# Bytes file

| ⏱ Created | @Jun 10, 2021 11:17 AM |
|---|---|
| ☰ Tags | |

## Exploratory Data Analysis

### Number of data points in each class



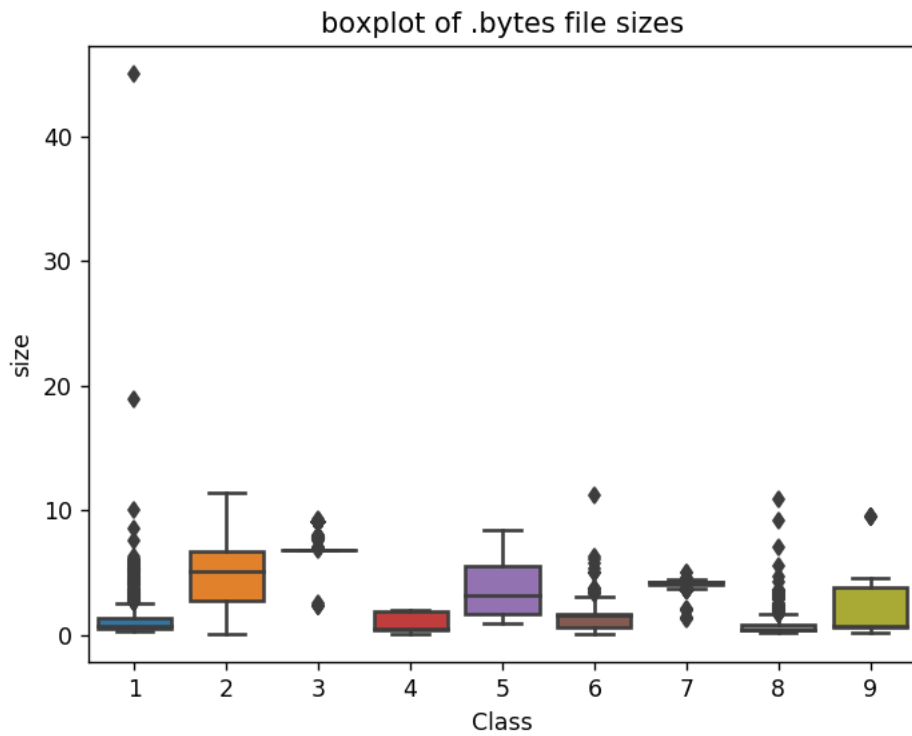**Observation**:- Class 5 has less number of data points, **Imbalance data problem.**

### File size as feature

| ID | File Name | Size |
|---|---|---|
| 0 | 01azqd4InC7m9JpocGv5 | 4.234863 |
| 1 | 01IsoiSMh5gxyDYTI4CB | 5.538818 |
| 2 | 01jsnpXSAlgw6aPeDxrU | 3.887939 |
| 3 | 01kcPWA9K2BOxQeS5Rju | 0.574219 |
| 4 | 01SuzwMJEIXsK7A8dQbl | 0.370850 |

### Box plot of file size as feature

boxplot of .bytes file sizes

**Observation:**- Class 2, 5 and 9 can be easily distinguished from other classes, using only the file size feature

**Copy of Bag of word as feature of the file**

| # ID | Aa File Name | # 0 | # 1 | # 2 | # 3 | # 4 | # 5 | # 6 | # 7 | # 8 | ... | # f9 | # fa | # fb | # fc |
|------|--------------|------|------|------|------|------|------|------|------|------|-----|------|------|------|------|
| 0 | 01azqd4InC7m9JpocGv5 | 601905 | 3905 | 2816 | 3832 | 3345 | 3242 | 3650 | 3201 | 2965 | ... | 3101 | 3211 | 3097 | 2758 |
| 1 | 01IsoiSMh5gxyDYTI4CB | 39755 | 8337 | 7249 | 7186 | 8663 | 6844 | 8420 | 7589 | 9291 | ... | 439 | 281 | 302 | 7639 |
| 2 | 01jsnpXSAlgw6aPeDxrU | 93506 | 9542 | 2568 | 2438 | 8925 | 9330 | 9007 | 2342 | 9107 | ... | 2242 | 2885 | 2863 | 2471 |
| 3 | 01kcPWA9K2BOxQeS5Rju | 21091 | 1213 | 726 | 817 | 1257 | 625 | 550 | 523 | 1078 | ... | 485 | 462 | 516 | 1133 |
| 4 | 01SuzwMJElXsK7A8dQbl | 19764 | 710 | 302 | 433 | 559 | 410 | 262 | 249 | 422 | ... | 350 | 209 | 239 | 653 |

**Copy of Combining Bag of Words and File size as Features**

| # ID | Aa File Name | # 0 | # 1 | # 2 | # 3 | # 4 | # 5 | # 6 | # 7 | # 8 | ... | # f9 | # fa | # fb | # fc |
|------|--------------|------|------|------|------|------|------|------|------|------|-----|------|------|------|------|
| 0 | 01azqd4InC7m9JpocGv5 | 601905 | 3905 | 2816 | 3832 | 3345 | 3242 | 3650 | 3201 | 2965 | ... | 3101 | 3211 | 3097 | 2758 |
| 1 | 01IsoiSMh5gxyDYTI4CB | 39755 | 8337 | 7249 | 7186 | 8663 | 6844 | 8420 | 7589 | 9291 | ... | 439 | 281 | 302 | 7639 |
| 2 | 01jsnpXSAlgw6aPeDxrU | 93506 | 9542 | 2568 | 2438 | 8925 | 9330 | 9007 | 2342 | 9107 | ... | 2242 | 2885 | 2863 | 2471 |
| 3 | 01kcPWA9K2BOxQeS5Rju | 21091 | 1213 | 726 | 817 | 1257 | 625 | 550 | 523 | 1078 | ... | 485 | 462 | 516 | 1133 |
| 4 | 01SuzwMJElXsK7A8dQbl | 19764 | 710 | 302 | 433 | 559 | 410 | 262 | 249 | 422 | ... | 350 | 209 | 239 | 653 |

