	or you can run the follown	that command-line switch used in this notebook anal files to get the data sets fro line: cd src python main.py	lysis simply from data			l already,
	perform analysis on the datransformer json file to csv  • pandas, numpy, seabo  1. Any major "gotchas" to to it takes more than 10 minut  1. Anything else you feel is  My data is about movie rate	relevant to the grading o	e stored in data subfoollowing packages: oup,sklearn lon't work, go slowly n websites.  f your project your person of the story of the	oroject.	oved, etc.)  vie genres, etc. I hope to	-line can
	1. What did you set out to so but if you switched gears  My project wants to predict  1. What did you Discover/weetc.?)	s or changed things, note	ome characteristics of point of your project it here.)  In g data from the past its (i.e. what were you	This should be three years.	close to your Milestone e your original assumpti	1 assignment,
	predict, it is found that the Therefore, the conclusion i In addition, I did some stat A. More than 2/3 of the I B. Action movies have fa C. Animation duration is D. Animation is easier to	s that these factors are indecistical analysis and concludate hotter movies are all action or more followers than other usually shorter than other	eed related to movie ed:  movies r types types of duration	ratings, and we ca	-	
	complicated. I failed many a lot of time.  1. What skills did you wish	st be scrape the imdb informatimes and I needed to wait  you had while you were demost for this project is to us	mation data from IMo for awhile every time loing the project?	. Encoding issues	in the process of crawling	
In [1]:	<pre>import pandas as pd from bs4 import Beautif import seaborn as sns import matplotlib.pyplo import numpy as np  # Read in the three dat movie_info_ = pd.read_c movie_info = movie_info</pre>	t <b>as</b> plt  a sets  sv('movie_imb_data.cs	v')			
<pre>In [3]: In [4]: Out[4]:</pre>	#Take a look of the col  movie_info.head()	n_user_for_reviews plot_keyo	words genres time travel Action 181.0			ritic_for_reviews 583.0 http
In [5]:	1 1.0 Avengers: Infinity War  Star Wars: Episode IX 2 2.0 - The Rise of Skywalker  4 4.0 No Time to Die  5 5.0 Tenet	4400.0 rec cha 7900.0 2700.0 sho	ath of urring Action 149.0 racter  evil Action 141.0  poting Action 163.0  NaN Action 150.0  the movie metadat	2018.0 2019.0 2021.0 2020.0	8.4 NaN 6.5 2.39 7.5 2.39 7.4 2.20	624.0 http 518.0 http 387.0 http 498.0 http
In [6]: Out[6]:	#like rate, and so on  #Get the average score mean_score_per_year = mmean_score_per_year.res  mean_score_per_year.hea  title_year Unnamed: 0 num  0 2018.0 317.000000  1 2019.0 313.000000  2 2020.0 320.666667	ovie_info.groupby(mov: et_index(inplace= <b>True</b> d()	ne imdb_score aspec 57 7.500000 2.3 00 7.450000 2.3		_for_reviews 459.333333 550.500000 359.666667	
In [7]:	#Trying to figure out w #Get the regression lin fig, axes = plt.subplot m1 = sns.distplot(movie)  C:\Users\49181\anaconda3 ated function and will b -level function with sim warnings.warn(msg, Fut	e s(1, 1) _info.imdb_score, bins \lib\site-packages\se se removed in a future silar flexibility) or	elates the scores s=15) aborn\distributio version. Please	adapt your cod	le to use either `dis	splot` (a figure
	0.8 - 0.6 - 0.6 - 0.2 - 0.0 - 0.5 in	7.0 7.5 8.0 8.5 s	9.0			
<pre>In [8]: Out[8]: In [9]:</pre>	Animation 98 Adventure 49 Name: genres, dtype: int  # More than 2/3 of the  vis1 = sns.lmplot(data=	64 hotter movies are all	er_for_reviews',		, \	
	C:\Users\49181\anaconda3 n renamed to `height`; p warnings.warn(msg, Use 8.5	lease update your cod			rning: The `size` pa	rameter has bee
		000 6000 80 um_user_for_reviews		on nation enture		
