Exploratory Data Analysis of IMDB Movies Dataset using Python In [1]: import warnings warnings.filterwarnings('ignore') In [2]: import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns In [12]: # Read the IMDB movies data movies = pd.read csv("Movie Assignment Data.csv") pd.set_option('display.max_columns',70) movies.head() Out[12]: Title title_year budget Gross actor_1_name actor_2_name actor_3_name actor_1_facebook_likes actor_2_facebook_likes La La Land 2016 30000000 151101803 Ryan Gosling Emma Stone Amiée Conn 14000 1900 Ginnifer Jason Zootopia 2016 150000000 341268248 Idris Elba 2800 2800 Goodwin Bateman 2 Lion 2016 12000000 51738905 Dev Patel Nicole Kidman Rooney Mara 33000 9600 Jeremy Forest 3 Arrival 47000000 100546139 35000 530 2016 Amy Adams Whitaker Renner Manchester Michelle 518 7100 2016 9000000 47695371 Casey Affleck Kyle Chandler by the Sea Williams #understanding the structure of the data In [4]: movies.shape Out[4]: (100, 62) In [5]: movies.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 100 entries, 0 to 99 Data columns (total 62 columns): # Column Non-Null Count Dtype Title 100 non-null object 100 non-null int64 title_year 100 non-null int64 100 non-null int64 budget Gross 4 actor_1_name 100 non-null object 5 actor_2_name 100 non-null object 6 actor_3_name 100 non-null object 7 actor 1 facebook likes 100 non-null int64 8 actor_2_facebook_likes 99 non-null float64 9 actor_3_facebook_likes 98 non-null float64 10 IMDb_rating 100 non-null float64 object 11 genre 1 100 non-null 12 genre_2 97 non-null object 13 genre_3 74 non-null object 14 MetaCritic 95 non-null float64 100 non-null int64 15 Runtime 16 CVotes10 100 non-null int64 17 CVotes09 100 non-null int64 18 CVotes08 int64 100 non-null 19 CVotes07 100 non-null int64 20 CVotes06 100 non-null int64 CVotes05 21 100 non-null int64 22 CVotes04 100 non-null int64 CVotes03 23 100 non-null CVotes02 100 non-null int64 25 CVotes01 100 non-null int.64 26 CVotesMale 100 non-null int64 27 CVotesFemale 100 non-null int64 100 non-null 28 CVotesU18 int64 100 non-null 29 CVotesU18M int64 100 non-null 30 CVotesU18F int64 31 CVotes1829 100 non-null 32 CVotes1829M 100 non-null int64 33 CVotes1829F 100 non-null int64 100 non-null 34 CVotes3044 int64 35 CVotes3044M 100 non-null int64 100 non-null CVotes3044F 36 int64 37 CVotes45A 100 non-null int64 100 non-null 38 CVotes45AM int64 100 non-null 39 CVotes45AF 40 CVotes1000 100 non-null int64 41 CVotesUS 100 non-null int64 42 CVotesnUS 100 non-null int64 43 VotesM 100 non-null float64 44 VotesF 100 non-null float64 45 VotesU18 100 non-null float64 46 VotesU18M 100 non-null float64 47 VotesU18F 100 non-null float64 48 Votes1829 100 non-null float64 49 Votes1829M 100 non-null float64 100 non-null 50 Votes1829F float64 51 Votes3044 100 non-null float64 52 Votes3044M 100 non-null float64 53 Votes3044F 100 non-null float64 54 Votes45A 100 non-null float64 55 Votes45AM 100 non-null float64 56 Votes45AF 100 non-null float64 57 Votes1000 100 non-null float64 58 VotesUS 100 non-null float64 59 VotesnUS 100 non-null float64 60 content_rating 100 non-null object 61 Country 100 non-null object dtypes: float64(21), int64(32), object(9) memory usage: 48.6+ KB In [8]: #summary statistics of numerical columns movies.describe(include="all") Out[8]: Title title_year budget Gross actor_1_name actor_2_name actor_3_name actor_1_facebook_likes actor_2_face 1.000000e+02 1.000000e+02 count 100 100.000000 100 100 100 100.000000 75 94 96 unique 100 NaN NaN NaN NaN La Leonardo Tom Hardy Mackenzie Foy La NaN NaN NaN NaN top DiCaprio Land freq NaN NaN NaN 5 2 NaN 2012.820000 NaN 7.838400e+07 1.468679e+08 NaN NaN NaN 13407.270000 mean std NaN 1.919491 7.445295e+07 1.454004e+08 NaN NaN NaN 10649.037862 13 2010.000000 3.000000e+06 2.238380e+05 NaN NaN NaN 39.000000 min NaN NaN NaN 1000.000000 25% NaN 2011.000000 1.575000e+07 4.199752e+07 NaN 2013.000000 50% NaN 4.225000e+07 1.070266e+08 NaN NaN NaN 13000.000000 1 1.500000e+08 2.107548e+08 NaN 20000.000000 75% NaN 2014.000000 NaN NaN 11 NaN 2016.000000 2.600000e+08 9.366622e+08 NaN NaN NaN 35000.000000 96 #converting budget and gross columns from \$ to million \$ for better readability In [13]: movies['budget'] = movies['budget'].floordiv(1000000) movies['Gross'] = movies['Gross'].floordiv(1000000) movies.head() In [14]: Out[14]: actor_2_name actor_3_name actor_1_facebook_likes actor 2 facebook likes Title title_year budget Gross actor_1_name La La Land 2016 30 151 Ryan Gosling Emma Stone Amiée Conn 14000 19000.0 Ginnifer Jason 1 Zootopia 2016 150 341 Idris Elba 2800 28000.0 Goodwin Bateman Dev Patel 2 2016 12 Nicole Kidman 33000 96000.0 Lion 51 Rooney Mara **Forest** Jeremy 3 35000 Arrival 2016 47 100 Amy Adams 5300.0 Whitaker Renner Manchester Michelle 71000.0 2016 47 Casey Affleck Kyle Chandler 518 Williams by the Sea **Profit** create a new column Profit = Gross - budget sort the data according to profit and finding the top 10 profitable movies In [15]: movies.insert(4, column = 'Profit', value = movies['Gross'] - movies['budget']) In [16]: movies.head() Out[16]: Title title_year budget **Gross Profit** actor_1_name actor_2_name actor_3_name actor_1_facebook_likes actor_2_facebook_lik La La Land Ryan Gosling 0 2016 30 151 121 Emma Stone Amiée Conn 14000 19000 Ginnifer Jason Zootopia 2016 150 341 191 Idris Elba 2800 28000 Goodwin Bateman Dev Patel 2 Lion 2016 12 51 39 Nicole Kidman Rooney Mara 33000 96000 Jeremy Forest 3 Arrival 2016 35000 5300 47 100 53 Amy Adams Renner Whitaker Manchester Michelle 2016 9 47 Kyle Chandler 518 71000 38 Casey Affleck by the Sea Williams movies.sort values(by=['Profit'], ascending=False, inplace = True, ignore index = True) In [17]: top10 movies = movies.loc[0:9,'Title'] top10 movies Out[17]: 0 Star Wars: Episode VII - The Force Awakens The Avengers 2 Deadpool 3 The Hunger Games: Catching Fire 4 Toy Story 3 5 The Dark Knight Rises The Lego Movie 6 7 Zootopia 8 Despicable Me Inside Out Name: Title, dtype: object In [20]: plt.figure(figsize=[8,6]) h = sns.jointplot(x='budget', y='Profit', data = movies) h.set_axis_labels('Budget','Profit',fontsize=12) plt.show() <Figure size 576x432 with 0 Axes> 700 600 500 400 Profit 300 200 100 -100100 200 250 Budget Inference Most of the movies are profitable, there is a slight positive trend between 'budget' and 'profit' There