

NULL CLASS INTERNSHIP REPORT

REAL TIME GOOGLE PLAY STORE DATA ANALYTICS - PYTHON

19/12/2024 – 19/2/2025

Live Project URL: <https://elaborate-marzipan-5b573c.netlify.app/>

GitHub Link: <https://github.com/sowmiyaM-1/GooglePlayStoreDataAnalysis>

Introduction:

During my internship, I had the opportunity to dive into data visualization and analysis using Python. The project focused on creating meaningful and interactive visualizations for the Google Play Store dataset. Each task challenged me to apply advanced filtering techniques and dynamic logic to extract insights and present them effectively through dashboards. This hands-on experience helped me enhance both my technical and analytical skills.

Background:

The Google Play Store dataset contains detailed information about apps, such as their categories, installs, revenue, size, and ratings. My internship revolved around understanding this data, identifying patterns, and visualizing insights that could assist decision-making. By working on this project, I explored how data can drive impactful storytelling when presented effectively.

Learning Objectives:

- To develop skills in data visualization and filtering using Python.
- To create time-based and dynamic dashboards.
- To learn how to work with real-world datasets and apply business logic to extract insights.
- To enhance creativity and problem-solving when faced with complex conditions in data analysis.

Activities and Tasks:

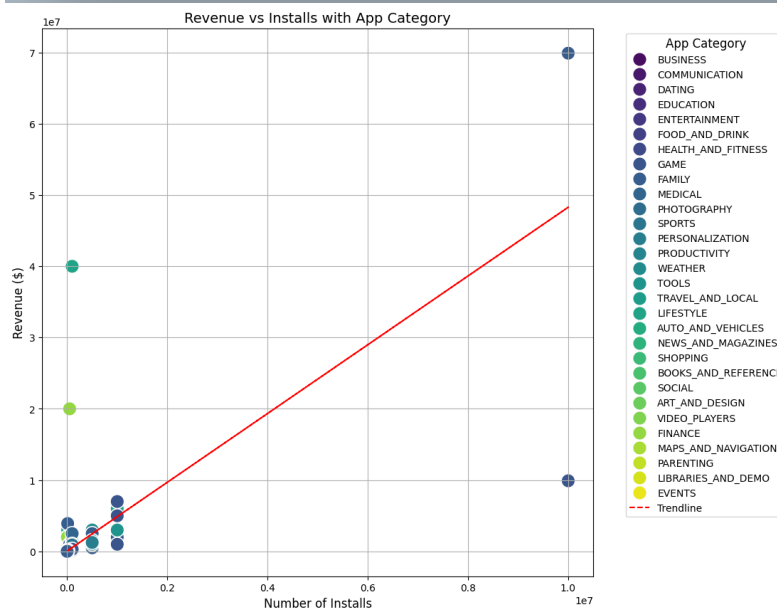
Here's a summary of the main tasks I worked on:

Task 1 Create a scatter plot to visualize the relationship between revenue and the number of installs for paid apps only. Add a trendline to show the correlation and color-coded the points based on app categories.

