

```
import numpy as np
import pandas as pd
import sklearn as sl
import sklearn.preprocessing as preprocessing
from sklearn.model_selection import train_test_split
import seaborn as sns
import matplotlib.pyplot as plt
pd.set_option('display.float_format', lambda x: '%.3f'%x)
pd.set_option('display.max_columns', None)
pd.set_option('display.max_rows', None)
```

```
TV=pd.read_table('/content/TVdata.txt',header=0,sep=',',lineterminator='\n')
print(TV.head())
```

```

0  video_id  cvt_per_day  weighted_categorical_position  \
0   385504   307127.606                                1
1   300175   270338.426                                1
2   361899   256165.867                                1
3   308314   196622.721                                3
4   307201   159841.652                                1

   weighted_horizontal_poitio  import_id  release_year  \
0                             3  lionsgate           2013
1                             3  lionsgate           2013
2                             3    other           2012
3                             4  lionsgate           2008
4                             3  lionsgate           2013

   genres  imdb_votes  budget  \
0  Action,Thriller,Drama    69614  15000000
1  Comedy,Crime,Thriller    46705  15000000
2  Crime,Drama            197596  26000000
3  Thriller,Drama,War,Documentary,Mystery,Action    356339  15000000
4  Crime,Thriller,Mystery,Documentary    46720  27220000

   boxoffice  imdb_rating  duration_in_mins  metacritic_score  awards  \
0  42930462      6.500      112.301          51  other award
1   3301046      6.500       94.983          41    no award
2  37397291      7.300      115.764          58  other award
3  15700000      7.600      130.704          94    Oscar
4   8551228      6.400      105.546          37  other award

   mpaa  star_category
0  PG-13      1.710
1     R      3.250
2     R      2.647
3     R      1.667
4     R      3.067
```

#Part1 Data Exploration

```
if TV['video_id'].duplicated().sum() == 0:
    print('no duplicated index')
```

no duplicated index

```
TV.info()
print(TV.drop(columns = ['video_id','release_year'],axis = 1).describe(percentiles = [0.1,0.25,0.5,0.75,0.95]))
(TV==0).sum(axis = 0)/TV.shape[0]
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4226 entries, 0 to 4225
Data columns (total 16 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   video_id                             4226 non-null   int64
1   cvt_per_day                           4226 non-null   float64
2   weighted_categorical_position          4226 non-null   int64
3   weighted_horizontal_poitio            4226 non-null   int64
4   import_id                             4226 non-null   object
5   release_year                           4226 non-null   int64
6   genres                                4226 non-null   object
7   imdb_votes                             4226 non-null   int64
8   budget                                4226 non-null   int64
9   boxoffice                             4226 non-null   int64
10  imdb_rating                           4226 non-null   float64
11  duration_in_mins                       4226 non-null   float64
12  metacritic_score                       4226 non-null   int64
```