are some big budget movies with negative profits. • between a budget of 0 - 100 million, there are a lot of profitable movies with upto 300 million profit. **Worst Performing Movies (Negative Profit)** negative profit = movies[movies['Profit']<0]</pre> In [21]: negative profit.reset index(drop = True) #Worst performing 5 movies negative_profit.sort_values(by='Profit')[0:5] Out[21]: title_year budget Gross Profit actor_1_name actor_2_name actor_3_name actor_1_facebook_likes actor_2_facebook_like Chloë Grace Christopher 99 2011 170 73 -97 Ray Winstone 17000 16000 Hugo Moretz The Little 2015 98 81 1 -80 Jeff Bridges Mackenzie Foy 12000 11000 James Franco Prince Edge of 97 2014 100 Tom Cruise 10000 178 -78 Lara Pulver Noah Taylor 854 Tomorrow Tangled 96 2010 260 200 -60 Brad Garrett Donna Murphy M.C. Gainey 799 553 Scott Anna Kendrick 1000 95 2010 60 31 10000 -29 Kieran Culkin Ellen Wong vs. the World **Movie Ratings Analysis** • MetaCritic (Critics opinion) rating and IMDB rating (public opinion) · Average overall rating of a movie will be the average of both the ratings # MetaCritic rating is out of 100, so to make it into the same scale as IMDB rating In [22]: movies['MetaCritic'] = movies['MetaCritic']/10 In [23]: movies.insert(18,column='Avg Rating', value = (movies['MetaCritic']+movies['IMDb_rating'])/2) In [24]: movies.head() Out [24]: Title title_year budget Gross Profit actor_1_name actor_2_name actor_3_name actor_1_facebook_likes actor_2_facebook_likes Star Wars: Episode Doug Walker 2015 245 936 691 Rob Walker 131 12.0 VII - The Force **Awakens** The Chris Robert Scarlett 2012 220 623 403 26000 21000.C Avengers Hemsworth Johansson Downey Jr. Ryan Stefan Kapicic 2 Deadpool 2016 58 363 305 Ed Skrein 16000 805.C Reynolds The Hunger Jennifer Josh Sandra Ellis 2013 130 424 294 34000 14000.C Games: Lawrence Hutcherson Lafferty Catching Fire Toy Story John 2010 200 15000 1000.C 414 214 Tom Hanks Don Rickles Ratzenberger In [25]: #sorting movies by their avg rating (descending) movies.sort values(by='Avg Rating', ascending=False, inplace=True) movies.head() Out[25]: Title title_year budget Gross Profit actor_1_name actor_2_name actor_3_name actor_1_facebook_likes actor_2_facebook_li Lorelei Ellar Coltrane 69 Boyhood 2014 4 25 21 Libby Villari 230 19 Linklater 12 Years a Quvenzhané 51 2013 20 56 36 Scoot McNairy Taran Killam 2000 66 Slave Wallis Inside Out 2015 175 356 181 Amy Poehler Phyllis Smith 1000 76 Mindy Kaling John Toy Story 3 2010 Don Rickles 15000 200 414 214 Tom Hanks 100 Ratzenberger Michelle Manchester 49 2016 9 47 38 Kyle Chandler 7100 Casey Affleck 518 by the Sea Williams **Most Popular Actor Trios** popoularity based on facebook likes of lead actors in the movie In [26]: movies.insert(11,column='Popularity', value=movies.loc[:,['actor 1 facebook likes','actor 2 facebook li kes' ,'actor 3 facebook likes']].sum(axis=1)) In [27]: most popular actors = movies.sort values(by='Popularity', ascending=False, ignore index=True).loc[0:4,:] most_popular_actors Out[27]: Title title_year budget Gross Profit actor_1_name actor_3_name actor_1_facebook_likes actor_2_facebook_lik actor_2_name 0 Lion 2016 12 51 39 