Task 1: Movie Rating Prediction Going to build a model that predicts the rating of a movie based on features like actors, directors, and genres In [4]: **import** numpy **as** bp import os import pandas as pd import matplotlib.pyplot as plt %matplotlib inline from matplotlib import style import seaborn as sns import heapq In [16]: #import csv file df = pd.read_csv('Movies.csv') **Data Cleaning Process** In [74]: #Check the overall data details df.info() <class 'pandas.core.frame.DataFrame'> Index: 5659 entries, 1 to 15508 Data columns (total 10 columns): # Column Non-Null Count Dtype --- -----0 Name 5659 non-null object 1 Year 5659 non-null int32 2 Duration 5659 non-null object 3 Genre 5659 non-null object 4 Rating 5659 non-null int32 5 Votes 5659 non-null object 6 Director 5659 non-null object 7 Actor 1 5659 non-null object 8 Actor 2 5659 non-null object 9 Actor 3 5659 non-null object dtypes: int32(2), object(8) memory usage: 442.1+ KB In [19]: #check the null values pd.isnull(df).sum() Out[19]: Name 528 8269 Duration 1877 Genre 7590 Rating Votes 7589 Director 525 Actor 1 1617 2384 Actor 2 Actor 3 dtype: int64 In [20]: df.shape Out[20]: (15509, 10) In [21]: #Remove the null values df.dropna(inplace=True) In [22]: df.shape Out[22]: (5659, 10) In [23]: #Null values removed pd.isnull(df).sum() Out[23]: Name Year Duration Genre Rating Votes Director Actor 1 Actor 2 Actor 3 dtype: int64 In [25]: #change from float to int df['Rating'] = df['Rating'].astype('int') In [145... | df['Duration'] = df['Duration'].astype('int') In [102... df['Year'] = df['Year'].astype(str).str.lstrip('<').str.rstrip('+')</pre> In [28]: df.columns Out[28]: Index(['Name', 'Year', 'Duration', 'Genre', 'Rating', 'Votes', 'Director', 'Actor 1', 'Actor 2', 'Actor 3'], dtype='object') In [106... | #Remove the special character in front of the years df['Year'] = [x.strip("-") for x in df['Year']] In [129... df['Votes'] = df['Votes'].str.replace(',','').astype(int) In [130... $df['Name'] = df['Name'].str.replace(r"[\"\'\|\=\.\...\:\<\>\@\#*\?\$\%\^\&\!_\-\;\'\~\,]", '')$ In [53]: df['Name'] = [x.strip(":",) for x in df['Name']] In [48]: df['Name'] = [x.strip("#",) for x in df['Name']] In [55]: df['Name'] = [x.strip("?",) for x in df['Name']] In [56]: | df['Name'] = [x.strip(":",) for x in df['Name']] In [61]: df['Name'] = [x.strip("@",) for x in df['Name']] In [40]: df['Name'] = [x.strip("...") for x in df['Name']]In [38]: df['Duration'] = [x.strip("min",) for x in df['Duration']] In [108... #Drop duplicate values df.drop_duplicates(subset = 'Name', keep = "first").head(2) Out[108... Director Actor 1 Actor 2 Actor 3 Name Year Duration Genre Rating Votes 8 Gaurav Bakshi Rasika Dugal Vivek Ghamande 1 Gadhvi (He thought he was Gandhi) 2019 Drama Arvind Jangid Ishita Raj Siddhant Kapoor Yaaram 2019 4 35 Ovais Khan Prateik 110 Comedy, Romance In [107... #change the name of column df.rename(columns={'Duration':'Duration (min)'}).head(2) Out[107... Name Year Duration (min) Genre Rating Votes Director Actor 1 Actor 2 