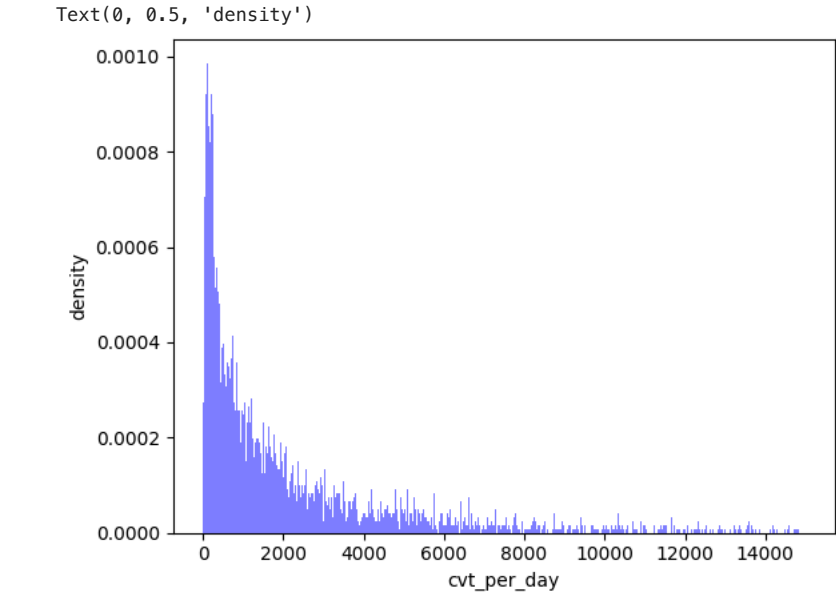
```
13 awards 4226 non-null object
14 mpaa 4226 non-null object
15 star_category 4226 non-null float64
dtypes: float64(4), int64(8), object(4)
memory usage: 528.4+ KB

count cvt_per_day weighted_categorical_position \
mean 4226.000 4226.000
std 4218.630 7.783
std 13036.080 6.134
min 2.188 1.000
10% 141.985 3.000
25% 351.169 4.000
50% 1193.500 6.000
75% 3356.789 9.000
95% 14692.834 22.000
max 307127.606 41.000

count weighted_horizontal_poiton imdb_votes budget boxoffice \
mean 28.104 6462.924 2150743.439 2536338.472
std 11.864 31596.007 7176604.483 8243516.266
min 1.000 0.000 0.000 0.000
10% 13.000 8.000 0.000 0.000
25% 20.000 81.000 0.000 0.000
50% 28.000 535.000 0.000 0.000
75% 36.000 3053.000 1500000.000 0.000
95% 48.000 26199.500 12000000.000 8551228.000
max 70.000 948630.000 107000000.000 184208848.000

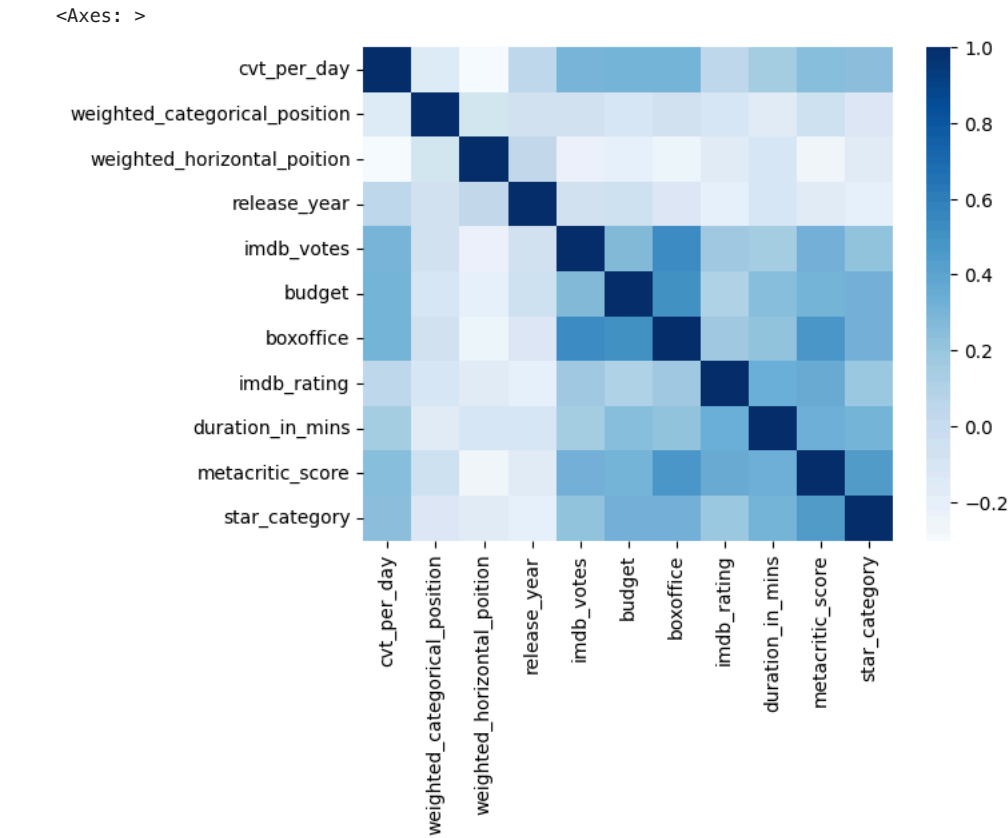
count imdb_rating duration_in_mins metacritic_score star_category
mean 5.257 89.556 15.974 0.955
std 2.123 21.086 26.205 0.955
min 0.000 4.037 0.000 0.000
10% 2.300 62.391 0.000 0.000
25% 4.300 82.602 0.000 0.000
50% 5.800 90.730 0.000 1.000
75% 6.800 99.500 41.000 1.667
95% 7.800 119.131 65.000 2.597
max 10.000 246.017 100.000 4.000
```

```
plt.hist(TV['cvt_per_day'],bins = range(0,15000,30),color = 'blue',label='cvt_per_day',density = True,alpha = 0.5)
plt.xlabel('cvt_per_day')
plt.ylabel('density')
```



```
corr = TV[['cvt_per_day','weighted_categorical_position','weighted_horizontal_poiton',
            'release_year', 'imdb_votes', 'budget', 'boxoffice', 'imdb_rating',
            'duration_in_mins', 'metacritic_score', 'star_category']].corr()

sns.heatmap(corr, cmap="Blues")
```