**Copy of Normalizing the Features**

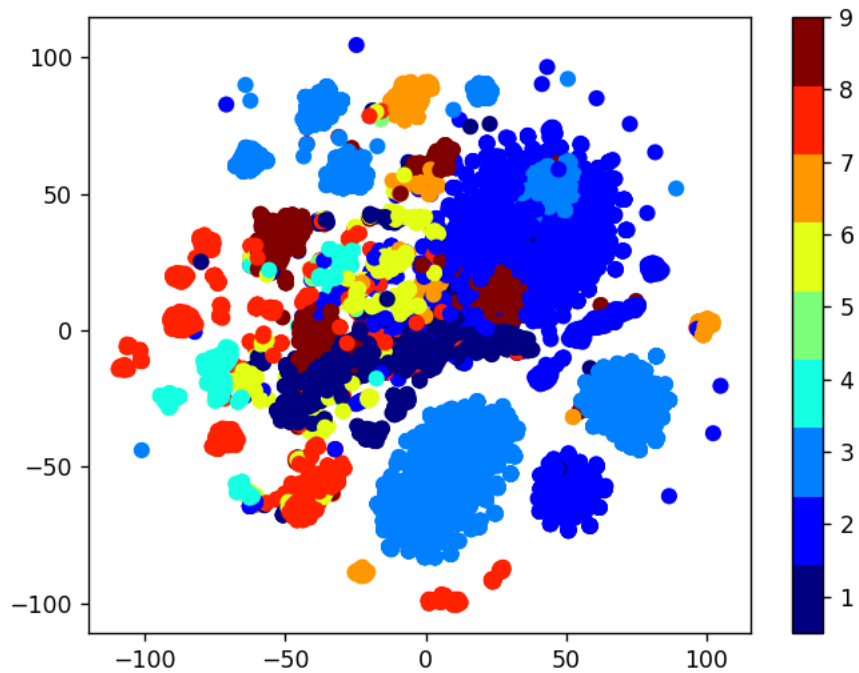| # ID | Aa File Name | # 0 | # 1 | # 2 | # 3 | # 4 | # 5 | # 6 | # 7 | # 8 |
|------|--------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0 | 01azqd4InC7m9JpocGv5 | 0.262806 | 0.005498 | 0.001567 | 0.002067 | 0.002048 | 0.001835 | 0.002058 | 0.002946 | 0.002638 |
| 1 | 01IsoiSMh5gxyDYTl4CB | 0.017358 | 0.011737 | 0.004033 | 0.003876 | 0.005303 | 0.003873 | 0.004747 | 0.006984 | 0.008267 |
| 2 | 01jsnpXSAlgw6aPeDxrU | 0.040827 | 0.013434 | 0.001429 | 0.001315 | 0.005464 | 0.00528 | 0.005078 | 0.002155 | 0.008104 |
| 3 | 01kcPWA9K2BOxQeS5Rju | 0.009209 | 0.001708 | 0.000404 | 0.000441 | 0.00077 | 0.000354 | 0.00031 | 0.000481 | 0.000959 |
| 4 | 01SuzwMJElXsK7A8dQbl | 0.008629 | 0.001 | 0.000168 | 0.000234 | 0.000342 | 0.000232 | 0.000148 | 0.000229 | 0.000376 |

# Multivariate analysis of the Features

## Perplexity = 50



**Observation**:- Class 2 and 3 are clearly separated whereas other classes have poor distinctions

## Perplexity = 30

**Observation**:- Class 2 and 3 are clearly separated whereas other classes have poor distinctions
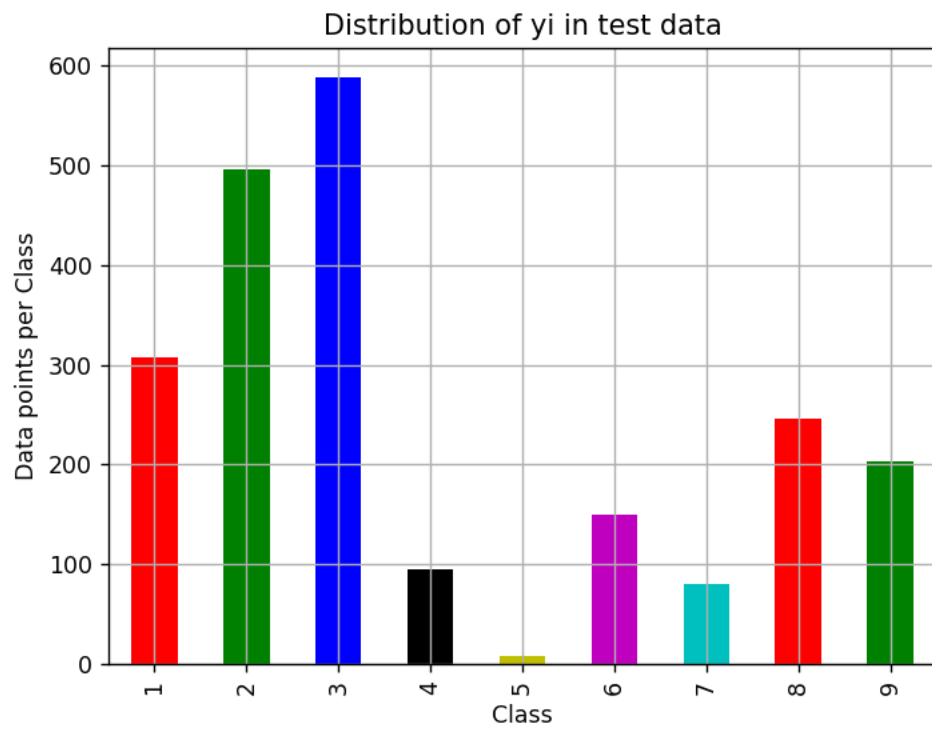
## Test Train Split

Number of data points in train data: 6955
Number of data points in test data: 2174
Number of data points in cross validation data: 1739

### Check for distribution of data

We check for the distribution of classes in each split by plotting a histogram.

Distribution of yi in train data



Distribution of yi in test data

Distribution of yi in cross validation data

## Machine Learning Model

### Random Model

Log loss on Cross Validation Data using Random Model 2.46
Log loss on Test Data using Random Model 2.48
Accuracy 11.49

### Confusion Matrix

**Precision Matrix**



**Recall Matrix**

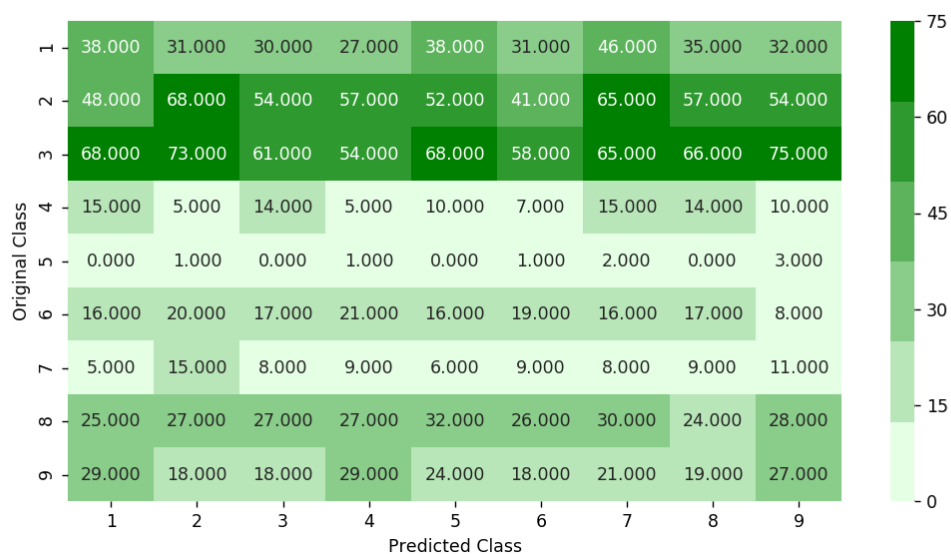|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.123 | 0.101 | 0.097 | 0.088 | 0.123 | 0.101 | 0.149 | 0.114 | 0.104 |
| 2 | 0.097 | 0.137 | 0.109 | 0.115 | 0.105 | 0.083 | 0.131 | 0.115 | 0.109 |
| 3 | 0.116 | 0.124 | 0.104 | 0.092 | 0.116 | 0.099 | 0.111 | 0.112 | 0.128 |
| 4 | 0.158 | 0.053 | 0.147 | 0.053 | 0.105 | 0.074 | 0.158 | 0.147 | 0.105 |
| 5 | 0.000 | 0.125 | 0.000 | 0.125 | 0.000 | 0.125 | 0.250 | 0.000 | 0.375 |
| 6 | 0.107 | 0.133 | 0.113 | 0.140 | 0.107 | 0.127 | 0.107 | 0.113 | 0.053 |
| 7 | 0.062 | 0.188 | 0.100 | 0.113 | 0.075 | 0.113 | 0.100 | 0.113 | 0.138 |
| 8 | 0.102 | 0.110 | 0.110 | 0.110 | 0.130 | 0.106 | 0.122 | 0.098 | 0.114 |
| 9 | 0.143 | 0.089 | 0.089 | 0.143 | 0.118 | 0.089 | 0.103 | 0.094 | 0.133 |

Original Class (rows) / Predicted Class (columns)
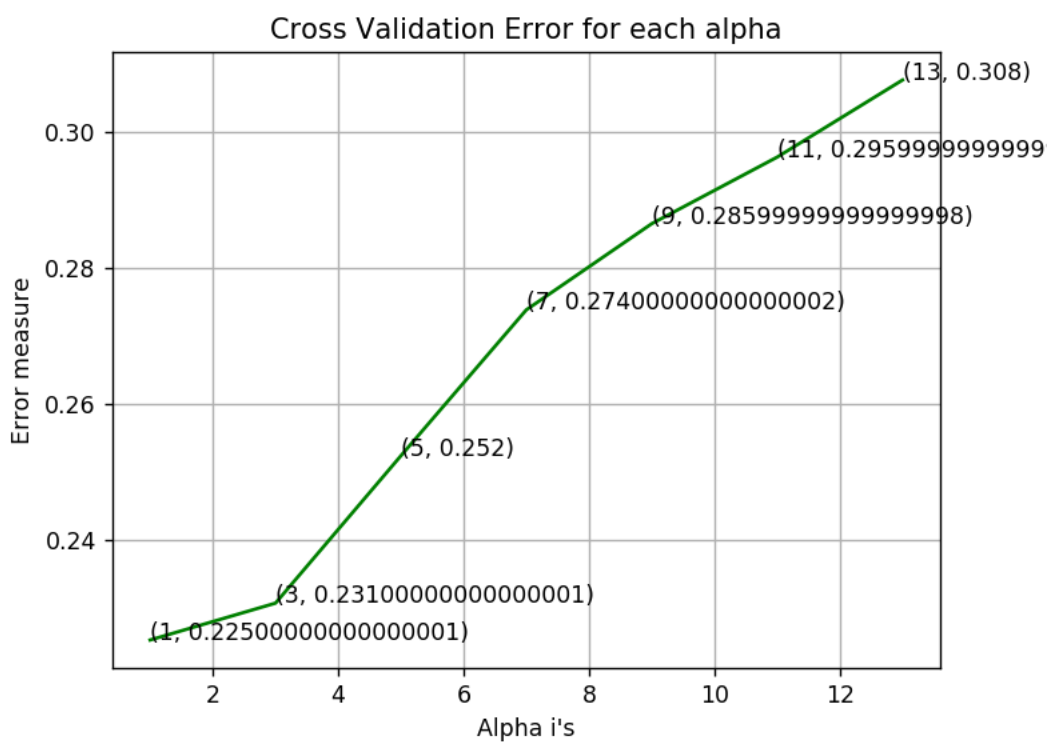
## K Nearest Neighbor Classification

### Hyperparameter Search

log_loss for k =  1 is 0.225386237304
log_loss for k =  3 is 0.230795229168
log_loss for k =  5 is 0.252421408646
log_loss for k =  7 is 0.273827486888
log_loss for k =  9 is 0.286469181555
log_loss for k =  11 is 0.29623391147
log_loss for k =  13 is 0.307551203154

## Cross Validation Error for each alpha



Error measure vs Alpha i's

(13, 0.308)
(11, 0.2959999999999
(9, 0.28599999999999998)
(7, 0.27400000000000002)
(5, 0.252)
(3, 0.23100000000000001)
(1, 0.22500000000000001)

## Results from the Best Model

For values of best alpha =  1 The train log loss is: 0.08
For values of best alpha =  1 The cross validation log loss is: 0.23
For values of best alpha =  1 The test log loss is: 0.24
Accuracy  95.49

