

LetsGrowMore (LGM) - Virtual Internship Program

Name :- KRUSHNA DEVENDRA HELGE

Task - Music Recommendation:

Music recommender systems can suggest songs to users based on their listening patterns.

Importing Packages

In []:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

Load the Dataset

In []:

```
1 ntr = 7000
2 nts = 3000
3 data_path = 'C:\\Users\\Shivv"'
4 train = pd.read_csv(data_path + 'train.csv',nrows=ntr , encoding = 'ISO-8859-1')
5 names=['msno','song_id','source_system_tab','source_screen_name',\
6        'source_type','target']
7 test1 = pd.read_csv(data_path+'train.csv',names=names,skiprows=ntr,nrows=nts)
8 songs = pd.read_csv(data_path + 'songs.csv')
9 members = pd.read_csv(data_path + 'members.csv')
```

In []:

```
1 train.head()
```

In []:

```
1 songs.head()
```

DATA VISUALIZATION

In []:

```
1 sns.countplot(x = train['source_system_tab'],hue=train['source_system_tab'])
```

In []:

```
1 sns.countplot(x = train['source_system_tab'],hue=train['target'])
```

In []:

```
1 sns.countplot(x = train['source_screen_name'],hue=train['target'],data = train,orient='v')
2 plt.xticks(rotation =90)
3 plt.show()
```

In []:

```
1 sns.countplot(x = train['source_type'],hue=train['source_type'],data = train,orient='v')
2 plt.xticks(rotation =90)
3 plt.show()
```

In []:

```
1 sns.countplot(x = train['source_type'],hue=train['target'],data = train,orient='v')
2 plt.xticks(rotation =90)
3 plt.show()
```

In []:

```
1 sns.countplot(x = songs['language'],data =train,hue=songs['language'],orient='v')
```

In []:

```
1 sns.countplot(x = members['registered_via'],hue=members['registered_via'],orient='v')
2 plt.xticks(rotation =90)
3 plt.show()
```

DATA PREPROCESSING AND CLEANING

In []:

```
1 test = test1.drop(['target'],axis=1)
2 ytr = np.array(test1['target'])
```

In []:

```
1 test_name = ['id','msno','song_id','source_system_tab',\
2             'source_screen_name','source_type']
3 test['id']=np.arange(nts)
4 test = test[test_name]
```

In []:

```
1 song_cols = ['song_id', 'artist_name', 'genre_ids', 'song_length', 'language']
2 train = train.merge(songs[song_cols], on='song_id', how='left')
3 test = test.merge(songs[song_cols], on='song_id', how='left')
```

In []:

```
1 members['registration_year'] = members['registration_init_time'].apply(lambda x: int(x))
2 members['registration_month'] = members['registration_init_time'].apply(lambda x: int(x))
3 members['registration_date'] = members['registration_init_time'].apply(lambda x: int(x))
```

In []:

```
1 members = members.drop(['registration_init_time'], axis=1)
2 members_cols = members.columns
3 train = train.merge(members[members_cols], on='msno', how='left')
4 test = test.merge(members[members_cols], on='msno', how='left')
```

In []:

```
1 members_cols = members.columns
2 train = train.merge(members[members_cols], on='msno', how='left')
3 test = test.merge(members[members_cols], on='msno', how='left')
```

In []:

```
1 train = train.fillna(-1)
2 test = test.fillna(-1)
```

In []:

```
1 import gc
2 del members, songs; gc.collect();
```

In []:

```
1 cols = list(train.columns)
2 cols.remove('target')
```

MODEL BUILDING

In []:

```
1 from tqdm import tqdm
2 from sklearn.preprocessing import LabelEncoder
3 for col in tqdm(cols):
4     if train[col].dtype == 'object':
5         train[col] = train[col].apply(str)
6         test[col] = test[col].apply(str)
7
8     le = LabelEncoder()
9     train_vals = list(train[col].unique())
10    test_vals = list(test[col].unique())
11    le.fit(train_vals + test_vals)
12    train[col] = le.transform(train[col])
13    test[col] = le.transform(test[col])
```

TRYING OUT BASIC CLASSIFICATION MODELS

In []:

```
1 unique_songs = range(max(train['song_id'].max(), test['song_id'].max()))
2 song_popularity = pd.DataFrame({'song_id': unique_songs, 'popularity':0})
3
4 train_sorted = train.sort_values('song_id')
5 train_sorted.reset_index(drop=True, inplace=True)
6 test_sorted = test.sort_values('song_id')
7 test_sorted.reset_index(drop=True, inplace=True)
```

In []:

```
1 !pip install lightgbm
```

In []:

```
1 from sklearn.model_selection import train_test_split
2 import lightgbm as lgb
3 X = np.array(train.drop(['target'], axis=1))
4 y = train['target'].values
5
6 X_test = np.array(test.drop(['id'], axis=1))
7 ids = test['id'].values
8
9 del train, test; gc.collect();
10
11 X_train, X_valid, y_train, y_valid = train_test_split(X, y, \
12     test_size=0.1, random_state = 12)
13
14 del X, y; gc.collect();
15
16 d_train = lgb.Dataset(X_train, label=y_train)
17 d_valid = lgb.Dataset(X_valid, label=y_valid)
18
19 watchlist = [d_train, d_valid]
```

In []:

```

1 def predict(m1_model):
2     model = m1_model.fit(X_train,y_train)
3     print('Training Score : {}'.format(model.score(X_train,y_train)))
4     y_pred = model.predict(X_valid)
5     #accuracy_score = m1_model.metrics.accuracy_score(y_valid,y_pred)
6     #print('Accuracy Score : {}'.format(accuracy_score))
7     v_test = model.predict(X_test)
8     yhat = (v_test>0.5).astype(int)
9     comp = (yhat==ytr).astype(int)
10    acc = comp.sum()/comp.size*100
11    print("Accuracy on test data for the model", acc)

```

In []:

```

1 from sklearn.linear_model import LogisticRegression
2 from sklearn.ensemble import RandomForestClassifier
3 predict(LogisticRegression())

```

In []:

```

1 predict(RandomForestClassifier())

```

PREDICTION USING LIGHTGBM

In []:

```

1 params = {}
2 params['learning_rate'] = 0.4
3 params['application'] = 'binary'
4 params['max_depth'] = 15
5 params['num_leaves'] = 2**8
6 params['verbosity'] = 0
7 params['metric'] = 'auc'
8
9 model1 = lgb.train(params, train_set=d_train, num_boost_round=200, valid_sets=watch1
10 early_stopping_rounds=10, verbose_eval=10)

```

In []:

```

1 p_test = model1.predict(X_test)

```

In []:

```

1 yhat = (p_test>0.5).astype(int)
2 comp = (yhat==ytr).astype(int)
3 acc = comp.sum()/comp.size*100
4 print('The accuracy of lgbm model on test data is: {0:f}%'.format(acc))

```

#Conclusion

This brings us to a conclusion that RandomForest and LGM perform very well in test data. But RandomForest in some cases tend to overfit the data. So LigthGBM is better at predicting the music the

