

```
In [15]: import numpy as np
import pandas as pd
from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import log_loss
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
import pickle as pkl
from scipy.stats import randint as sp_randint
```

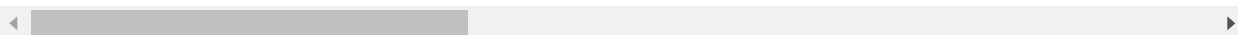
```
In [16]: train_df = pd.read_csv('X_train2.csv')
# train_df = pd.read_csv('X_train1.csv')
```

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

Out[17]:

	id	age	year	month	day	tfa_year	tfa_month	tfa_day	timediff	gender_ unknown-	...	view
0	d1mm9tcy42	62.0	2014	1	1	2014	1	1	0	0	...	
1	yo8nz8bqcq	-1.0	2014	1	1	2014	1	1	0	1	...	
2	4grx6yxey	-1.0	2014	1	1	2014	1	1	0	1	...	
3	ncf87guaf0	-1.0	2014	1	1	2014	1	1	0	1	...	
4	4rvqpxoh3h	-1.0	2014	1	1	2014	1	1	0	1	...	

5 rows × 321 columns



```
In [19]: train_df.set_index('id',inplace=True)
```

```
In [20]: with open('labels.pkl','rb') as f:
Y = pkl.load(f)
print(Y.shape)
```

(73812,)

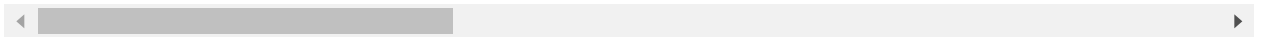
```
In [21]: # with open('labels1.pkl','rb') as f:
# Y = pkl.load(f)
# print(Y.shape)
```

```
In [22]: train_df.head()
```

```
Out[22]:
```

	age	year	month	day	tfa_year	tfa_month	tfa_day	timediff	gender_ unknown-	gender_FE
id										
d1mm9tcy42	62.0	2014	1	1	2014	1	1	0	0	
yo8nz8bqcq	-1.0	2014	1	1	2014	1	1	0	1	
4grx6yxeby	-1.0	2014	1	1	2014	1	1	0	1	
ncf87guaf0	-1.0	2014	1	1	2014	1	1	0	1	
4rvqpxoh3h	-1.0	2014	1	1	2014	1	1	0	1	

5 rows × 320 columns



In [23]:

```
x_cfl= RandomForestClassifier()

prams={
    'min_samples_split':[2,20],
    'n_estimators':[100,200,500,1000,2000],
    'max_depth':[3,5,10]
}
random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,cv= 3,verbose=10,n
random_cfl.fit(train_df,Y)
```

Fitting 3 folds for each of 10 candidates, totalling 30 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done   5 tasks      | elapsed:   19.1s
[Parallel(n_jobs=-1)]: Done  10 tasks      | elapsed:   50.3s
[Parallel(n_jobs=-1)]: Done  17 tasks      | elapsed:   2.9min
[Parallel(n_jobs=-1)]: Done  27 out of  30 | elapsed:   4.7min remaining:   31.3
s
[Parallel(n_jobs=-1)]: Done  30 out of  30 | elapsed:   4.8min finished
```

```
Out[23]: RandomizedSearchCV(cv=3, error_score='raise-deprecating',
    estimator=RandomForestClassifier(bootstrap=True, class_weight=None, c
riterion='gini',
    max_depth=None, max_features='auto', max_leaf_nodes=None,
    min_impurity_decrease=0.0, min_impurity_split=None,
    min_samples_leaf=1, min_samples_split=2,
    min_weight_fraction_leaf=0.0, n_estimators='warn', n_jobs=None,
    oob_score=False, random_state=None, verbose=0,
    warm_start=False),
    fit_params=None, iid='warn', n_iter=10, n_jobs=-1,
    param_distributions={'min_samples_split': [2, 20], 'n_estimators': [1
00, 200, 500, 1000, 2000], 'max_depth': [3, 5, 10]},
    pre_dispatch='2*n_jobs', random_state=None, refit=True,
    return_train_score='warn', scoring=None, verbose=10)
```

```
In [24]: # displaying the best parameters
random_cfl.best_params_
```

```
Out[24]: {'n_estimators': 500, 'min_samples_split': 2, 'max_depth': 10}
```

```
In [27]: #Using the best parameters to train the model
x_cfl=RandomForestClassifier(n_estimators=500,min_samples_split=2,max_depth=10)
x_cfl.fit(train_df,Y)
```