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5
6
7 data = pd.read_csv('Play Store Data.csv')
8
9
10 # 1. Filtered only paid apps
11 data = data[data['Type'] == 'Paid']
12
13 # 2. Handling of strings like $ in the price column
14 data['Price'] = data['Price'].apply(lambda x: str(x).replace('$', '') if isinstance(x, str) else x).astype(float)
15
16 # 3. Handling of the characters like + into int
17 data['Installs'] = data['Installs'].str.replace(r'[^\d]', '', regex=True).astype(int)
18
19 # 4. calculate revenue= Price* Installs to show the scatterplot
20 data['Revenue'] = data['Price'] * data['Installs']
21
22 # 5. Change the data frame based on our need
23 data = data[['Category', 'Price', 'Installs', 'Revenue']]
24
25 # 6. Scatterplot of Revenue vs Installs
26 plt.figure(figsize=(10, 10))
27 sns.scatterplot(
28     data=data,
29     x='Installs',
30     y='Revenue',
31     hue='Category',
32     palette='viridis',
33     s=100
34 )
35
36 # 7. Trendline for correlation
37 z = np.polyfit(data['Installs'], data['Revenue'], 1) # Linear fit
38 p = np.poly1d(z)
39 plt.plot(data['Installs'], p(data['Installs']), color='red', linestyle='--', label='Trendline')
40
41
42
43 plt.title('Revenue vs Installs with App Category', fontsize=14)
44 plt.xlabel('Number of Installs', fontsize=12)
45 plt.ylabel('Revenue ($)', fontsize=12)
46
47 plt.legend(
48     title='App Category',
49     fontsize=12,
50     title_fontsize=12,
51     bbox_to_anchor=(1.05, 1),
52     loc='upper left'
53 )
54
55 plt.grid(True)
56 plt.show()

```

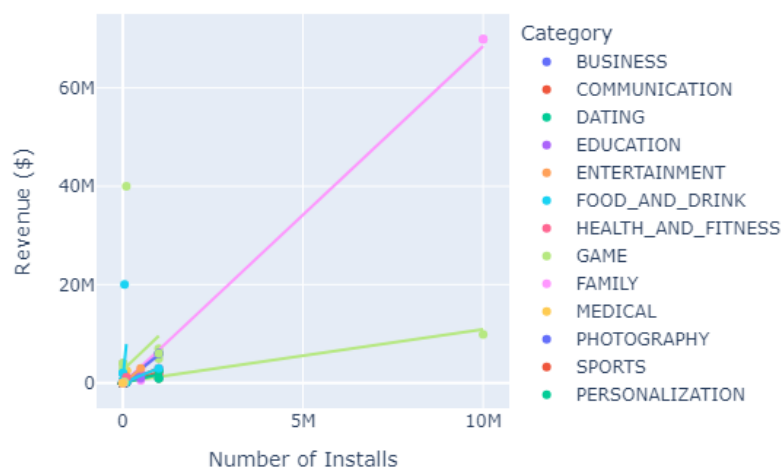


