rating'].value_counts(), 'Rating')])
```

Bestseller rating counts



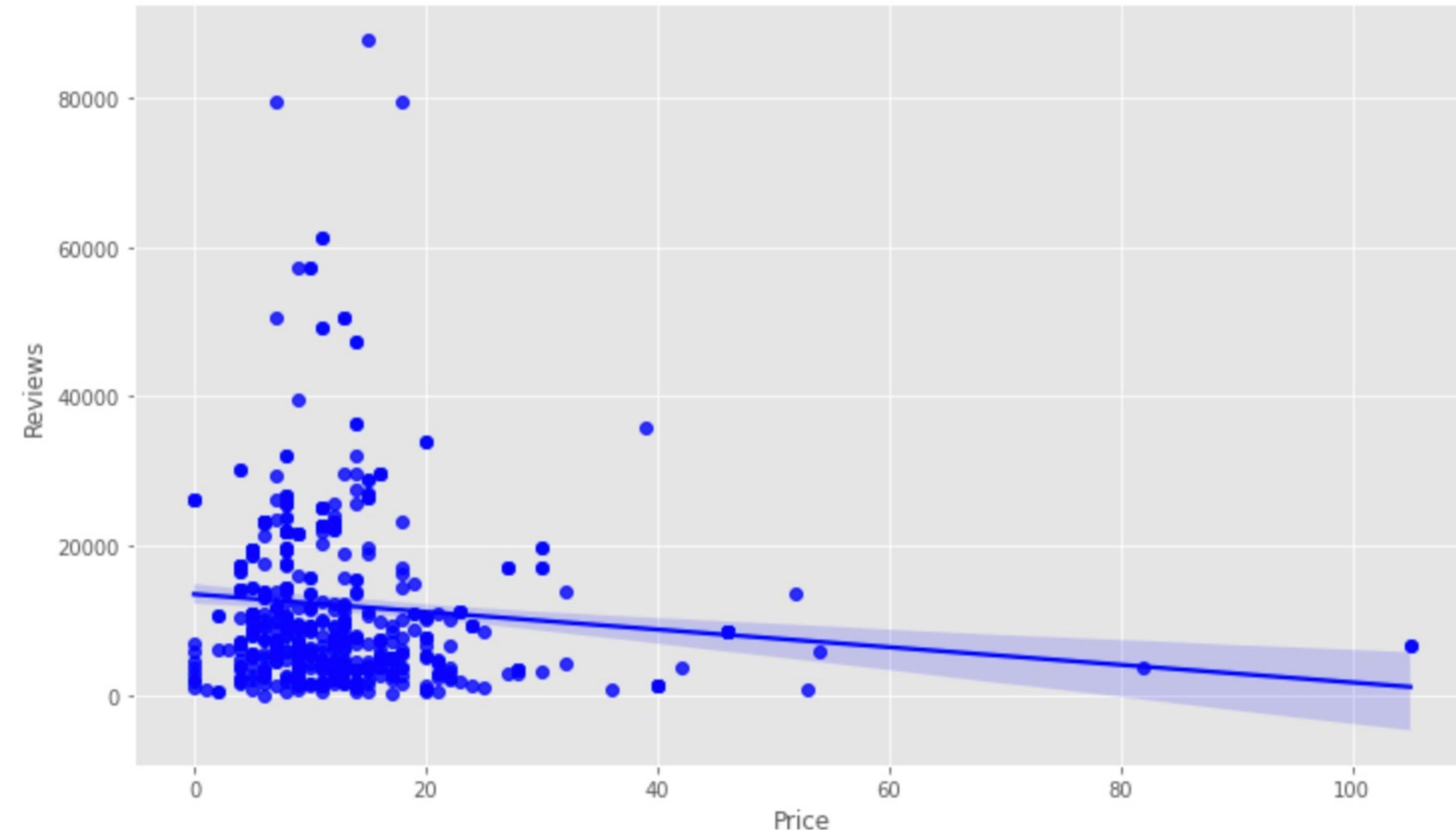
Relationship between Price and User Rating

```
1 fig, ax = plt.subplots(1,1, figsize=(12, 7), dpi=72)
2 sns.regplot(data=data, x='Price', y='User Rating', ax=ax)
3 plt.show()
```



Relationship between Price and Reviews

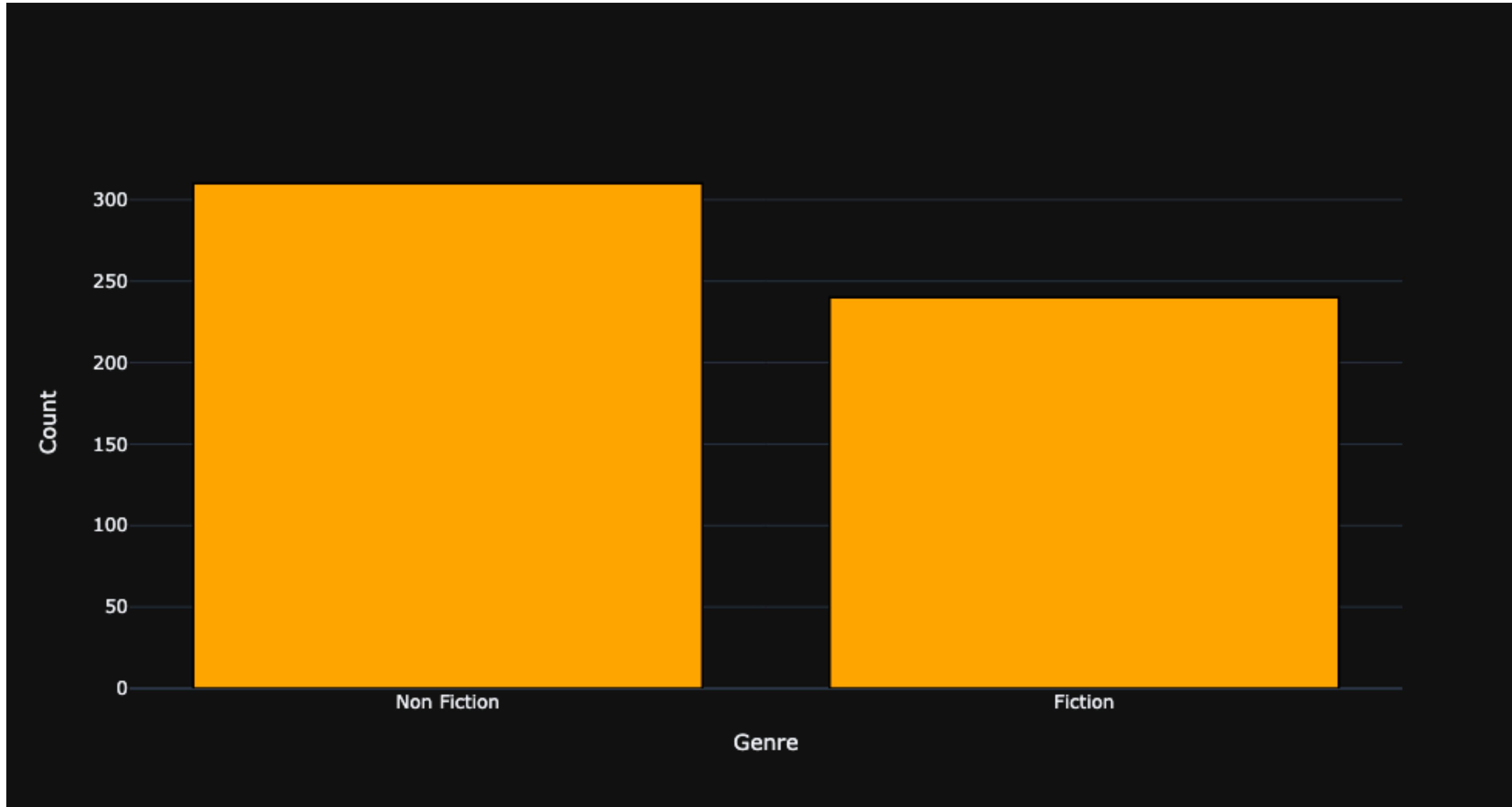
```
1 fig, ax = plt.subplots(1,1, figsize=(12, 7), dpi=72)
2 sns.regplot(data=data, x='Price', y='Reviews', color='b', ax=ax)
3 plt.show()
```



Fiction and Non-fiction Bestseller

```
1 temp_df = data['Genre'].value_counts().reset_index()
2
3 trace1 = go.Bar(
4     x = temp_df['index'],
5     y = temp_df['Genre'],
6     marker = dict(color = 'rgb(255,165,0)',
7                     line=dict(color='rgb(0,0,0)',width=1.5)))
8
9 layout = go.Layout(template= "plotly_dark", xaxis = dict(title = 'Genre'),
10                     yaxis = dict(title = 'Count'))
11 fig = go.Figure(data = [trace1], layout = layout)
12 fig.show()
```

Fiction and Non-fiction Bestseller