```
Out[27]: RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
    max_depth=10, max_features='auto', max_leaf_nodes=None,
    min_impurity_decrease=0.0, min_impurity_split=None,
    min_samples_leaf=1, min_samples_split=2,
    min_weight_fraction_leaf=0.0, n_estimators=500, n_jobs=None,
    oob_score=False, random_state=None, verbose=0,
    warm_start=False)
```

```
In [28]: #storing the model in a pickle file
import pickle
pickle.dump(x_cfl,open('RF.pickle.dat','wb'))
```

```
In [29]: classifier = pickle.load(open('RF.pickle.dat','rb'))
```

```
In [30]: test_df = pd.read_csv('X_test2.csv')
```

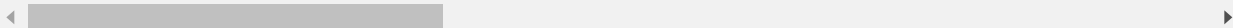
```
In [37]: # test_df.drop('Unnamed: 0',axis=1,inplace=True)
```

```
In [38]: test_df.head(15)
```

Out[38]:

	age	year	month	day	tfa_year	tfa_month	tfa_day	timediff	gender_ unknown-	gender_FEM
id										
5uwns89zht	35.0	2014	7	1	2014	7	1	0	0	
jtl0dijy2j	-1.0	2014	7	1	2014	7	1	0	1	
xx0ulgorjt	-1.0	2014	7	1	2014	7	1	0	1	
6c6puo6ix0	-1.0	2014	7	1	2014	7	1	0	1	
czqhjk3yfe	-1.0	2014	7	1	2014	7	1	0	1	
szx28ujmhf	28.0	2014	7	1	2014	7	1	0	0	
guenkfjcbq	48.0	2014	7	1	2014	7	1	0	0	
tkpq0mlugk	-1.0	2014	7	1	2014	7	1	0	1	
3xtgd5p9dn	-1.0	2014	7	1	2014	7	1	0	1	
md9aj22l5a	-1.0	2014	7	1	2014	7	1	0	1	
gg3eswjxdf	-1.0	2014	7	1	2014	7	1	0	1	
fyomoivygn	30.0	2014	7	1	2014	7	1	0	0	
iq4kkd5oan	24.0	2014	7	1	2014	7	1	0	0	
6k1xls6x5j	-1.0	2014	7	1	2014	7	1	0	1	
jodmb2ok1f	-1.0	2014	7	1	2014	7	1	0	1	

15 rows × 320 columns



```
In [40]: test_df.set_index('id',inplace=True)
```

In [41]: `test_df.head()`

Out[41]:

	age	year	month	day	tfa_year	tfa_month	tfa_day	timediff	gender_- unknown-	gender_FEM
id										
5uwns89zht	35.0	2014	7	1	2014	7	1	0	0	
jtl0dijy2j	-1.0	2014	7	1	2014	7	1	0	1	
xx0ulgorjt	-1.0	2014	7	1	2014	7	1	0	1	
6c6puo6ix0	-1.0	2014	7	1	2014	7	1	0	1	
czqjhjk3yfe	-1.0	2014	7	1	2014	7	1	0	1	

5 rows × 320 columns

In [42]: `# since in the problem statement it is mentioned that the we need to predict the
pred_probab = classifier.predict_proba(test_df)`

In [43]: `# storing the predictions of each user_id in a dataframe with user_id as the index
pred_probab_df = pd.DataFrame(pred_probab, index=test_df.index)`

In [44]: `pred_probab_df.head()`

Out[44]:

	0	1	2	3	4	5	6	7	
id									
5uwns89zht	0.001502	0.003722	0.002181	0.007116	0.012815	0.006783	0.007760	0.699447	0.002
jtl0dijy2j	0.000831	0.002190	0.001067	0.003871	0.008386	0.003826	0.005727	0.850855	0.001
xx0ulgorjt	0.000902	0.003111	0.001440	0.005313	0.011489	0.004990	0.007966	0.824988	0.001
6c6puo6ix0	0.000877	0.003087	0.001389	0.005552	0.011619	0.004953	0.007641	0.826830	0.001
czqjhjk3yfe	0.001580	0.013982	0.003514	0.019310	0.045565	0.020012	0.037435	0.212242	0.009

```
In [45]: # The dictionary is the label encoding of the countries feature
output_classes = {'AU': 0,
                  'CA': 1,
                  'DE': 2,
                  'ES': 3,
                  'FR': 4,
                  'GB': 5,
                  'IT': 6,
                  'NDF': 7,
                  'NL': 8,
                  'PT': 9,
                  'US': 10,
                  'other': 11}
```

```
In [46]: # inverting the dictionary
inv_classes = {v:k for k,v in output_classes.items()}
```