Dev Patel Nicole Kidman Rooney Mara 33000 96000 Leonardo Joseph 1 Inception 2010 160 292 132 Tom Hardy 29000 27000 Gordon-Levitt DiCaprio X-Men: Jennifer Peter Dinklage Days of 2014 200 233 33 34000 22000 Hugh Jackman Lawrence **Future Past** Manchester Michelle 9 2016 47 38 Casey Affleck Kyle Chandler 518 71000 Williams by the Sea The Dark Joseph 27000 23000 2012 250 448 198 Christian Bale Knight Tom Hardy Gordon-Levitt Rises In [29]: most popular actor trios = movies.sort values(by='Popularity', ascending=False, ignore index=True).loc[0] :4,['actor 1 name','actor 2 name','actor 3 name']].values.tolist() most popular actor trios Out[29]: [['Dev Patel', 'Nicole Kidman', 'Rooney Mara'], ['Leonardo DiCaprio', 'Tom Hardy', 'Joseph Gordon-Levitt'], ['Jennifer Lawrence', 'Peter Dinklage', 'Hugh Jackman'], ['Casey Affleck', 'Michelle Williams ', 'Kyle Chandler'], ['Tom Hardy', 'Christian Bale', 'Joseph Gordon-Levitt']] In [30]: #none of the three actors' Facebook likes should be less than half of the other two individual_popular_trio = movies[~(((movies['actor_1_facebook_likes'] < movies['actor_2_facebook_like</pre> **s'**]/2)| (movies['actor_1_facebook_likes'] < movies['actor_3_facebook_like</pre> **s'**]/2))==**True**) |(((movies['actor_2_facebook_likes'] < movies['actor_1_facebook_like</pre> **s'**]/2)| (movies['actor_2_facebook_likes'] < movies['actor_3_facebook_like</pre> **s'**]/2))==**True**) | (((movies['actor_3_facebook_likes'] < movies['actor_1_facebook_lik</pre> es']/2)| (movies['actor_3_facebook_likes'] < movies['actor_2_facebook_like</pre> **s'**]/2))==**True**))] In [31]: | #common in the two dfs common_actors = most_popular_actors[most_popular_actors['Title'].isin(individual_popular_trio['Title']) common_actors.loc[:,['actor_1_name','actor_2_name','actor_3_name']].values.tolist() Out[31]: [['Leonardo DiCaprio', 'Tom Hardy', 'Joseph Gordon-Levitt'], ['Jennifer Lawrence', 'Peter Dinklage', 'Hugh Jackman'], ['Tom Hardy', 'Christian Bale', 'Joseph Gordon-Levitt']] In [32]: top_popular = individual_popular_trio.sort_values(by='Popularity', ascending=False)[0:5].reset_index(dro p=**True**) top_popular.loc[:,['actor_1_name','actor_2_name','actor_3_name']].values.tolist() Out[32]: [['Leonardo DiCaprio', 'Tom Hardy', 'Joseph Gordon-Levitt'], ['Jennifer Lawrence', 'Peter Dinklage', 'Hugh Jackman'], ['Tom Hardy', 'Christian Bale', 'Joseph Gordon-Levitt'], ['Chris Hemsworth', 'Robert Downey Jr.', 'Scarlett Johansson'], ['Philip Seymour Hoffman', 'Robin Wright', 'Brad Pitt']] Movie Run-time Analysis In [33]: #movie run time histogram density plot plt.figure(figsize=[8,4]) runtime = sns.distplot(movies['Runtime']) runtime.axes.set_title("Movie Run-time distribution", fontsize=20) runtime.set xlabel("Runtime", fontsize=15) plt.show() Movie Run-time distribution 0.030 0.025 0.020 0.015 0.010 0.005 0.000 80 100 120 140 160 180 200 Runtime Most movies have a run-time between 120 to 130 minutes **Movie Genre Analysis** In [35]: #creating dataframe by genre columns = []for i in movies.columns: if i.startswith('CV')|i.startswith('V')|i.startswith("ge"): columns.append(i) df_genre = movies.loc[:,columns] df genre['cnt'] = 