Fiction and Non-fiction Bestseller

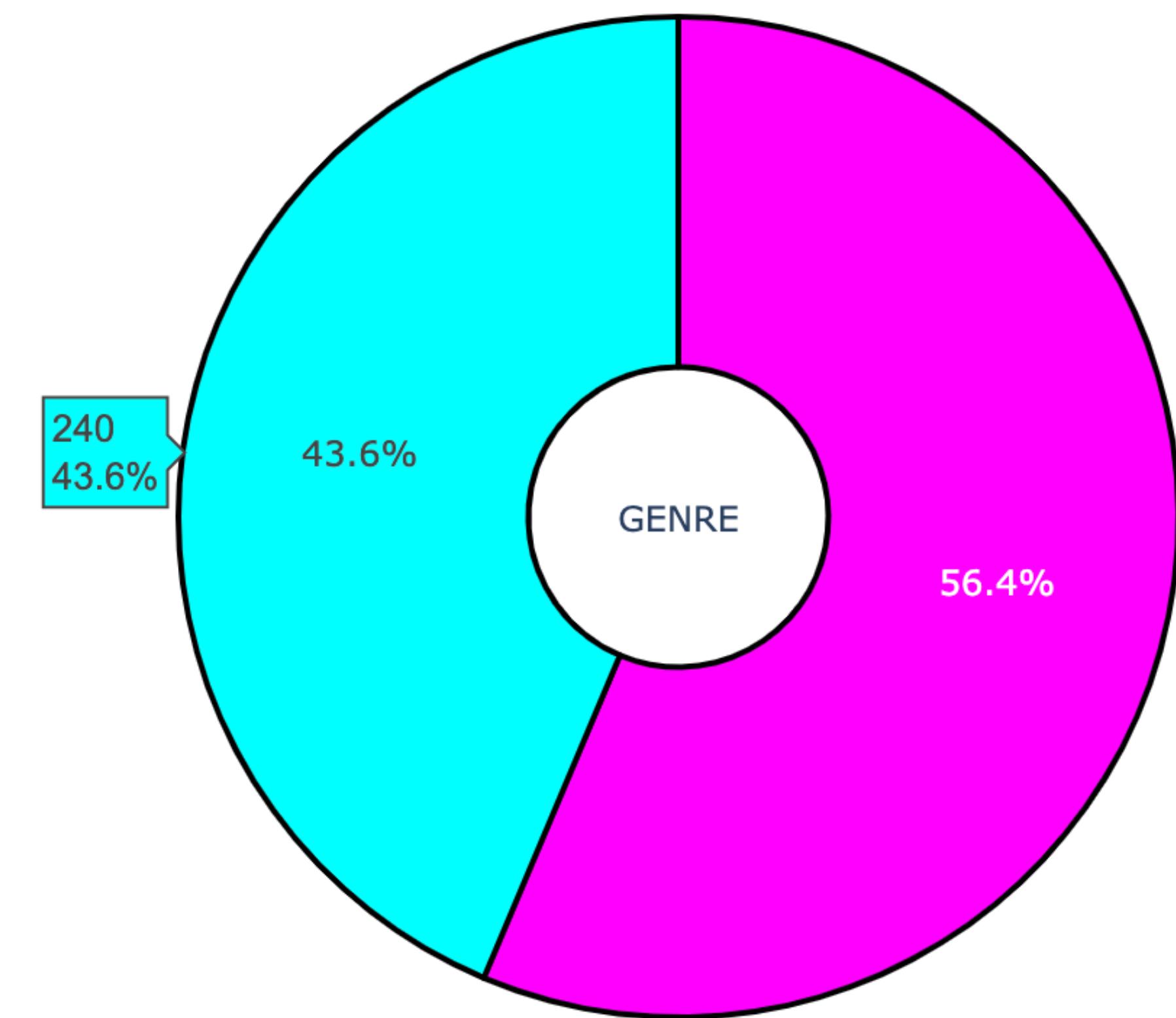
```
1 def pie_plot(cnt_srs, colors, title):
2     '''A Function To Plot Pie Plot using Plotly'''
3     labels=cnt_srs.index
4     values=cnt_srs.values
5     trace = go.Pie(labels=labels,
6                      values=values,
7                      title=title,
8                      hoverinfo='percent+value',
9                      textinfo='percent',
10                     textposition='inside',
11                     hole=0.3,
12                     showlegend=True,
13                     marker=dict(colors=colors,
14                                 line=dict(color='#000000',
15                                           width=2), ))
16
17     return trace
18
19 bold("**NON-FICTION BESTSELLERS ARE MORE THAN FICTION**")
20 py.iplot([pie_plot(data['Genre'].value_counts(), ['magenta', 'cyan'], 'GENRE')])
```

Fiction and Non-fiction Bestseller

NON-FICTION BESTSELLERS ARE MORE THAN FICTION



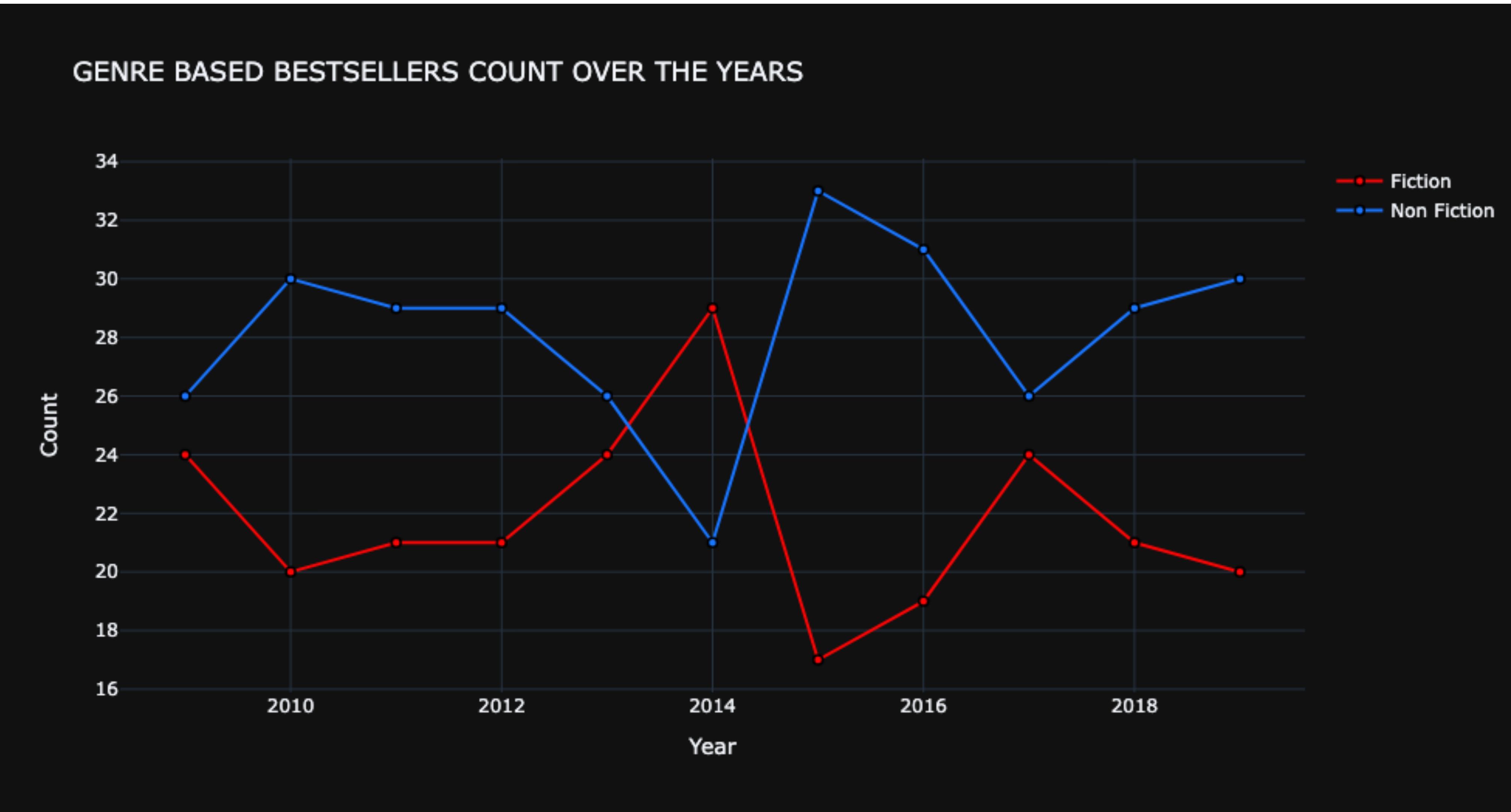
