

Data Description

GroupLens Research has collected and made available rating data sets from the MovieLens web site (<http://movielens.org>). The data sets were collected over various periods of time, depending on the size of the set.

100,000 ratings and 3,600 tag applications applied to 9,000 movies by 600 users. Last updated 9/2018.

```
In [34]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import statmodels.api as sm
import datetime

%matplotlib inline
sns.set_style('dark')
sns.set(font_scale=1.2)

from sklearn.linear_model import LinearRegression
from sklearn.model_selection import cross_val_score, train_test_split, GridSearchCV, RandomizedSearchCV
from sklearn.preprocessing import LabelEncoder, StandardScaler, MinMaxScaler, OneHotEncoder
from sklearn.metrics import confusion_matrix, classification_report, mean_absolute_error, mean_squared_error, r2_score
from sklearn.metrics import plot_confusion_matrix, plot_precision_recall_curve, plot_roc_curve, accuracy_score
from sklearn.metrics import auc, f1_score, precision_score, recall_score, roc_auc_score

import warnings
warnings.filterwarnings('ignore')

pd.set_option('display.max_columns', None)
#pd.set_option('display.max_rows', None)
np.random.seed(0)
np.set_printoptions(suppress=True)
```

```
In [2]: movies = pd.read_csv("movies.csv")
```

```
In [3]: ratings = pd.read_csv("ratings.csv")
```

```
In [4]: movies.head()
```

```
Out[4]:
```

movieId	title	genres
0	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	Jumanji (1995)	Adventure Children Fantasy
2	Grimpier Old Men (1995)	Comedy Romance
3	Waiting to Exhale (1995)	Comedy Drama Romance
4	Father of the Bride Part II (1995)	Comedy

```
In [5]: ratings.head()
```

```
Out[5]:
```

userId	movieId	rating	timestamp
0	1	4.0	964982703
1	1	3.0	964981247
2	1	6.0	964982224
3	1	4.7	5.0
4	1	5.0	5.0

```
In [6]: df = pd.merge(ratings,movies, on='movieId')
```

```
In [7]: df
```

```
Out[7]:
```

userId	movieId	rating	timestamp	title	genres
0	1	4.0	964982703	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	5	1.0	964981247	Jumanji (1995)	Adventure Children Fantasy
2	7	1.0	4.5	Grimpier Old Men (1995)	Comedy Romance
3	15	1.0	10.0635946	Waiting to Exhale (1995)	Comedy Drama Romance
4	17	1.0	4.5	Father of the Bride Part II (1995)	Comedy
...
100831	610	160347	4.5	Bloodmoon (1997)	Action Thriller
100832	610	160527	4.5	Sympathy for the Underdog (1991)	Action Crime Drama
100833	610	160836	3.0	Hazard (2005)	Action Drama Thriller
100834	610	163037	3.5	Blair Witch (2016)	Horror Thriller
100835	610	163981	3.5	31 (2016)	Horror

```
100836 rows × 6 columns
```

```
In [8]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 100836 entries, 0 to 100835
Data columns (total: 6 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   userId      100836 non-null  int64  
 1   movieId     100836 non-null  int64  
 2   rating      100836 non-null  float64 
 3   timestamp   100836 non-null  int64  
 4   title       100836 non-null  object  
 5   genres      100836 non-null  object  
dtypes: float64(4), int64(3), object(2)
memory usage: 5.4+ MB
```

```
In [9]: df.groupby(['title'])['rating'].mean().sort_values(ascending=False).head()
```

```
Out[9]:
```

title	rating
Winter Returns (1970)	5.0
Winter in Prostokvashino (1984)	5.0
My Love (2006)	5.0
Sacrifice House Massacre II (1990)	5.0
Winnie the Pooh and the Day of Concern (1972)	5.0

```
Name: rating, dtype: float64
```

```
In [10]: df.groupby(['title'])['rating'].count().sort_values(ascending=False).head()
```

```
Out[10]:
```

title	rating
Forrest Gump (1994)	329
Shawshank Redemption, The (1994)	317
Pulp Fiction (1994)	307
Silence of the Lambs, The (1991)	279
Matrix, The (1999)	278

```
Name: rating, dtype: int64
```

```
In [11]: rating = pd.DataFrame(df.groupby(['title'])['rating'].mean())
```

```
rating
```

```
Out[11]:
```

rating
71 (2014) 4.000000
'Hellboy': The Seeds of Creation (2004) 4.000000
'Round Midnight (1986) 3.500000
'Salem's Lot (2004) 5.000000
'Til There Was You (1997) 4.000000
...
eXistenZ (1999) 3.863636
xXx (2002) 2.770833
xXx: State of the Union (2005) 2.000000
Three Amigos! (1986) 3.134615
À nous la liberté (Freedom for Us) (1931) 1.000000

```
9719 rows × 1 columns
```

```
In [12]: rating['n_ratings'] = pd.DataFrame(df.groupby(['title'])['rating'].count())
```

```
rating
```

```
Out[12]:
```

rating	n_ratings
71 (2014) 4.000000	1
'Hellboy': The Seeds of Creation (2004) 4.000000	1
'Round Midnight (1986) 3.500000	2
'Salem's Lot (2004) 5.000000	1
'Til There Was You (1997) 4.000000	2
...	...
eXistenZ (1999) 3.863636	22
xXx (2002) 2.770833	24
xXx: State of the Union (2005) 2.000000	5
Three Amigos! (1986) 3.134615	26
À nous la liberté (Freedom for Us) (1931) 1.000000	1

```
9719 rows × 2 columns
```

```
In [13]: rating.hist(figsize=(20,5), bins=30)
```

```
plt.show()
```

```
9719 rows × 2 columns
```

```
In [14]: sns.jointplot(x='rating', y='n_ratings', data=rating);
```

```
9719 rows × 2 columns
```

```
In [15]: df.head()
```

```
Out[15]:
```

userId	movieId	rating	timestamp	title	genres
0	1	1	4.0	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	5	1	4.5	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
2	7	1	4.5	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
3	15	1	2.5	Til 105077970	Toy Story (1995)
4	17	1	4.5	130566483	Toy Story (1995)

```
In [16]: rating_matrix = df.pivot_table(values='rating', index='userId', columns='title')
```

```
In [17]: rating_matrix.head()
```

```
Out[17]:
```

title	'Hellboy': The Seeds of Creation (2004)	'Round Midnight (1986)	'Salem's Lot (2004)	'Til There Was You (1997)	'Burbs, The (1989)	'Burbs, The (1995)	'Burbs, The (1998)	'Burbs, The (2009)	'Burbs, The (2015)	'Burbs, The (2016)	'Burbs, The (2017)	'Burbs, The (2018)	'Burbs, The (2019)	'Burbs, The (2020)	'Burbs, The (2021)	'Burbs, The (2022)	'Burbs, The (2023)	'Burbs, The (2024)	'Burbs, The (2025)	'Burbs, The (2026)	'Burbs, The (2027)	'Burbs, The (2028)	'Burbs, The (2029)	'Burbs, The (2030)	'Burbs, The (2031)	'Burbs, The (2032)	'Burbs, The (2033)	'Burbs, The (2034)	'Burbs, The (2035)	'Burbs, The (2036)	'Burbs, The (2037)	'Burbs, The (2038)	'Burbs, The (2039)	'Burbs, The (2040)	'Burbs, The (2041)	'Burbs, The (2042)	'Burbs, The (2043)	'Burbs, The (2044)	'Burbs, The (2045)	'Burbs, The (2046)	'Burbs, The (2047)	'Burbs, The (2048)	'Burbs, The (2049)	'Burbs, The (2050)	'Burbs, The (2051)	'Burbs, The (2052)	'Burbs, The (2053)	'Burbs, The (2054)	'Burbs, The (2055)	'Burbs, The (2056)	'Burbs, The (2057)	'Burbs, The (2058)	'Burbs, The (2059)	'Burbs, The (2060)	'Burbs, The (2061)	'Burbs, The (2062)	'Burbs, The (2063)	'Burbs, The (2064)	'Burbs, The (2065)	'Burbs, The (2066)	'Burbs, The (2067)	'Burbs, The (2068)	'Burbs, The (2069)	'Burbs, The (2070)	'Burbs, The (2071)	'Burbs, The (2072)	'Burbs, The (2073)	'Burbs, The (2074)	'Burbs, The (2075)	'Burbs, The (2076)	'Burbs, The (2077)	'Burbs, The (2078)	'Burbs, The (2079)	'Burbs, The (2080)	'Burbs, The (2081)	'Burbs, The (2082)	'Burbs, The (2083)	'Burbs, The (2084)	'Burbs, The (2085)	'Burbs, The (2086)	'Burbs, The (2087)	'Burbs, The (2088)	'Burbs, The (2089)	'Burbs, The (2090)	'Burbs, The (2091)	'Burbs, The (2092)	'Burbs, The (2093)	'Burbs, The (2094)	'Burbs, The (2095)	'Burbs, The (2096)	'Burbs, The (2097)	'Burbs, The (2098)	'Burbs, The (2099)	'Burbs, The (2010)	'Burbs, The (2011)	'Burbs, The (2012)	'Burbs, The (2013)	'Burbs, The (2014)	'Burbs, The (2015)	'Burbs, The (2016)	'Burbs, The (2017)	'Burbs, The (2018)	'Burbs, The (2019)	'Burbs, The (2020)	'Burbs, The (2021)	'Burbs, The (2022)	'Burbs, The (2023)	'Burbs, The (2024)	'Burbs, The (2025)	'Burbs, The (2026)	'Burbs, The (2027)	'Burbs, The (2028)	'Burbs, The (2029)	'Burbs, The (2030)	'Burbs, The (2031)	'Burbs, The (2032)	'Burbs, The (2033)	'Burbs, The (2034)	'Burbs, The (2035)	'Burbs, The (2036)	'Burbs, The (2037)	'Burbs, The (2038)	'Burbs, The (2039)	'Burbs, The (2040)	'Burbs, The (2041)	'Burbs, The (2042)	'Burbs, The (2043)	'Burbs, The (2044)	'Burbs, The (2045)	'Burbs, The (2046)	'Burbs, The (2047)	'Burbs, The (2048)	'Burbs, The (2049)	'Burbs, The (2050)	'Burbs, The (2051)	'Burbs, The (2052)	'Burbs, The (2053)	'Burbs, The (2054)	'Burbs, The (2055)	'Burbs, The (2056)	'Burbs, The (2057)	'Burbs, The (2058)	'Burbs, The (2059)	'Burbs, The (2060)	'Burbs, The (2061)	'Burbs, The (2062)	'Burbs, The (2063)	'Burbs, The (2064)	'Burbs, The (2065)	'Burbs, The (2066)	'Burbs, The (2067)	'Burbs, The (2068)	'Burbs, The (2069)	'Burbs, The (2070)	'Burbs, The (2071)	'Burbs, The (2072)	'Burbs, The (2073)	'Burbs, The (2074)	'Burbs, The (2075)	'Burbs, The (2076)	'Burbs, The (2077)	'Burbs, The (2078)	'Burbs, The (2079)	'Burbs, The (2080)	'Burbs, The (2081)	'Burbs, The (2082)	'Burbs, The (2083)	'Burbs, The (2084)	'Burbs, The (2085)	'Burbs, The (2086)	'Burbs, The (2087)	'Burbs, The (2088)	'Burbs, The (2089)	'Burbs, The (2090)	'Burbs, The (2091)	'Burbs, The (2092)	'Burbs, The (2093)	'Burbs, The (2094)	'Burbs, The (2095)	'Burbs, The (2096)	'Burbs, The (2097)	'Burbs, The (2098)	'Burbs, The (2099)	'Burbs, The (2010)	'Burbs, The (2011)	'Burbs, The (2012)	'Burbs, The (2013)	'Burbs, The (2014)	'Burbs, The (2015)	'Burbs, The (2016)	'Burbs, The (2017)	'Burbs, The (2018)	'Burbs, The (2019)	'Burbs, The (2020)	'Burbs, The (2021)	'Burbs, The (2022)	'Burbs, The (2023)	'Burbs, The (2024)	'Burbs, The (2025)	'Burbs, The (2026)	'Burbs, The (2027)	'Burbs, The (2028)	'Burbs, The (2029)	'Burbs, The (2030)	'Burbs, The (2031)	'Burbs, The (2032)	'Burbs, The (2033)	'Burbs, The (2034)	'Burbs, The (2035)	'Burbs, The (2036)	'Burbs, The (2037)	'Burbs, The (2038)	'Burbs, The (2039)	'Burbs, The (2040)	'Burbs, The (2041)	'Burbs, The (2042)	'Burbs, The (2043)	'Burbs, The (2044)	'Burbs, The (2045)	'Burbs, The (2046)	'Burbs, The (2047)	'Burbs, The (2048)	'Burbs, The (2049)	'Burbs, The (2050)	'Burbs, The (2051)	'Burbs, The (2052)	'Burbs, The (2053)	'Burbs, The (2054)	'Burbs, The (2055)	'Burbs, The (2056)	'Burbs, The (2057)	'Burbs, The (2058)	'Burbs, The (2059)	'Burbs, The (2060)	'Burbs, The (2061)	'Burbs, The (2062)	'Burbs, The (2063)	'Burbs, The (2064)	'Burbs, The (2065)	'Burbs, The (2066)	'Burbs, The (2067)	'Burbs, The (2068)	'Burbs, The (2069)	'Burbs, The (2070)	'Burbs, The (2071)	'Burbs, The (2072)	'Burbs, The (2073)	'Burbs, The (2074)	'Burbs, The (2075)	'Burbs, The (2076)	'Burbs, The (2077)	'Burbs, The (2078)	'Burbs, The (2079)	'Burbs, The (2080)	'Burbs, The (2081)	'Burbs, The (2082)	'Burbs, The (2083)	'Burbs, The (2084)	'Burbs, The (2085)	'Burbs, The (2086)	'Burbs, The (2087)	'Burbs, The (2088)	'Burbs, The (2089)	'Burbs, The (2090)	'Burbs, The (2091)	'Burbs, The (2092)	'Burbs, The (2093)	'Burbs, The (2094)	'Burbs, The (2095)	'Burbs, The (2096)	'Burbs, The (2097)	'Burbs, The (2098)	'Burbs, The (2099)	'Burbs, The (2010)	'Burbs, The (2011)	'Burbs, The (2012)	'Burbs, The (2013)	'Burbs, The (2014)	'Burbs, The (2015)	'Burbs, The (2016)	'Burbs, The (2017)	'Burbs, The (2018)	'Burbs, The (2019)	'Burbs, The (2020)	'Burbs, The (2021)	'Burbs, The (2022)	'Burbs, The (2023)	'Burbs, The (2024)	'Burbs, The (2025)	'
-------	---	------------------------	---------------------	---------------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	--------------------	---