## Confusion Matrix



| Original Class \ Predicted Class | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 298.000 | 0.000 | 0.000 | 0.000 | 2.000 | 3.000 | 1.000 | 1.000 | 3.000 |
| 2 | 18.000 | 463.000 | 0.000 | 0.000 | 0.000 | 6.000 | 0.000 | 1.000 | 8.000 |
| 3 | 0.000 | 0.000 | 588.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 4 | 1.000 | 0.000 | 0.000 | 92.000 | 0.000 | 1.000 | 0.000 | 0.000 | 1.000 |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 6.000 | 1.000 | 0.000 | 1.000 | 0.000 |
| 6 | 5.000 | 1.000 | 0.000 | 1.000 | 0.000 | 138.000 | 3.000 | 2.000 | 0.000 |
| 7 | 0.000 | 2.000 | 0.000 | 0.000 | 0.000 | 0.000 | 73.000 | 1.000 | 4.000 |
| 8 | 10.000 | 0.000 | 0.000 | 1.000 | 0.000 | 3.000 | 0.000 | 230.000 | 2.000 |
| 9 | 4.000 | 7.000 | 0.000 | 0.000 | 0.000 | 2.000 | 0.000 | 2.000 | 188.000 |

**Precision Matrix**



**Recall Matrix**



## Logistic Regression
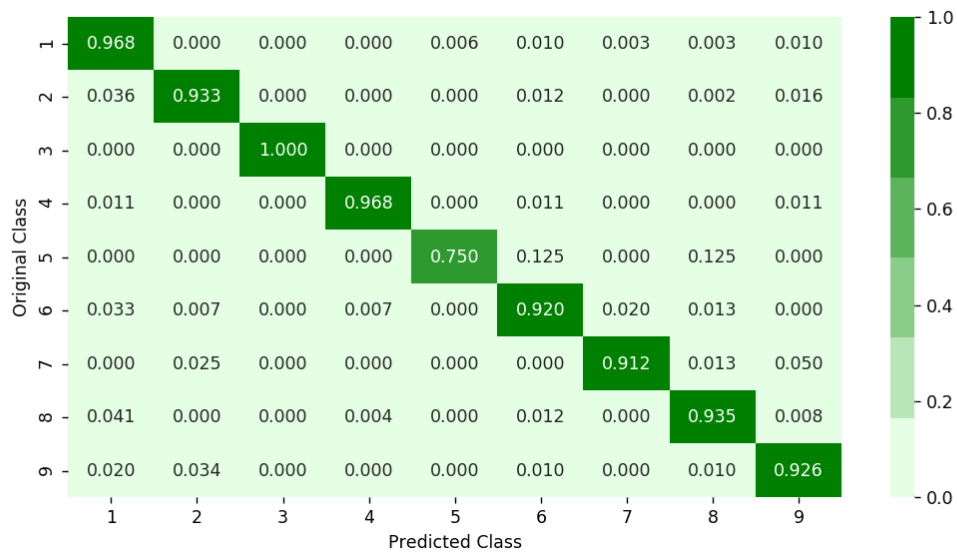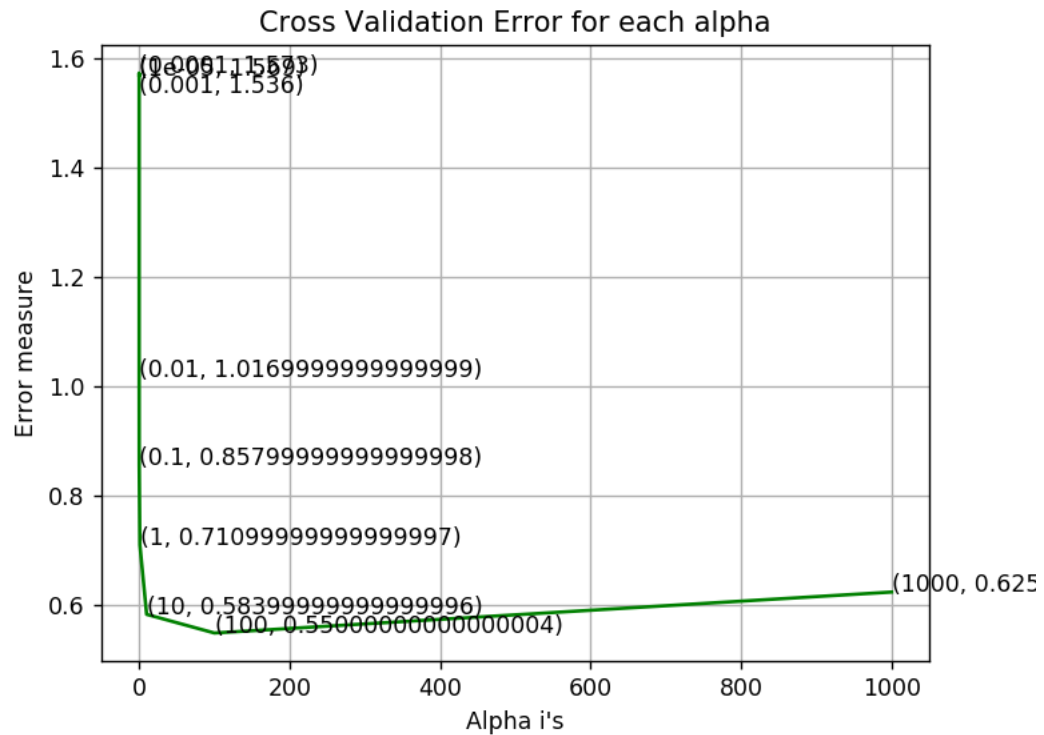
### Hyperparameter Search

log_loss for c =  1e-05 is 1.56916911178
log_loss for c =  0.0001 is 1.57336384417
log_loss for c =  0.001 is 1.53598598273
log_loss for c =  0.01 is 1.01720972418
log_loss for c =  0.1 is 0.857766083873

