

```
[3]: import warnings
warnings.filterwarnings("ignore")
import shutil
import os
import pandas as pd
import matplotlib
matplotlib.use("nbAgg")
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
import pickle
from sklearn.manifold import TSNE
from sklearn import preprocessing
import pandas as pd
from multiprocessing import Process # this is used for multithreading
import multiprocessing
import codecs # this is used for file operations
import random as r
from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.calibration import CalibratedClassifierCV
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

In [14]: df=pd.read_csv("microsoft_final_file.csv")
df.shape

Out[14]: (18866, 2460)

In [15]: y=df['Class_y']

In [16]: x=df.drop(['Class_y','ID'],axis=1)

Train_CV_test_Split
```

```
[47]: X_train,X_test,y_train,y_test=train_test_split(X,Y,stratify=y,test_size=0.20)
      X_train,X_cv,y_train,y_cv=train_test_split(X_train,y_train,stratify=y_train,test_size=0.20)

In [48]: print('Number of data points in train data:', X_train.shape[0])
      print('Number of data points in test data:', X_test.shape[0])
      print('Number of data points in cross validation data:', X_cv.shape[0])

      Number of data points in train data: 6953
      Number of data points in test data: 2274
      Number of data points in cross validation data: 1739

In [49]: def plot_confusion_matrix(test_y, predict_y):
      C = confusion_matrix(test_y, predict_y)
      print("Number of misclassified points ",(len(test_y)-np.trace(C))/len(test_y)*100)
      # C = 9,9 matrix, each cell (i,j) represents number of points of class i are predicted class j
```

```
A = ((C.T)/(C.sum(axis=1))).T
#divide each element of the confusion matrix with the sum of elements in that column

# C = [[1, 2],
#       [3, 4]]
# C.T = [[1, 3],
#         [2, 4]]
# C.sum(axis = 1) axis=0 corresponds to columns and axis=1 corresponds to rows in two dimensional array
# C.sum(axis = 1) = [[3, 7]]
# ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7]
#                             [2/3, 4/7]]

# ((C.T)/(C.sum(axis=1))) T = [[1/3, 2/3]
#                               [3/7, 4/7]]
```

```

# (C[:,0] + C[:,1] + C[:,2]) / 3
# sum of row elements = 1 (3/3, 4/3)

B = (C/C.sum(axis=0))
#divides each element of the confusion matrix with the sum of elements in that row
# C = [[2, 2],
#      [3, 4]]
# C.sum(axis = 0) axis=0 corresponds to columns and axis=1 corresponds to rows in two dimensional array
# C.sum(axis=0) = [4, 6]
# (C/C.sum(axis=0)) = [[1/4, 2/6],
#                      [3/4, 4/6]]

labels = [1, 2, 3, 4, 5, 6, 7, 8, 9]

```

```

cmap=sns.light_palette("green")
# representing A in heatmap format
print("---", "Confusion matrix", "\n---")
plt.figure(figsize=(10,5))
sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.show()

print("---", "Precision matrix", "\n---")
plt.figure(figsize=(10,5))
sns.heatmap(P, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')

```

```
plt.xlabel('Original Class')
plt.show()
print('Sum of columns in precision matrix', B.sum(axis=0))

# representing B in heatmap format
print("n=","50", "Recall matrix" , "n=","50")
plt.figure(figsize=(10,5))
sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.show()
print('Sum of rows in precision matrix', A.sum(axis=1))
```

```
In [21]: alpha = [x for x in range(1, 15, 2)]
cv_log_error_array=[]
for i in alpha:
    k_cflf=KNeighborsClassifier(n_neighbors=1)
    k_cflf.fit(X_train,y_train)
    sig_cflf = CalibratedClassifierCV(k_cflf, method="sigmoid")
    sig_cflf.fit(X_train, y_train)
    predict_y = sig_cflf.predict_proba(X_cv)
    cv_log_error_array.append(log_loss(y_cv, predict_y, labels=k_cflf.classes_, eps=1e-15))

for i in range(len(cv_log_error_array)):
    print("alpha: %d cv_log_error: %f" % (alpha[i], cv_log_error_array[i]))
```

```
print (log_loss for k = ,alpha[i], is_cv_log_error_array[i])

best_alpha = np.argmax(cv_log_error_array)

fig, ax = plt.subplots()
ax.plot(alpha, cv_log_error_array, c='g')
for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate([alpha[i],np.round(txt,3)], (alpha[i],cv_log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()
```

```
k_clf = KNeighborsClassifier(n_neighbors=alpha[best_alpha])
k_clf.fit(X_train, y_train)
sig_clf = CalibratedClassifierCV(k_clf, method="sigmoid")
sig_clf.fit(X_train, y_train)

predict_y = sig_clf.predict_proba(X_train)
print("For values of best alpha = ", alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y))
predict_y = sig_clf.predict_proba(X_cv)
print("For values of best alpha = ", alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y))
predict_y = sig_clf.predict_proba(X_test)
print("For values of best alpha = ", alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y))
plot_confusion_matrix(x_test, sig_clf.predict(X_test))
```

```
log_loss for k = 1 is 0.33104582364681093
log_loss for k = 3 is 0.31847154913664727
log_loss for k = 5 is 0.3315458954549312
log_loss for k = 7 is 0.34846444337334
log_loss for k = 9 is 0.3618248768628205
log_loss for k = 11 is 0.37421866440751402
log_loss for k = 13 is 0.3861672771048972
```