Non Fiction
Fiction



Bestseller count based on genre over the years

```
1 d1 = data[data["Genre"] == "Fiction"]
2 d2 = data[data["Genre"] == "Non Fiction"]
3
4 # Fiction
5 vc1 = pd.DataFrame(data={'count' : d1.Year.value_counts() })
6 vc1['percent'] = vc1['count'].apply(lambda x : 100*x/sum(vc1['count']))
7 vc1 = vc1.sort_index()
8
9 # Non-Fiction
10 vc2 = pd.DataFrame(data={'count' : d2.Year.value_counts() })
11 vc2['percent'] = vc2['count'].apply(lambda x : 100*x/sum(vc2['count']))
12 vc2 = vc2.sort_index()
13
14 trace1 = go.Scatter(x=vc1.index, y=vc1["count"], name="Fiction",
15                      marker=dict(color = 'rgb(249, 6, 6)',
16                                 line=dict(color='rgb(0,0,0)',width=1.5)))
17
18 trace2 = go.Scatter(x=vc2.index, y=vc2["count"], name="Non Fiction",
19                      marker= dict(color = 'rgb(26, 118, 255)',
20                                 line=dict(color='rgb(0,0,0)',width=1.5)))
21 layout = go.Layout(hovermode= 'closest',
22                     title = 'GENRE BASED BESTSELLERS COUNT OVER THE YEARS' ,
23                     xaxis = dict(title = 'Year'), yaxis = dict(title = 'Count'),
24                     template= "plotly_dark")
25 fig = go.Figure(data = [trace1, trace2], layout=layout)
26 fig.show()
```

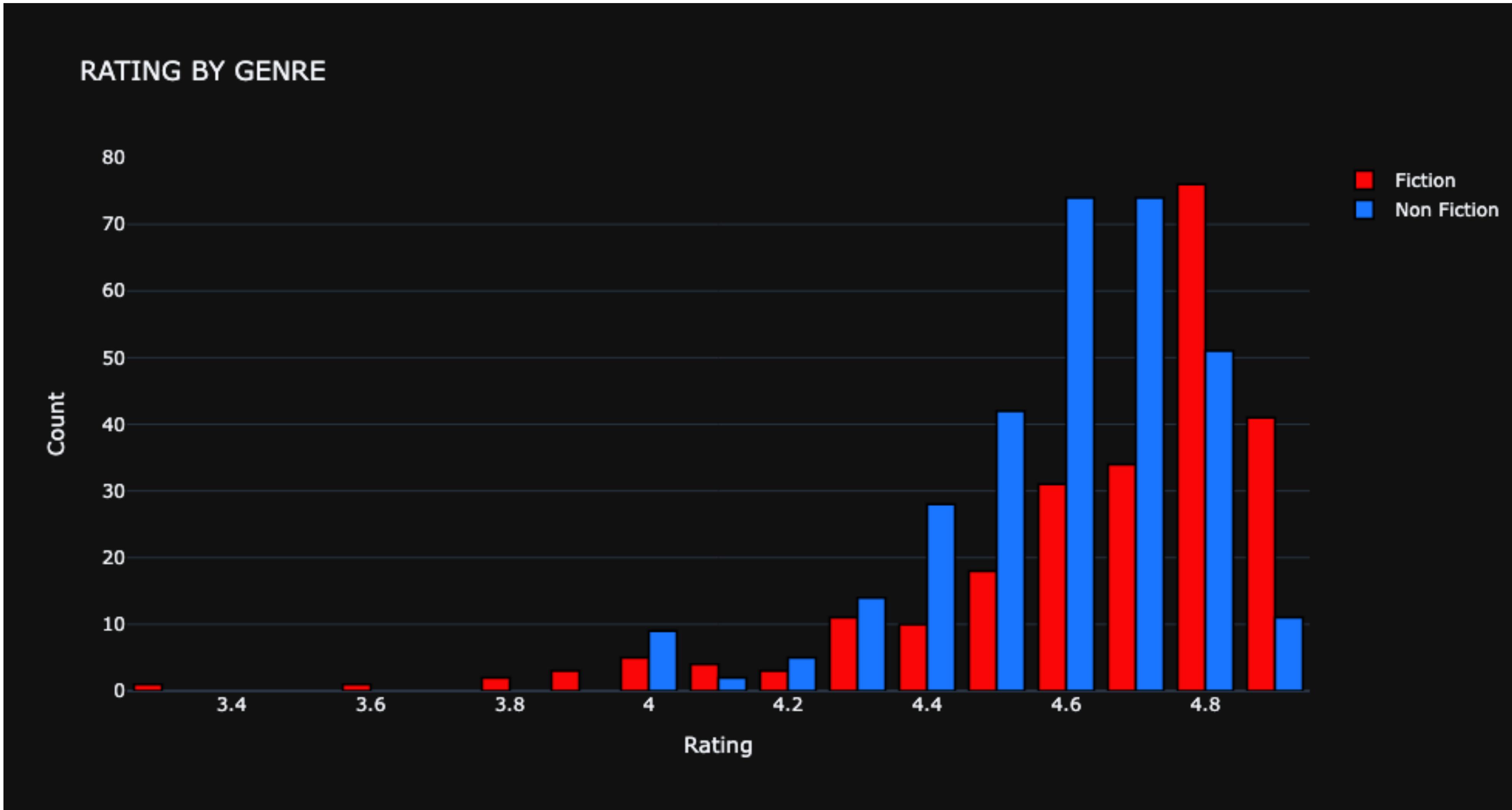
Bestseller count based on genre over the years



Rating on genre

```
1 df1 = data[data["Genre"] == "Fiction"]
2 df2 = data[data["Genre"] == "Non Fiction"]
3
4 temp_df1 = pd.DataFrame(data={'Rating':df1['User Rating'].value_counts()}).sort_index()
5 temp_df2 = pd.DataFrame(data={'Rating':df2['User Rating'].value_counts()}).sort_index()
6
7 trace1 = go.Bar(x = temp_df1.index, y = temp_df1['Rating'],
8                  name="Fiction", marker = dict(color = 'rgb(249, 6, 6)',
9                  line=dict(color='rgb(0,0,0)',width=1.5)))
10 # create trace2
11 trace2 = go.Bar(x = temp_df2.index, y = temp_df2['Rating'],
12                  name = "Non Fiction", marker = dict(color = 'rgb(26, 118, 255)',