```

1 import pandas as pd
2 import plotly.express as px
3
4 data = pd.read_csv('Play Store Data.csv')
5
6
7 data = data[data['Type'] == 'Paid']
8
9 data['Price'] = data['Price'].apply(lambda x: str(x).replace('$', '') if isinstance(x, str) else x).astype(float)
10
11 data['Installs'] = data['Installs'].str.replace(r'^[d]', '', regex=True).astype(int)
12
13
14 data['Revenue'] = data['Price'] * data['Installs']
15
16
17 data = data[['Category', 'Price', 'Installs', 'Revenue']]
18
19
20 fig = px.scatter(
21     data,
22     x='Installs',
23     y='Revenue',
24     color='Category',
25     title='Revenue vs Installs with App Category',
26     labels={'Installs': 'Number of Installs', 'Revenue': 'Revenue ($)'},
27     hover_data=['Price'],
28     template='plotly'
29 )
30
31
32 fig = px.scatter(
33     data,
34     x='Installs',
35     y='Revenue',
36     color='Category',
37     trendline='ols',
38     title='Revenue vs Installs with App Category',
39     labels={'Installs': 'Number of Installs', 'Revenue': 'Revenue ($)'},
40     hover_data=['Price'],
41     template='plotly'
42 )
43
44
45 fig.show()
46

```

Revenue vs Installs with App Categories



Task 2 Create a dual-axis chart comparing the average installs and revenue for free vs. paid apps within the top 3 app categories.

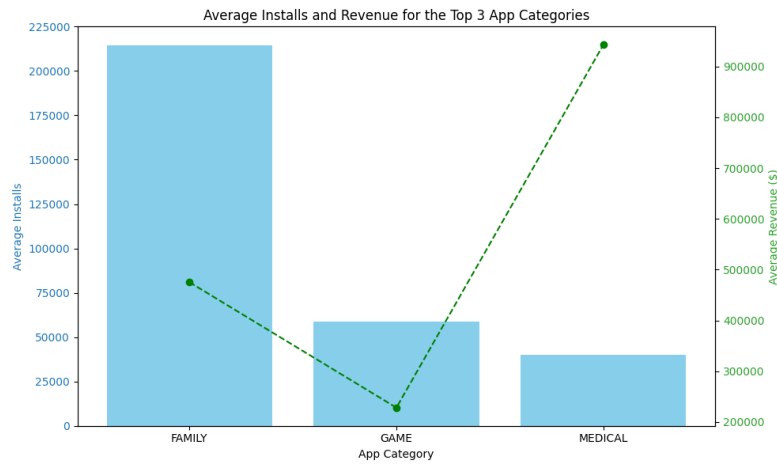
Apply filters - Exclude apps with fewer than 10,000 installs and revenue below \$10,000, android version should be more than 4.0 as well as size should be more than 15M and content rating should be Everyone and app name should not have more than 30 characters including space and special character.

This graph should work only between 1 PM IST to 2 PM IST apart from that time we should not show this graph in dashboard itself.

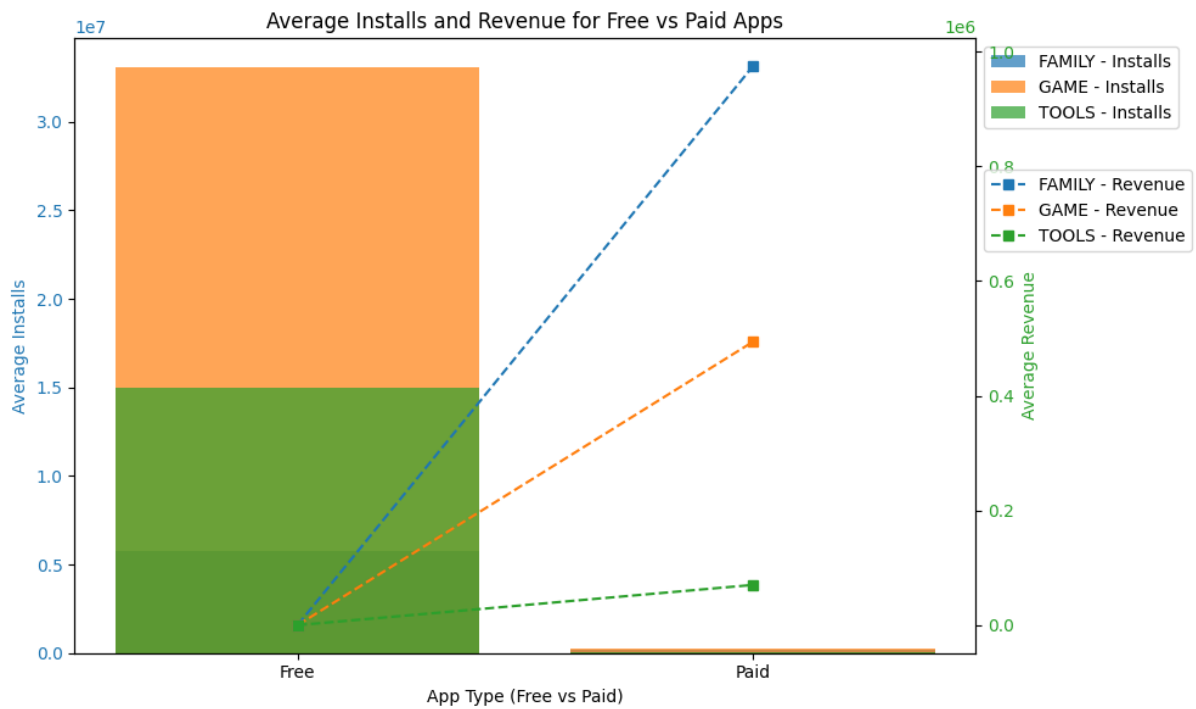
1.Filters and cleaning

