

## Import required libraries

```
In [198]: 1 import warnings  
2 warnings.filterwarnings('ignore')
```

```
In [199]: 1 from surprise.model_selection import RandomizedSearchCV, GridSearchCV, KFold  
2 from surprise import BaselineOnly, KNNBaseline, SVD, SVDpp  
3 from surprise import Reader, Dataset  
4 import xgboost as xgb  
5 import pandas as pd  
6 import numpy as np  
7 from prettytable import PrettyTable  
8 from sklearn.metrics import confusion_matrix  
9 import matplotlib.pyplot as plt  
10 from sklearn.metrics import make_scorer, mean_squared_error  
11 from sklearn.model_selection import GridSearchCV as GridSearch, RandomizedSearchCV as RandomSearch, KFold as KF  
12 from collections import OrderedDict
```

### NOTE:

Must have to use aliasing for sklearn GridSearchCV, RandomizedSearchCV and KFold because we are using same method from surprise library also.

## Reading Train and Test data to dataframe :

### 1.Train\_dataframe

```
In [223]: 1 train_reg = pd.read_csv('reg_train.csv', names = ['user', 'movie', 'GAvg', 'sur1', 'sur2', 'sur3', 'sur4', 'sur5',
2                                         'smr1', 'smr2', 'smr3', 'smr4', 'smr5',
3                                         'UAvg', 'MAvg', 'rating'], header=None)
4 train_reg.head()
```

Out[223]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MAvg	rating
0	174683	10	3.587581	5.0	5.0	3.0	4.0	4.0	3.0	5.0	4.0	3.0	2.0	3.882353	3.611111	5
1	233949	10	3.587581	4.0	4.0	5.0	1.0	3.0	2.0	3.0	2.0	3.0	3.0	2.692308	3.611111	3
2	555770	10	3.587581	4.0	5.0	4.0	4.0	5.0	4.0	2.0	5.0	4.0	4.0	3.795455	3.611111	4
3	767518	10	3.587581	2.0	5.0	4.0	4.0	3.0	5.0	5.0	4.0	4.0	3.0	3.884615	3.611111	5
4	894393	10	3.587581	3.0	5.0	4.0	4.0	3.0	4.0	4.0	4.0	4.0	4.0	4.000000	3.611111	4

- **GAvg** : Average rating of all the ratings
- **Similar users rating of this movie:**
  - sur1, sur2, sur3, sur4, sur5 ( top 5 similar users who rated that movie.. )
- **Similar movies rated by this user:**
  - smr1, smr2, smr3, smr4, smr5 ( top 5 similar movies rated by this movie.. )
- **UAvg** : User's Average rating
- **MAvg** : Average rating of this movie
- **rating** : Rating of this movie by this user.

## 2.Test\_dataframe

```
In [224]: 1 test_reg = pd.read_csv('reg_test.csv', names = ['user', 'movie', 'GAvg', 'sur1', 'sur2', 'sur3', 'sur4', 'sur5',
2                                         'smr1', 'smr2', 'smr3', 'smr4', 'smr5',
3                                         'UAvg', 'MAvg', 'rating'], header=None)
4 test_reg.head(4)
```

Out[224]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MAvg
0	808635	71	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581
1	941866	71	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581
2	1737912	71	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581
3	1849204	71	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581

- **GAvg** : Average rating of all the ratings
- **Similar users rating of this movie:**
  - sur1, sur2, sur3, sur4, sur5 ( top 5 similiar users who rated that movie.. )
- **Similar movies rated by this user:**
  - smr1, smr2, smr3, smr4, smr5 ( top 5 similiar movies rated by this movie.. )
- **UAvg** : User AVerage rating
- **MAvg** : Average rating of this movie
- **rating** : Rating of this movie by this user.

## Transforming data for Surprise models

### 1.Transforming train data

- We can't give raw data (movie, user, rating) to train the model in Surprise library.
- They have a separate format for TRAIN and TEST data, which will be useful for training the models like SVD, KNNBaselineOnly....etc..,in Surprise.
- We can form the trainset from a file, or from a Pandas DataFrame. [http://surprise.readthedocs.io/en/stable/getting\\_started.html#load-dom-dataframe-py](http://surprise.readthedocs.io/en/stable/getting_started.html#load-dom-dataframe-py) ([http://surprise.readthedocs.io/en/stable/getting\\_started.html#load-dom-dataframe-py](http://surprise.readthedocs.io/en/stable/getting_started.html#load-dom-dataframe-py))

In [225]:

```

1 # It is to specify how to read the dataframe.
2 # for our dataframe, we don't have to specify anything extra..
3 reader = Reader(rating_scale=(1,5))
4
5 # create the traindata from the dataframe...
6 train_data = Dataset.load_from_df(train_reg[['user', 'movie', 'rating']], reader)
7
8 # build the trainset from traindata.., It is of dataset format from surprise library..
9 trainset = train_data.build_full_trainset()

```

## 2.Transforming test data

- Testset is just a list of (user, movie, rating) tuples. (Order in the tuple is important)

In [226]:

```

1 testset = list(zip(test_reg.user.values, test_reg.movie.values, test_reg.rating.values))
2 testset[:3]

```

Out[226]:

## Applying Machine Learning models

- Global dictionary that stores rmse and mape for all the models....
  - It stores the metrics in a dictionary of dictionaries

**keys** : model names(string)

```
value: dict(key : metric, value : value )
```

In [227]:

```
1 models_evaluation_train = dict()
2 models_evaluation_test = dict()
3
4 models_evaluation_train, models_evaluation_test
```

Out[227]: ({}, {})

### Function to print parameters summary:

In [205]:

```
1 def param_list(params, para_summ):
2     para_summ.clear_rows()
3     for name, val in zip(params.keys(), params.values()):
4         para_summ.add_row([name, val])
5     print(para_summ)
```

### Function for feature importance:

In [206]:

```
1 def plot_importance(model, clf):
2     fig = plt.figure(figsize = (8, 6))
3     ax = fig.add_axes([0,0,1,1])
4     model.plot_importance(clf, ax = ax, height = 0.3)
5     plt.xlabel("F Score", fontsize = 20)
6     plt.ylabel("Features", fontsize = 20)
7     plt.title("Feature Importance", fontsize = 20)
8     plt.tick_params(labelsize = 15)
9
10    plt.show()
```

### Function required for Surprise Models:

## SURPRISE:

## 1. Function for Hyperparameter Tunning :

```
In [207]: 1 def hyperparamTUNNING(clf, data, params, CV, name, searchMethod='random'):
2     if searchMethod=='random':
3         model=RandomizedSearchCV(clf,\n4                         param_distributions=params,\n5                         n_iter=20,\n6                         n_jobs=16,\n7                         measures=['rmse', 'mae'],\n8                         return_train_measures=True)
9     elif searchMethod=='grid':
10        model=GridSearchCV(clf,\n11                         param_grid=params,\n12                         n_jobs=-1,\n13                         measures=['rmse', 'mae'],\n14                         cv=CV,\n15                         return_train_measures=True)
16    model.fit(data)
17    tr_rmse = model.cv_results['mean_train_rmse']
18    cv_rmse = model.cv_results['mean_test_rmse']
19    for i in range(len(tr_rmse)):
20        print('Train RMSE: %.4f'%tr_rmse[i], 'CV RMSE: %.4f'%cv_rmse[i])
21    # BEST RMSE AND COMBINATION OF OPTIMAL PARAMETERS VALUES TO GET THAT RMSE
22    print('\n\nBest RMSE: %.5f'%model.best_score['rmse'], 'For Optimal Params Combination Value: ')
23    if name=='knn':
24        print('\t[1.]K: %d'%model.best_params['rmse']['k'],\n25              '\n\t[2.]shrinkage: %d'%model.best_params['rmse']['sim_options']['shrinkage'],\n26              '\n\t[3.]lerning_rate: %.4f'%model.best_params['rmse']['bsl_options']['learning_rate'],\n27              '\n\t[4.]Epochs: %d'%model.best_params['rmse']['bsl_options']['n_epochs'])
28    elif name=='bsl':
29        #print(model.best_params)
30        print('\t[1.]Learning-Rate: %.4f'%model.best_params['rmse']['bsl_options']['learning_rate'],\n31              '\n\t[2.]Epochs: %d'%model.best_params['rmse']['bsl_options']['n_epochs'])
32    elif name=='svd':
33        print('\t[1.]N-Factors: %d'%model.best_params['rmse']['n_factors'],\n34              '\n\t[2.]Epochs: %d'%model.best_params['rmse']['n_epochs'],\n35              '\n\t[3.]Learning-Rate: %.4f'%model.best_params['rmse']['lr_all'],\n36              '\n\t[4.]Reg-All: %.4f'%model.best_params['rmse']['reg_all'])
37    return model
```

**2. function for measure rating:**

```
In [208]: 1 def get_ratings(predictions):
2     actual = np.array([pred.r_ui for pred in predictions])
3     predicted = np.array([pred.est for pred in predictions])
4     return actual, predicted
5 #in surprise prediction of every data point is returned as dictionary like this:
6 #{"user": 196, "item": 302, "r_ui": 4.0, "est": 4.06, "actual_k": 40, "was_impossible": False}
7 #In this dictionary, "r_ui" is a key for actual rating and "est" is a key for predicted rating
```

**3. function for error measurement:**

```
In [209]: 1 def get_error(predictions):
2     actual, predicted = get_ratings(predictions)
3     rmse = np.sqrt(mean_squared_error(actual, predicted))
4     mape = np.mean(np.abs((actual - predicted) / actual)) * 100
5     return rmse, mape
```

**4. function for measure performance on test data:**

In [210]:

```
1 my_seed = 15
2 np.random.seed(my_seed)
3
4 def run_surprise(algo, trainset, testset, model_name, model_summ_local, model_summary=None):
5
6     train = dict()
7     test = dict()
8
9     algo.fit(trainset)
10
11    train_pred = algo.test(trainset.build_testset())
12
13    train_actual, train_predicted = get_ratings(train_pred)
14    train_rmse, train_mape = get_error(train_pred)
15    #print("TRAIN RMSE(WITH OPTIMAL VALUE OF PARAMS) = %.6f"%train_rmse)
16    #print("TRAIN MAPE(WITH OPTIMAL VALUE OF PARAMS) = %.6f"%train_mape)
17    #print("-"*50)
18    train = {"RMSE": train_rmse, "MAPE": train_mape, "prediction": train_predicted}
19
20 #-----Evaluating Test Data-----
21 #print("TEST DATA")
22 #PREDICTIONS ON TEST DATA
23 test_pred = algo.test(testset)
24 test_actual, test_predicted = get_ratings(test_pred)
25 test_rmse, test_mape = get_error(test_pred)
26 #print("TEST RMSE(WITH OPTIMAL VALUE OF PARAMS) = %.6f"%test_rmse)
27 #print("TEST MAPE(WITH OPTIMAL VALUE OF PARAMS) = %.6f"%test_mape)
28 #print("-"*50)
29 test = {"RMSE": test_rmse, "MAPE": test_mape, "prediction": test_predicted}
30
31 #PRINT PERFORMANCE WITH BEST PARAMS
32 model_summ_local.add_row([model_name, '%.5f'%train_rmse, '%.5f'%train_mape, '%.5f'%test_rmse, '%.5f'%test_mape])
33 print('Performance Summary (WITH OPTIMAL VALUE OF PARAMETERS):')
34
35 if model_summary !=None:
36     model_summary.add_row([model_name, '%.5f'%train_rmse, '%.5f'%train_mape, '%.5f'%test_rmse, '%.5f'%test_mape])
37     #make_table(model_name, train_rmse, train_mape, test_rmse, test_mape)
38
39 return train, test, model_summ_local
```

```
In [211]: 1 def testPerformance(clf, trainset, testset, train_reg, test_reg, model_name, model_summary=None, add_to_xgboost=False)
2     model_summ_local=PrettyTable()
3     model_summ_local.field_names=['Model', 'Train(RMSE)', 'Train(MAPE)', 'Test(RMSE)', 'Test(MAPE)']
4     #APPLYING MODEL WITH OPTIMAL VALUE OF HYPERPARAMS
5     train_result, test_result, model_summ_local = run_surprise(clf, trainset, testset, model_name, model_summary_local,
6     #ADD PREDICTED DATA AS NEW COLUMN FROM SURPRISE MODEL TO XGBOOST DATA = XGBOOST_FEATS + SURPRISE_O/P
7     if add_to_xgboost:
8         train_reg[model_name]=train_result['prediction']
9         test_reg[model_name]=test_result['prediction']
10    return train_result, test_result, model_summ_local
```



## Functions required for XGBOOST models :

### XGBOOST (GBDT)

#### 1. function for scorer for Grid search and Random search:

```
In [212]: 1 def cust_scorer(y_true, y_pred):
2     rmse = np.sqrt(mean_squared_error(y_true, y_pred))
3     #mape = np.mean(abs((y_true - y_pred)/y_true))*100
4     return rmse
```

#### 2. function for Tunning hyperparam:

In [213]:

```
1 def Ensemble_Regressor(x_train,y_train,CV,params_, tune_param, searchMethod):
2     #INITIALIZE GBDT CLASSIFIER
3     reg=xgb.XGBRegressor(n_estimators=params_['n_estimators'],\
4                           max_depth=params_['max_depth'],\
5                           eta=.02,\
6                           reg_alpha=params_['reg_alpha'],\
7                           min_child_weight=params_['min_child_weight'],\
8                           gamma=params_['gamma'],\
9                           subsample=params_['subsample'],\
10                          colsample_bytree=params_['colsample_bytree'],\
11                          booster='gbtree')
12     # APPLY RANDOM OR GRID SEARCH FOR HYPERPARAMETER TUNNING
13     if searchMethod=='random':
14         model=RandomSearch(reg,\
15                             n_jobs=-1,\
16                             cv=CV,\
17                             param_distributions=tune_param,\
18                             n_iter=3,\
19                             return_train_score=True,\
20                             scoring=make_scorer(cust_scoring, greater_is_better=False))
21     elif searchMethod=='grid':
22         model= GridSearch(estimator=reg,\
23                            param_grid=tune_param,\
24                            scoring=make_scorer(cust_scoring, greater_is_better=False),\
25                            n_jobs=-1,\
26                            cv=CV,\
27                            return_train_score=True
28        )
29     model.fit(x_train,y_train)
30     #PLOT THE PERFORMANCE OF MODEL ON CROSSVALIDATION DATA FOR EACH HYPERPARAM VALUE
31     train_rmse = model.cv_results_['mean_train_score']
32     cv_rmse = model.cv_results_['mean_test_score']
33     #if len(param_name.split())==1:
34     for i in range(len(train_rmse)):
35         print('Train RMSE: %.4f'%abs(train_rmse[i]),'CV RMSE: %.4f'%abs(cv_rmse[i]))
36     print()
37     return model
```

In [214]:

```
1 def tuneALL_PARAM_XGB(train, y_train, CV, params_range, params_, searchMethod):
2     for param_name, param_list in zip(params_range.keys(), params_range.values()):
3         tune_param={}
4         tune_param[param_name]=param_list
5         #TUNNING HYPERPARAM
6         print('Tunning {}'.format(param_name.upper()))
7         model=Ensemble_Regressor(train, y_train, CV, params_, tune_param, searchMethod)
8         #UPDATE OTIMAL VALUE OF PARAMETER
9         params_[param_name] = model.best_params_[param_name]
10    return model, params_
```

### 3. function for measuring performance on test data:

In [215]:

```

1 def test_performance_xgb(x_train,y_train,x_test,y_test,params_, model_summary,model_use=None,summery=False):
2     '''FUNCTION FOR TEST PERFORMANCE(PLOT ROC CURVE FOR BOTH TRAIN AND TEST) WITH OPTIMAL HYPERPARAM'''
3     train_result={}
4     test_result={}
5     #INITIALIZE GBDT WITH OPTIMAL VALUE OF HYPERPARAMS
6     clf=xgb.XGBRegressor(n_estimators=params_['n_estimators'],\
7                           max_depth=params_['max_depth'],\
8                           eta=.02,\
9                           reg_alpha=params_['reg_alpha'],\
10                          min_child_weight=params_['min_child_weight'],\
11                          gamma=params_['gamma'],\
12                          subsample=params_['subsample'],\
13                          colsample_bytree=params_['colsample_bytree'],\
14                          booster='gbtree',\
15                          verbose=1)
16
17     clf.fit(x_train, y_train)
18     #PREDICTION FOR TRAIN AND TEST
19     y_pred=clf.predict(x_test)
20     y_pred_tr=clf.predict(x_train)
21
22     #TEST RMSE AND MAPE
23     test_rmse=abs(np.sqrt(mean_squared_error(y_test, y_pred)))
24     test_mape=np.mean(np.abs((y_test - y_pred) / y_test)) * 100
25     #TRAIN RMSE AND MAPE
26     train_rmse=abs(np.sqrt(mean_squared_error(y_train, y_pred_tr)))
27     train_mape=np.mean(np.abs((y_train - y_pred_tr) / y_train)) * 100
28
29     #SAVE RESULT GLOBALLY
30     train_result['RMSE']=train_rmse; train_result['MAPE']=train_mape; train_result['prediction']=y_pred_tr
31     test_result['RMSE']=test_rmse; test_result['MAPE']=test_mape; test_result['prediction']=y_pred
32
33     print('FOR OPTIMAL PARAMETERS, TRAIN RMSE: %.5f, TEST RMSE: %.5f'%(train_rmse, test_rmse))
34     model_summ_local=PrettyTable()
35     model_summ_local.field_names=['Model', 'Train(RMSE)', 'Train(MAPE)', 'Test(RMSE)', 'Test(MAPE)']
36     model_summ_local.add_row([model_use, '%.5f'%train_rmse, '%.5f'%train_mape, '%.5f'%test_rmse, '%.5f'%test_mape])
37     plot_importance(xgb, clf)
38     #print(model_local_summ)
39     if summary:
40         model_summary.add_row([model_use, '%.5f'%train_rmse, '%.5f'%train_mape, '%.5f'%test_rmse, '%.5f'%test_mape])
41

```



## MODELS:

**Initialization of common objects for Baseline and Matrix Factorization models:**

In [228]:

```

1 #DATA
2 data = Dataset.load_from_df(train_reg[['user', 'movie', 'rating']], reader)
3
4 #TUNNING ALGORITHM TO BE USED
5 searchMethod=['random', 'grid']
6
7 #HYPERPARAM RANGE
8 min_k=[1]
9 k=[20,30,40,60]
10 bsl_options = {'method' : ['sgd'], 'learning_rate':[.001,.01,.05,.1], 'n_epochs':[10,20]}
11 sim_options= {'name': ['pearson_baseline'] , 'user_based': [True], 'shrinkage':[80,100,120,140]}
12
13 #DICT OF HYPERPARAM
14 params={}
15
16 #DICT HYPERPARAMS FOR SVD AND SVDPP
17 params_svd=OrderedDict([('n_factors',[10,30,50,70,90]),\
18                         ('n_epochs',[10,20,30,40]),\
19                         ('lr_all',[.003,.005,.008,.1]),\
20                         ('reg_all',[0.0,.02,.05,.07,.09])])
21 #MODEL USED
22 model_name=['XGBOOST-13-FEATS', 'BaselineOnly', 'XGBOOST-BSL', 'KNNBaseline-U-U',\
23             'KNNBaseline-I-I', 'XGBOOST-KNN', 'SVD', 'SVD++', 'XGBOOST-13-BSL-KNN-MF', 'XGBOOST-BSL-KNN-MF', ]
24
25
26 #GLOBAL SUMMARY OF ALL THE MODELS
27 model_summary=PrettyTable()
28 model_summary.field_names=['Model', 'Train(RMSE)', 'Train(MAPE)', 'Test(RMSE)', 'Test(MAPE)']
29
30 #CROSS VALIDATION SPLIT OBJECT
31 CV_=KFold(n_splits=5)

```

## Initialization of common objects for XGBOOST model:

In [229]:

```

1 #DICT OF PARAMETERS, INITIALLY SET REASONABLE VALUES FOR PARAMETER, AND AFTER TUNNING UPDATE VALUE WITH OPT. VALUE
2 params_=OrderedDict([
3     ('n_estimators',64),
4     ('max_depth',5),
5     ('min_child_weight',1),
6     ('gamma',0),
7     ('subsample',.8),
8     ('colsample_bytree',.8),
9     ('reg_alpha',.1)
10 )
11
12
13 # DICT OF HYPERPARAMETER FOR XGBOOST WITH RANGE OF VALUES
14 params_range=OrderedDict([('n_estimators', [32,64,80,100]),\
15     ('max_depth', [5,7,9]),\
16     ('min_child_weight', [1,3,5]),\
17     ('gamma', [i/10.0 for i in range(0,5)]),\
18     ('subsample', [.6,.7,.8,.9]),\
19     ('colsample_bytree', [.6,.7,.8,.9]),\
20     ('reg_alpha',[0, 0.001, 0.005, 0.01, 0.05])])
21
22 #CV SPLIT (HERE KF STANDS/ALIAS FOR KFold() method of sklearn)
23 CV = KF(n_splits=5)
24
25 #SEARCHING ALGO FOR HYPERPARAM
26 searchMethod=['random', 'grid']
27
28 #
29 param_summ=PrettyTable()
30 param_summ.field_names=['Parameter', 'Value']

```

## 1. XGBOOST-13:

[1.1]Hyperparameter Tunnnning:

In [230]:

```

1  %%time
2  #TRAIN AND TEST DATA FOR XGBOOST MODELS
3  train, test = train_reg.drop(["user", "movie", "rating"], axis = 1), test_reg.drop(["user", "movie", "rating"], axis = 1)
4  y_train, y_test = train_reg["rating"], test_reg["rating"]
5  print('HYPERPARAMETER:\n')
6  param_list(params_range, param_summ)
7  print()
8  model, params_ = tuneALL_PARAM_XGB(train, y_train, CV, params_range, params_, searchMethod[1])

```

HYPERPARAMETER:

Parameter	Value
n_estimators	[32, 64, 80, 100]
max_depth	[5, 7, 9]
min_child_weight	[1, 3, 5]
gamma	[0.0, 0.1, 0.2, 0.3, 0.4]
subsample	[0.6, 0.7, 0.8, 0.9]
colsample_bytree	[0.6, 0.7, 0.8, 0.9]
reg_alpha	[0, 0.001, 0.005, 0.01, 0.05]

Tunning N\_ESTIMATORS:

Train RMSE: 0.8421 CV RMSE: 0.8624  
 Train RMSE: 0.8235 CV RMSE: 0.8507  
 Train RMSE: 0.8197 CV RMSE: 0.8507  
 Train RMSE: 0.8157 CV RMSE: 0.8511

Tunning MAX\_DEPTH:

Train RMSE: 0.8197 CV RMSE: 0.8507  
 Train RMSE: 0.7778 CV RMSE: 0.8551  
 Train RMSE: 0.7043 CV RMSE: 0.8640

Tunning MIN\_CHILD\_WEIGHT:

Train RMSE: 0.8197 CV RMSE: 0.8507  
 Train RMSE: 0.8200 CV RMSE: 0.8505  
 Train RMSE: 0.8203 CV RMSE: 0.8503

Tunning GAMMA:

Train RMSE: 0.8203 CV RMSE: 0.8503

```
Train RMSE: 0.8203 CV RMSE: 0.8502
Train RMSE: 0.8203 CV RMSE: 0.8501
Train RMSE: 0.8202 CV RMSE: 0.8500
Train RMSE: 0.8203 CV RMSE: 0.8501
```

Tunning SUBSAMPLE:

```
Train RMSE: 0.8212 CV RMSE: 0.8500
Train RMSE: 0.8208 CV RMSE: 0.8503
Train RMSE: 0.8202 CV RMSE: 0.8500
Train RMSE: 0.8199 CV RMSE: 0.8503
```

Tunning COLSAMPLE\_BYTREE:

```
Train RMSE: 0.8230 CV RMSE: 0.8492
Train RMSE: 0.8208 CV RMSE: 0.8502
Train RMSE: 0.8202 CV RMSE: 0.8500
Train RMSE: 0.8194 CV RMSE: 0.8508
```

Tunning REG\_ALPHA:

```
Train RMSE: 0.8230 CV RMSE: 0.8494
Train RMSE: 0.8230 CV RMSE: 0.8494
Train RMSE: 0.8229 CV RMSE: 0.8503
Train RMSE: 0.8228 CV RMSE: 0.8502
Train RMSE: 0.8228 CV RMSE: 0.8491
```

CPU times: user 27.1 s, sys: 172 ms, total: 27.3 s

Wall time: 1min 11s

**[1.2]Optimal value of parameters after tuning:**

```
In [231]: 1 print('Optimal Value of Hyperparameters after Tuning:\n')
2 param_list(params_,param_summ)
3
```

Optimal Value of Hyperparameters after Tuning:

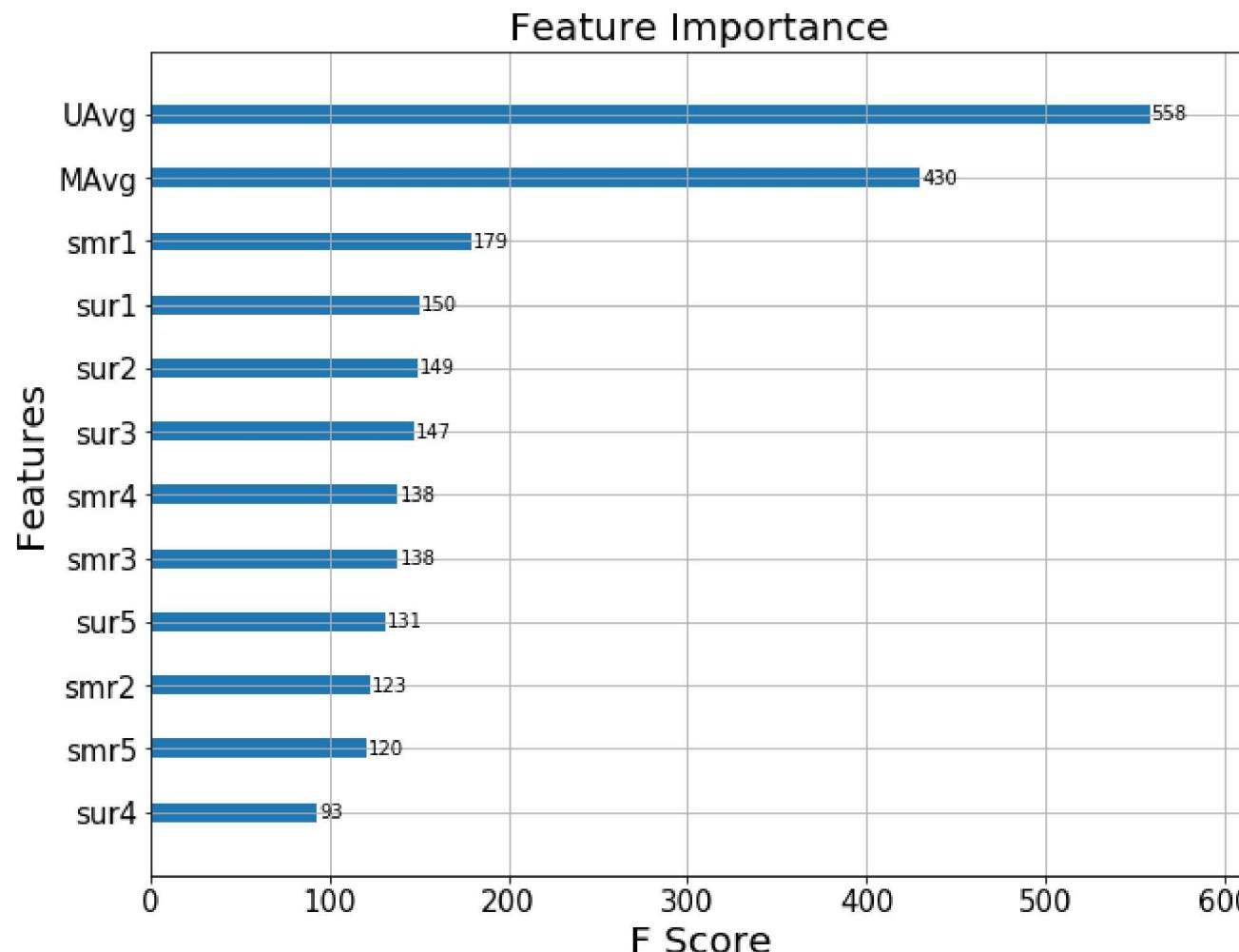
Parameter	Value
n_estimators	80
max_depth	5
min_child_weight	5
gamma	0.3
subsample	0.8
colsample_bytree	0.6
reg_alpha	0.05

[1.3]Test performance :

In [232]:

```
1 %%time
2 train_result, test_result, model_summ_local=test_performance_xgb(train, y_train, test, y_test,\n3                                         params_, model_summary, model_name[0], summary=True
4
5 #SAVE THE TEST RESULTS TO DICT OF DICT
6 models_evaluation_train[model_name[0]] = train_result
7 models_evaluation_test[model_name[0]] = test_result
```

FOR OPTIMAL PARAMETERS, TRAIN RMSE: 0.82606, TEST RMSE: 1.07701



CPU times: user 4.62 s, sys: 1.86 s, total: 6.48 s

Wall time: 3.83 s

**[1.4] Model Summary:**

In [233]:

```

1 #MODEL_SUMMARY
2 print(model_summ_local)

```

Model	Train(RMSE)	Train(MAPE)	Test(RMSE)	Test(MAPE)
XGBOOST-13-FEATS	0.82606	24.16566	1.07701	33.59137

**2. BaselineOnly:**

$$\hat{r}_{ui} = b_{ui} = \mu + b_u + b_i$$

- $\mu$  : Average of all trainings in training data.
- $b_u$  : User bias
- $b_i$  : Item bias (movie biases)

Optimization function:

$$\sum_{r_{ui} \in R_{train}} (r_{ui} - (\mu + b_u + b_i))^2 + \lambda (b_u^2 + b_i^2) . \text{ [minimize } b_u, b_i]$$

**[a.] Hyperparam Tuning:**

In [234]:

```

1  %%time
2  #MODEL USED
3  clf = BaselineOnly
4
5  #DICTIONARY OF HYPERPARAM TO BE TUNED
6  params['bsl_options']=bsl_options
7
8  print('HYPERPARAMETER:\n')
9  param_list(params, param_summ)
10
11 #HYPERPARAM TUNNING
12 model = hyperparamTUNNING(clf, data, params, CV_, 'bsl', searchMethod[1])

```

HYPERPARAMETER:

Parameter	Value
bsl_options	{'method': ['sgd'], 'n_epochs': [10, 20], 'learning_rate': [0.001, 0.01, 0.05, 0.1]}

Train RMSE: 0.9780 CV RMSE: 0.9899  
Train RMSE: 0.8758 CV RMSE: 0.9525  
Train RMSE: 0.7894 CV RMSE: 0.9767  
Train RMSE: 0.7777 CV RMSE: 1.0077  
Train RMSE: 0.9564 CV RMSE: 0.9779  
Train RMSE: 0.8309 CV RMSE: 0.9510  
Train RMSE: 0.7647 CV RMSE: 0.9968  
Train RMSE: 0.7687 CV RMSE: 1.0247

Best RMSE: 0.95096 For Optimal Params Combination Value:

```

[1.]Learning-Rate: 0.0100
[2.]Epochs: 20

```

CPU times: user 44.4 s, sys: 84 ms, total: 44.5 s

Wall time: 45.2 s

**[b.]Performance on test data and summary:**

In [235]:

```
1 %%time
2 #INITIALIZE CLASSIFIER WITH BEST PARAMS VALUE
3 clf=BaselineOnly(bsl_options=model.best_params['rmse']['bsl_options'])
4 #MEASURE PERFORMANCE ON TEST DATA
5 train_result, test_result, model_summ_local=testPerformance(clf,trainset, testset, train_reg, test_reg,\n                model_name[1],model_summary, add_to_xgboost=True)
6 #SAVE THE TEST RESULTS TO DICT OF DICT
7 models_evaluation_train[model_name[1]] = train_result
8 models_evaluation_test[model_name[1]] = test_result
9
10 #PRINT PERFROMANCE SUMMARY
11 print(model_summ_local)
```

Estimating biases using sgd...

Performance Summary (WITH OPTIMAL VALUE OF PARAMETERS):

Model	Train(RMSE)	Train(MAPE)	Test(RMSE)	Test(MAPE)
BaselineOnly	0.83552	25.25685	1.07429	35.12019

CPU times: user 828 ms, sys: 4 ms, total: 832 ms

Wall time: 830 ms

### 3. XGBOOST(13 + BSL):

#### [3.1.] Hyperparameter Tuning:

In [236]:

```

1  %%time
2  #TRAIN AND TEST DATA FOR XGBOOST MODELS
3  train, test = train_reg.drop(["user", "movie", "rating"], axis = 1), test_reg.drop(["user", "movie", "rating"], axis = 1)
4  y_train, y_test = train_reg["rating"], test_reg["rating"]
5  print('HYPERPARAMETER:\n')
6  param_list(params_range, param_summ)
7  print()
8  model, params_ = tuneALL_PARAM_XGB(train, y_train, CV, params_range, params_, searchMethod[1])

```

HYPERPARAMETER:

Parameter	Value
n_estimators	[32, 64, 80, 100]
max_depth	[5, 7, 9]
min_child_weight	[1, 3, 5]
gamma	[0.0, 0.1, 0.2, 0.3, 0.4]
subsample	[0.6, 0.7, 0.8, 0.9]
colsample_bytree	[0.6, 0.7, 0.8, 0.9]
reg_alpha	[0, 0.001, 0.005, 0.01, 0.05]

Tunning N\_ESTIMATORS:

Train RMSE: 0.8435 CV RMSE: 0.8628  
 Train RMSE: 0.8255 CV RMSE: 0.8502  
 Train RMSE: 0.8220 CV RMSE: 0.8502  
 Train RMSE: 0.8179 CV RMSE: 0.8506

Tunning MAX\_DEPTH:

Train RMSE: 0.8220 CV RMSE: 0.8502  
 Train RMSE: 0.7852 CV RMSE: 0.8533  
 Train RMSE: 0.7253 CV RMSE: 0.8579

Tunning MIN\_CHILD\_WEIGHT:

Train RMSE: 0.8216 CV RMSE: 0.8505  
 Train RMSE: 0.8217 CV RMSE: 0.8504  
 Train RMSE: 0.8220 CV RMSE: 0.8502

Tunning GAMMA:

Train RMSE: 0.8220 CV RMSE: 0.8503

```
Train RMSE: 0.8219 CV RMSE: 0.8503
Train RMSE: 0.8219 CV RMSE: 0.8502
Train RMSE: 0.8220 CV RMSE: 0.8502
Train RMSE: 0.8220 CV RMSE: 0.8501
```

Tunning SUBSAMPLE:

```
Train RMSE: 0.8227 CV RMSE: 0.8505
Train RMSE: 0.8223 CV RMSE: 0.8510
Train RMSE: 0.8220 CV RMSE: 0.8501
Train RMSE: 0.8215 CV RMSE: 0.8504
```

Tunning COLSAMPLE\_BYTREE:

```
Train RMSE: 0.8220 CV RMSE: 0.8501
Train RMSE: 0.8207 CV RMSE: 0.8507
Train RMSE: 0.8190 CV RMSE: 0.8508
Train RMSE: 0.8182 CV RMSE: 0.8505
```

Tunning REG\_ALPHA:

```
Train RMSE: 0.8221 CV RMSE: 0.8504
Train RMSE: 0.8221 CV RMSE: 0.8505
Train RMSE: 0.8220 CV RMSE: 0.8507
Train RMSE: 0.8220 CV RMSE: 0.8507
Train RMSE: 0.8220 CV RMSE: 0.8501
```

CPU times: user 25.8 s, sys: 164 ms, total: 26 s

Wall time: 1min 8s

### [3.2]Optimal value of parameters after tuning:

In [237]:

```
1 print('Optimal Value of Hyperparameters after Tuning:\n')
2 param_list(params_,param_summ)
3
```

Optimal Value of Hyperparameters after Tuning:

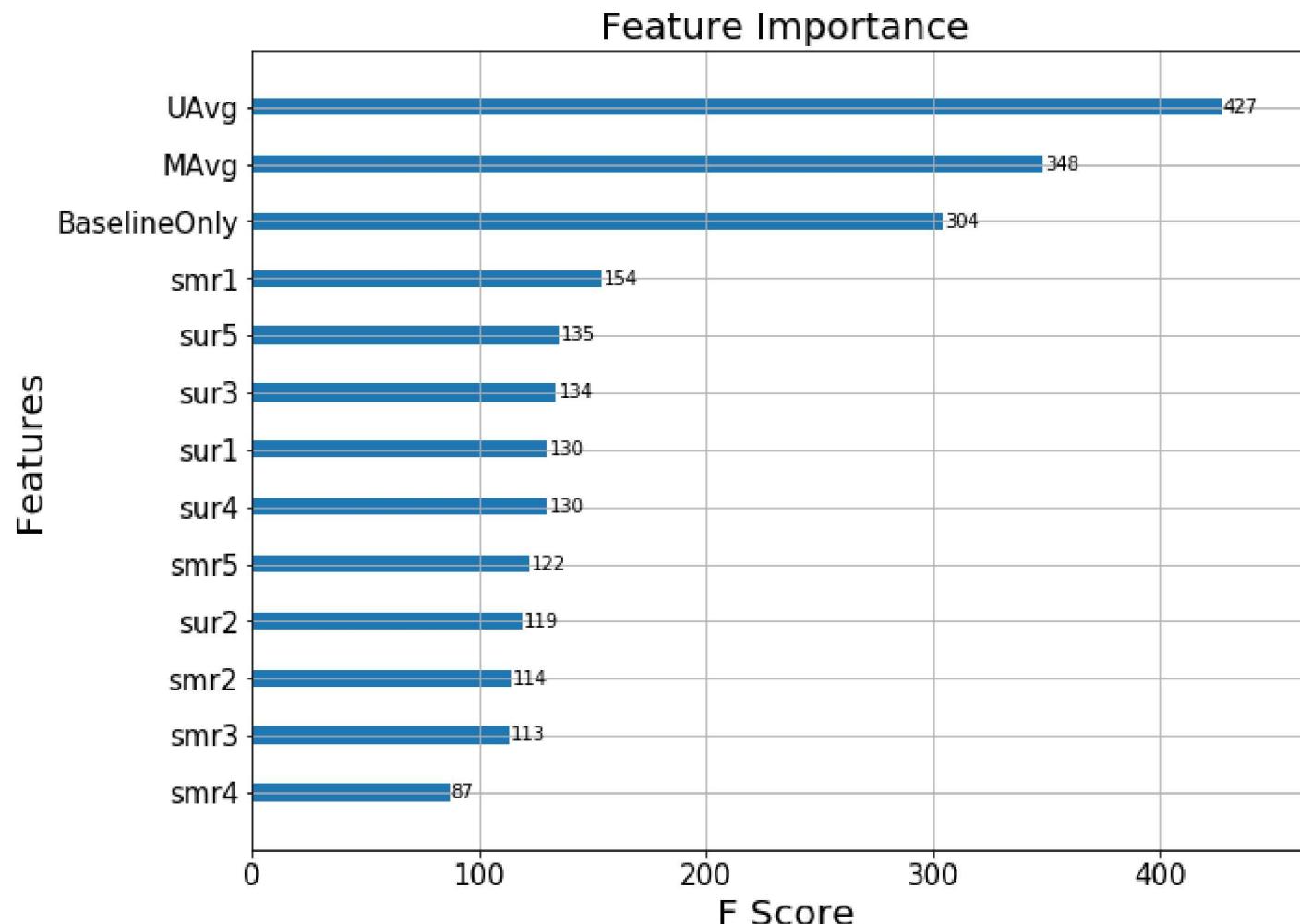
Parameter	Value
n_estimators	80
max_depth	5
min_child_weight	5
gamma	0.4
subsample	0.8
colsample_bytree	0.6
reg_alpha	0.05

[3.3]Test performance :

In [238]:

```
1 %%time
2 train_result, test_result, model_local_summ=test_performance_xgb(train, y_train, test, y_test,\n3                                         params_, model_summary, model_name[2], summary=True
4
5 #SAVE THE TEST RESULTS TO DICT OF DICT
6 models_evaluation_train[model_name[2]] = train_result
7 models_evaluation_test[model_name[2]] = test_result
```

FOR OPTIMAL PARAMETERS, TRAIN RMSE: 0.82509, TEST RMSE: 1.06767



CPU times: user 4.33 s, sys: 300 ms, total: 4.63 s

Wall time: 4.18 s

**[3.4]Model Summary:**

In [239]:

```

1 #MODEL SUMMARY
2 print(model_local_summ)

```

Model	Train(RMSE)	Train(MAPE)	Test(RMSE)	Test(MAPE)
XGBOOST-BSL	0.82509	24.14721	1.06767	34.57024

**4.KNNBaseline-U-U:**

Based on User-User similarity

$$\hat{r}_{ui} = b_{ui} + \frac{\sum_{v \in N_i^k(u)} \text{sim}(u, v) \cdot (r_{vi} - b_{vi})}{\sum_{v \in N_i^k(u)} \text{sim}(u, v)}$$

Based on Item-Item similarity

$$\hat{r}_{ui} = b_{ui} + \frac{\sum_{j \in N_u^k(i)} \text{sim}(i, j) \cdot (r_{uj} - b_{uj})}{\sum_{j \in N_u^k(i)} \text{sim}(i, j)}$$

- $b_{ui}$  - Baseline prediction of (user,movie) rating
- $N_i^k(u)$  - Set of K similar users (neighbours) of user (u) who rated movie(i)
- $\text{sim}(u, v)$  - Similarity between users u and v
  - Generally, it will be cosine similarity or Pearson correlation coefficient.

- But we use **shrunk Pearson-baseline correlation coefficient**, which is based on the pearsonBaseline similarity ( we take base line predictions instead of mean rating of user/item) Based on Item-Item similarity

#### [a.] Hyperparam Tuning:

In [240]:

```

1  %%time
2  #DICT HYPERPARAM
3  params['k']=k
4  params['min_k']=min_k
5  params['sim_options']=sim_options
6  print('HYPERPARAMETER:\n')
7  param_list(params, param_summ)
8  print()
9  #HYPERPARAM TUNNING
10 clf=KNNBaseline
11 model = hyperparamTUNNING(clf, data, params, CV_, 'knn', searchMethod[0])

```

HYPERPARAMETER:

Parameter	Value
k	[20, 30, 40, 60]
bsl_options	{'method': ['sgd'], 'n_epochs': [10, 20], 'learning_rate': [0.001, 0.01, 0.05, 0.1]}
sim_options	{'user_based': [True], 'shrinkage': [80, 100, 120, 140], 'name': ['pearson_baseline']}
min_k	[1]

Train RMSE: 0.1190 CV RMSE: 0.9583  
 Train RMSE: 0.1356 CV RMSE: 0.9664  
 Train RMSE: 0.1151 CV RMSE: 0.9549  
 Train RMSE: 0.2085 CV RMSE: 0.9655  
 Train RMSE: 0.1289 CV RMSE: 0.9688  
 Train RMSE: 0.1174 CV RMSE: 0.9708  
 Train RMSE: 0.1117 CV RMSE: 1.0057  
 Train RMSE: 0.1370 CV RMSE: 1.0026  
 Train RMSE: 0.1308 CV RMSE: 0.9640  
 Train RMSE: 0.1510 CV RMSE: 0.9725  
 Train RMSE: 0.1647 CV RMSE: 1.0179  
 Train RMSE: 0.1210 CV RMSE: 0.9603  
 Train RMSE: 0.1301 CV RMSE: 1.0015  
 Train RMSE: 0.0965 CV RMSE: 0.9770  
 Train RMSE: 0.1618 CV RMSE: 0.9655  
 Train RMSE: 0.1789 CV RMSE: 1.0179  
 Train RMSE: 0.1276 CV RMSE: 1.0188  
 Train RMSE: 0.1083 CV RMSE: 0.9770

```
Train RMSE: 0.1438 CV RMSE: 0.9549
```

```
Train RMSE: 0.1252 CV RMSE: 1.0230
```

Best RMSE: 0.95492 For Optimal Params Combination Value:

[1.]K: 30

[2.]shrinkage: 120

[3.]lerning\_rate: 0.0100

[4.]Epochs: 20

```
CPU times: user 2min 48s, sys: 132 ms, total: 2min 48s
```

```
Wall time: 13min 7s
```

**[b.]Performance on test data:**

In [241]:

```
1 %%time
2 #INITIALIZE CLASSIFIER WITH BEST PARAMS VALUE
3 clf=KNNBaseline(k=model.best_params['rmse']['k'],\
4                  bsl_options=model.best_params['rmse']['bsl_options'],\
5                  sim_options=model.best_params['rmse']['sim_options'],\
6                  min_k=model.best_params['rmse']['min_k'])
7
8 #MEASURE PERFORMANCE ON TEST DATA
9 train_result, test_result, model_summ_local=testPerformance(clf, trainset, testset, train_reg, test_reg, model_name[3],\
10                                                            model_summary=model_summary, add_to_xgboost=True)
11 #SAVE THE TEST RESULTS TO DICT OF DICT
12 models_evaluation_train[model_name[3]] = train_result
13 models_evaluation_test[model_name[3]] = test_result
14
15 #PRINT PERFORMANCE SUMMARY
16 print(model_summ_local)
```

Estimating biases using sgd...

Computing the pearson\_baseline similarity matrix...

Done computing similarity matrix.

Performance Summary (WITH OPTIMAL VALUE OF PARAMETERS):

Model	Train(RMSE)	Train(MAPE)	Test(RMSE)	Test(MAPE)
KNNBaseline-U-U	0.14034	3.41399	1.07426	35.11784

CPU times: user 2min, sys: 7.64 s, total: 2min 8s

Wall time: 2min 8s

## 5.KNNBaseline-I-I:

[a.] Hyperparam Tunning:

In [242]:

```

1 %%time
2 #DICT OF HYPERPARAM
3 params['sim_options']['user_based'][0]=False
4 print('HYPERPARAMETER:\n')
5 param_list(params, param_summ)
6 print()
7 #HYPERPARAM TUNNING
8 clf=KNNBaseline
9 model = hyperparamTUNNING(clf, data, params, CV, 'knn', searchMethod[0])

```

HYPERPARAMETER:

Parameter	Value
k	[20, 30, 40, 60]
bsl_options	{'method': ['sgd'], 'n_epochs': [10, 20], 'learning_rate': [0.001, 0.01, 0.05, 0.1]}
sim_options	{'user_based': [False], 'shrinkage': [80, 100, 120, 140], 'name': ['pearson_baseline']}
min_k	[1]

Train RMSE: 0.1022 CV RMSE: 1.0509  
 Train RMSE: 0.0856 CV RMSE: 1.0512  
 Train RMSE: 0.0498 CV RMSE: 1.0467  
 Train RMSE: 0.0955 CV RMSE: 1.0510  
 Train RMSE: 0.0496 CV RMSE: 1.0867  
 Train RMSE: 0.2770 CV RMSE: 1.0339  
 Train RMSE: 0.0458 CV RMSE: 1.0559  
 Train RMSE: 0.0421 CV RMSE: 1.0559  
 Train RMSE: 0.0590 CV RMSE: 1.0758  
 Train RMSE: 0.0545 CV RMSE: 1.0468  
 Train RMSE: 0.0604 CV RMSE: 1.0468  
 Train RMSE: 0.0654 CV RMSE: 1.0758  
 Train RMSE: 0.2341 CV RMSE: 1.0331  
 Train RMSE: 0.0460 CV RMSE: 1.0467  
 Train RMSE: 0.2863 CV RMSE: 1.0335  
 Train RMSE: 0.2770 CV RMSE: 1.0339  
 Train RMSE: 0.0440 CV RMSE: 1.0689  
 Train RMSE: 0.0901 CV RMSE: 1.0511  
 Train RMSE: 0.0856 CV RMSE: 1.0512  
 Train RMSE: 0.1022 CV RMSE: 1.0509

Best RMSE: 1.03314 For Optimal Params Combination Value:

- [1.]K: 60
- [2.]shrinkage: 140
- [3.]lerning\_rate: 0.0010
- [4.]Epochs: 20

CPU times: user 1min 46s, sys: 256 ms, total: 1min 46s

Wall time: 1min 48s

**[b.]Performance on test data and summary:**

In [244]:

```

1 %%time
2 #INITIALIZE CLASSIFIER WITH BEST PARAMS VALUE
3 clf=KNNBaseline(k=model.best_params['rmse']['k'],\
4                  bsl_options=model.best_params['rmse']['bsl_options'],\
5                  sim_options=model.best_params['rmse']['sim_options'],\
6                  min_k=model.best_params['rmse']['min_k'])
7
8 #MEASURE PERFORMANCE ON TEST DATA
9 train_result, test_result, model_summ_local=testPerformance(clf, trainset, testset, train_reg, test_reg, model_name[4],\
10                                                               model_summary=model_summary, add_to_xgboost=True)
11
12 #SAVE THE TEST RESULTS TO DICT OF DICT
13 models_evaluation_train[model_name[4]] = train_result
14 models_evaluation_test[model_name[4]] = test_result
15
16 #PRINT PERFORMANCE SUMMARY
17 print(model_summ_local)

```

Estimating biases using sgd...

Computing the pearson\_baseline similarity matrix...

Done computing similarity matrix.

Performance Summary (WITH OPTIMAL VALUE OF PARAMETERS):

Model	Train(RMSE)	Train(MAPE)	Test(RMSE)	Test(MAPE)
KNNBaseline-I-I	0.27398	6.45966	1.07461	35.19819

CPU times: user 2.29 s, sys: 8 ms, total: 2.3 s

Wall time: 2.3 s

## 6.XGBOOST(13 + BSL + KNN)

### [6.1.] Hyperparameter Tuning:

In [245]:

```

1  %%time
2  #TRAIN AND TEST DATA FOR XGBOOST MODELS
3  train, test = train_reg.drop(["user", "movie", "rating"], axis = 1), test_reg.drop(["user", "movie", "rating"], axis = 1)
4  y_train, y_test = train_reg["rating"], test_reg["rating"]
5  print('HYPERPARAMETER:\n')
6  param_list(params_range, param_summ)
7  print()
8  model, params_ = tuneALL_PARAM_XGB(train, y_train, CV, params_range, params_, searchMethod[1])

```

HYPERPARAMETER:

Parameter	Value
n_estimators	[32, 64, 80, 100]
max_depth	[5, 7, 9]
min_child_weight	[1, 3, 5]
gamma	[0.0, 0.1, 0.2, 0.3, 0.4]
subsample	[0.6, 0.7, 0.8, 0.9]
colsample_bytree	[0.6, 0.7, 0.8, 0.9]
reg_alpha	[0, 0.001, 0.005, 0.01, 0.05]

Tunning N\_ESTIMATORS:

Train RMSE: 0.8430 CV RMSE: 0.8606  
 Train RMSE: 0.8240 CV RMSE: 0.8498  
 Train RMSE: 0.8201 CV RMSE: 0.8497  
 Train RMSE: 0.8153 CV RMSE: 0.8500

Tunning MAX\_DEPTH:

Train RMSE: 0.8201 CV RMSE: 0.8497  
 Train RMSE: 0.7793 CV RMSE: 0.8536  
 Train RMSE: 0.7095 CV RMSE: 0.8596

Tunning MIN\_CHILD\_WEIGHT:

Train RMSE: 0.8193 CV RMSE: 0.8503  
 Train RMSE: 0.8195 CV RMSE: 0.8501  
 Train RMSE: 0.8201 CV RMSE: 0.8497

Tunning GAMMA:

Train RMSE: 0.8200 CV RMSE: 0.8497

```
Train RMSE: 0.8200 CV RMSE: 0.8498
Train RMSE: 0.8202 CV RMSE: 0.8502
Train RMSE: 0.8202 CV RMSE: 0.8500
Train RMSE: 0.8201 CV RMSE: 0.8497
```

Tunning SUBSAMPLE:

```
Train RMSE: 0.8214 CV RMSE: 0.8514
Train RMSE: 0.8204 CV RMSE: 0.8501
Train RMSE: 0.8200 CV RMSE: 0.8497
Train RMSE: 0.8192 CV RMSE: 0.8500
```

Tunning COLSAMPLE\_BYTREE:

```
Train RMSE: 0.8200 CV RMSE: 0.8497
Train RMSE: 0.8181 CV RMSE: 0.8504
Train RMSE: 0.8178 CV RMSE: 0.8500
Train RMSE: 0.8169 CV RMSE: 0.8505
```

Tunning REG\_ALPHA:

```
Train RMSE: 0.8200 CV RMSE: 0.8500
Train RMSE: 0.8200 CV RMSE: 0.8500
Train RMSE: 0.8199 CV RMSE: 0.8496
Train RMSE: 0.8199 CV RMSE: 0.8497
Train RMSE: 0.8200 CV RMSE: 0.8497
```

CPU times: user 28.8 s, sys: 548 ms, total: 29.4 s

Wall time: 1min 19s

## [6.2]Optimal value of parameters after tuning:

```
In [246]: 1 print('Optimal Value of Hyperparameters after Tuning:\n')
2 param_list(params_,param_summ)
```

Optimal Value of Hyperparameters after Tuning:

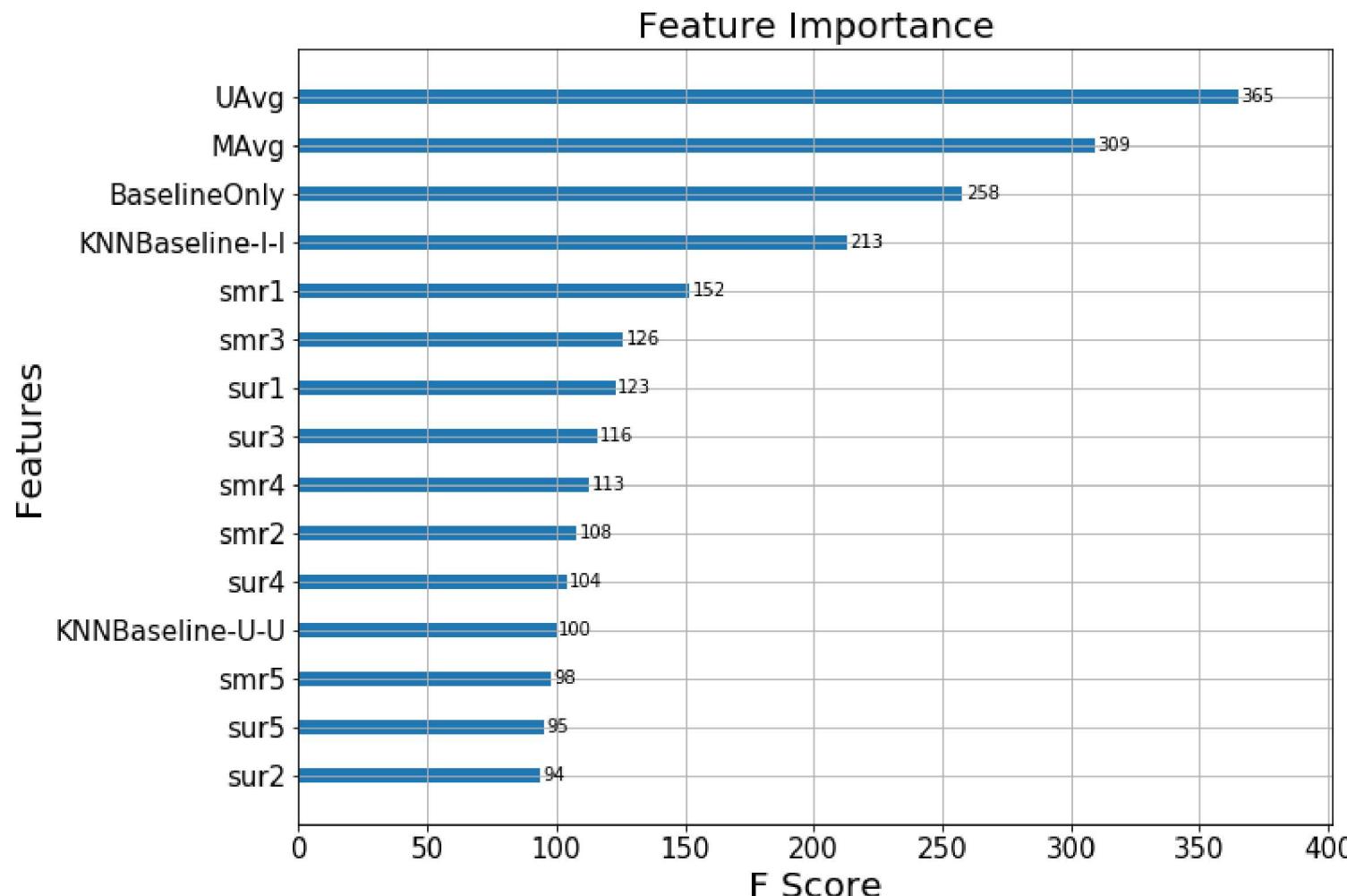
Parameter	Value
n_estimators	80
max_depth	5
min_child_weight	5
gamma	0.0
subsample	0.8
colsample_bytree	0.6
reg_alpha	0.005

### [6.3]Test Performance:

In [252]:

```
1 %time
2 train_result, test_result, model_summ_local=test_performance_xgb(train, y_train, test, y_test,\n3                                         params_, model_summary, model_name[5], summary=True
4
5 #SAVE THE TEST RESULTS TO DICT OF DICT
6 models_evaluation_train[model_name[5]] = train_result
7 models_evaluation_test[model_name[5]] = test_result
```

FOR OPTIMAL PARAMETERS, TRAIN RMSE: 0.82392, TEST RMSE: 1.06919



CPU times: user 4.86 s, sys: 320 ms, total: 5.18 s

Wall time: 4.72 s

## [6.4] Model Summary

In [253]:

```
1 #MODEL_SUMMARY  
2 print(model_summ_local)
```

Model	Train(RMSE)	Train(MAPE)	Test(RMSE)	Test(MAPE)
XGBOOST-KNN	0.82392	24.12600	1.06919	34.61556

## 7.SVD:

Predicted Rating:

- $\hat{r}_{ui} = \mu + b_u + b_i + q_i^T p_u$
- $q_i$  - Representation of item(movie) in latent factor space
- $p_u$  - Representation of user in new latent factor space

Optimization problem:

- $\sum_{r_{ui} \in R_{train}} (r_{ui} - \hat{r}_{ui})^2 + \lambda (b_i^2 + b_u^2 + ||q_i||^2 + ||p_u||^2)$

### [a.] Hyperparam tunning:

In [260]:

```
1 %%time
2 print('HYPERPARAMETER:')
3 param_list(params, param_summ)
4 print()
5 #HYPERPARAM TUNNING
6 clf=SVD
7 model = hyperparamTUNNING(clf, data, params_svd, CV_, 'svd', searchMethod[1])
```

HYPERPARAMETER:

Parameter	Value
k	[20, 30, 40, 60]
bsl_options	{'method': ['sgd'], 'n_epochs': [10, 20], 'learning_rate': [0.001, 0.01, 0.05, 0.1]}
sim_options	{'user_based': [False], 'shrinkage': [80, 100, 120, 140], 'name': ['pearson_baseline']}
min_k	[1]

```
Train RMSE: 0.9336 CV RMSE: 0.9719
Train RMSE: 0.9372 CV RMSE: 0.9707
Train RMSE: 0.9386 CV RMSE: 0.9710
Train RMSE: 0.9392 CV RMSE: 0.9711
Train RMSE: 0.9395 CV RMSE: 0.9711
Train RMSE: 0.8963 CV RMSE: 0.9646
Train RMSE: 0.9088 CV RMSE: 0.9621
Train RMSE: 0.9125 CV RMSE: 0.9616
Train RMSE: 0.9136 CV RMSE: 0.9614
Train RMSE: 0.9144 CV RMSE: 0.9617
```

**[b.] Performance on test data and summary:**

In [261]:

```

1  %%time
2  #INITIALIZE CLASSIFIER WITH BEST PARAMS VALUE
3  clf=SVD(n_factors=model.best_params['rmse']['n_factors'],\
4            n_epochs=model.best_params['rmse']['n_epochs'],\
5            lr_all=model.best_params['rmse']['lr_all'],\
6            reg_all=model.best_params['rmse']['reg_all'])
7
8  #MEASURE PERFORMANCE ON TEST DATA
9  train_result, test_result, model_summ_local=testPerformance(clf, trainset, testset, train_reg, test_reg, model_name[6],\
10                                         model_summary=model_summary, add_to_xgboost=True)
11
12 #SAVE THE TEST RESULTS TO DICT OF DICT
13 models_evaluation_train[model_name[6]] = train_result
14 models_evaluation_test[model_name[6]] = test_result
15
16 #PRINT PERFORMANCE SUMMARY
17 print(model_summ_local)

```

Performance Summary (WITH OPTIMAL VALUE OF PARAMETERS):

Model	Train(RMSE)	Train(MAPE)	Test(RMSE)	Test(MAPE)
SVD	0.84508	25.80980	1.07418	35.13581

CPU times: user 2.76 s, sys: 4 ms, total: 2.76 s

Wall time: 2.76 s

## 8.SVDpp:

Predicted Rating:

- $\hat{r}_{ui} = \mu + b_u + b_i + q_i^T \left( p_u + |I_u|^{-\frac{1}{2}} \sum_{j \in I_u} y_j \right)$
- $I_u$  --- the set of all items rated by user u

- $y_j$  --- Our new set of item factors that capture implicit ratings.

Optimization Problem:

- $\sum_{r_{ui} \in R_{train}} (r_{ui} - \hat{r}_{ui})^2 + \lambda (b_i^2 + b_u^2 + ||q_i||^2 + ||p_u||^2 + ||y_j||^2)$

[a.] Hyperparam tunning:

In [263]:

```

1 %time
2 print('HYPERPARAMETER:')
3 param_list(params, param_summ)
4 print()
5 #HYPERPARAM TUNNING
6 clf=SVDpp
7 model = hyperparamTUNNING(clf, data, params_svd, CV_, 'svd', searchMethod[0])

```

HYPERPARAMETER:

Parameter	Value
k	[20, 30, 40, 60]
bsl_options	{'method': ['sgd'], 'n_epochs': [10, 20], 'learning_rate': [0.001, 0.01, 0.05, 0.1]}
sim_options	{'user_based': [False], 'shrinkage': [80, 100, 120, 140], 'name': ['pearson_baseline']}
min_k	[1]

Train RMSE: 0.4572 CV RMSE: 1.0431  
 Train RMSE: 0.2714 CV RMSE: 0.9722  
 Train RMSE: 0.6510 CV RMSE: 1.0004  
 Train RMSE: 0.5286 CV RMSE: 0.9497  
 Train RMSE: 0.3022 CV RMSE: 0.9912  
 Train RMSE: 0.3605 CV RMSE: 0.9758  
 Train RMSE: 0.8703 CV RMSE: 0.9507  
 Train RMSE: 0.3622 CV RMSE: 0.9776  
 Train RMSE: 0.6502 CV RMSE: 0.9481  
 Train RMSE: 0.5682 CV RMSE: 1.2116  
 Train RMSE: 0.5508 CV RMSE: 0.9826  
 Train RMSE: 0.9181 CV RMSE: 0.9698  
 Train RMSE: 0.8819 CV RMSE: 0.9524  
 Train RMSE: 0.6068 CV RMSE: 0.9569  
 Train RMSE: 0.1751 CV RMSE: 1.0765  
 Train RMSE: 0.4871 CV RMSE: 0.9751  
 Train RMSE: 0.9065 CV RMSE: 0.9575  
 Train RMSE: 0.8455 CV RMSE: 0.9505  
 Train RMSE: 0.8960 CV RMSE: 0.9576  
 Train RMSE: 0.4957 CV RMSE: 0.9790

Best RMSE: 0.94813 For Optimal Params Combination Value:

```
[1.]N-Factors: 50
[2.]Epochs: 40
[3.]Learning-Rate: 0.0080
[4.]Reg-All: 0.0900
CPU times: user 2min 2s, sys: 352 ms, total: 2min 3s
Wall time: 2min 50s
```

### [b.] Performance on test data and summary:

In [264]:

```
1 %%time
2 #INITIALIZE CLASSIFIER WITH BEST PARAMS VALUE
3 clf=SVDpp(n_factors=model.best_params['rmse']['n_factors'],\
4             n_epochs=model.best_params['rmse']['n_epochs'],\
5             lr_all=model.best_params['rmse']['lr_all'],\
6             reg_all=model.best_params['rmse']['reg_all'])
7
8 #MEASURE PERFORMANCE ON TEST DATA
9 train_result, test_result, model_summ_local=testPerformance(clf, trainset, testset, train_reg, test_reg, model_name[
10                                         model_summary=model_summary, add_to_xgboost=True])
11
12 #SAVE THE TEST RESULTS TO DICT OF DICT
13 models_evaluation_train[model_name[7]] = train_result
14 models_evaluation_test[model_name[7]] = test_result
15
16 #PRINT PERFORMANCE SUMMARY
17 print(model_summ_local)
```

Performance Summary (WITH OPTIMAL VALUE OF PARAMETERS):

Model	Train(RMSE)	Train(MAPE)	Test(RMSE)	Test(MAPE)
SVD++	0.67078	19.84641	1.07450	35.11816

CPU times: user 38.5 s, sys: 8 ms, total: 38.6 s  
Wall time: 38.6 s



## 9. XGBOOST( 13 + BSL + KNN + MF[SVD+SVDpp] ):

[9.1.] Hyperparameter Tuning:

In [265]:

```

1 %time
2 #TRAIN AND TEST DATA FOR XGBOOST MODELS
3 train, test = train_reg.drop(["user", "movie", "rating"], axis = 1), test_reg.drop(["user", "movie", "rating"], axis = 1)
4 y_train, y_test = train_reg["rating"], test_reg["rating"]
5 print('HYPERPARAMETER:\n')
6 param_list(params_range, param_summ)
7 print()
8 model, params_ = tuneALL_PARAM_XGB(train, y_train, CV, params_range, params_, searchMethod[1])

```

HYPERPARAMETER:

Parameter	Value
n_estimators	[32, 64, 80, 100]
max_depth	[5, 7, 9]
min_child_weight	[1, 3, 5]
gamma	[0.0, 0.1, 0.2, 0.3, 0.4]
subsample	[0.6, 0.7, 0.8, 0.9]
colsample_bytree	[0.6, 0.7, 0.8, 0.9]
reg_alpha	[0, 0.001, 0.005, 0.01, 0.05]

Tunning N\_ESTIMATORS:

Train RMSE: 0.8430 CV RMSE: 0.8622  
 Train RMSE: 0.8236 CV RMSE: 0.8497  
 Train RMSE: 0.8194 CV RMSE: 0.8500  
 Train RMSE: 0.8148 CV RMSE: 0.8501

Tunning MAX\_DEPTH:

Train RMSE: 0.8236 CV RMSE: 0.8497  
 Train RMSE: 0.7857 CV RMSE: 0.8529  
 Train RMSE: 0.7179 CV RMSE: 0.8581

Tunning MIN\_CHILD\_WEIGHT:

Train RMSE: 0.8232 CV RMSE: 0.8495  
 Train RMSE: 0.8233 CV RMSE: 0.8500  
 Train RMSE: 0.8236 CV RMSE: 0.8497

Tunning GAMMA:

Train RMSE: 0.8232 CV RMSE: 0.8495

```
Train RMSE: 0.8232 CV RMSE: 0.8495
Train RMSE: 0.8232 CV RMSE: 0.8496
Train RMSE: 0.8232 CV RMSE: 0.8498
Train RMSE: 0.8232 CV RMSE: 0.8498
```

Tunning SUBSAMPLE:

```
Train RMSE: 0.8240 CV RMSE: 0.8508
Train RMSE: 0.8235 CV RMSE: 0.8506
Train RMSE: 0.8232 CV RMSE: 0.8495
Train RMSE: 0.8233 CV RMSE: 0.8496
```

Tunning COLSAMPLE\_BYTREE:

```
Train RMSE: 0.8232 CV RMSE: 0.8495
Train RMSE: 0.8220 CV RMSE: 0.8501
Train RMSE: 0.8213 CV RMSE: 0.8506
Train RMSE: 0.8205 CV RMSE: 0.8510
```

Tunning REG\_ALPHA:

```
Train RMSE: 0.8232 CV RMSE: 0.8496
Train RMSE: 0.8232 CV RMSE: 0.8495
Train RMSE: 0.8232 CV RMSE: 0.8495
Train RMSE: 0.8232 CV RMSE: 0.8494
Train RMSE: 0.8232 CV RMSE: 0.8495
```

CPU times: user 27.9 s, sys: 572 ms, total: 28.5 s

Wall time: 1min 15s

## [9.2]Optimal value of parameters after tuning:

```
In [266]: 1 print('Optimal Value of Hyperparameters after Tuning:\n')
2 param_list(params_,param_summ)
3
```

Optimal Value of Hyperparameters after Tuning:

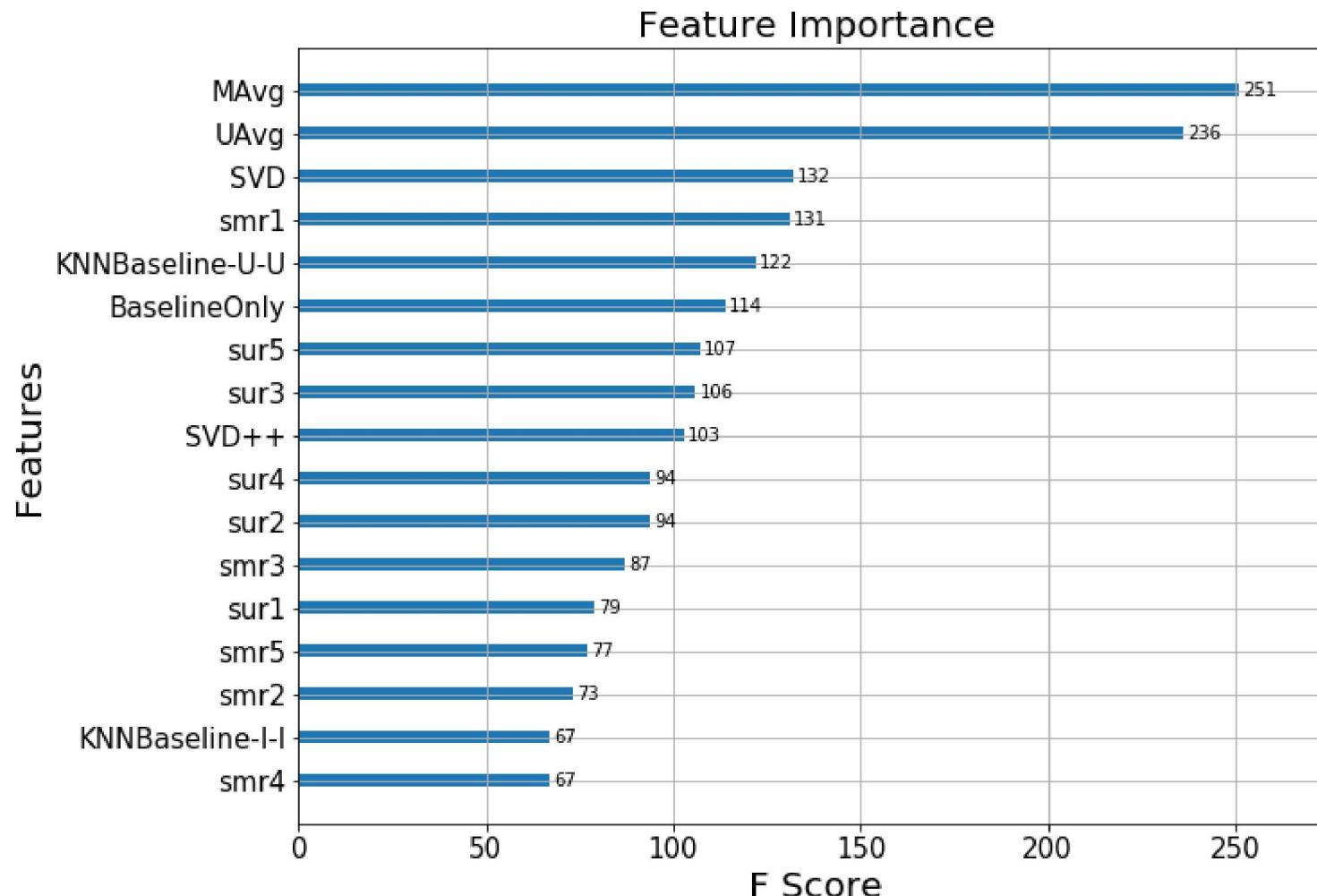
Parameter	Value
n_estimators	64
max_depth	5
min_child_weight	1
gamma	0.1
subsample	0.8
colsample_bytree	0.6
reg_alpha	0.01

[9.3]Test performance :

In [269]:

```
1 %%time
2 train_result, test_result, model_local_summ=test_performance_xgb(train, y_train, test, y_test,\n3                                         params_, model_summary, model_name[8], summary=True
4
5 #SAVE THE TEST RESULTS TO DICT OF DICT
6 models_evaluation_train[model_name[8]] = train_result
7 models_evaluation_test[model_name[8]] = test_result
8
```

FOR OPTIMAL PARAMETERS, TRAIN RMSE: 0.82766, TEST RMSE: 1.06981



CPU times: user 4.58 s, sys: 332 ms, total: 4.92 s

Wall time: 4.46 s

#### [9.4] model summary:

In [270]:

```
1 #MODEL_SUMMARY  
2 print(model_summ_local)
```

Model	Train(RMSE)	Train(MAPE)	Test(RMSE)	Test(MAPE)
SVD++	0.67078	19.84641	1.07450	35.11816

## 10. XGBOOST( BSL + KNN + MF[SVD+SVDpp] ):

#### [10.1] Hyperparameter Tuning:

In [274]:

```

1  %%time
2  #SELECT COLUMNS
3  cols=list(train_reg.columns)
4  del cols[0:16]
5
6  #TRAIN AND TEST DATA FOR XGBOOST MODELS
7  train, test = train_reg[cols], test_reg[cols]
8  y_train, y_test = train_reg['rating'], test_reg['rating']
9
10 print('HYPERPARAMETER:\n')
11 param_list(params_range, param_summ)
12 print()
13 #HYPERPARAM TUNNING
14 model, params_ = tuneALL_PARAM_XGB(train, y_train, CV, params_range, params_, searchMethod[0])

```

HYPERPARAMETER:

Parameter	Value
n_estimators	[32, 64, 80, 100]
max_depth	[5, 7, 9]
min_child_weight	[1, 3, 5]
gamma	[0.0, 0.1, 0.2, 0.3, 0.4]
subsample	[0.6, 0.7, 0.8, 0.9]
colsample_bytree	[0.6, 0.7, 0.8, 0.9]
reg_alpha	[0, 0.001, 0.005, 0.01, 0.05]

Tunning N\_ESTIMATORS:

Train RMSE: 1.0447 CV RMSE: 1.0695  
 Train RMSE: 1.0571 CV RMSE: 1.0739  
 Train RMSE: 1.0415 CV RMSE: 1.0698

Tunning MAX\_DEPTH:

Train RMSE: 1.0447 CV RMSE: 1.0695  
 Train RMSE: 1.0292 CV RMSE: 1.0710  
 Train RMSE: 1.0060 CV RMSE: 1.0727

Tunning MIN\_CHILD\_WEIGHT:

Train RMSE: 1.0439 CV RMSE: 1.0701

```
Train RMSE: 1.0442 CV RMSE: 1.0699
Train RMSE: 1.0447 CV RMSE: 1.0695
```

Tunning GAMMA:

```
Train RMSE: 1.0447 CV RMSE: 1.0695
Train RMSE: 1.0448 CV RMSE: 1.0697
Train RMSE: 1.0447 CV RMSE: 1.0694
```

Tunning SUBSAMPLE:

```
Train RMSE: 1.0445 CV RMSE: 1.0697
Train RMSE: 1.0454 CV RMSE: 1.0702
Train RMSE: 1.0447 CV RMSE: 1.0694
```

Tunning COLSAMPLE\_BYTREE:

```
Train RMSE: 1.0453 CV RMSE: 1.0698
Train RMSE: 1.0447 CV RMSE: 1.0694
Train RMSE: 1.0447 CV RMSE: 1.0694
```

Tunning REG\_ALPHA:

```
Train RMSE: 1.0448 CV RMSE: 1.0695
Train RMSE: 1.0447 CV RMSE: 1.0694
Train RMSE: 1.0450 CV RMSE: 1.0696
```

```
CPU times: user 18.6 s, sys: 112 ms, total: 18.7 s
Wall time: 43.6 s
```

## [10.2]Optimal value of parameters after tuning:

In [276]:

```
1 print('Optimal Value of Hyperparameters after Tuning:\n')
2 param_list(params_,param_summ)
3
```

Optimal Value of Hyperparameters after Tuning:

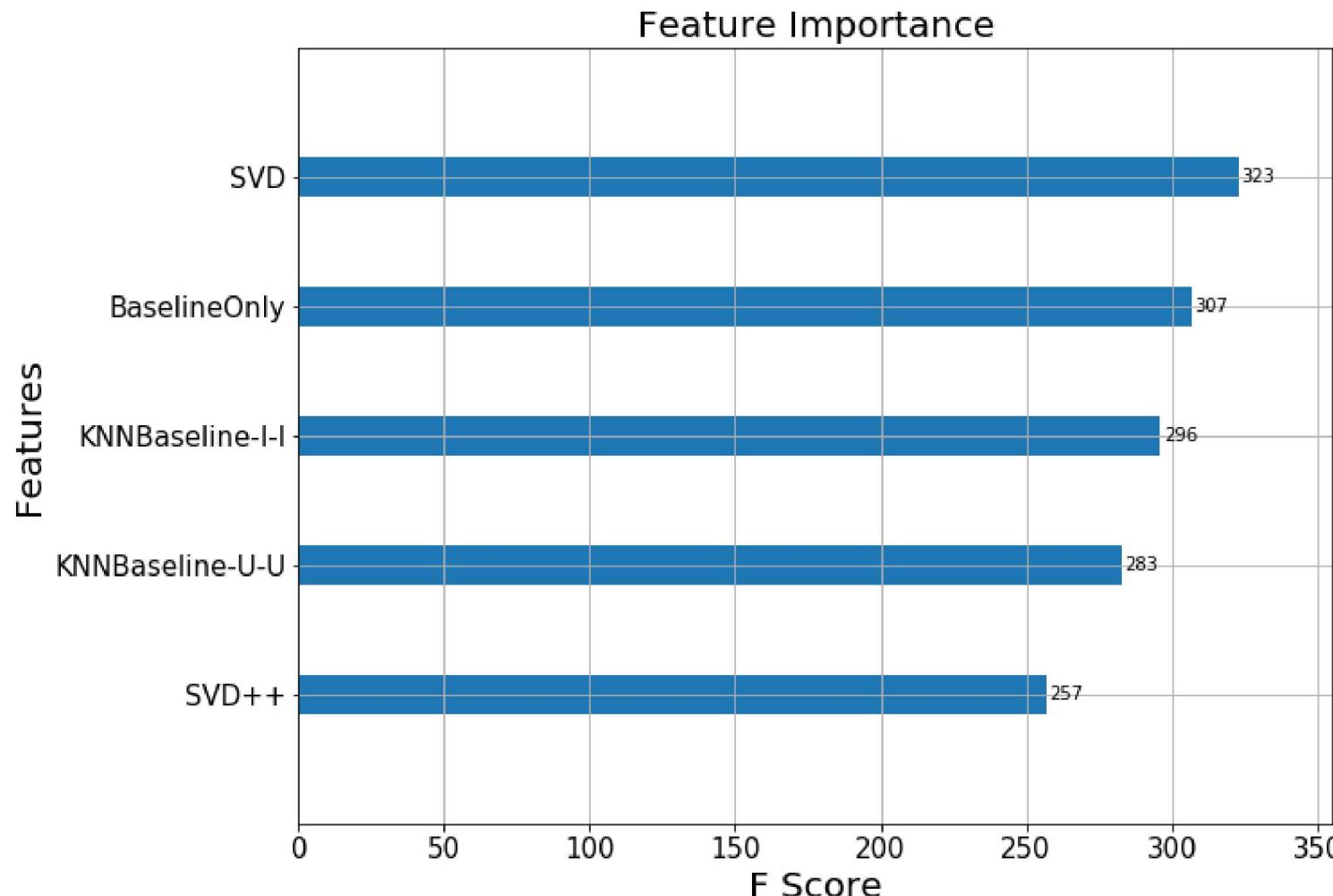
Parameter	Value
n_estimators	64
max_depth	5
min_child_weight	5
gamma	0.3
subsample	0.9
colsample_bytree	0.9
reg_alpha	0.005

[10.3]Test performance :

In [277]:

```
1 %%time
2 train_result, test_result, model_local_summ=test_performance_xgb(train, y_train, test, y_test,\n3                                         params_, model_summary, model_name[9], summary=True
4
5 #SAVE THE TEST RESULTS TO DICT OF DICT
6 models_evaluation_train[model_name[9]] = train_result
7 models_evaluation_test[model_name[9]] = test_result
8
```

FOR OPTIMAL PARAMETERS, TRAIN RMSE: 1.04781, TEST RMSE: 1.07510



CPU times: user 3.42 s, sys: 996 ms, total: 4.42 s

Wall time: 3.03 s

#### [10.4] Model Summary:

In [280]:

```
1 #MODEL_SUMMARY  
2 print(model_local_summ)
```

Model	Train(RMSE)	Train(MAPE)	Test(RMSE)	Test(MAPE)
XGBOOST-BSL-KNN-MF	1.04781	33.96307	1.07510	35.23055



## Conclusion:

In [281]:

```
1 print(model_summary)
```

Model	Train(RMSE)	Train(MAPE)	Test(RMSE)	Test(MAPE)
XGBOOST-13-FEATS	0.82606	24.16566	1.07701	33.59137
BaselineOnly	0.83552	25.25685	1.07429	35.12019
XGBOOST-BSL	0.82509	24.14721	1.06767	34.57024
KNNBaseline-U-U	0.14034	3.41399	1.07426	35.11784
KNNBaseline-I-I	0.27398	6.45966	1.07461	35.19819
XGBOOST-KNN	0.82392	24.12600	1.06919	34.61556
SVD	0.84508	25.80980	1.07418	35.13581
SVD++	0.67078	19.84641	1.07450	35.11816
XGBOOST-13-BSL-KNN-MF	0.82766	24.27682	1.06981	34.19354
XGBOOST-BSL-KNN-MF	1.04781	33.96307	1.07510	35.23055

## Best performance given by model:

1. XGBOOST+BSL with RMSE=1.06767
2. XGBOOST+KNN with RMSE=1.06919

## Procedure:

1. We have to solve a problem, where we have to predict a rating given by user to movie not yet watched by that user.
2. We can map this problem to Machine Learning problem in 2 ways:
  - a. Recommendation system.
  - b. Regression problem.
3. We proceed with loading the dataset and with basic details of the dataset:
  - a. How many users are present in dataset?
  - b. How many movies are present in dataset?
  - c. How many ratings are present?
4. As the time stamps are given in dataset so we sorted the data in ascending order of time and splitted it into train and test data for further analysis(Temporal Fashion).
5. We did EDA and Feature Extraction:
  - a. Ratings:
    - >ratings per movie.
    - >how many movies a user rated?
    - >PDF and CDF of ratings given by a user.
  - b. Build movie-movie similarity matrix.
  - c. As we have huge amount of data we build sample train data for train data and sample test data from test data .
  - d. We feturize sample train data and sample test data for regression:
    - >global avg, user avg and movie avg.
    - >sur = top 5 similar users rating for a movie.
    - >smr = top 5 similar movies rating by a user.
    - >handled cold start problem(new user / new movie) by taking global average.

6. We convert data so that accepted by surprise library.

7. We Build ML models:

a. Regression(XGBOOST) model.

b. Surprise library models:

->BaselineOnly

->KNNBaseline

->SVD

->SVDpp

**Note:**

**-After applying each of the surprise library model we add output of that model as a new feature to our regression data and build regression model on that data.**

**-Ways to handle cold start problem:**

1. Recommend globally top movies initially to a new user.
2. From ip address we can find information like user belongs to which region, and based on that region we recommend movies which are popular in that region to a new user.
3. Recommend movies/items based on the few of the items/movies watched by the user.