13                  line=dict(color='rgb(0,0,0)',width=1.5)))
14
15 layout = go.Layout(template= "plotly_dark",title = 'RATING BY GENRE' ,
16                     xaxis = dict(title = 'Rating'), yaxis = dict(title = 'Count'))
17 fig = go.Figure(data = [trace1, trace2], layout = layout)
18 fig.show()
```

Rating on genre



Top 9 multiple time bestseller

```
1 common_books = data['Name'].value_counts()[:9].rename_axis('Common Books')\
2 .reset_index(name='count')
3
4 fig = px.treemap(common_books, path=['Common Books'],
5                   values='count', title='TOP 9 MULTIPLE TIMES BESTSELLERS')
6
7 fig.show()
```

Top 9 multiple time bestseller

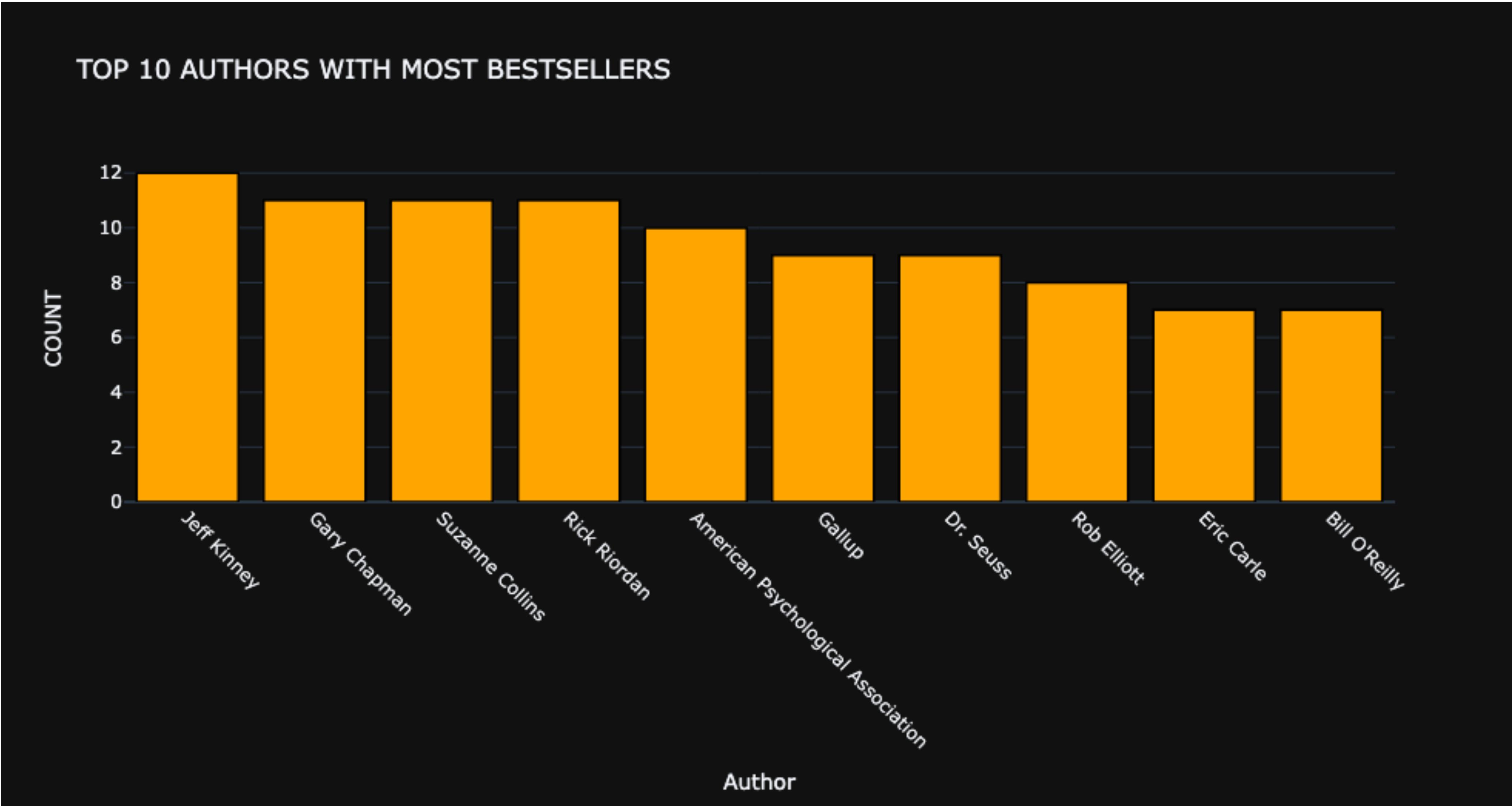
TOP 9 MULTIPLE TIMES BESTSELLERS



TOP 10 AUTHORS WITH MOST BESTSELLERS

```
1 temp_df1 = data.groupby('Author').count().reset_index().sort_values('Name', ascending=False).head(10)
2
3 trace1 = go.Bar(
4     x = temp_df1['Author'],
5     y = temp_df1['Name'],
6     marker = dict(color = 'rgb(255,165,0)',
7                     line=dict(color='rgb(0,0,0)',width=1.5)))