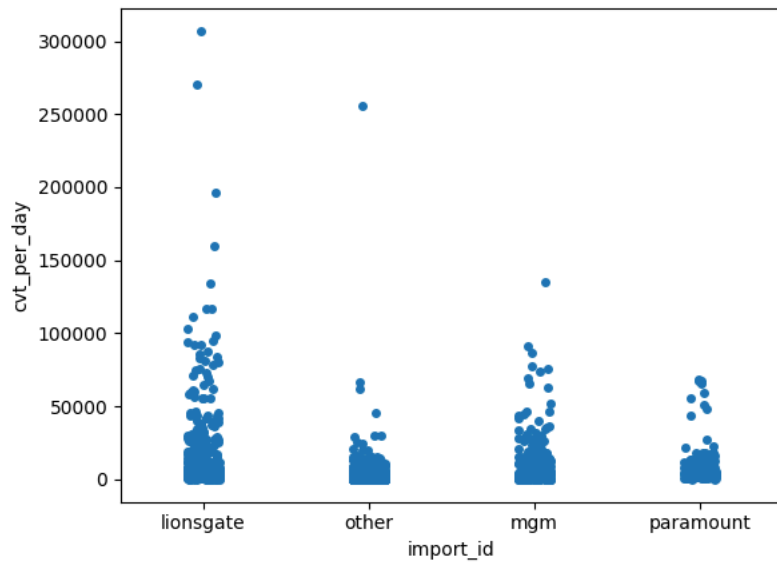
corr

	cvt_per_day	weighted_categorical_position	weighted_horizontal_poition	release_year	imdb_votes	
cvt_per_day	1.000	-0.148	-0.302	0.046	0.298	
weighted_categorical_position	-0.148	1.000	-0.084	-0.069	-0.064	
weighted_horizontal_poition	-0.302	-0.084	1.000	0.027	-0.221	
release_year	0.046	-0.069	0.027	1.000	-0.083	
imdb_votes	0.298	-0.064	-0.221	-0.083	1.000	
budget	0.316	-0.090	-0.211	-0.040	0.273	
boxoffice	0.312	-0.074	-0.245	-0.128	0.531	
imdb_rating	0.059	-0.116	-0.162	-0.199	0.163	
duration_in_mins	0.152	-0.174	-0.102	-0.097	0.156	
metacritic_score	0.249	-0.044	-0.255	-0.167	0.322	
star_category	0.247	-0.123	-0.168	-0.198	0.225	

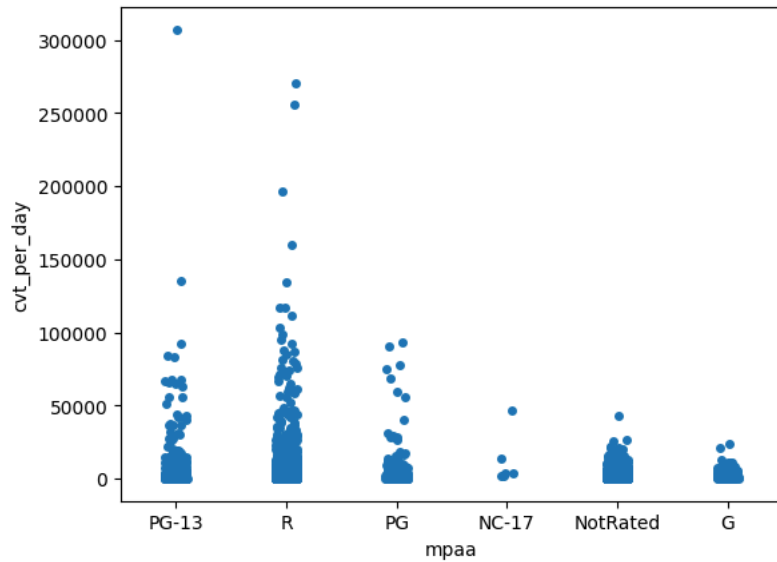
```
sns.stripplot(x='import_id',y='cvt_per_day',data = TV,jitter = True)
plt.show()
print(TV['import_id'].value_counts())

sns.stripplot(x='mpaa',y='cvt_per_day',data = TV,jitter = True)
plt.show()
print(TV['mpaa'].value_counts())

sns.stripplot(x='awards',y='cvt_per_day',data = TV,jitter = True)
plt.show()
print(TV['awards'].value_counts())
```



```
import_id
other      2963
lionsgate   677
mgm         445
paramount   141
Name: count, dtype: int64
```