```
1 data = pd.read_csv('Play Store Data.csv', encoding='utf-8', on_bad_lines='skip')
2
3 data = data[data['Type'].isin(['Free', 'Paid'])]
4
5 # 2. Handling of string in price value to float
6 data['Price'] = data['Price'].apply(lambda x: str(x).replace('$', '' if isinstance(x, str) else x).astype(float))
7
8 # 3. Handling of the characters like + into int
9 data['Installs'] = data['Installs'].str.replace(r'[^\d]', '', regex=True).astype(int)
10
11 # 4. Calculate revenue = Price * Installs
12 data['Revenue'] = data['Price'] * data['Installs']
13 data['Revenue'] = data['Revenue'].apply(lambda x: int(round(x)))
14
15 # Filter out apps with fewer than 10,000 installs and revenue below 10,000
16 data = data[(data['Installs'] >= 10000) & (data['Revenue'] >= 10000)]
17
18 # 5. Exclude all except M
19 data = data[~data['Size'].isin(['Varies with device'])]
20 data = data[~data['Size'].str.contains('k')]
21 data['Size'] = data['Size'].str.replace('M', '').astype(float)
22
23 # 6. Filter size greater than 15M
24 data = data[data['Size'] > 15]
25
26 # 7. Filter Content Rating = 'Everyone'
27 data = data[data['Content Rating'] == 'Everyone']
28
29 # 8. Define a function to extract the number
30 def extract_version(version_str):
31     try:
32         version = version_str.split(' ')[0]
33         return float(version) if '.' in version else int(version)
34     except:
35         return None
36
37 # 9. Apply the function to the 'Android Ver' column
38 data['Android Ver'] = data['Android Ver'].apply(extract_version)
39
40 # Filter for Android version >= 4.0
41 data = data[data['Android Ver'] >= 4.0]
42
43 # Trimming to get app names < 30 chars
44 data['App'] = data['App'].apply(lambda x: x[:30] if len(x) > 30 else x)
45
46 # 10. Select needed columns for the final output
47 data = data[['Category', 'Type', 'Installs', 'Revenue', 'Size', 'App']]
48
49 # Remove duplicates
50 data = data.drop_duplicates()
51
52 # Print the cleaned data
53 print(data)
```

1.Task2.py



2.Task21.py

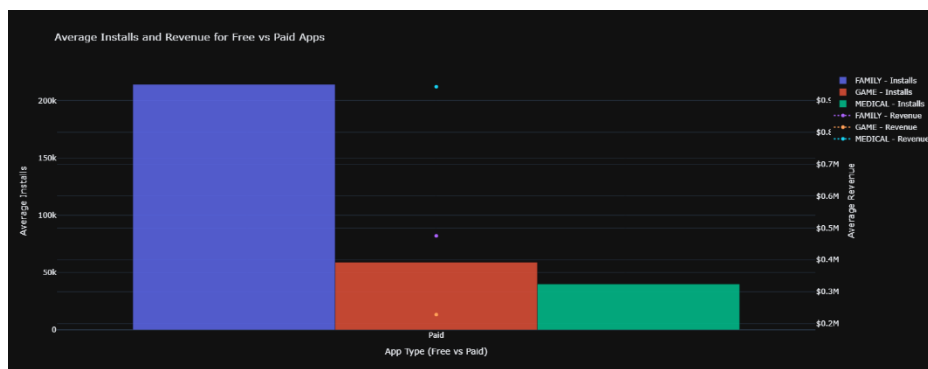


3.Task21Dashboard.py

Dashboard is visible during the timeframe

4.SaveHTMLFile.py → Final code

The plot will be saved as dashboard.html if within the mentioned time frame, else plot is not saved.