8 layout = go.Layout(template= "plotly_dark",title = 'TOP 10 AUTHORS WITH MOST BESTSELLERS ' ,
9                     xaxis = dict(title = 'Author',tickangle=45), yaxis = dict(title = 'COUNT'))
10 fig = go.Figure(data = [trace1], layout = layout)
11 fig.show()
```

TOP 10 AUTHORS WITH MOST BESTSELLERS



Authors with bestsellers

Authors with bestsellers in both the genre: fiction and non-fiction

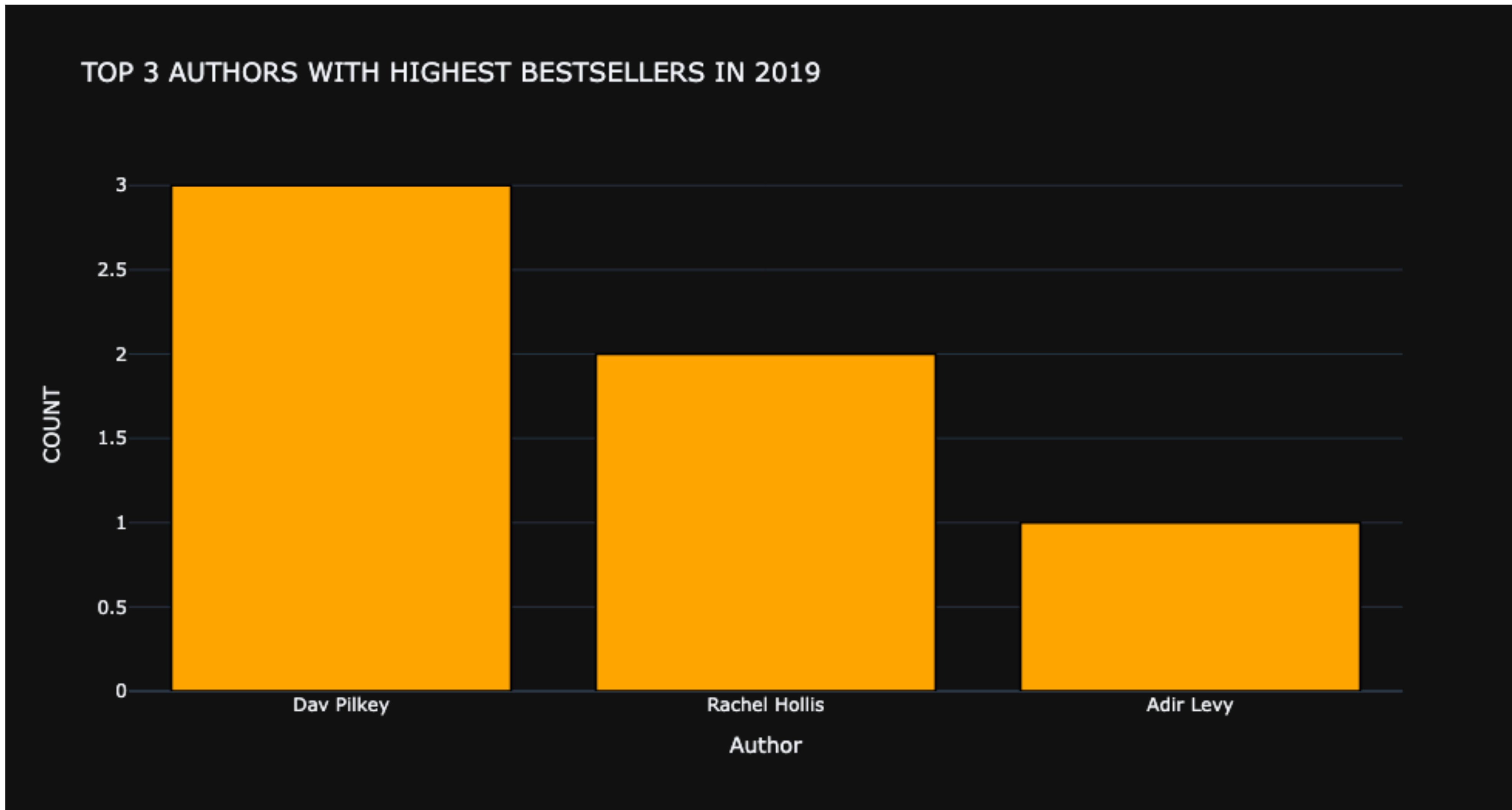
```
1 data[data['Author'].isin(list(data.groupby(['Author','Genre']).count().reset_index()\n2 ['Author'].value_counts()[:2].index))]
```

		Name	Author	User Rating	Reviews	Price	Year	Genre
28	Baby Touch and Feel: Animals		DK	4.6	5360	5	2015	Non Fiction
158	Harry Potter Coloring Book	Scholastic		4.7	3564	9	2015	Non Fiction
268	Pokémon Deluxe Essential Handbook: The Need-to... The Need-to... The Need-to...	Scholastic		4.7	3503	9	2016	Fiction
514	Ultimate Sticker Book: Frozen: More Than 60 Re...		DK	4.5	2586	5	2014	Fiction

Top 3 authors with highest bestsellers in 2019