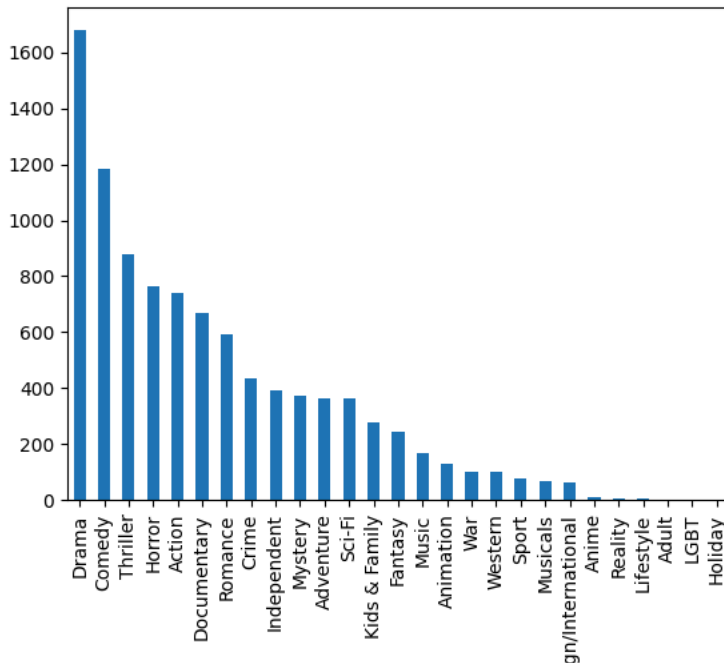
```
mpaa
NotRated   2152
```

```
gen_split = TV['genres'].str.get_dummies(sep = ',').sum()
print(gen_split)
gen_split.sort_values(ascending = False).plot.bar()
```

```

Action          739
Adult           3
Adventure       363
Animation       129
Anime           11
Comedy          1184
Crime           437
Documentary     671
Drama           1677
Fantasy         243
Foreign/International 64
Holiday         1
Horror          762
Independent     393
Kids & Family   280
LGBT            2
Lifestyle       7
Music           171
Musicals        68
Mystery         375
Reality         9
Romance         591
Sci-Fi          363
Sport           77
Thriller        879
War             102
Western         102
dtype: int64
<Axes: >

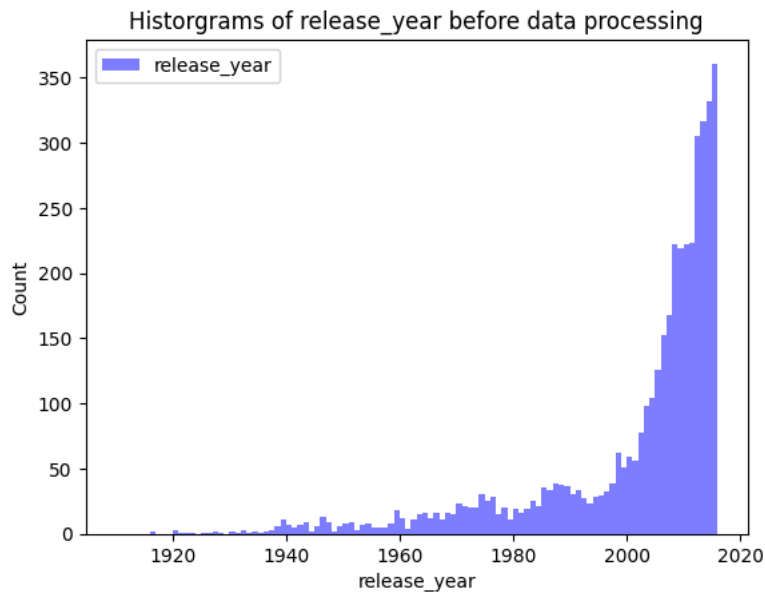
```



```

plt.hist(TV['release_year'].values, bins = range(1910, 2017, 1), alpha = 0.5, color='blue', label = 'release_year')
plt.legend(loc = 'upper left')
plt.title('Histograms of release_year before data processing')
plt.xlabel('release_year')
plt.ylabel('Count')
plt.show()

```



#Part2: Data Preprocessing

```
d_import_id = pd.get_dummies(TV['import_id']).astype(np.int64)
d_mpaa = pd.get_dummies(TV['mpaa']).astype(np.int64)
d_awards = pd.get_dummies(TV['awards']).astype(np.int64)

d_genres=TV['genres'].str.get_dummies(sep=',').astype(np.int64)
d_genres['Misc_genres']=d_genres['Anime']|d_genres['Reality']|d_genres['Lifestyle']|d_genres['Adult']|d_genres['LGBT']|d_genres['Holiday']
d_genres.drop(['Anime', 'Reality','Lifestyle', 'Adult','LGBT','Holiday'], inplace=True, axis=1)
```

```
TV['release_year'].quantile([0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9])
```

```
0.100    1974.000
0.200    1991.000
0.300    2001.000
0.400    2006.000
0.500    2008.000
0.600    2010.000
0.700    2012.000
0.800    2013.000
0.900    2014.000
Name: release_year, dtype: float64
```

```
bin_year = [1916, 1974, 1991, 2001, 2006, 2008, 2010, 2012, 2013, 2014, 2017]
year_range = ['1916-1974', '1974-1991', '1991-2001', '2001-2006', '2006-2008', '2008-2010', '2010-2012', '2012-2013',
              '2013-2014', '2014-2017']
```

```
year_bin = pd.cut(TV['release_year'], bin_year, labels=year_range)
d_year = pd.get_dummies(year_bin).astype(np.int64)
```

```
temp_tv=TV.drop(['import_id', 'mpaa','awards','genres', 'release_year'], axis=1)
```

```
newTV = pd.concat([temp_tv, d_import_id, d_mpaa, d_awards, d_genres, d_year], axis=1)
print(newTV.head())
```

	video_id	cvt_per_day	weighted_categorical_position	\
0	385504	307127.606		1
1	300175	270338.426		1
2	361899	256165.867		1
3	308314	196622.721		3
4	307201	159841.652		1

	weighted_horizontal_poition	imdb_votes	budget	boxoffice	imdb_rating	\
0		3	69614	15000000	42930462	6.500
1		3	46705	15000000	3301046	6.500
2		3	197596	26000000	37397291	7.300
3		4	356339	15000000	15700000	7.600
4		3	46720	27220000	8551228	6.400

	duration_in_mins	metacritic_score	star_category	lionsgate	mgm	other	\
--	------------------	------------------	---------------	-----------	-----	-------	---