```

1  def within_time_window():
2      ist = pytz.timezone('Asia/Kolkata')
3      current_time = datetime.now(ist)
4      return 1 <= current_time.hour < 24
5
6  if within_time_window():
7
8      fig = go.Figure()
9
10     # bar chart for average installs
11     for category in top_categories:
12         category_data = filtered_data[filtered_data['Category'] == category]
13         fig.add_trace(go.Bar(
14             x=category_data['Type'],
15             y=category_data['Installs'],
16             name=f'{category} - Installs',
17             opacity=0.8,
18             yaxis='y1'
19         ))
20
21     # line chart for average revenue
22     for category in top_categories:
23         category_data = filtered_data[filtered_data['Category'] == category]
24         fig.add_trace(go.Scatter(
25             x=category_data['Type'],
26             y=category_data['Revenue'],
27             mode='lines+markers',
28             name=f'{category} - Revenue',
29             line=dict(dash='dot'),
30             yaxis='y2'
31         ))
32
33     fig.update_layout(
34         title='Average Installs and Revenue for Free vs Paid Apps',
35         xaxis_title='App Type (Free vs Paid)',
36         yaxis_title='Average Installs',
37         yaxis2=dict(
38             title='Average Revenue',
39             overlaying='y',
40             side='right',
41             tickprefix='$',
42         ),
43         barmode='group',
44         template='plotly_dark',
45         height=600,
46         showlegend=True,
47     )
48
49     fig.write_html('dashboard.html')
50
51     print("Graph generated and saved as 'dashboard.html'.")
52
53 else:
54     print("This graph is only available between 1 PM and 2 PM IST.")
55

```

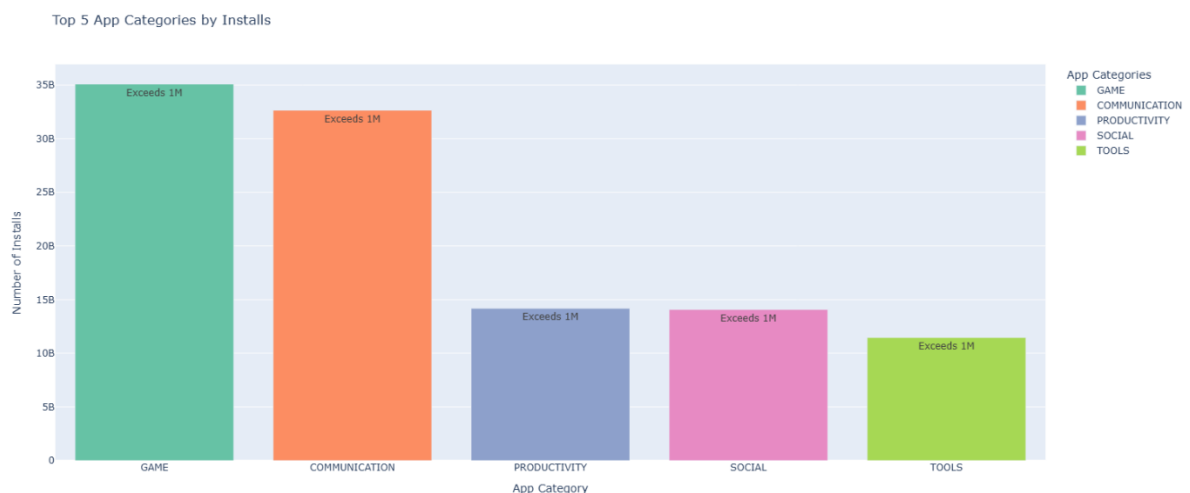
Task 3 Create an interactive choropleth map using Plotly to visualize global installs by categories . Apply filters to show data for only the top 5 app categories and highlight categories where the number of installs exceeds 1 million and App category should not start with character “A” , “C” , “G” and “S” . this graph should work only between 6PM IST to 8 PM IST apart from that time we should not show this graph in dashboard itself.

task3Barchart.py