In [11]: In [12]:	<pre>#Action movies have far  fig, axes = plt.subplot fig.set_size_inches(11. plt.subplots_adjust(wsp)  sns.boxplot(data=movie_ sns.boxplot(data=movie_ sns.boxplot(data=movie_ sns.boxplot(data=movie_ plt.suptitle('Some data plt.show()</pre>	s(2, 2) 7, 5.27) ace=0.2, hspace=0.4) info, x='genres', y='s info, x='genres', y='s info, x='genres', y='s	<pre>imdb_score', ax=a num_user_for_revi time', ax=axes[1, num_critic_for_re</pre>	ews', ax=axes[ 0])		
			ding to their genre	Animation genres	Adventure	
In [13]:		movie_info['imdb_score		Animation genres	Adventure	
Out[13]:	Action Action	Animation				
In [14]:	<pre>#df.head() for feature in columnsT     try:         df[feature] = 1     except:</pre>	<pre>movie_info[['title_year','genres','num_r ,'imdb_score','num_cr</pre>	ar','genres']]) user_for_reviews' itic_for_reviews'	_	s',	
Out[14]:	['title_year', 'genres'] <ipython-input-14-125820 .loc[row_index="" a="" be="" caveats="" d="" df[feature]="le.fit_t&lt;/th" in="" is="" see="" the="" to="" try="" trying="" urning-a-view-versus-a-c="" using="" value=""><th>17a042&gt;:11: SettingWi set on a copy of a sl er,col_indexer] = val documentation: https:/ copy ransform(df[feature])</th><th>ice from a DataFrue instead  /pandas.pydata.or  ot_keywords time im  time travel 181.0</th><th>g/pandas-docs/</th><th></th><th>ndexing.html#ret</th></ipython-input-14-125820>	17a042>:11: SettingWi set on a copy of a sl er,col_indexer] = val documentation: https:/ copy ransform(df[feature])	ice from a DataFrue instead  /pandas.pydata.or  ot_keywords time im  time travel 181.0	g/pandas-docs/		ndexing.html#ret
In [15]:	<pre># get x y X = df.dropna() y = X['imdb_score'] X = X.drop(['imdb_score  #split data from sklearn.model sele</pre>		shooting 163.0  NaN 150.0	7.5 7.4	387.0 498.0	
	from sklearn.preprocess scaler=StandardScaler()  X = scaler.fit_transform y = np.array(y*10).asty X_train, X_test, y_train #print(y_test)  [67 67 74 75 84 84 84 74 69 65 65 67 76 84 84 74 69 69 84 74 76 57 73 75 69 84 84 65 75 76 84 84 57 57 67 76 84 74 74 67 65 67 84 65 74 74 65 74 74 75 67 65 74 74 65 74 57 74 84 76 76 69 76	ing import StandardScam(X) pe(int) n, y_test = train_test 74 84 84 65 74 67 84 74 65 84 76 69 84 84 76 69 84 84 67 84 84 57 74 65 76 65 65 84 65 65 57 74 65 73 57 74 74 67 65 84 76 65 84 84 65 69 74 73 76 84 84 65 69 74 73 76	57 75 75 84 84 8 73 57 76 57 76 7 69 74 75 74 84 6 65 73 75 69 73 5 65 73 74 74 67 7 76 69 69 65 75 6 74 84 65 69 74 6 73 69 74 69 65 8	4 73 75 84 3 67 84 74 7 74 57 74 7 75 65 84 4 65 65 75 9 65 65 65 5 73 74 84 4 84 73 84	ndom_state=33)	
	74       73       84       65       73       84       57       76         67       73       84       65       75       74       74       67         84       67       75       84       74       65       65       67         73       65       65       75       65       73       84       74         74       74       74       84       73       76       57       73         76       74       74       57       57       75       65       76         74       57       84       65       75       67       69       57         74       57       65       57       74       67       69       69         73       65       76       65       57       69       84       74         84       65       74       67       73       69       69       74         84       65       74       67       73       69       69       74         85       75       74       84       73       67       69       84         65       75       74       84 <td>74       65       57       84       74       65       76         74       73       74       69       74       84       73         84       74       76       74       69       74       67         67       74       75       73       84       76         65       84       73       67       73       74       67         75       76       74       65       57       57       69         73       74       84       65       74       65       65         65       74       73       74       57       84         84       65       57       73       84       69         74       76       65       73       84       