0	112.301	51	1.710	1	0	0
1	94.983	41	3.250	1	0	0
2	115.764	58	2.647	0	0	1
3	130.704	94	1.667	1	0	0
4	105.546	37	3.067	1	0	0

	paramount	G	NC-17	NotRated	PG	PG-13	R	BAFTA	Golden	Globe	Oscar	\
0	0	0	0	0	0	1	0	0		0	0	
1	0	0	0	0	0	0	1	0		0	0	
2	0	0	0	0	0	0	1	0		0	0	
3	0	0	0	0	0	0	1	0		0	1	
4	0	0	0	0	0	0	1	0		0	0	

	no award	other award	Action	Adventure	Animation	Comedy	Crime	\
0	0		1	1	0	0	0	
1	1		0	0	0	0	1	1
2	0		1	0	0	0	0	1
3	0		0	1	0	0	0	0
4	0		1	0	0	0	0	1

	Documentary	Drama	Fantasy	Foreign/International	Horror	Independent	\
0	0	1	0		0	0	0
1	0	0	0		0	0	0
2	0	1	0		0	0	0
3	1	1	0		0	0	0
4	1	0	0		0	0	0

	Kids & Family	Music	Musicals	Mystery	Romance	Sci-Fi	Sport	Thriller	\
0	0	0	0	0	0	0	0	1	
1	0	0	0	0	0	0	0	1	
2	0	0	0	0	0	0	0	0	
3	0	0	0	1	0	0	0	1	
4	0	0	0	1	0	0	0	1	

	War	Western	Misc_genres	1916-1974	1974-1991	1991-2001	2001-2006	\
0	0	0	0	0	0	0	0	
1	0	0	0	0	0	0	0	
2	0	0	0	0	0	0	0	
3	1	0	0	0	0	0	0	
4	0	0	0	0	0	0	0	

	2006-2008	2008-2010	2010-2012	2012-2013	2013-2014	2014-2017
0	0	0	0	1	0	0

```
newTV[['budget','boxoffice','metacritic_score', 'star_category','imdb_votes', 'imdb_rating']] = newTV[['budget','boxoffice','met
print(newTV.info())
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 4226 entries, 0 to 4225
```

```
Data columns (total 58 columns):
```

#	Column	Non-Null Count	Dtype
0	video_id	4226 non-null	int64
1	cvt_per_day	4226 non-null	float64
2	weighted_categorical_position	4226 non-null	int64
3	weighted_horizontal_poition	4226 non-null	int64
4	imdb_votes	3882 non-null	float64
5	budget	1772 non-null	float64
6	boxoffice	1032 non-null	float64
7	imdb_rating	3882 non-null	float64
8	duration_in_mins	4226 non-null	float64
9	metacritic_score	1214 non-null	float64
10	star_category	2380 non-null	float64
11	lionsgate	4226 non-null	int64
12	mgm	4226 non-null	int64
13	other	4226 non-null	int64
14	paramount	4226 non-null	int64
15	G	4226 non-null	int64
16	NC-17	4226 non-null	int64
17	NotRated	4226 non-null	int64
18	PG	4226 non-null	int64
19	PG-13	4226 non-null	int64
20	R	4226 non-null	int64
21	BAFTA	4226 non-null	int64
22	Golden Globe	4226 non-null	int64
23	Oscar	4226 non-null	int64
24	no award	4226 non-null	int64
25	other award	4226 non-null	int64
26	Action	4226 non-null	int64
27	Adventure	4226 non-null	int64
28	Animation	4226 non-null	int64
29	Comedy	4226 non-null	int64
30	Crime	4226 non-null	int64
31	Documentary	4226 non-null	int64

```

32 Drama 4226 non-null int64
33 Fantasy 4226 non-null int64
34 Foreign/International 4226 non-null int64
35 Horror 4226 non-null int64
36 Independent 4226 non-null int64
37 Kids & Family 4226 non-null int64
38 Music 4226 non-null int64
39 Musicals 4226 non-null int64
40 Mystery 4226 non-null int64
41 Romance 4226 non-null int64
42 Sci-Fi 4226 non-null int64
43 Sport 4226 non-null int64
44 Thriller 4226 non-null int64
45 War 4226 non-null int64
46 Western 4226 non-null int64
47 Misc_genres 4226 non-null int64
48 1916-1974 4226 non-null int64
49 1974-1991 4226 non-null int64
50 1991-2001 4226 non-null int64
51 2001-2006 4226 non-null int64
52 2006-2008 4226 non-null int64