```
In [47]: inv_classes
```

```
Out[47]: {0: 'AU',
          1: 'CA',
          2: 'DE',
          3: 'ES',
          4: 'FR',
          5: 'GB',
          6: 'IT',
          7: 'NDF',
          8: 'NL',
          9: 'PT',
          10: 'US',
          11: 'other'}
```

```
In [48]: # taking the indices from 0-11
indices = np.arange(0,12)
```

```
In [49]: #prediction values of the first user_id
pred_probab[0]
```

```
Out[49]: array([1.50174261e-03, 3.72211676e-03, 2.18083470e-03, 7.11581193e-03,
                1.28147245e-02, 6.78309589e-03, 7.76009020e-03, 6.99446630e-01,
                2.33953031e-03, 6.59914619e-04, 2.20036945e-01, 3.56385638e-02])
```

```
In [50]: # creating a dictionary of the predictio and indices value
pred_dict = dict(zip(indices,pred_probab[0]))
```

```
In [51]: # sorting the dictionary and taking only the top 5 values
sorted_abc = sorted(pred_dict.items(),key=lambda x:x[1],reverse=True)[:5]
```

```
In [52]: sorted_abc
```

```
Out[52]: [(7, 0.6994466295747693),
          (10, 0.22003694510947955),
          (11, 0.035638563773071855),
          (4, 0.01281472452314817),
          (6, 0.007760090202513086)]
```

```
In [53]: # taking only the index value of the tuple sorted_abc
row_indices = [x[0] for x in sorted_abc]
```

```
In [54]: row_indices
```

```
Out[54]: [7, 10, 11, 4, 6]
```

```
In [55]: # taking the indices and giving the country names
top_five = [inv_classes[i] for i in row_indices]
```

```
In [56]: top_five
```

```
Out[56]: ['NDF', 'US', 'other', 'FR', 'IT']
```

```
In [57]: type(top_five)
```

```
Out[57]: list
```

```
In [58]: # Combining the above steps into a fuction so that it can be applied to the predi
def top_5_countries(s):
    """
    This function takes the probability values of each id, sorts the top 5 values
    """
    indices = np.arange(0,12)
    pred_dict = dict(zip(indices,s))
    sorted_abc = sorted(pred_dict.items(),key=lambda x:x[1],reverse=True)[:5]
    row_indices = [x[0] for x in sorted_abc]
    top_five = [inv_classes[i] for i in row_indices]
    return top_five
```

```
In [59]: # here we apply the above function on each row of the dataframe to get the top 5
pred_probab_df['top_five'] = pred_probab_df.apply(top_5_countries,axis=1)
```

In [60]: `pred_probab_df.head()`

Out[60]:

	0	1	2	3	4	5	6	7	
id									
5uwns89zht	0.001502	0.003722	0.002181	0.007116	0.012815	0.006783	0.007760	0.699447	0.002
jtl0dijy2j	0.000831	0.002190	0.001067	0.003871	0.008386	0.003826	0.005727	0.850855	0.001
xx0ulgorjt	0.000902	0.003111	0.001440	0.005313	0.011489	0.004990	0.007966	0.824988	0.001
6c6puo6ix0	0.000877	0.003087	0.001389	0.005552	0.011619	0.004953	0.007641	0.826830	0.001
czqjhjk3yfe	0.001580	0.013982	0.003514	0.019310	0.045565	0.020012	0.037435	0.212242	0.009

In [61]: `# ungrouping the list values of the top_five column`
`s = pred_probab_df.apply(lambda x: pd.Series(x['top_five']),axis=1).stack().reset`
`s.name = 'country'`

In [62]: `submission = pred_probab_df.drop([i for i in range(0,12)] + ['top_five'],axis=1).`
`submission.head()`

Out[62]:

	country
id	
0010k6l0om	NDF
0010k6l0om	US
0010k6l0om	other
0010k6l0om	FR
0010k6l0om	IT

In [63]: `submission.to_csv('RFsubmission.csv')`

The final Public Score(ndcg)

Your most recent submission				
Name	Submitted	Wait time	Execution time	Score
RFsubmission.csv	a few seconds ago	0 seconds	5 seconds	0.87304
Complete				
Jump to your position on the leaderboard ▼				

The final public and private score

Submission and Description	Private Score	Public Score	Use for Final Score
RFsubmission.csv a few seconds ago by AdityaBantwal add submission details	0.87817	0.87304	<input type="checkbox"/>

Conclusion

The official Kaggle score is 0.87304 for a Random Forest model which can be further improved by using some text features from the train and session data and doing more hyper parameter tuning. I would like to conclude this notebook here.

Thank You!

In []: