<ul> <li>IMDB Top 250 Movies Dataset</li> <li>IMDB (Internet Movie Database) is one of the largest online databases for movies and television shows, providing comprehensive information about movies, including r</li> <li>This dataset contains the top 250 rated movies on IMDB as of 2021, providing a snapshot of the most popular and highly rated movies of recent times. By analyzing this</li> </ul>	
<pre>import pandas as pd import numpy as np import seaborn as sns import os import warnings import matplotlib.pyplot as plt import plotly.graph_objects as go import plotly.express as px import plotly.graph_objs as go  warnings.filterwarnings('ignore') sns.set(style="whitegrid")</pre>	
Loading data  In [2]: folder = './datasets' for file in os.listdir(folder):     print(file)  IMDB-Top-250-Movies.csv  In [3]: df = pd.read_csv('datasets/IMDB-Top-250-Movies.csv') df.head()	
value[3]:         rank         mame         year         rating         genre         certificate         run_time         fear can hold you prisoner. Hope can set you f         25000000         28884504           1         2         The Godfather         1972         9.2         Crime,Drama         R         2h 32m         An offer you can't refuse.         6000000         250341816           2         3         4         The Godfather Part II         1974         9.0         Crime,Drama         R         3h 22m         All the power on earth can't change destiny.         13000000         47961919	Tim Robbins,Morgan Freeman,Bob Gunton,William Frank Darabont Marlon Brando,Al Pacino,James Caan,Diane Keato Francis Ford Coppola Christian Bale,Heath Ledger,Aaron Eckhart,Mich Christopher Nolan Jonathan Nolan,Christopher Nolan,David S. Goyer
<pre><class 'pandas.core.frame.dataframe'=""> RangeIndex: 250 entries, 0 to 249 Data columns (total 13 columns): # Column Non-Null Count Dtype</class></pre>	
7 tagline 250 non-null object 8 budget 250 non-null object 9 box_office 250 non-null object 10 casts 250 non-null object 11 directors 250 non-null object 12 writers 250 non-null object dtypes: float64(1), int64(2), object(10) memory usage: 25.5+ KB  These are what the column names represent:  • rank: The rank of the movie based on its popularity or any other criteria.	
<ul> <li>name: The name of the movie.</li> <li>year: The year in which the movie was released.</li> <li>rating: The rating of the movie on IMDb's scale of 1 to 10.</li> <li>genre: The category or type of movie, such as action, comedy, drama, etc.</li> <li>certificate: The rating given to the movie by the relevant certification board, which indicates the age appropriateness of the movie and may contain information all run_time: The duration of the movie.</li> <li>tagline: The phrase or sentence used to promote the movie.</li> <li>budget: The estimated cost of making the movie.</li> <li>box_office: The amount of revenue generated by the movie at the box office.</li> <li>casts: The actors who appear in the movie.</li> <li>directors: The person who directed the movie.</li> <li>writers: The people who wrote the screenplay of the movie.</li> </ul>	about the content.
count 250.00000 250.00000 250.00000 250.00000 mean 125.50000 1986.360000 8.307200  std 72.312977 25.125356 0.229081  min 1.000000 1921.000000 8.000000  25% 63.250000 1966.250000 8.100000  50% 125.500000 1994.000000 8.200000  75% 187.750000 2006.000000 8.400000	
<pre>max 250.00000 2022.00000 9.300000  in [7]: df.isnull().sum()  put[7]: rank</pre>	
box_office 0 casts 0 directors 0 writers 0 dtype: int64  In [8]: # splits each cell with genres by separator genres = df['genre'].str.split(',', expand=True).stack().reset_index(level=1, drop=True).rename('director')  # the list of genres with their respective counts list_of_genres = genres.value_counts() list_of_genres  dut[8]: director	
Drama       177         Adventure       60         Crime       51         Action       50         Comedy       45         Mystery       31         Thriller       30         Biography       29         Romance       23         War       23         Animation       23         Sci-Fi       20         Fantasy       14         Family       13	
History 10 Western 7 Horror 5 Sport 5 Music 4 Film-Noir 4 Musical 1 Name: count, dtype: int64  Data Visualization	