```

```

newTV1=newTV.copy()
newTV1['boxoffice']=newTV1['boxoffice'].fillna(newTV1['boxoffice'].mean())
newTV1['metacritic_score']=newTV1['metacritic_score'].fillna(newTV1['metacritic_score'].mean())
newTV1['star_category']=newTV1['star_category'].fillna(newTV1['star_category'].mean())
newTV1['imdb_votes']=newTV1['imdb_votes'].fillna(newTV1['imdb_votes'].mean())
newTV1['imdb_rating']=newTV1['imdb_rating'].fillna(newTV1['imdb_rating'].mean())
newTV1['budget']=newTV1['budget'].fillna(newTV1['budget'].mean())
print(newTV1.info())

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4226 entries, 0 to 4225
Data columns (total 58 columns):
#   Column                Non-Null Count  Dtype
---  -
0   video_id              4226 non-null   int64
1   cvt_per_day           4226 non-null   float64
2   weighted_categorical_position 4226 non-null   int64
3   weighted_horizontal_poition 4226 non-null   int64
4   imdb_votes            4226 non-null   float64
5   budget                4226 non-null   float64
6   boxoffice             4226 non-null   float64
7   imdb_rating           4226 non-null   float64
8   duration_in_mins      4226 non-null   float64
9   metacritic_score      4226 non-null   float64
10  star_category          4226 non-null   float64
11  lionsgate              4226 non-null   int64
12  mgm                    4226 non-null   int64
13  other                  4226 non-null   int64
14  paramount              4226 non-null   int64
15  G                      4226 non-null   int64
16  NC-17                  4226 non-null   int64
17  NotRated               4226 non-null   int64
18  PG                     4226 non-null   int64
19  PG-13                  4226 non-null   int64
20  R                      4226 non-null   int64
21  BAFTA                  4226 non-null   int64
22  Golden Globe           4226 non-null   int64
23  Oscar                  4226 non-null   int64
24  no award               4226 non-null   int64
25  other award            4226 non-null   int64
26  Action                 4226 non-null   int64
27  Adventure               4226 non-null   int64
28  Animation              4226 non-null   int64
29  Comedy                 4226 non-null   int64
30  Crime                  4226 non-null   int64
31  Documentary             4226 non-null   int64
32  Drama                  4226 non-null   int64
33  Fantasy                 4226 non-null   int64
34  Foreign/International  4226 non-null   int64
35  Horror                  4226 non-null   int64
36  Independent             4226 non-null   int64
37  Kids & Family           4226 non-null   int64
38  Music                  4226 non-null   int64
39  Musicals               4226 non-null   int64
40  Mystery                 4226 non-null   int64
41  Romance                 4226 non-null   int64
42  Sci-Fi                 4226 non-null   int64
43  Sport                  4226 non-null   int64
44  Thriller                4226 non-null   int64
45  War                    4226 non-null   int64
46  Western                 4226 non-null   int64
47  Misc_genres            4226 non-null   int64

```




```

48 1916-1974      4226 non-null  int64
49 1974-1991      4226 non-null  int64
50 1991-2001      4226 non-null  int64
51 2001-2006      4226 non-null  int64
52 2006-2008      4226 non-null  int64
-- -- -- -- --

```

```

scale_lst = ['weighted_categorical_position', 'weighted_horizontal_position', 'budget', 'boxoffice',
             'imdb_votes', 'imdb_rating', 'duration_in_mins', 'metacritic_score', 'star_category']
newTV_sc = newTV1.copy()
sc_scale = preprocessing.StandardScaler().fit(newTV_sc[scale_lst])
newTV_sc[scale_lst] = sc_scale.transform(newTV_sc[scale_lst])
newTV_sc.head()

```

	video_id	cvt_per_day	weighted_categorical_position	weighted_horizontal_position	imdb_votes	budget	boxoffice	imdb_rating
0	385504	307127.606	-1.106	-2.116	1.984	1.470	4.696	0.5
1	300175	270338.426	-1.106	-2.116	1.258	1.470	-1.022	0.5
2	361899	256165.867	-1.106	-2.116	6.043	3.108	3.898	1.1
3	308314	196622.721	-0.780	-2.032	11.077	1.470	0.767	1.3
4	307201	159841.652	-1.106	-2.116	1.258	3.290	-0.265	0.4

#Part3:Model Training

```

train,test = train_test_split(newTV_sc, test_size=0.15, random_state = 3)
model_train_x = train.drop(['video_id', 'cvt_per_day'], axis = 1)
model_test_x = test.drop(['video_id', 'cvt_per_day'], axis = 1)
model_train_y = train['cvt_per_day']
model_test_y = test['cvt_per_day']

```

```

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, Lasso, Ridge
from sklearn.metrics import mean_squared_error, r2_score #not used
from math import sqrt