69         74       76       67       74       73       74       84         69       75       74       65       73       84       69         74       76       67       74       73       74       84         76       76       84       67       74       74       65         74       57       65       73&lt;</td> <td>57       57       84       65       75       73       6         67       57       65       75       73       6         65       75       76       76       76       7         73       76       76       76       7         74       57       69       69       76       7         73       65       76       57       65       84       6         76       69       57       74       69       7         73       57       57       65       84       6         69       84       75       84       84       6         76       65       65       69       69       8         57       57       74       75       75       6         75       84       84       74       8         74       69       67       65       75       75       7</td> <td>7 65 65 84 5 65 67 75 5 73 74 84 4 67 84 67 3 73 69 76 4 84 84 65 9 84 74 67 4 84 65 76 5 69 67 73 7 74 84 65 7 65 75 75 7 65 74 69 4 69 57 74 5 84 84 57</td> <td></td> <td></td>	74       65       57       84       74       65       76         74       73       74       69       74       84       73         84       74       76       74       69       74       67         67       74       75       73       84       76         65       84       73       67       73       74       67         75       76       74       65       57       57       69         73       74       84       65       74       65       65         65       74       73       74       57       84         84       65       57       73       84       69         74       76       65       73       84       69         74       76       67       74       73       74       84         69       75       74       65       73       84       69         74       76       67       74       73       74       84         76       76       84       67       74       74       65         74       57       65       73<	57       57       84       65       75       73       6         67       57       65       75       73       6         65       75       76       76       76       7         73       76       76       76       7         74       57       69       69       76       7         73       65       76       57       65       84       6         76       69       57       74       69       7         73       57       57       65       84       6         69       84       75       84       84       6         76       65       65       69       69       8         57       57       74       75       75       6         75       84       84       74       8         74       69       67       65       75       75       7	7 65 65 84 5 65 67 75 5 73 74 84 4 67 84 67 3 73 69 76 4 84 84 65 9 84 74 67 4 84 65 76 5 69 67 73 7 74 84 65 7 65 75 75 7 65 74 69 4 69 57 74 5 84 84 57		
In [123	# predict randomForest from sklearn.ensemble in from sklearn.metrics im clf = RandomForestClass print(X_train) clf = clf.fit(X_train, clf_y_predict = clf.pre #print(accuracy_score(y)  [[ 1.13137085 -0.5488213 [-1.41421356 2.0855209 [-0.56568542 -0.5488213 [-1.41421356 -0.5488213 [ 1.13137085 -0.5488213	<pre>port accuracy_score ifier(n_estimators=20  y_train) dict(X_test) _test, clf_y_predict)  4 -0.34840069 -1.0690 4 -0.97002445 0.9354 1.67187652 -1.4699 0.31207455 -1.8708</pre>	4497 1.11058519 1435 -1.02363546 3683 0.06718843 2869 0.44660543	-0.04617783] 0.94087319] 1.85865746]		
	[ 1.13137085 -0.5488213 [ 1.13137085 -0.5488213 [ 0.28284271	-0.34840069 -1.0690 4 -1.10173098 0.9354 4 -1.10173098 0.9354 1.67187652 -1.4699 2 -0.81461851 0.5345	4497 1.11058519 1435 -1.78246947 1435 -1.78246947 3683 0.06718843 2248 -0.26480145	-0.19336964] -0.72152852] -0.72152852] 0.94087319] -0.5137283]		