log_loss for c =  1 is 0.711154393309
log_loss for c =  10 is 0.583929522635
log_loss for c =  100 is 0.549929846589
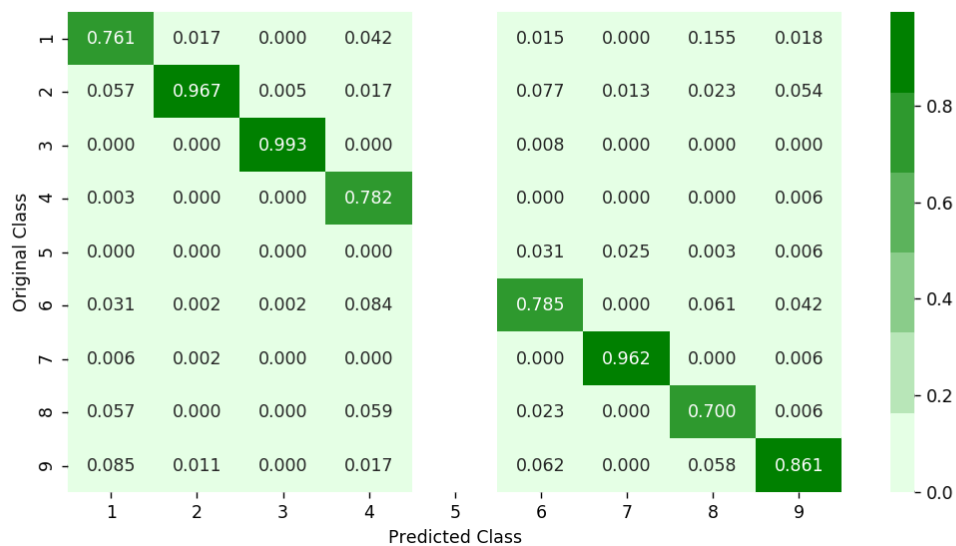log_loss for c =  1000 is 0.624746769121

## Cross Validation Error for each alpha



## Results from the Best Model

log loss for train data 0.50
log loss for cv data 0.55
log loss for test data 0.53
Number of misclassified points  87.67

## Confusion Matrix

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 242.000 | 8.000 | 0.000 | 5.000 | 0.000 | 2.000 | 0.000 | 48.000 | 3.000 |
| 2 | 18.000 | 446.000 | 3.000 | 2.000 | 0.000 | 10.000 | 1.000 | 7.000 | 9.000 |
| 3 | 0.000 | 0.000 | 587.000 | 0.000 | 0.000 | 1.000 | 0.000 | 0.000 | 0.000 |
| 4 | 1.000 | 0.000 | 0.000 | 93.000 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 4.000 | 2.000 | 1.000 | 1.000 |
| 6 | 10.000 | 1.000 | 1.000 | 10.000 | 0.000 | 102.000 | 0.000 | 19.000 | 7.000 |
| 7 | 2.000 | 1.000 | 0.000 | 0.000 | 0.000 | 0.000 | 76.000 | 0.000 | 1.000 |
| 8 | 18.000 | 0.000 | 0.000 | 7.000 | 0.000 | 3.000 | 0.000 | 217.000 | 1.000 |
| 9 | 27.000 | 5.000 | 0.000 | 2.000 | 0.000 | 8.000 | 0.000 | 18.000 | 143.000 |

**Precision Matrix**



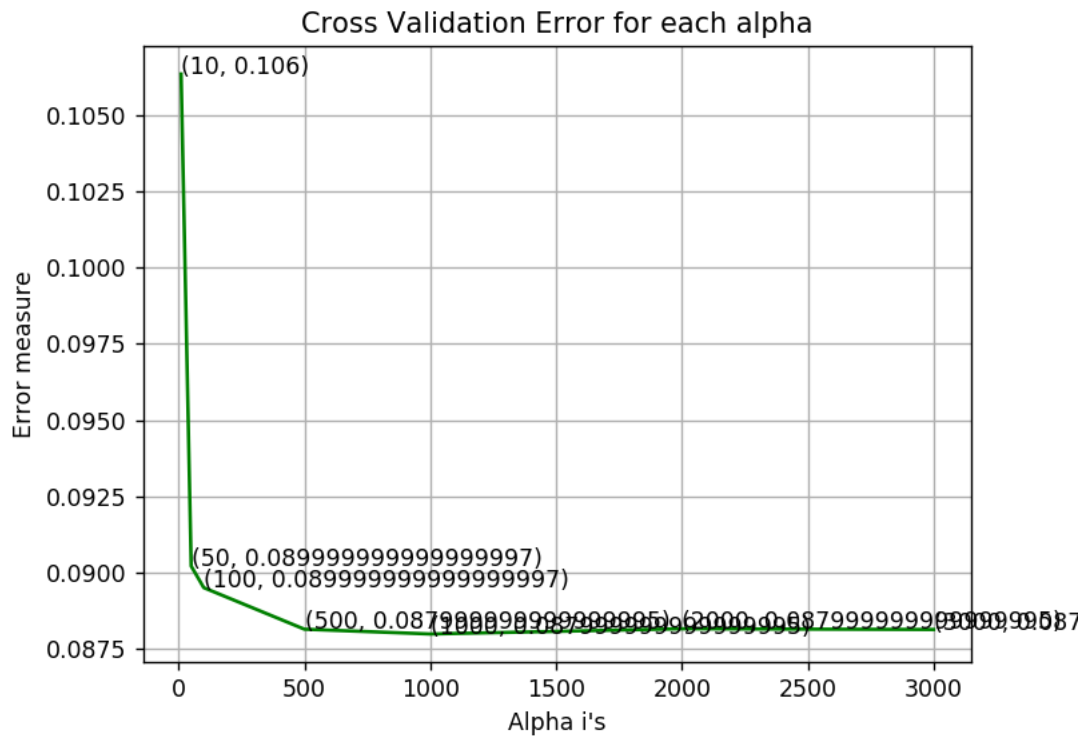| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.761 | 0.017 | 0.000 | 0.042 | 0.015 | 0.000 | 0.155 | 0.018 | |
| 2 | 0.057 | 0.967 | 0.005 | 0.017 | 0.077 | 0.013 | 0.023 | 0.054 | |
| 3 | 0.000 | 0.000 | 0.993 | 0.000 | 0.008 | 0.000 | 0.000 | 0.000 | |
| 4 | 0.003 | 0.000 | 0.000 | 0.782 | 0.000 | 0.000 | 0.000 | 0.006 | |
| 5 | 0.000 | 0.000 | 0.000 | 0.000 | 0.031 | 0.025 | 0.003 | 0.006 | |
| 6 | 0.031 | 0.002 | 0.002 | 0.084 | 0.785 | 0.000 | 0.061 | 0.042 | |
| 7 | 0.006 | 0.002 | 0.000 | 0.000 | 0.000 | 0.962 | 0.000 | 0.006 | |
| 8 | 0.057 | 0.000 | 0.000 | 0.059 | 0.023 | 0.000 | 0.700 | 0.006 | |
| 9 | 0.085 | 0.011 | 0.000 | 0.017 | 0.062 | 0.000 | 0.058 | 0.861 | |

