

## Step 1, Data Preprocessing and Exploratory Data Analysis:

We import the 'movie\_metadata.csv' data.

There are 28 columns. We need to select the most useful features to do the visualization and prediction.

First, keep the numerical columns, drop the non-numerical and non-related columns, or we should assume the value to those columns, for example color (0/1), language (0/1), 'plot\_keywords' and country.

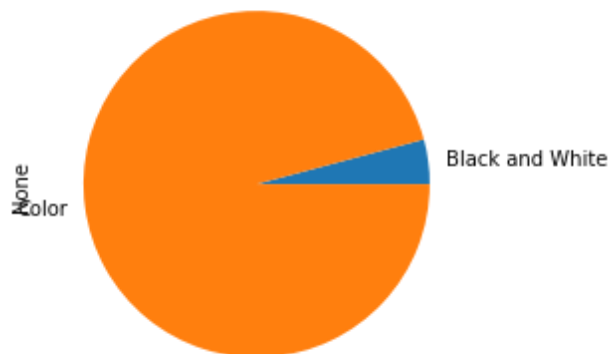
Second, some columns like director\_name, actor\_2\_name, actor\_1\_name are either difficult to get the patterns or related to other columns like 'director\_facebook\_likes', 'actor\_2\_facebook\_likes' and 'actor\_1\_facebook\_likes'. So we could remove those columns.

Third, the movie\_imdb\_link, movie\_title should be removed for the purpose of predicting the movie rating.

Finally, some columns are numerical but useless in prediction, like 'duration', 'facenumber\_in\_poster', 'title\_year'.

## Step 2, visualization:

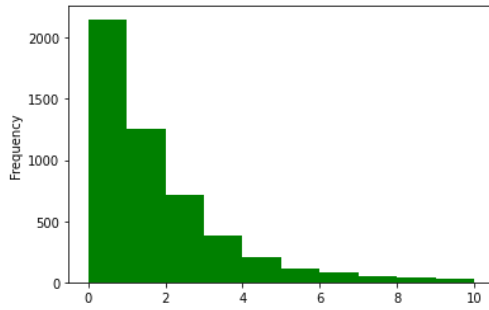
Pie chart for the color:



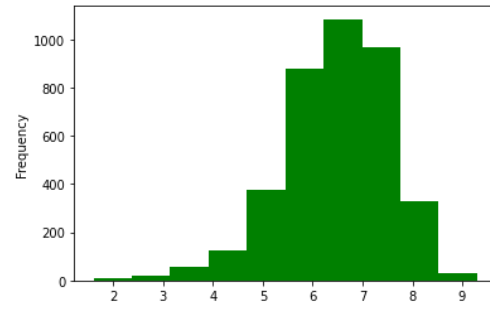
Pie chart for language:



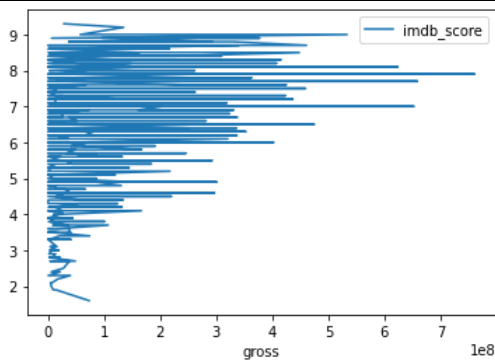
Histogram for facenumber in poster:



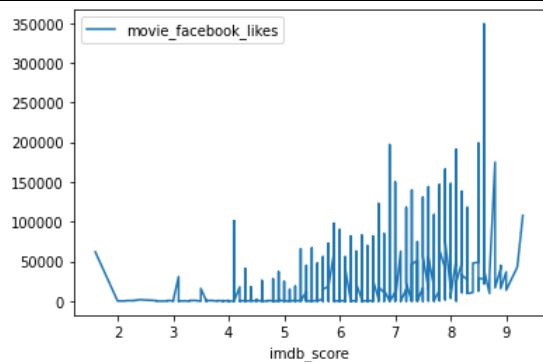
Histogram for the imdb score



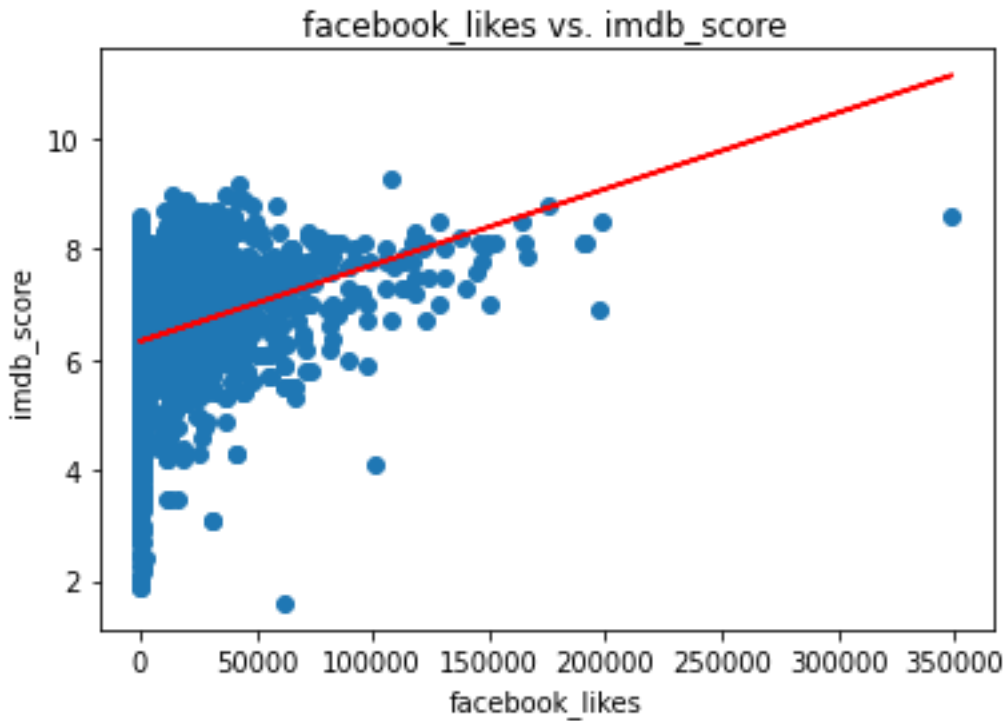
Line chart 'gross' and 'imdb\_score'



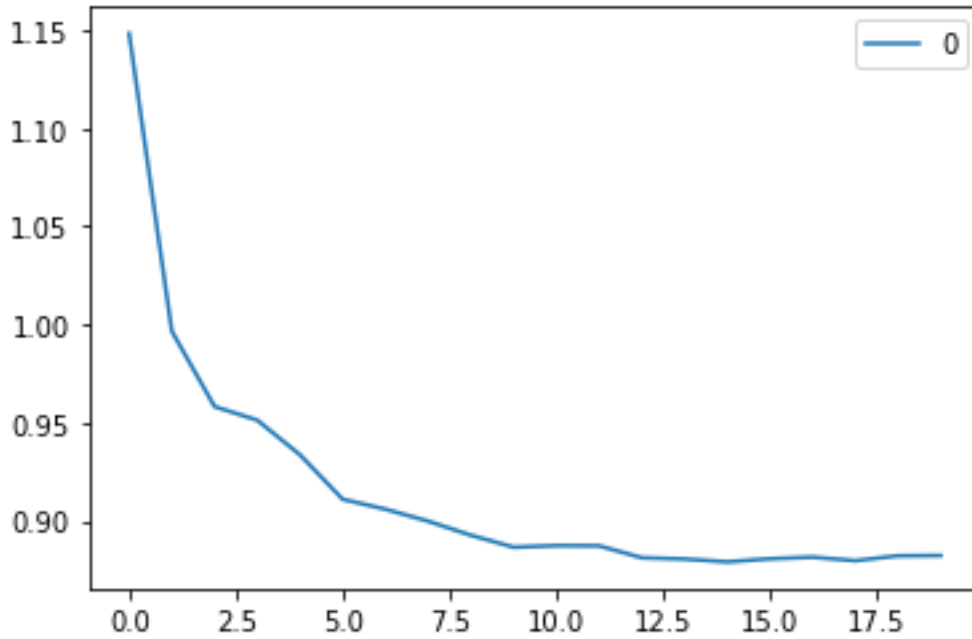
Line chart 'facebook\_likes' and 'imdb\_score'



Step 3, Simple Linear Regression:



Step 4, K Nearest Neighbor:



So K = 9 is a better choice.

Step 5, Neural Network:

	precision	recall	f1-score	support
0	0.84	0.91	0.87	689
1	0.71	0.58	0.64	281
accuracy			0.81	970
macro avg	0.78	0.74	0.75	970
weighted avg	0.80	0.81	0.80	970

The precision, recall and f1-score seem good, Neural Network is a good model for prediction.