```

1 import pandas as pd
2 import plotly.express as px
3
4 data = pd.read_csv('C:/Users/sowmi/OneDrive/Documents/AI/Nullclass/Play Store Data.csv', encoding='utf-8', on_bad_lines='skip')
5
6 data = data[data['Price'] != "Everyone"]
7
8 data['Price'] = data['Price'].apply(lambda x: str(x).replace('$', '') if isinstance(x, str) else x).astype(float)
9
10 data['Installs'] = data['Installs'].str.replace(r'[^\d]', '', regex=True).astype(int)
11
12 category_installs = data.groupby('Category')['Installs'].sum().reset_index()
13
14 top_5_categories = category_installs.nlargest(5, 'Installs')
15
16 filtered_data = data[data['Category'].isin(top_5_categories['Category'])]
17
18 top_5_categories['Highlight'] = top_5_categories['Installs'].apply(lambda x: 'Exceeds 1M' if x > 1_000_000 else 'Below 1M')
19
20 fig = px.bar(
21     top_5_categories,
22     x='Category',
23     y='Installs',
24     color='Category',
25     text='Highlight',
26     title='Top 5 App Categories by Installs',
27     color_discrete_sequence=px.colors.qualitative.Set2
28 )
29
30 fig.update_layout(
31     legend_title_text='App Categories',
32     xaxis_title='App Category',
33     yaxis_title='Number of Installs',
34     showlegend=True
35 )
36
37 fig.write_html('top_5_categories.html')
38 print("Bar chart generated and saved as 'top_5_categories.html'.")
39

```

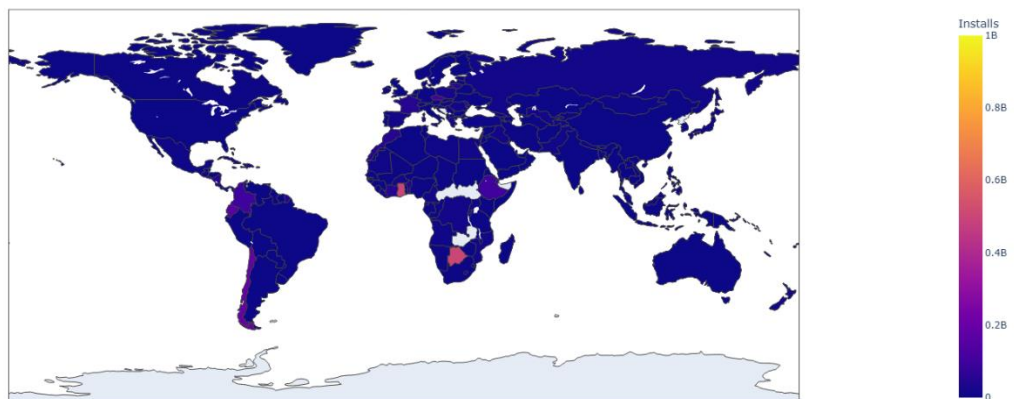


```

1 import pandas as pd
2 import plotly.express as px
3 from datetime import datetime
4 import pytz
5
6 data = pd.read_csv('C:/Users/sowmi/OneDrive/Documents/AI/Nullclass/Play Store Data.csv', encoding='utf-8', on_bad_lines='skip')
7
8 data = data[data['Price'] != "Everyone"]
9
10 data['Price'] = data['Price'].apply(lambda x: str(x).replace('$', '') if isinstance(x, str) else x).astype(float)
11
12
13 data['Installs'] = data['Installs'].str.replace(r'^\d+', '', regex=True).astype(int)
14
15 data['Revenue'] = data['Price'] * data['Installs']
16 data['Revenue'] = data['Revenue'].apply(lambda x: int(round(x)))
17
18 excluded_categories = ['A', 'C', 'G', 'S']
19 data = data[~data['Category'].str.startswith(tuple(excluded_categories))]
20
21 data = data[~data['Size'].isin(['Varies with device'])]
22 data = data[~data['Size'].str.contains('k')]
23 data['Size'] = data['Size'].str.replace('M', '').astype(float)
24
25 data = data[data['Size'] > 1]
26
27 top_5_categories = data.nlargest(5, 'Installs')
28
29 filtered_data = data[data['Category'].isin(top_5_categories['Category'])]
30
31 final_grouped_data = filtered_data.groupby(['Category', 'Country'])['Installs'].sum().reset_index()
32
33 def within_time_window():
34     ist = pytz.timezone('Asia/Kolkata')
35     current_time = datetime.now(ist)
36     return 1 <= current_time.hour < 24
37
38 if within_time_window():
39
40     fig = px.choropleth(
41         final_grouped_data,
42         locations='Country',
43         locationmode='country names',
44         color='Installs',
45         hover_name='Category',
46         title='Global Installs by Top 5 App Categories',
47         color_continuous_scale=px.colors.sequential.Plasma
48     )
49
50     fig.write_html('dashboard3.html')
51     print("Graph generated and saved as 'dashboard3.html'.")
52 else:
53     print("This graph is only available between 1 PM and 2 PM IST.")
54
55

```

Global Installs by Top 5 App Categories

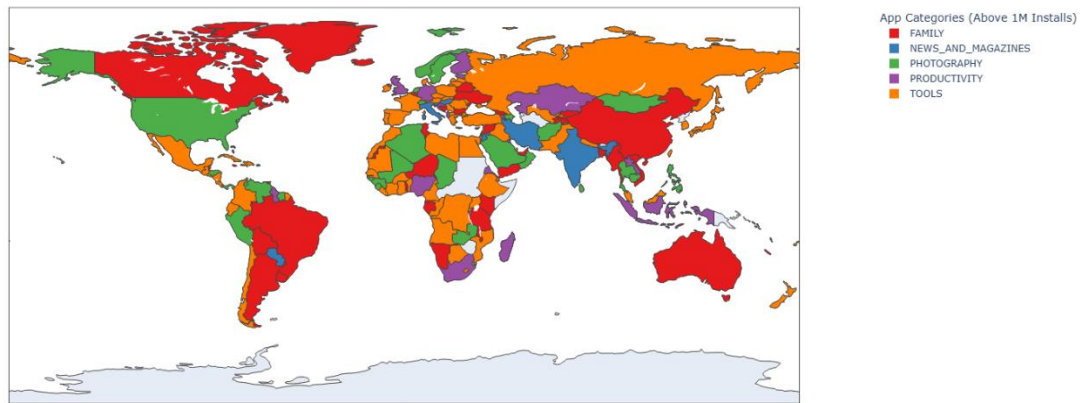