**Recall Matrix**

## Random Forest Classifier

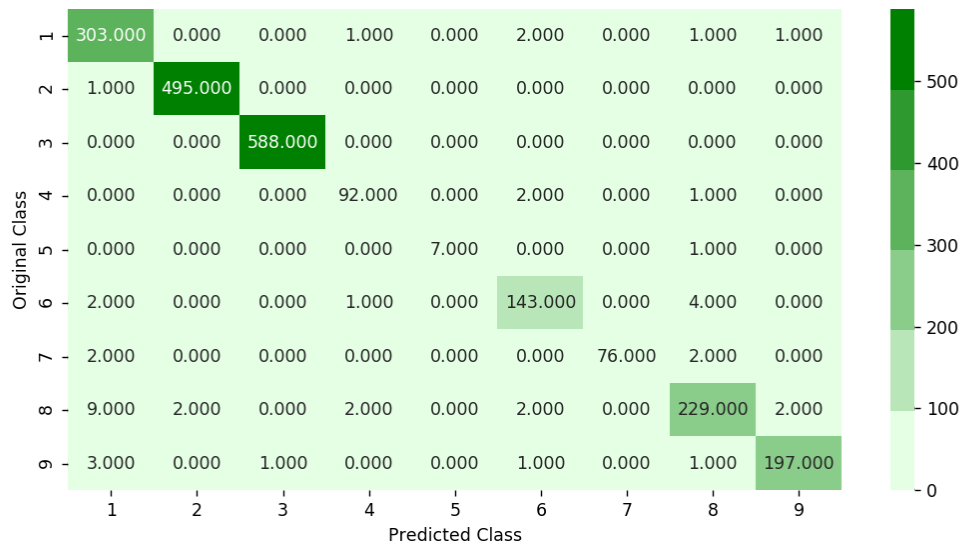### Hyperparameter Search

log_loss for c =  10 is 0.106357709164
log_loss for c =  50 is 0.0902124124145
log_loss for c =  100 is 0.0895043339776
log_loss for c =  500 is 0.0881420869288
log_loss for c =  1000 is 0.0879849524621
log_loss for c =  2000 is 0.0881566647295
log_loss for c =  3000 is 0.0881318948443

## Cross Validation Error for each alpha



(10, 0.106)

(50, 0.089999999999999997)
(100, 0.089999999999999997)
(500, 0.0879999999999999995) (2000, 0.0879999999999999587)
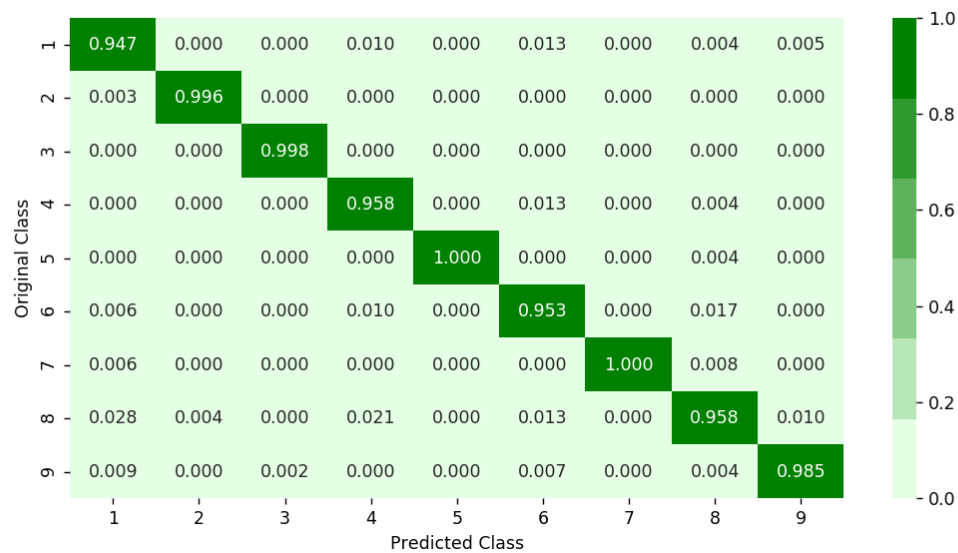(1000, 0.0879999999999999995)

## Results from the Best model

For values of best alpha = 1000 The train log loss is: 0.031
For values of best alpha = 1000 The cross validation log loss is: 0.09
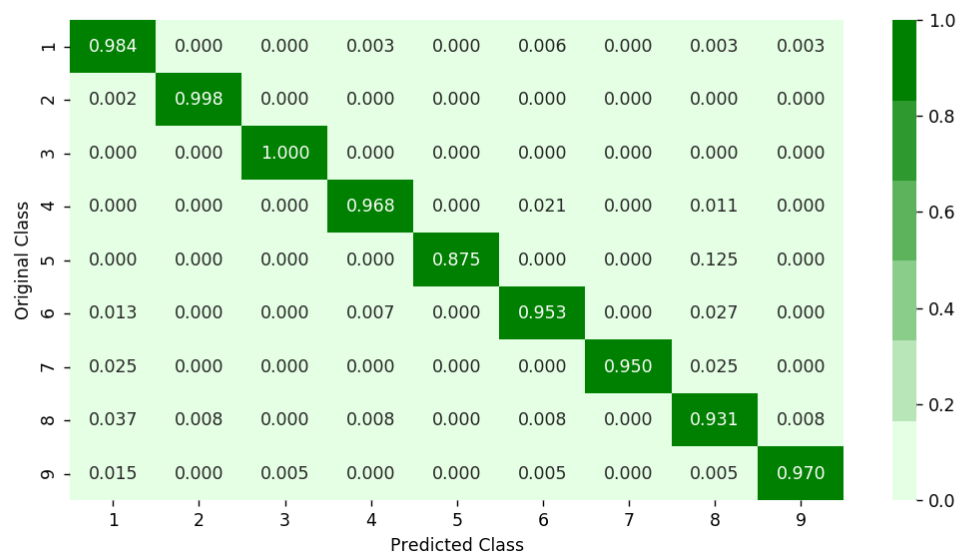For values of best alpha = 1000 The test log loss is: 0.08
Accuracy 96.76