Most popular movies across different genres  plt.figure(facecolor = "c", figsize = (6, 6))  ax = sns.barplot(y = genres.value_counts().index, x = genres.value_counts().values)  ax.set(xlabel='Number of Movies',ylabel='Genre', title='Most Popular Movie Genres')  ax.set_facecolor("xkcd:eggshell")  ax.bar_label(ax.containers[0], fontsize=10)  plt.show()  Most Popular Movie Genres	
Adventure	
plt.xlabel('Year') plt.ylabel('Number of Movies Released') plt.title('Average Number of Movies over Time')  plt.show()  Average Number of Movies over Time  8  7	
Average Rating over Time  [ii]: rating_year_counts = df.groupby('year')['rating'].mean() plt.figure(facecolor = "xkcd:eggshell", figsize = (6, 6))	
plt.ylabel('Average raiting') plt.ylabel('Average Rating over Time') plt.show()  Average Rating over Time  9.2  9.0  8.8  8.8  8.8  8.4	
Top 10 Movies with the Longest Run Times  For now, the runtime format is not suitable for visualization:  [12]: df[['run_time']]  [12]: run_time 0 2h 22m 1 2h 55m 2 2h 32m 3 3h 22m 4 1h 36m	
245 2h 26m 246 2h 22m 247 1h 30m 248 3h 11m 249 3h 1m 250 rows × 1 columns  Let's apply the function bellow and convert it to minutes:  [13]: def convert_to_mins(time_string):     time_list = time_string.split()     hours = 0     minutes = 0     for time_segment in time_list:         if 'h' in time_segment:	
<pre>hours = int(time_segment[:-1]) elif 'm' in time_segment:     minutes = int(time_segment[:-1]) total_minutes = hours * 60 + minutes     return total_minutes  df['run_time'] = df['run_time'].apply(convert_to_mins)  [14]:</pre>	
1 175 2 152 3 202 4 96 245 146 246 142 247 90 248 191 249 181	
top_10_budget = df.nlargest(10, 'run_time') top_10_budget[['name', 'run_time']]  158	
3 The Godfather Part II 202 6 The Lord of the Rings: The Return of the King 201 5 Schindler's List 195 248 Gandhi 191 26 The Green Mile 189  [16]: plt.figure(facecolor = 'c', figsize = (8, 4)) ax = sns.barplot(x='run_time', y='name', data=top_10_budget, palette='mako') plt.title('Top 10 Movies by Writers') plt.xlabel('Duraction in Minutes')	
## Pit.ylabel('Name of the Movie') ax.set_facecolor("xkcd:eggshell") ax.se	
Top 10 Movies by Rating  top_10_rating = df.nlargest(10, 'rating') top_10_rating[['name', 'rating']]  7	
plt.ylabel('Movie Title') ax.set_Tacecolor('mkcd:eggshell'') ax.bar_label(ax.containers[0], fontsize=10) plt.show()  Top 10 Movies by Rating  The Shawshank Redemption The Godfather The Dark Knight The Godfather Part II 12 Angry Men Schindler's List The Lord of the Rings: The Return of the King Pulp Fiction The Lord of the Rings: The Fellowship of the Ring The Good, the Bad and the Ugly  Top 10 Moviee by Writers:	
top_10_writers = df.nlargest(10, 'rating') top_10_writers[('writers', 'rating')]    Writers	
plt.ylabel('Writers Name') ax. set_facecolor("xkcd:eggshel1") ax. bar_label(ax.containers[0], fontsize=12) plt.show()  Top 10 Movies by Writers  Stephen King,Frank Darabont Mario Puzo,Francis Ford Coppola Jonathan Nolan,Christopher Nolan,David S. Goyer Francis Ford Coppola,Mario Puzo Reginald Rose Reginald Rose Thomas Keneally,Steven Zaillian J.R.R. Tolkien,Fran Walsh,Philippa Boyens Quentin Tarantino,Roger Avary Luciano Vincenzoni,Sergio Leone,Agenore Incrocci  8.8  0 2 4 6 8  Rating	
Directors with the highest number of movies  df_directors = df['directors'].str.split(',', expand=True).stack().reset_index(level=1, drop=True).rename('Director Name') df_top_directors = df_directors.value_counts().nlargest(10)  Director Name Christopher Nolan	
plt. ylabel('Number of Movies') plt. xticks(rotation=45) # Rotates x-axis labels by 90 degrees plt. show()  Top Directors by Number of Movies  7  6  5  9  10  10  10  10  10  10  10  10  10	
Oldest movies on the charts  Oldest_movies=df.sort_values(by='year')[0:10]  fig = go.Figure(data=[go.Table(header=dict(values=['Movie Name', 'Release Year'],fill_color='cyan'),	
Movie Name  The Kid Sherlock Jr.  The Gold Rush The General Metropolis The Passion of Joan of Arc M City Lights It Happened One Night Modern Times	Release Year  1921  1924  1925  1926  1927  1931  1931  1934  1936
Latest movies on the chart  oldest_movies=df.sort_values(by='year', ascending=False)[0:15]  fig = go.Figure(data=[go.Table(header=dict(values=['Movie Name', 'Release Year'],fill_color='cyan'),	
Movie Name  Top Gun: Maverick  Spider-Man: No Way Home  Jai Bhim  The Father  Hamilton  1917  Parasite	Release Year  2022 2021 2021 2021 2020 2020 2019 2019
Parasite Ford v Ferrari  Avengers: Endgame Klaus Joker  Capernaum Green Book Avengers: Infinity War  Spider-Man: Into the Spider-Verse	2019 2019 2019 2019 2019 2019 2018 2018 2018 2018