1 #grouping movies by individual genre df_g1 = df_genre.groupby('genre_1') df_g2 = df_genre.groupby('genre_2') df g3 = df genre.groupby('genre 3') #finding the sum for genre1 df g1 = pd.DataFrame(df g1.sum()) df g2 = pd.DataFrame(df g2.sum()) df g3 = pd.DataFrame(df g3.sum()) In [36]: add_df1_df2 = df_g1.add(df_g2,fill_value=0) df add=add df1 df2.add(df g3,fill value=0) genre_top10 = df_add[df_add['cnt'] >=10] genre top10.head() Out[36]: CVotes10 CVotes09 CVotes08 CVotes07 CVotes06 CVotes05 CVotes04 CVotes03 CVotes02 CVotes01 CVotesMale CVot **Action** 3166467.0 3547429.0 4677755.0 2922126.0 1075354.0 393484.0 166970.0 95004.0 65573.0 171247.0 10837034.0 **Adventure** 3594659.0 4014192.0 5262328.0 3281981.0 1212075.0 11759815.0 438970.0 183070.0 103318.0 69737.0 173858.0 **Animation** 681562.0 798227.0 1153214.0 722782.0 251076.0 83069.0 30718.0 15733.0 10026.0 25193.0 2282985.0 **Biography** 852003.0 1401608.0 2231078.0 1332980.0 425595.0 138648.0 53718.0 29510.0 20613.0 51297.0 4329471.0 Comedy 1383616.0 1774987.0 2506851.0 1591069.0 600287.0 226852.0 97469.0 56218.0 39391.0 88367.0 5223033.0 In [37]: #extracting count column count column= genre top10['cnt'] #taking the mean of every column genre top10=genre top10.div(genre top10['cnt'],axis='index') genre_top10=genre_top10.apply(lambda x:round(x,2),axis=0) genre top10.drop(columns='cnt',inplace=True) In [38]: #convert CVotes into int col cvotes=[] for i in genre top10.columns: if i.startswith('CVotes'): col_cvotes.append(i) genre_top10[col_cvotes]=genre_top10[col_cvotes].astype('int') genre_top10.head() Out[38]: CVotes10 CVotes09 CVotes08 CVotes07 CVotes06 CVotes05 CVotes04 CVotes03 CVotes02 CVotes01 CVotes01 CVotesMale CVotes 12693 349581 Action 102144 114433 150895 94262 34688 5386 3064 2115 5524 **Adventure** 94596 105636 138482 86367 31896 11551 4817 2718 1835 4575 309468 **Animation** 61960 104837 65707 22825 7551 2792 1430 2290 207544 72566 911 **Biography** 47333 77867 123948 74054 23644 7702 2984 1639 1145 2849 240526 60157 69176 26099 4237 2444 3842 227088 Comedy 77173 108993 9863 1712 In [39]: cnt = count column.tolist() genre top10['cnt'] = cnt In [41]: plt.figure(figsize=[10,6]) sns.barplot(x=genre top10.index, y = genre top10.cnt) plt.title('Genres with count', fontsize = 20) plt.xlabel('Genre', fontsize=15) plt.ylabel('Count', fontsize=15) plt.show() Genres with count 60 50 30 20 10 0 Adventure Animation Biography Comedy Crime Drama Romance Sci-Fi Genre Top 1000 Voters in IMDB vs Genre genre top10 cvotes = genre top10.sort values(by='CVotes1000',ascending = False) In [42]: plt.figure(figsize=[10,6]) In [43]: sns.barplot(x=genre_top10_cvotes.index, y = genre_top10_cvotes.CVotes1000) plt.title('Top 1000 voters vs Genre', fontsize = 20) plt.xlabel('Movie Genre', fontsize=15) plt.ylabel('CV1000 Votes', fontsize=15) plt.show() Top 1000 voters vs Genre 700 600 500 400 200 100 0 Thriller Adventure Sci-Fi Action Crime Comedy Biography Drama Animation Romance Movie Genre **Inferences** • Romance genre was voted the least by top 1000 voters · Sci-Fi is the most popular · Action, Thriller and Adventure are similar, as they are genres that goes hand in hand