## Confusion Matrix

**Precision Matrix**

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| **1** | 0.947 | 0.000 | 0.000 | 0.010 | 0.000 | 0.013 | 0.000 | 0.004 | 0.005 |
| **2** | 0.003 | 0.996 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| **3** | 0.000 | 0.000 | 0.998 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| **4** | 0.000 | 0.000 | 0.000 | 0.958 | 0.000 | 0.013 | 0.000 | 0.004 | 0.000 |
| **5** | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 | 0.000 | 0.000 | 0.004 | 0.000 |
| **6** | 0.006 | 0.000 | 0.000 | 0.010 | 0.000 | 0.953 | 0.000 | 0.017 | 0.000 |
| **7** | 0.006 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 1.000 | 0.008 | 0.000 |
| **8** | 0.028 | 0.004 | 0.000 | 0.021 | 0.000 | 0.013 | 0.000 | 0.958 | 0.010 |
| **9** | 0.009 | 0.000 | 0.002 | 0.000 | 0.000 | 0.007 | 0.000 | 0.004 | 0.985 |

Original Class (rows) · Predicted Class (columns)

**Recall Matrix**

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| **1** | 0.984 | 0.000 | 0.000 | 0.003 | 0.000 | 0.006 | 0.000 | 0.003 | 0.003 |
| **2** | 0.002 | 0.998 | 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 | 0.968 | 0.000 | 0.021 | 0.000 | 0.011 | 0.000 |
| **5** | 0.000 | 0.000 | 0.000 | 0.000 | 0.875 | 0.000 | 0.000 | 0.125 | 0.000 |
| **6** | 0.013 | 0.000 | 0.000 | 0.007 | 0.000 | 0.953 | 0.000 | 0.027 | 0.000 |
| **7** | 0.025 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.950 | 0.025 | 0.000 |
| **8** | 0.037 | 0.008 | 0.000 | 0.008 | 0.000 | 0.008 | 0.000 | 0.931 | 0.008 |
| **9** | 0.015 | 0.000 | 0.005 | 0.000 | 0.000 | 0.005 | 0.000 | 0.005 | 0.970 |

Original Class (rows) · Predicted Class (columns)

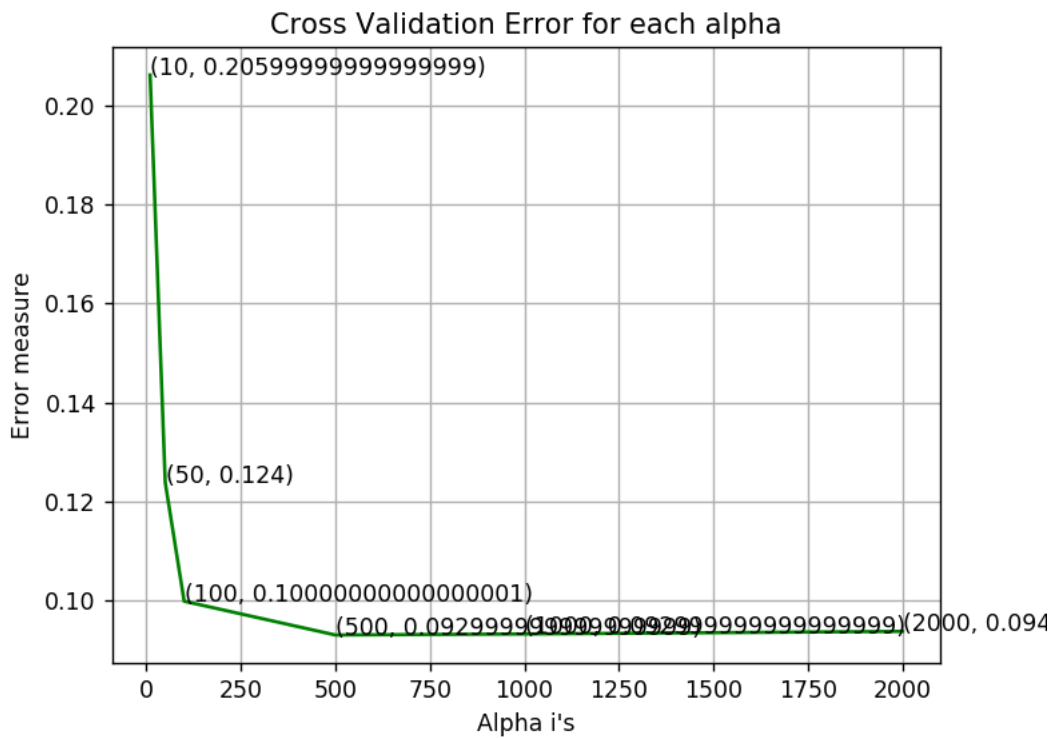# XgBoost Classification

## Hyperparameter Search

log_loss for c =  10 is 0.20615980494
log_loss for c =  50 is 0.123888382365
log_loss for c =  100 is 0.099919437112
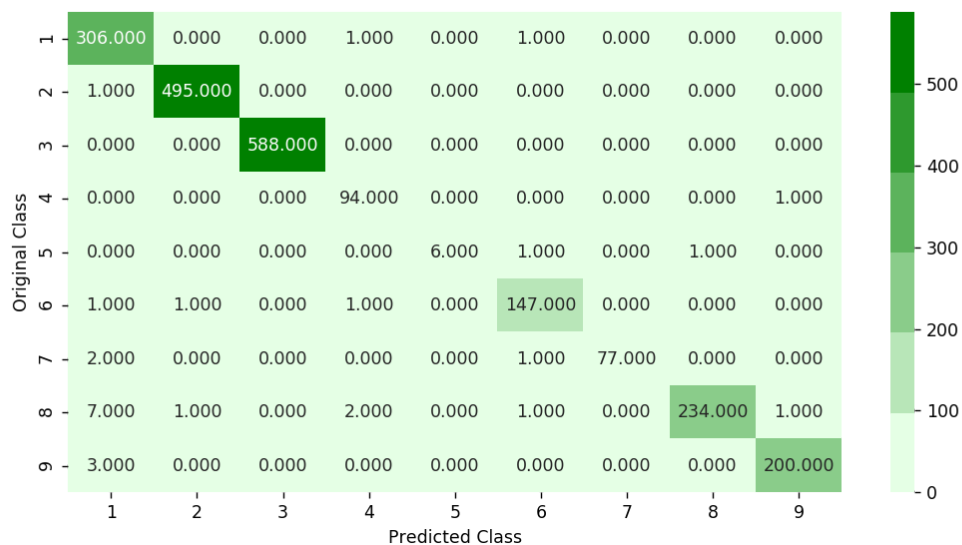log_loss for c =  500 is 0.0931035681289

log_loss for c =  1000 is 0.0933084876012
log_loss for c =  2000 is 0.0938395690309