```

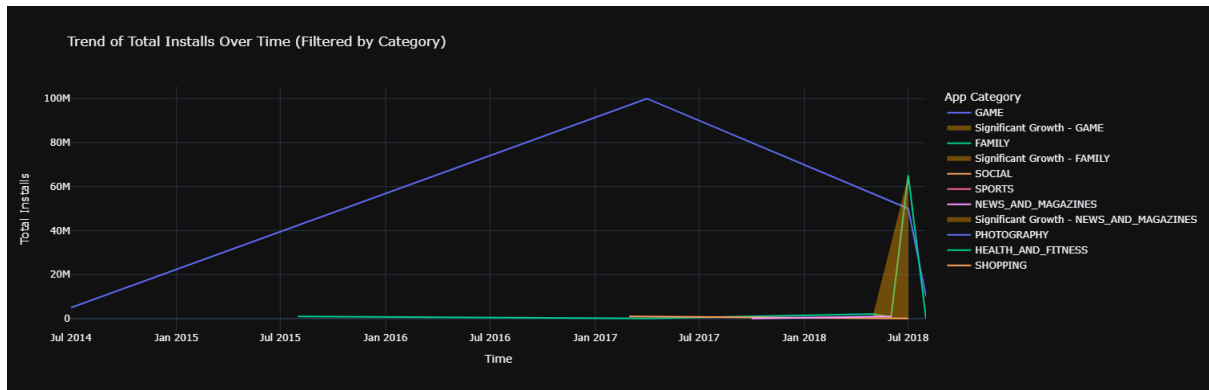
1 import pandas as pd
2 import plotly.express as px
3 from datetime import datetime
4 import pytz
5
6 # Load the data
7 data = pd.read_csv('C:/Users/sowmi/OneDrive/Documents/AI/Nullclass/Play Store Data.csv', encoding='utf-8', on_bad_lines='skip')
8
9 # Data Cleaning
10 data = data[data['Price'] != "Everyone"]
11
12 data['Price'] = data['Price'].apply(lambda x: str(x).replace('$', '') if isinstance(x, str) else x).astype(float)
13
14
15 data['Installs'] = data['Installs'].str.replace(r'[^\d]', '', regex=True).astype(int)
16
17
18 data['Revenue'] = data['Price'] * data['Installs']
19 data['Revenue'] = data['Revenue'].apply(lambda x: int(round(x)))
20
21
22 ex_categories = ['A', 'C', 'G', 'S']
23 data = data[~data['Category'].str.startswith(tuple(ex_categories))]
24
25 data = data[~data['Size'].isin(['Varies with device'])]
26 data = data[~data['Size'].str.contains('k')]
27 data['Size'] = data['Size'].str.replace('M', '').astype(float)
28
29
30 data = data[data['Size'] > 15]
31
32
33 top_5_categories = data.groupby('Category')['Installs'].sum().nlargest(5).reset_index()
34
35
36 data = data[data['Category'].isin(top_5_categories['Category'])]
37
38
39 data['Highlight'] = data['Installs'].apply(lambda x: 'Above 1M Installs' if x > 1_000_000 else 'Other')
40
41
42 data1 = data[data['Highlight'] == 'Above 1M Installs']
43
44
45 final_data = data1.groupby(['Category', 'Country'])['Installs'].sum().reset_index()
46
47 def within_time_window():
48     ist = pytz.timezone('Asia/Kolkata')
49     current_time = datetime.now(ist)
50     return 11 <= current_time.hour < 20
51
52 if within_time_window():
53
54     fig = px.choropleth(
55         final_data,
56         locations='Country',
57         locationmode='country names',
58         color='Category',
59         hover_name='Category',
60         hover_data={
61             'Category': True,
62             'Installs': True
63         },
64         title='Global Installs by Top Categories Above 1 Million',
65         color_discrete_sequence=px.colors.qualitative.Set1
66     )
67
68     # Add a legend and save the graph
69     fig.update_layout(
70         legend_title_text='App Categories (Above 1M Installs)',
71         showlegend=True