n	error measure
12	0.35
19	0.362
21	0.374
25	0.386

For values of best alpha = 3 The train log loss is: 0.16367879563562591

```

For values of best alpha = 3 The cross validation log loss is: 0.31847154913664727
For values of best alpha = 3 The test log loss is: 0.28938572637564644
Number of misclassified points: 7.311681572680893
----- Confusion matrix -----

```

[illegible]

Confusion matrix

	1	2	3	4	5	6	7	8	9
1	11,000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
2	8,000	0.000	0.000	1,000	0.000	9,000	0.000	0.000	0.000
3	31,000	6,000	0.000	1,000	0.000	3,000	0.000	1,000	160,000
4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
5	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
6	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
7	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
8	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
9	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000

Predicted Class

Precision matrix

[illegible]

Original City	1	2	3	4	5	6	7	8	9	10
1	0.003	0.000	0.002	0.000	0.500	0.000	0.012	0.016	0.000	0.000
2	0.032	0.000	0.000	0.030	0.500	0.815	0.000	0.012	0.000	0.000
3	0.000	0.000	0.000	0.000	0.000	0.000	0.855	0.012	0.033	0.000
4	0.023	0.000	0.000	0.010	0.000	0.056	0.000	0.908	0.011	0.000
5	0.089	0.013	0.000	0.010	0.000	0.019	0.000	0.004	0.874	0.000

Sum of columns in precision matrix [1. 1. 1. 1. 1. 1. 1. 1. 1.]

----- Recall matrix -----

	1	2	3	4	5	6	7	8	9
1	0.912	0.006	0.003	0.003	0.000	0.036	0.010	0.019	0.010

Original Class	0	1	2	3	4	5	6	7	8	9
0	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
1	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
2	0.030	0.909	0.000	0.000	0.000	0.014	0.016	0.006	0.024	0.000
3	0.000	0.000	0.997	0.000	0.000	0.000	0.000	0.003	0.000	0.000
4	0.011	0.000	0.000	0.979	0.000	0.000	0.000	0.011	0.000	0.000
5	0.125	0.000	0.125	0.000	0.125	0.000	0.125	0.500	0.000	0.000
6	0.073	0.000	0.000	0.020	0.007	0.880	0.000	0.020	0.000	0.000

	1	2	3	4	5	6	7	8	9
1	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
2	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
3	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000
4	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000
5	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000
6	0.000	0.000	0.000	0.000	0.000	1.000	0.897	0.037	0.075
7	0.000	0.000	0.000	0.000	0.000	0.897	1.000	0.000	0.008
8	0.000	0.000	0.000	0.000	0.000	0.037	0.000	1.000	0.792
9	0.000	0.000	0.000	0.000	0.000	0.075	0.008	0.792	1.000

Sum of rows in precision matrix [1. 1. 1. 1. 1. 1. 1. 1. 1.]

```
[In]: # Training a hyper-parameter tuned Xg-Boost regressor on our train data

# find more about XGBClassifier function here http://xgboost.readthedocs.io/en/latest/python/python_api.html?xgb=XGBClassifier
# -----
# default parameters
# class xgboost.XGBClassifier(max_depth=6, learning_rate=0.1, n_estimators=100, silent=True,
#                             verbose=False, warm_start=None, booster='gbtree', n_jobs=-1, ntrees=None, gamma=0, min_child_weight=1,
#                             max_delta_step=0, subsample=1, colsample_bytree=1, colsample_bynode=1, reg_alpha=0, reg_lambda=1,
#                             scale_pos_weight=1, base_score=0.5, random_state=None, seed=None, missing=None, **kwargs)
# some of methods of RandomForestRegressor()
# -----
# Note: Some of the parameters are deprecated and will be removed in future releases. Please check documentation for details.
```

```

#-----
# Get parameters for this estimator.
# predict(data, output_margin=False, ntree_line=0) : Predict with data. NOTE: This function is not thread safe.
# get_score(importance_type='weight') -> get the feature importance
#-----
# video link: https://www.appliedaiacourse.com/course/applai-ai-course-online/lessons/regression-using-decision-trees-2/
# video link: https://www.appliedaiacourse.com/course/applai-ai-course-online/lessons/what-are-ensembles/
#-----

alpha=[10,50,100,500,1000,2000]
cv_log_error_array=[]
for i in alpha:
    x_cfl=XGBClassifier(n_estimators=i,nthread=-1)
    x_cfl.fit(X_train,y_train)

```

```
sig_clf = CalibratedClassifierCV(x_cv, method="sigmoid")
sig_clf.fit(X_train, y_train)
predict_y = sig_clf.predict_proba(X_test)
cv_log_error_array.append(log_loss(y_cv, predict_y, labels=x_cv.classes_, eps=1e-15))

for i in range(len(cv_log_error_array)):
    print('log_loss for c = ', alpha[i], 'is', cv_log_error_array[i])

best_alpha = np.argmax(cv_log_error_array)
fig, ax = plt.subplots()
check_labels = fig.get_labels()
fig.set_xlabel('log_loss array: %s' % check_labels)
```

```

for i, txt in enumerate(np.round(cv_log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],cv_log_error_array[i]))
    grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

x_cfl=XGBClassifier(n_estimators=alpha[best_alpha],nthread=-1)
x_cfl.fit(X_train,y_train)
sig_cfl = CalibratedClassifierCV(x_cfl, method='sigmoid')
sig_cfl.fit(X_train,y_train)

```

```

predict_y = sig_clf.predict_proba(X_train)
print('For values of best_alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y))
predict_y = sig_clf.predict_proba(X_cv)
print('For values of best_alpha = ', alpha[best_alpha], "The cross validation log loss is:", log_loss(y_cv, predict_y))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best_alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y))
plot_confusion_matrix(x_test, sig_clf.predict(X_test))

```

[11:12:02] WARNING: ~\src\learner\c1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'multi:softmax' was changed from 'merror' to 'mlogloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[11:12:09] WARNING: ~\src\learner\c1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'multi:softmax' was changed from 'merror' to 'mlogloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

```
[12:12:14] WARNING: ~$rsccleaner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'multi:softmax' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[12:12:20] WARNING: ~$rsccleaner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'multi:softmax' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[12:12:26] WARNING: ~$rsccleaner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'multi:softmax' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[12:12:32] WARNING: ~$rsccleaner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'multi:softmax' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[12:12:38] WARNING: ~$rsccleaner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'multi:softmax' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[12:12:44] WARNING: ~$rsccleaner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'multi:softmax' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[12:12:50] WARNING: ~$rsccleaner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'multi:softmax' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[12:12:56] WARNING: ~$rsccleaner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'multi:softmax' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[12:13:02] WARNING: ~$rsccleaner.c:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'multi:softmax' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
```

[illegible][illegible][illegible][illegible]

```
Explicitly set validation metric to use to restore the validation
[11:58:13] WARNING: --xvrclearer.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective "multi:softmax" was changed from "merror" to "mlogloss". Explicitly set eval_metric if you'd like to restore the old behavior.
log_loss for c = 10 is 0.01216248490432462
log_loss for c = 50 is 0.00789971482454154
log_loss for c = 100 is 0.00790745321256206
log_loss for c = 500 is 0.007906980892597702
log_loss for c = 1000 is 0.007907356314128268
log_loss for c = 2000 is 0.007907501293380684
```

[illegible][illegible][illegible]

Confusion matrix visualization showing counts for predicted classes (1-9) against actual classes (1-9). The matrix is a 9x9 grid. The diagonal elements (top-left to bottom-right) are all 202,000, indicating perfect classification. The off-diagonal elements are all 0.000. The color scale ranges from 0 (light green) to 100 (dark green).

	1	2	3	4	5	6	7	8	9
1	202,000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
2	0.000	202,000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
3	0.000	0.000	202,000	0.000	0.000	0.000	0.000	0.000	0.000
4	0.000	0.000	0.000	202,000	0.000	0.000	0.000	0.000	0.000
5	0.000	0.000	0.000	0.000	202,000	0.000	0.000	0.000	0.000
6	0.000	0.000	0.000	0.000	0.000	202,000	0.000	0.000	0.000
7	0.000	0.000	0.000	0.000	0.000	0.000	202,000	0.000	0.000
8	0.000	0.000	0.000	0.000	0.000	0.000	0.000	202,000	0.000
9	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	202,000

Precision matrix

[illegible][illegible][illegible][illegible]

Sum of rows in precision matrix [1. 1. 1. 1. 1. 1. 1. 1. 1.]

For values of best alpha = 50 The train log loss is: 0.0077342595303011  
For values of best alpha = 50 The cross validation log loss is: 0.007899714824544154  
For values of best alpha = 50 The test log loss is: 0.00859675577434954  
Number of misclassified points 0.0

## Summary

IN this assignment i have directly used XGBoost for modeling

By using this xgboost model and calibration i have got a log loss of 0.008 which is has improved a lot as compared to other models

```
In [ ]:
```