

```
2306 AML117 MAHESH BABU M ASSIGNMENT 15 OR ML8
```

```
In [1]: import numpy as np
import pandas as pd
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.metrics import classification_report
import matplotlib.pyplot as plt
import seaborn as sns

In [5]: data=pd.read_csv('data.csv')
data

Out[5]:
```

Table with 18 columns: valence, year, acousticness, artists, danceability, duration\_ms, energy, explicit, id, instrumentalness, key, liveness, loudness, mode, popularity, r. It shows the first 5 rows of data.

```
In [6]: data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 176653 entries, 0 to 176652
Data columns (total 19 columns):
 #   Column                               Non-Null Count  Dtype
---  --
 0   valence                               176653 non-null float64
 1   year                                 176653 non-null int64
 2   acousticness                         176653 non-null float64
 3   artists                              176653 non-null object
 4   danceability                         176653 non-null float64
 5   duration_ms                          176653 non-null int64
 6   energy                               176653 non-null float64
 7   explicit                             176653 non-null int64
 8   id                                    176653 non-null object
 9   instrumentalness                    176653 non-null float64
10   key                                   176653 non-null int64
11   liveness                            176653 non-null float64
12   loudness                            176653 non-null float64
13   mode                                 176653 non-null int64
14   name                                 176653 non-null object
15   popularity                           176653 non-null int64
16   release_date                       176653 non-null object
17   speechiness                         176653 non-null float64
18   tempo                               176653 non-null float64
dtypes: float64(9), int64(6), object(4)
memory usage: 24.7+ MB

In [7]: import seaborn as sns
In [8]: sns.heatmap(data.corr(),annot=True)
```



```
In [9]: artist=pd.read_csv('data_by_artist.csv')
artist

Out[9]:
```

Table with 15 columns: mode, count, acousticness, artists, danceability, duration\_ms, energy, instrumentalness, liveness, loudness, speechiness, tempo, valence, popularity, key. It shows the first 5 rows of data for artists.

```
In [10]: year=pd.read_csv('data_by_year.csv')
year

Out[10]:
```

Table with 14 columns: mode, year, acousticness, danceability, duration\_ms, energy, instrumentalness, liveness, loudness, speechiness, tempo, valence, popularity, key. It shows the first 5 rows of data for years.

```
In [11]: genres=pd.read_csv('data_w_genres.csv')
genres

Out[11]:
```

Table with 16 columns: genres, artists, acousticness, danceability, duration\_ms, energy, instrumentalness, liveness, loudness, speechiness, tempo, valence, popularity, key, mode, count. It shows the first 5 rows of data for genres.

```
In [12]: print(artist.size)
```

430200

```
In [14]: print(artist,artist.size)
print(artist,artist.shape)
print(artist,artist.columns)
print('year',year.size)
print('year',year.shape)
print('year',year.columns)
print('genres',genres.size)
print('genres',genres.shape)
print('genres',genres.columns)
```

artist 430200
artist Index(['mode', 'count', 'acousticness', 'artists', 'danceability', 'duration\_ms', 'energy', 'instrumentalness', 'liveness', 'loudness', 'speechiness', 'tempo', 'valence', 'popularity', 'key'], dtype='object')

year 100
year Index(['mode', 'year', 'acousticness', 'danceability', 'duration\_ms', 'energy', 'instrumentalness', 'liveness', 'loudness', 'speechiness', 'tempo', 'valence', 'popularity', 'key'], dtype='object')

genres 458800
genres Index(['genres', 'artists', 'acousticness', 'danceability', 'duration\_ms', 'energy', 'instrumentalness', 'liveness', 'loudness', 'speechiness', 'tempo', 'valence', 'popularity', 'key', 'mode', 'count'], dtype='object')

```
In [15]: artist.info()
year.info()
genres.info()
```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 28680 entries, 0 to 28679
Data columns (total 15 columns):
 # Column Non-Null Count Dtype
--- --
 0 mode 28680 non-null int64
 1 count 28680 non-null object
 2 acousticness 28680 non-null float64
 3 artists 28680 non-null object
 4 danceability 28680 non-null float64
 5 duration\_ms 28680 non-null int64
 6 energy 28680 non-null float64
 7 instrumentalness 28680 non-null float64
 8 liveness 28680 non-null float64
 9 loudness 28680 non-null float64
10 speechiness 28680 non-null float64
11 tempo 28680 non-null float64
12 valence 28680 non-null float64
13 popularity 28680 non-null float64
14 key 28680 non-null int64
dtypes: float64(11), int64(3), object(1)
memory usage: 3.3+ MB

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 14 columns):
 # Column Non-Null Count Dtype
--- --
 0 mode 100 non-null int64
 1 year 100 non-null int64
 2 acousticness 100 non-null float64
 3 danceability 100 non-null float64
 4 duration\_ms 100 non-null int64
 5 energy 100 non-null float64
 6 instrumentalness 100 non-null float64
 7 liveness 100 non-null float64
 8 loudness 100 non-null float64
 9 speechiness 100 non-null float64
10 tempo 100 non-null float64
11 valence 100 non-null float64
12 popularity 100 non-null float64
13 key 100 non-null int64
dtypes: float64(11), int64(3), object(2)
memory usage: 3.5+ MB

```
In [19]: artist.dtypes
year.dtypes
genres.dtypes
```

genres object
artists object
acousticness float64
danceability float64
duration\_ms float64
energy float64
instrumentalness float64
liveness float64
loudness float64
speechiness float64
tempo float64
popularity float64
key int64
mode int64
dtype: object

```
In [21]: artist['artists'].unique()
```

array(['Cats' 1981 Original London Cast', 'Cats' 1983 Broadway Cast', 'Fiddler On The Roof' Motion Picture Chorus', ...], dtype=object)

```
In [22]: year['year'].unique()
```

array([1921, 1922, 1923, 1924, 1925, 1926, 1927, 1928, 1929, 1930, 1931, 1932, 1933, 1934, 1935, 1936, 1937, 1938, 1939, 1940, 1941, 1942, 1943, 1944, 1945, 1946, 1947, 1948, 1949, 1950, 1951, 1952, 1953, 1954, 1955, 1956, 1957, 1958, 1959, 1960, 1961, 1962, 1963, 1964, 1965, 1966, 1967, 1968, 1969, 1970, 1971, 1972, 1973, 1974, 1975, 1976, 1977, 1978, 1979, 1980, 1981, 1982, 1983, 1984, 1985, 1986, 1987, 1988, 1989, 1990, 1991, 1992, 1993, 1994, 1995, 1996, 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020], dtype=int64)

```
In [23]: genres['genres'].unique()
```

array(['show tunes', ''], ['comedy rock', 'comic', 'parody'], ..., ['mainland chinese pop', 'zhongguo feng'], ['c-pop', 'classic mandopop', 'mainland chinese pop', 'mandopop'], ['chinese indie', 'chinese indie rock'], dtype=object)

```
In [33]: features=['acousticness', 'danceability', 'duration_ms', 'energy', 'instrumentalness', 'liveness', 'loudness', 'speechiness', 'tempo', 'valence', 'key', 'mode']
target = artist['popularity']
labels = artist['popularity']
```

```
In [34]: from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import GradientBoostingClassifier
```

```
In [35]: models = []
models.append(('KNN', KNeighborsClassifier()))
models.append(('SVC', SVC()))
models.append(('LR', LogisticRegression()))
models.append(('DT', DecisionTreeClassifier()))
models.append(('NB', GaussianNB()))
models.append(('RF', RandomForestClassifier()))
models.append(('GB', GradientBoostingClassifier()))
```

```
In [58]: from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
```

```
In [59]: from sklearn.datasets import make_classification
from sklearn.overfitting import SMOTE
```

```
# Create a toy imbalanced dataset for demonstration
X, y = make_classification(n_classes=2, class_sep=2, weights=[0.1, 0.9], n_informative=3, n_redundant=1, flip_y=0, n_features=20, n_clusters_per_class=1, n_samples=1000, random_state=42)
```

```
# Instantiate the SMOTE object
smote = SMOTE(sampling_strategy='auto', random_state=42)
```

```
# Fit and transform the dataset to generate synthetic samples
X_resampled, y_resampled = smote.fit_resample(X, y)
```

```
# Check the class distribution after SMOTE
import numpy as np
unique, counts = np.unique(y_resampled, return_counts=True)
print(dict(zip(unique, counts)))
```

0: 900, 1: 100

```
In [60]: X1=X_resampled
y1=y_resampled
```

```
In [61]: X_train, X_test, y_train, y_test = train_test_split(X1, y1, stratify = y1, random_state=0)
```

```
In [62]: names = []
for name, model in models:
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    scores.append(accuracy_score(y_test, y_pred))
names.append(name)
tr_split = pd.DataFrame({'Name': names, 'Score': scores})
print(tr_split)
```

Table with 2 columns: Name, Score. It shows the accuracy scores for each model.

```
In [63]: acc = []
model = []
```

```
In [64]: !pip install xgboost
```

Requirement already satisfied: xgboost in c:\anaconda\lib\site-packages (1.7.6)
Requirement already satisfied: scipy in c:\anaconda\lib\site-packages (from xgboost) (1.10.0)
Requirement already satisfied: numpy in c:\anaconda\lib\site-packages (from xgboost) (1.23.5)

```
In [65]: from sklearn import metrics
```

```
In [68]: import xgboost as xgb
XB = xgb.XGBClassifier()
XB.fit(X_train,y_train)
```

predicted\_values = XB.predict(X\_test)

```
x = metrics.accuracy_score(y_test, predicted_values)
acc.append(x)
model.append('XGBoost')
print("XGBoost's Accuracy is: ", x)
```

print(classification\_report(y\_test,predicted\_values))

Classification report for XGBoost showing precision, recall, f1-score, and support for each class.

```
In [75]: from sklearn.ensemble import RandomForestClassifier
RF = RandomForestClassifier(n_estimators=20, random_state=0)
RF.fit(X_train,y_train)
```

predicted\_values = RF.predict(X\_test)

```
x = metrics.accuracy_score(y_test, predicted_values)
acc.append(x)
model.append('RF')
print("RF's Accuracy is: ", x)
```

print(classification\_report(y\_test,predicted\_values))

Classification report for RandomForestClassifier showing precision, recall, f1-score, and support for each class.

```
In [77]: accuracy_models = dict(zip(model, acc))
for k, v in accuracy_models.items():
    print(k, '-->', v)
```

XGBoost --> 1.0
RF --> 1.0

```
In [78]: data = np.array([[1, 9, 0.5981, 0.46, 290318.55, 0.394063333, 0.811399851, 0.29083333, -14.448, 0.210388889, 117.518111, 0.3895, 38.333, 5, 20, 3, 5, 12, 0.5]])
prediction = RF.predict(data)
print(prediction)
```

[1]

```
In [ ]:
```