Actor 3 1 Gadhvi (He thought he was Gandhi) 2019 8 Gaurav Bakshi Rasika Dugal Vivek Ghamande Arvind Jangid 3 Yaaram 2019 110 Comedy, Romance Ishita Raj Siddhant Kapoor Ovais Khan Prateik In [109... #Rename the column of duration df.rename(columns={'Duration':'Duration (min)'}).head(2) Out[109... Name Year Duration (min) Actor 2 Actor 3 Genre Rating Votes Director Actor 1 Arvind Jangid 1 Gadhvi (He thought he was Gandhi) 2019 Gaurav Bakshi Rasika Dugal Vivek Ghamande Drama Yaaram 2019 110 Comedy, Romance Prateik Ishita Raj Siddhant Kapoor 35 Ovais Khan **Exploratory Data Analysis** In [110... import warnings warnings.filterwarnings('ignore') **Rating Counts** In [111... plt.figure(figsize = (19,10)) ax = sns.countplot(x = 'Rating', data = df)for bars in ax. containers: ax.bar_label(bars) plt.title('Rating Count') plt.show() Rating Count 1661 1600 1400 1249 1200 1106 1000 count 818 800 600 412 400 237 200 135 18 10 Rating In [136... df['Rating'].max() Out[136... **10** In [135... df['Rating'].min() Out[135... **1** Maximum Duration of movie In [149... df['Duration'].max() Out[149... **321** In [280... df['Duration'].max() df['Name'][df['Duration'].idxmax()] Out[280... 'Gangs of Wasseypur' Minimum Voting In [133... df['Votes'].min() Out[133... 5 Highest Vote to which one movie? In [142... print(df['Votes'].max()) df['Name'][df['Votes'].idxmax()] 591417 Out[142... 'Life of Pi' Which year has highest Rating? In [138... plt.figure(figsize = (50,19)) ax = sns.barplot(x = 'Year', y = 'Rating', data = df)**for** bars **in** ax. containers: ax.bar_label(bars) plt.title('Yearly status') plt.show() Top 10 movies with highest voting df.nlargest(10, columns=['Votes'], keep='first') Out[210... Name Year Duration Genre Rating Votes Director Actor 1 Actor 2 Actor 3 8219 127 Adventure, Drama, Fantasy Suraj Sharma Life of Pi 2012 7 591417 Ang Lee Irrfan Khan Adil Hussain Rajkumar Hirani 75 3 Idiots 2009 170 Comedy, Drama 8 357889 Aamir Khan Madhavan Mona Singh 8233 Garth Davis 8 220526 Dev Patel Nicole Kidman Lion 2016 118 Biography, Drama Rooney Mara Gandhi 1982 4848 Biography, Drama, History 8 220118 Richard Attenborough Ben Kingsley John Gielgud Rohini Hattangadi 191 The Darjeeling Limited 2007 91 Adventure, Comedy, Drama 7 185127 Adrien Brody 14038 Wes Anderson Owen Wilson Jason Schwartzman 165 8 175810 Aamir Khan 8228 Like Stars on Earth 2007 Drama, Family Amole Gupte **Darsheel Safary** Aamir Khan 10882 PK 2014 153 Comedy, Drama, Musical 8 168150 Rajkumar Hirani Aamir Khan Anushka Sharma Sanjay Dutt 3410 Dangal 2016 Fatima Sana Shaikh 161 Action, Biography, Drama 8 165074 Nitesh Tiwari Aamir Khan Sakshi Tanwar Radhe 2021 11463 135 Action, Crime, Thriller 1 162455 Prabhu Deva Salman Khan Disha Patani Randeep Hooda 6 117377 Mukesh Chhabra Sushant Singh Rajput 3829 Dil Bechara 2020 101 Comedy, Drama, Romance Sanjana Sanghi Sahil Vaid Movie's duration having less than 100 minutes. df[(df.Duration <= 100)]</pre> Out[201... Genre Rating Votes Name Year Duration Director Actor 1 Actor 2 Actor 