## Cross Validation Error for each alpha



(10, 0.205999999999999)

(50, 0.124)

(100, 0.10000000000000001)

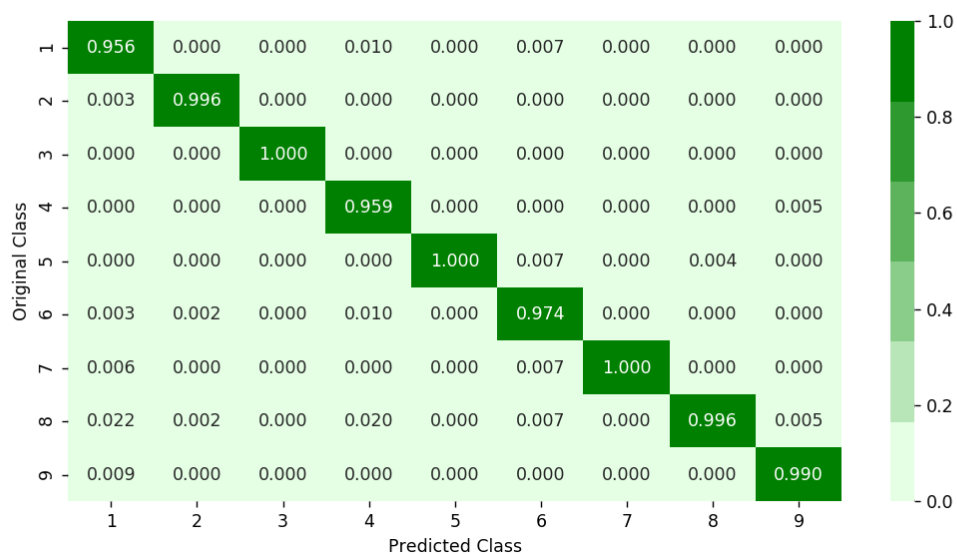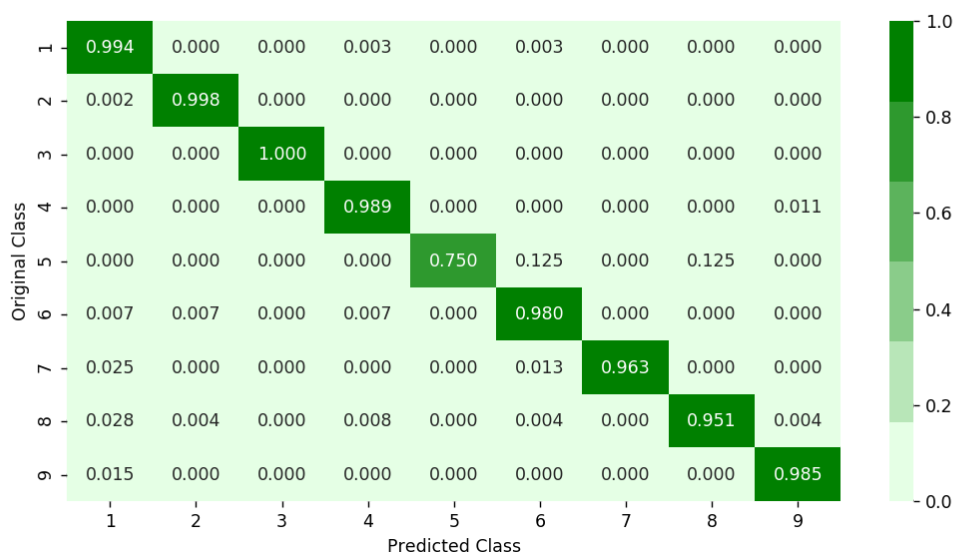(500, 0.0929999... (1000, 0.0932999999999999)(2000, 0.094

### Results from the Best Model

For values of best alpha =  500 The train log loss is: 0.022
For values of best alpha =  500 The cross validation log loss is: 0.09
For values of best alpha =  500 The test log loss is: 0.08
Accuracy 98.67

### Confusion Matrix

**Precision Matrix**



**Recall Matrix**

## XgBoost Classification with best hyper parameters using Random Search

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


[Parallel(n_jobs=-1)]: Done   2 tasks      | elapsed:   26.5s
[Parallel(n_jobs=-1)]: Done   9 tasks      | elapsed:   5.8min
[Parallel(n_jobs=-1)]: Done  19 out of  30 | elapsed:  9.3min remaining:  5.4min
[Parallel(n_jobs=-1)]: Done  23 out of  30 | elapsed: 10.1min remaining:  3.1min
[Parallel(n_jobs=-1)]: Done  27 out of  30 | elapsed: 14.0min remaining:  1.6min
[Parallel(n_jobs=-1)]: Done  30 out of  30 | elapsed: 14.2min finished


RandomizedSearchCV(cv=None, error_score='raise',
          estimator=XGBClassifier(base_score=0.5, colsample_bylevel=1, colsample_bytree=1,
       gamma=0, learning_rate=0.1, max_delta_step=0, max_depth=3,
       min_child_weight=1, missing=None, n_estimators=100, nthread=-1,
       objective='binary:logistic', reg_alpha=0, reg_lambda=1,
       scale_pos_weight=1, seed=0, silent=True, subsample=1),
          fit_params=None, iid=True, n_iter=10, n_jobs=-1,
          param_distributions={'learning_rate': [0.01, 0.03, 0.05, 0.1, 0.15, 0.2], 'n_estimators': [100, 200, 500, 1000, 2000], 'max_dep
th': [3, 5, 10], 'colsample_bytree': [0.1, 0.3, 0.5, 1], 'subsample': [0.1, 0.3, 0.5, 1]},
          pre_dispatch='2*n_jobs', random_state=None, refit=True,
          return_train_score=True, scoring=None, verbose=10)
```

### Best Parameters

{'subsample': 1, 'n_estimators': 500, 'max_depth': 5, 'learning_rate': 0.05, 'colsample_bytree': 0.5}

### Results from the Best Parameter Model

train loss 0.022
cv loss 0.09
test loss 0.08                                                          Accuracy 98.67