```

```
lr_train, lr_validate = train_test_split(train, test_size=0.15, random_state = 0)
```

```

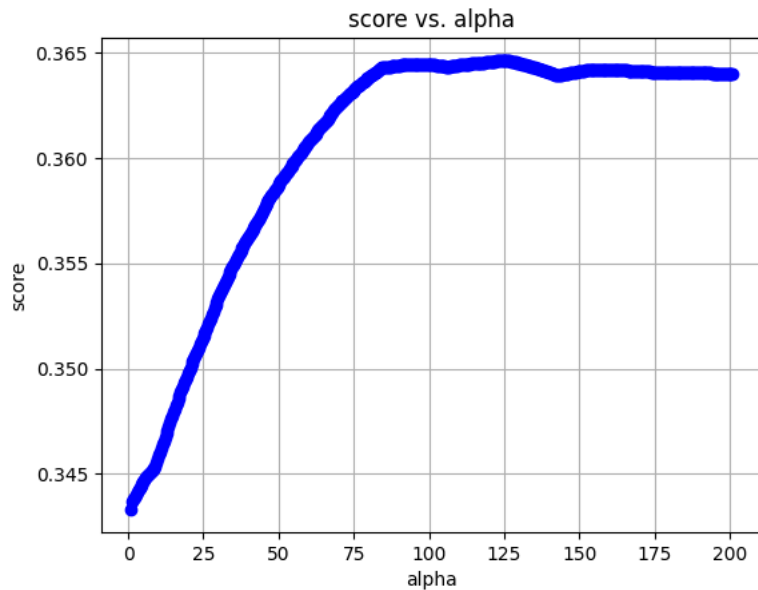
lr_train_x = lr_train.drop(['video_id', 'cvt_per_day'], axis = 1)
lr_validate_x = lr_validate.drop(['video_id', 'cvt_per_day'], axis = 1)
lr_train_y = lr_train['cvt_per_day']
lr_validate_y = lr_validate['cvt_per_day']

```

```

alphas = np.linspace(1,201,num=500)
scores = np.empty_like(alphas)
opt_a = float('-inf')
max_score = float('-inf')
for i,a in enumerate(alphas):
    lasso = Lasso()
    lasso.set_params(alpha = a)
    lasso.fit(lr_train_x,lr_train_y)
    scores[i] = lasso.score(lr_validate_x,lr_validate_y)
    if scores[i] > max_score:
        max_score = scores[i]
        opt_a = a
        lasso_save = lasso
plt.plot(alphas, scores, color='b', linestyle='dashed', marker='o',markerfacecolor='blue', markersize=6)
plt.xlabel('alpha')
plt.ylabel('score')
plt.grid(True)
plt.title('score vs. alpha')
plt.show()
model1_para = opt_a
print ('The optimized alpha and score of Lasso linear is: '), opt_a, max_score

```



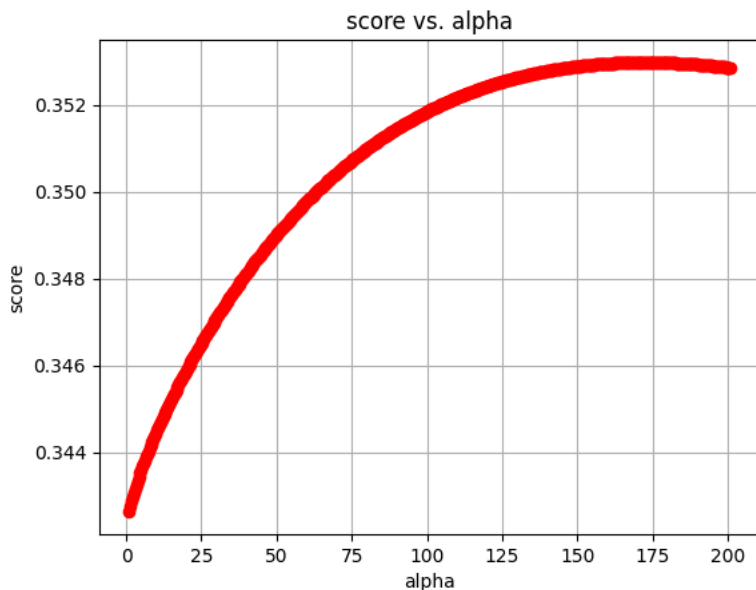
The optimaized alpha and score of Lasso linear is:
(None, 125.64929859719439, 0.3646194368519905)

```
lasso_f = Lasso()
lasso_f.set_params(alpha = opt_a)
lasso_f.fit(model_train_x,model_train_y)
```

```
▼ Lasso
Lasso(alpha=125.64929859719439)
```

```
lr_train, lr_validate = train_test_split(train, test_size=0.15, random_state = 0)
```

```
alphas = np.linspace (1, 201, num=500)
scores = np.empty_like(alphas)
opt_a = float('-inf')
max_score = float('-inf')
for i, a in enumerate(alphas):
    ridge = Ridge()
    ridge.set_params(alpha = a)
    ridge.fit(lr_train_x, lr_train_y)
    scores[i] = ridge.score(lr_validate_x, lr_validate_y)
    if scores[i] > max_score:
        max_score = scores[i]
        opt_a = a
    ridge_save = ridge #optional code
plt.plot(alphas, scores, color='r', linestyle='dashed', marker='o',markerfacecolor='r', markersize=6)
plt.xlabel('alpha')
plt.ylabel('score')
plt.grid(True)
plt.title('score vs. alpha')
plt.show()
model3_para = opt_a
print ('The optimaized alpha and score of Ridge linear is: '), opt_a, max_score
```



The optimized alpha and score of Ridge linear is:
(None, 172.5430861723447, 0.3529726994066543)

```
ridge_f = Ridge()
ridge_f.set_params(alpha = opt_a)
ridge_f.fit(model_train_x, model_train_y)
```

```
▼ Ridge
Ridge(alpha=172.5430861723447)
```

#Random Forest

```
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import GridSearchCV
rf=RandomForestRegressor(random_state=2,max_features='sqrt')
param_grid={'n_estimators': [55,56,57,58,59,60,61,62,63,64,65], 'max_depth': [15,16,17,18,19,20,21]}
clf=GridSearchCV(estimator=rf,param_grid=param_grid,cv=5)
clf.fit(model_train_x,model_train_y)
```

```
► GridSearchCV
► estimator: RandomForestRegressor
  ► RandomForestRegressor
```

```
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import GridSearchCV
rf=RandomForestRegressor(random_state=2,max_features='sqrt')
param_grid={'n_estimators': [55,56,57,58,59,60,61,62,63,64,65], 'max_depth': [15,16,17,18,19,20,21]}
clf=GridSearchCV(estimator=rf,param_grid=param_grid,cv=5)
clf.fit(model_train_x,model_train_y)
```

```
► GridSearchCV
► estimator: RandomForestRegressor
  ► RandomForestRegressor
```



```

from sklearn.linear_model import Lasso
from sklearn.metrics import mean_squared_error
lasso = Lasso(alpha = model1_para)
lasso.fit(train_x,train_y)
pred_y = lasso.predict(test_x)
lasso_score = lasso.score(test_x,test_y)
MSE_lasso = mean_squared_error(test_y,pred_y)
RMSE_lasso=np.sqrt(MSE_lasso)
print ('lasso score: ', ridge_score)
print ('Mean square error of ridge: ', MSE_ridge)
print ('Root mean squared error of ridge:', RMSE_ridge)

lasso score: 0.1138238752321652
Mean square error of ridge: 235165131.2526933
Root mean squared error of ridge: 15335.094758516925

False. False. False. False. False. False. False. False.

```

```

from sklearn.metrics import mean_squared_error
ridge=Ridge(alpha=model3_para)
ridge.fit(train_x,train_y)
pred_y=ridge.predict(test_x)
ridge_score=ridge.score(test_x,test_y)
MSE_ridge=mean_squared_error(test_y,pred_y)
RMSE_ridge=np.sqrt(MSE_ridge)
print ('ridge score: ', ridge_score)
print ('Mean square error of ridge: ', MSE_ridge)
print ('Root mean squared error of ridge:', RMSE_ridge)

```

```

ridge score: 0.1138238752321652
Mean square error of ridge: 235165131.2526933
Root mean squared error of ridge: 15335.094758516925

```

```

0.12125232 0.12004206 0.12085472 0.12110635 0.12277452

```

```

from sklearn.ensemble import RandomForestRegressor
rf=RandomForestRegressor(n_estimators=clf.best_params_['n_estimators'],max_depth=clf.best_params_['max_depth'],max_features='sqr
rf.fit(train_x,train_y)
pred_y=rf.predict(test_x)
rf_score=rf.score(test_x,test_y)
MSE_rf=mean_squared_error(test_y,pred_y)
RMSE_rf=np.sqrt(MSE_rf)
print ('rf score: ', rf_score)
print ('Mean square error of rf: ', MSE_rf)
print ('Root mean squared error of rf:', RMSE_rf)

```

```

rf score: 0.5169719005665712
Mean square error of rf: 128181479.08437566
Root mean squared error of rf: 11321.725976386095

```

```

0.50000000, 0.50000000, 0.50000000, 0.50000000, 0.50000000,
0.50000000, 0.50000000, 0.50000000, 0.50000000, 0.50000000

```

```
#Model COMparison
```

```

0.50000000, 0.50000000, 0.50000000, 0.50000000, 0.50000000,
0.50000000, 0.50000000, 0.50000000, 0.50000000, 0.50000000

```

```

lst_score = [lasso_score, ridge_score, rf_score]
MSE_lst = [MSE_lasso, MSE_ridge, MSE_rf]
RMSE_lst = [RMSE_lasso, RMSE_ridge, RMSE_rf]
model_lst = ['Lasso_linear', 'Ridge linear', 'Random forest']

```

```

plt.figure(1)
plt.plot(model_lst, lst_score, 'ro')
plt.legend(['r-squre / score'])
plt.xlabel('model names', fontsize =16)
plt.ylabel('score / r square', fontsize =16)
plt.grid(True)
plt.show()

```

```

plt.figure(2)
plt.plot(model_lst, MSE_lst, 'g^')
plt.legend(['mean square error (MSE)'])
plt.xlabel('model names', fontsize =16)
plt.ylabel('mean square error', fontsize =16)
plt.grid(True)
plt.show()

```

```

plt.figure(3)
plt.plot(model_lst, RMSE_lst, 'bs')
plt.legend(['root mean square error (RMSE)'])
plt.xlabel('model names', fontsize =16)
plt.ylabel('root mean square error', fontsize =16)
plt.grid(True)
plt.show()

```