3 8 A Question Mark 2012 82 Allyson Patel Muntazir Ahmad Kiran Bhatia Horror, Mystery, Thriller 5 326 Yash Dave 10 96 Madhu Ambat 1:1.6 An Ode to Lost Love 2004 Drama 6 17 Rati Agnihotri Gulshan Grover Atul Kulkarni 18 10ml LOVE 2010 87 Comedy, Drama, Romance Anusha Bose Manu Rishi Chadha 6 162 Sharat Katariya Neil Bhoopalam 36 19 Revolutions 2004 Sridhar Reddy Tarun Arora Vinay Pandey Drama 16 Gulshan Grover 50 2 Nights in Soul Valley 2012 80 Adventure, Horror, Mystery Sumeet Sharma 6 21 Harish Sharma **Hemant Pandey** Sumeet Sharma 9 Sridhar Rangayan 15229 Yeh Hai Chakkad Bakkad Bumbe Bo 2003 90 Tom Alter Mona Ambegaonkar Adventure Aardra Athalye 15288 Yeh Suhaagraat Impossible 2019 Aarav Mavi Comedy Abhinav Thakur Preetika Chauhan Aloknath Pathak 15312 Yours Truly 2018 84 Soni Razdan Aahana Kumra Pankaj Tripathi Drama 5 102 Sanjoy Nag 15320 Yug the law of karma 2021 75 Action, Crime, Drama 10 Dilip Kumar Vinod Kumar Dilip Kumar Saurav Bagga 15488 Zoo 2018 100 78 Shlok Sharma Shashank Arora Prince Daniel Shatakshi Gupta Drama 5 596 rows × 10 columns Movie's duration having less than 50 minutes In [279... df['Duration'].loc[lambda x:x < 50]</pre> Out[279... 280 1513 48 2723 49 4504 45 5756 47 8046 48 8709 45 10504 49 10974 21 12002 45 12658 47 13217 45 14071 48 Name: Duration, dtype: int32 Top 5 movies with highest-rated In [287... Top5_movies=df.nlargest(5,'Rating')[['Name', 'Votes', 'Rating', 'Director']]\ .set_index('Name') In [288... Top5_movies Out[288... Votes Rating Director Name Love Qubool Hai 10 Saif Ali Sayeed Ashok Vatika Rahul Mallick Baikunth 29 9 Vishwa Bhanu Consequence Karma 9 Shadab Ahmad Gho Gho Rani 47 9 Munni Pankaj In [234... ax = sns.barplot(x='Rating', y=Top5_movies.index,data=Top5_movies, hue='Director') for bars in ax. containers: ax.bar_label(bars) Love Qubool Hai Ashok Vatika Director Saif Ali Sayeed Name Rahul Mallick Baikunth Vishwa Bhanu Shadab Ahmad Munni Pankaj Consequence Karma Gho Gho Rani 10 Rating Top 5 movies with highest-voted In [285... Top5_movies=df.nlargest(5,'Votes')[['Name', 'Year', 'Rating','Votes', 'Director','Actor 1','Actor 2', 'Actor 3']]\ .set_index('Name') In [286... Top5_movies Out[286... Year Rating Director Actor 1 Actor 2 Actor 3 Votes Name 7 591417 Adil Hussain Life of Pi 2012 Ang Lee Suraj Sharma Irrfan Khan **3 Idiots** 2009 8 357889 Rajkumar Hirani Aamir Khan Madhavan Mona Singh Lion 2016 8 220526 Garth Davis Dev Patel Nicole Kidman Rooney Mara Gandhi 1982 8 220118 Richard Attenborough John Gielgud Rohini Hattangadi Ben Kingsley The Darjeeling Limited 2007 7 185127 Wes Anderson Owen Wilson Adrien Brody Jason Schwartzman In [241... ax = sns.barplot(x='Votes', y=Top5_movies.index,data=Top5_movies, hue='Director') for bars in ax. containers: ax.bar_label(bars) 591417 Life of Pi 357889 3 Idiots Name 220526 Lion Director Gandhi Ang Lee 220118 Rajkumar Hirani Garth Davis Richard Attenborough The Darjeeling Limited Wes Anderson 185127 100000 200000 300000 400000 500000 600000 Votes Lower_rated movies In [243... df.nsmallest(5, columns=['Rating'], keep='first') Out[243... Genre Rating Votes Name Year Duration Director Actor 1 Actor 2 Actor 3 150 Action, Adventure, Comedy 2918 Chatur Singh Two Star 2011 624 Ajay Chandhok Sanjay Dutt Ameesha Patel Anupam Kher 3618 Desh Drohi 2008 140 **Gracy Singh** Action, Thriller 1 3899 Jagdish A. Sharma Kamal Rashid Khan Hrishitaa Bhatt 4170 Dracula 2012 2013 137 Horror 1 128 Vinayan Thilakan Shraddha Das Monal Gajjar 5523 Hari Puttar: A Comedy of Terrors 2008 Comedy, Drama, Family Rajesh Bajaj 90 1 314 Lucky Kohli Jackie Shroff Sarika 5711 Himmatwala 2013 150 Action, Comedy, Drama 1 8186 Sajid Khan Ajay Devgn Tamannaah Bhatia Mahesh Manjrekar Least_voted movies In [284... df.nsmallest(50, columns=['Votes'], keep='first') Out[284... Genre Rating Votes Name Year Duration Director Actor 1 Actor 2 Actor 3 1116 Anmol Sitaare 1982 153 Drama 5 Geethapriya Master Baboo Rakesh Bedi Ramesh Deo 1341 Atal Faisla 2018 127 Abdul Sattar Drama Sahil Akhtar Himayat Ali Aman Jain 131 1434 Awesome Mausam 2016 3 5 Yogesh Bharadwaj Vaarssh Bhatnagar Sunil Chaurasiyaa Suhasini Mulay Romance 1469 B for Bundelkhand 2017 117 8 Bharat Chawla Moumita Nandi Drama Vishal Mourya Nemi Chandra Jha 2118 Bhai-Bahen 1959 136 Drama, Family 6 G.P. Sippy Kathana 5 Daisy Irani Rajan Kapoor Chipku 2016 3061 135 Drama Priyanka Gouri Ash Rakesh Bedi Neha Pal Chowdhury Sanjay Khan 3107 Chori Chori 1972 115 5 5 Keval P. Kashyap Action, Drama Radha Saluja Kamran Daku Kali Bhawani 2000 120 Dharmendra 3361 Action 3 S.R. Pratap Milind Gunaji Mohan Joshi 7 A.K. Hangal 4046 Do Nambar Ke Amir 1974 136 5 P.D. Shenoy Asha Sachdev Sajid Khan Crime, Drama 4384 Ek Daku Saher Mein 1985 133 Kalidas Action Ashok Kumar Pradeep Kumar Suresh Oberoi 4523 120 Ashwini Kalsekar Ayush Mahesh Khedekar Ek Tha Hero 2018 Drama, Family 4 5 Yogesh Pagare Asrani 5087 Ghulam Begum Badshah 1956 124 Malika Drama, History Jugal Kishore Daljeet Nishi 5353 Gyara Hazar Ladkian 1962 152 Bharat Bhushan 5 Khwaja Ahmad Abbas Mala Sinha Drama Helen 6160 Intezar 1973 116 Drama Mohan Singh Kavia Rinku Jaiswal Padmini Kapila Baldev Khosa 6284 Jaan Lada Denge 1990 137 6 5 Dilip Gulati Hemant Birje Sahila Chaddha Natasha Action 6457 Jai Mahakali 1951 162 Fantasy Dhirubhai Desai Nirupa Roy Shahu Modak Ulhas 6573 Jawab Ayega 1968 82 6 5 Ismat Chughtai Shaheed Latif Yogesh Bali Family Meena Rai 7082 Kadke Kamal Ke 2019 136 Laxman Singh Comedy Aaryan Adhikari Neeta Dhungana Rajpal Yadav 7262 Kaneez 1949 140 Drama, History 6 5 Krishna Kumar Shyam Urmila Munawwar Sultana 8288 Lorni - The Flaneur 2019 107 Wanphrang Diengdoh Elizer Bareh Deihokhlang Basaiawmoit Crime, Mystery Poonam Agarwal 8339 Love Qubool Hai 2020 94 Drama, Romance 10 Saif Ali Sayeed Ahaan Jha Mahesh Narayan Rajasree Rajakumari Priyanshu Chatterjee 8726 Majaz: Ae Gham-e-Dil Kya Karun 2017 178 Neelima Azim Biography Ravindra Singh Rashmi Mishra Manchala 1999 8801 137 Jayprakkash Shaw Vivek Mushran Rakesh Bedi 5 Gauri Khopkar Romance Mashaal 1950 Ruma Guha Thakurta 8923 136 Drama Nitin