72     )
73
74     # Save the graph as an HTML file
75     fig.write_html('dashboard31.html')
76     print("Graph generated and saved as 'dashboard31.html'.")
77 else:
78     print("This graph is only available between 1 PM and 2 PM IST.")
79

```

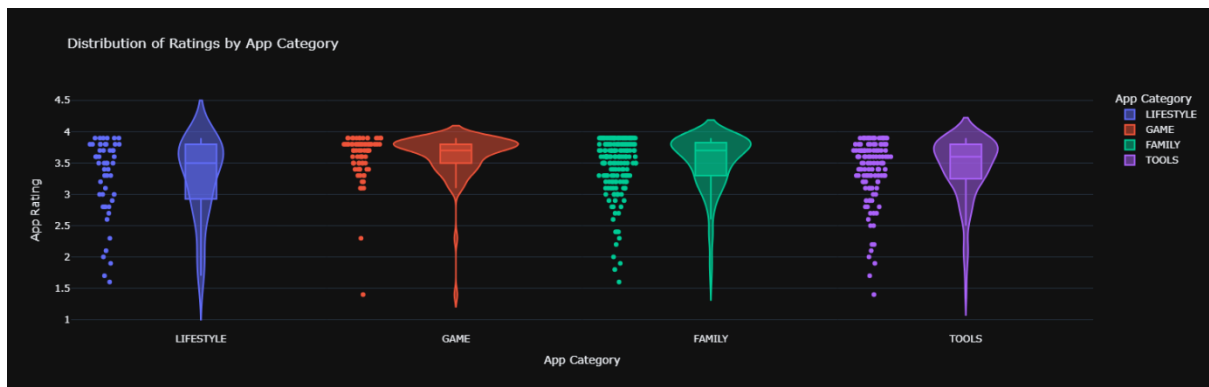
Global Installs by Top Categories Above 1 Million



Task 4 Plot a time series line chart to show the trend of total installs over time, segmented by app category. Highlight periods of significant growth by shading the areas under the curve where the increase in installs exceeds 20% month-over-month and content rating should be teen and app name should start with letter 'E' and installs should be more than 10k as well as this graph should work only between 6 PM IST to 9 PM IST apart from that time we should not show this graph in dashboard itself.



Task 5 Create a violin plot to visualize the distribution of ratings for each app category, but only include categories with more than 50 apps and app name should contain letter "C" and exclude apps with fewer than 10 reviews and rating should be less 4.0. this graph should work only between 4 PM IST to 6 PM IST apart from that time we should not show this graph in dashboard itself.



Challenges and Solutions:

1: Filtering data based on multiple conditions.

- I used pandas and advanced filtering techniques, combining multiple logical conditions seamlessly.

2: Restricting graphs to specific time windows.

- By integrating Python's datetime module, I ensured that graphs were displayed only during specified hours.

Outcomes and Impact:

This internship gave me a deeper understanding of how data visualization can drive meaningful insights. I successfully delivered interactive and visually engaging dashboards that met all business requirements. Beyond the technical skills, I learned to think critically and creatively when solving data-related challenges.

Conclusion:

Overall, this internship was a learning experience. It pushed me to explore new tools, apply advanced concepts, and find solutions to complex problems. I leave this internship feeling more confident in my ability to analyse and visualize data effectively, skills that I am excited to carry forward in future projects.