```
1 temp_df1=data[data['Year']==2019].groupby('Author').count().reset_index()\n2 .sort_values('Name',ascending=False).head(3)\n3\n4 trace1 = go.Bar(\n5     x = temp_df1['Author'],\n6     y = temp_df1['Name'],\n7     marker = dict(color = 'rgb(255,165,0)',\n8                     line=dict(color='rgb(0,0,0)',width=1.5)))\n9 layout = go.Layout(template= "plotly_dark",title = 'TOP 3 AUTHORS WITH HIGHEST BESTSELLERS IN 2019 ' ,\n10                      xaxis = dict(title = 'Author'), yaxis = dict(title = 'COUNT'))\n11 fig = go.Figure(data = [trace1], layout = layout)\n12 fig.show()
```

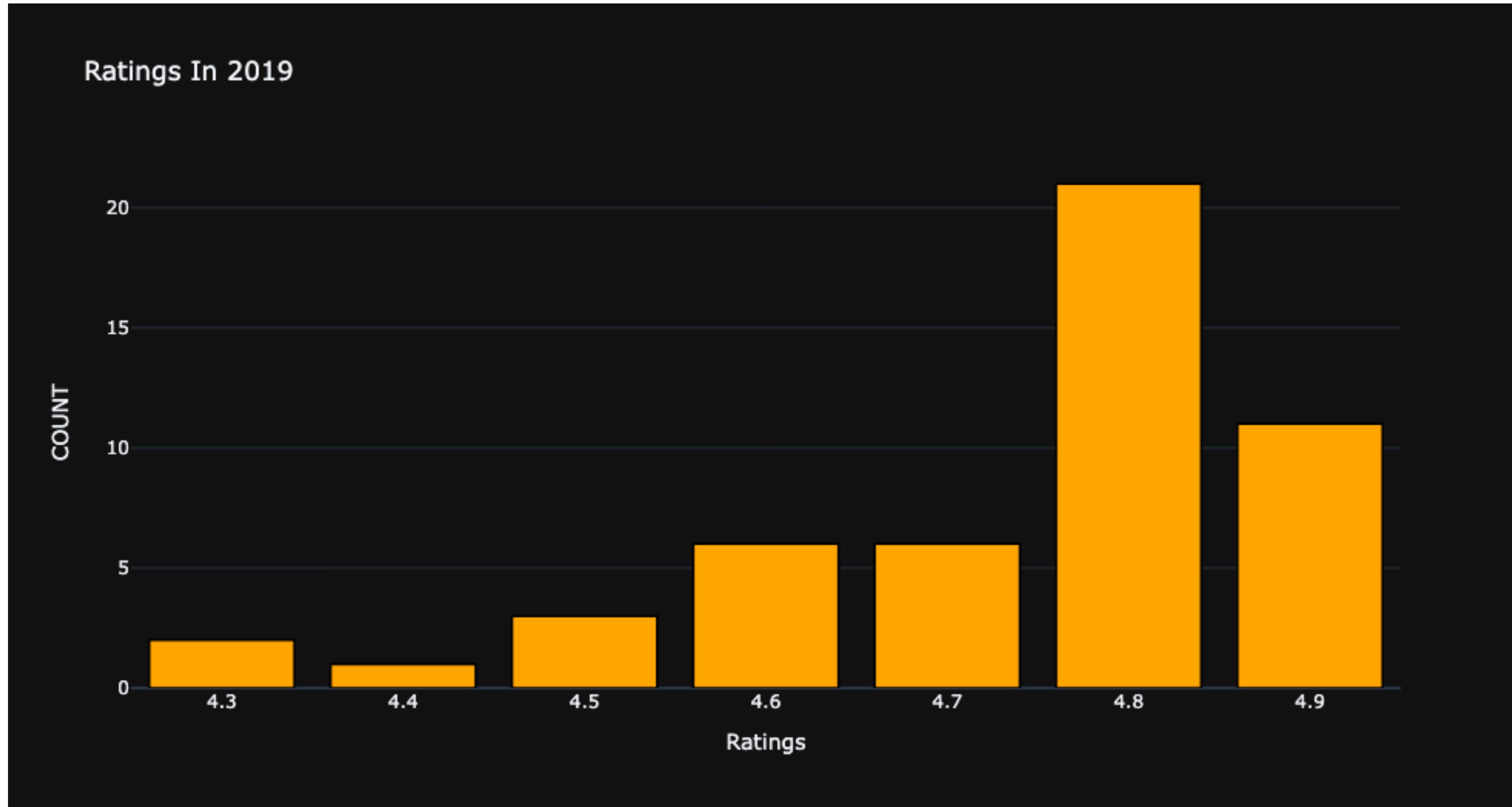
Top 3 authors with highest bestsellers in 2019



Rating in 2019

```
1 temp_df1=data[data['Year']==2019].groupby('User Rating').count().reset_index()\n2 .sort_values('Name',ascending=False)\n3\n4 trace1 = go.Bar(\n5     x = temp_df1['User Rating'],\n6     y = temp_df1['Name'],\n7     marker = dict(color = 'rgb(255,165,0)',\n8                     line=dict(color='rgb(0,0,0)',width=1.5)))\n9 layout = go.Layout(template= "plotly_dark",title = 'Ratings In 2019 ' ,\n10                      xaxis = dict(title = 'Ratings'), yaxis = dict(title = 'COUNT' ))\n11 fig = go.Figure(data = [trace1], layout = layout)\n12 fig.show()
```

Rating in 2019



Chapter Wrap Up

Open Discussion:

- What is the key of the bestseller
- What is the genre and content shall the customer like
- What ideas shall we get from rating and review
- What further data we shall get for better estimation
- What books would you suggest to sell



Reference & Resources

Official Website:

<https://plotly.com/python/>

Plotly Graph Objects:

<https://plotly.com/python/graph-objects/>

Kaggle dataset:

<https://www.kaggle.com/datasets>

WorldBank API:

- <https://blogs.worldbank.org/opendata/introducing-wbgapi-new-python-package-accessing-world-bank-data>
- <https://nbviewer.org/github/tgherzog/wbgapi/blob/master/examples/wbgapi-cookbook.ipynb>
- <https://pypi.org/project/wbgapi/>

GitHub Open Source Code:

<https://github.com/plotly/plotly.py>