Bose Ashok Kumar Sumitra Devi 9525 Mounto 1975 114 5 Action, Crime Jambulingam Navin Nischol Saira Banu 9568 Mr. X 1987 120 Fantasy 5 Khwaja Ahmad Abbas Amol Palekar Shabana Azmi Tom Alter Naami Chor 1977 127 Biswajit Chatterjee Shatrughan Sinha 9829 5 Manmohan Desai Action Kamal Mehra 10179 Night Club 1958 137 6 Crime, Drama Naresh Saigal Ashok Kumar Kamini Kaushal Mubarak 11179 116 5 Rudra Patla Venugopal Pyaar Karle 2019 3 Faheem UI Haq Drama Chitram Basha Udita Goswami 11352 Qatil Chudail 2001 Jhony Nirmal Horror Kanti Shah Anil Nagrath Amit Pachori 11628 Rakshaa Bandhan 1977 90 5 Shantilal Soni Jairaj Pallavi Joshi Lalita Pawar Drama 12058 Vijaya Choudhury Rustom-E-Rome 1964 Action, Adventure, Drama Radhakant Dara Singh Azaad Irani Dara Singh 12885 Shankar Khan 1966 140 5 Nanabhai Bhatt Action, Drama, Sport 5 Prithviraj Kapoor Randhawa 12969 Shehnai 1947 133 P.L. Santoshi Romance 5 Rehana Indumati Nasir Khan 13759 Tangewali 1955 120 8 5 Lekhraj Bhakri Shammi Kapoor Balraj Sahni Anita Guha Drama, Family 14825 Veeru Ustaad 1975 127 Drama 5 Jagdish Nirula Aruna Irani Shakti Kapoor Satish Kaul 15466 Zindagi Aur Maut 1965 134 6 Bela Bose Drama 5 Nisar Ahmad Ansari Nisar Ahmad Ansari Chandrima Bhaduri 79 99 3.0 Megapixel 2015 Adventure, Drama Divakar Ghodake Seema Azmi Meenal Chirankar Pooja Kasekar 91 3rd EYE 2019 145 6 6 Deepak Bhatia Sakshi Dwivedi Sanjay Niranjan Mukul Dev Horror 172 Khwaja Ahmad Abbas 493 Aasmaan Mahal 1965 Drama Prithviraj Kapoor Dilip Roy Surekha 530 Ab Hoga Dharna Unlimited 2012 78 4 6 Deepak Tanwar Saurabh Malik Rekha Rana Omkar Das Manikpuri Comedy 539 127 Arvind Joshi Ab To Aaja Saajan Mere 1994 Drama, Romance Rakesh Nahata Keshto Iqbal Shalini Kapoor 1189 Apaatkaal 1993 137 5 6 V. Subba Rao Aatish Devgan Sabah Rita Bhaduri Action Apradhi 1974 122 7 Yogeeta Bali 1236 Drama Jugal Kishore Kiran Kumar **Imtiaz** 1531 Baazigar 1972 133 Adventure, Comedy, Crime 7 Karunesh Thakur Vijayalalitha Bindu 6 Roopesh Kumar Bad Boys 2003 1575 118 Thriller 3 6 Vicky S. Kumar Rakesh Bedi Suresh Chatwal Avtar Gill 1684 Bagpat Ka Dulha 2019 137 5 6 Puneet Vashisht Lalit Parimoo Drama Karan Kashyap Raza Murad 148 G. Rakesh 1861 Banwra 1950 Drama Raj Kapoor Lalita Pawar Nimmi 1932 Bazaar 1949 148 6 Nigar Sultana Shyam Drama, Family K. Amarnath Gope 2207 Bhayaanak 1998 105 Tina Ghai Horror R. Mittal Anil Dhawan Bharat Kapoor Insights 1. Out of 10 the 6 ratings are mostly given to movies at a maximum of 1661 counts. 2. Gangs of Wasseypur have taken a maximum duration, of 321 minutes. Action, comedy, and crime movie made in 2012 by Anurag Kashyap director. 3. The highest voting goes to life of pi movie, 591417. Life of pi was released in 2012 by Ang lee director. 4. Least_voted movies are like, Anmol sitaare, Atal Faisla, Awesome Mausam (Flop), B for Bundelkhand, Bhai-behan, and many more. Overall 37 movies got the least voting. 5. There is one movie, Love Qubool Hai got 10 out of 10 rating (highest_rated) directed by Saif Ali Sayeed but got the least number of votes.

