


```
In [52]: from sklearn.linear_model import LogisticRegression
logmodel = LogisticRegression()
logmodel.fit(X_train,y_train)
logpred = logmodel.predict(X_test)

print(confusion_matrix(y_test, logpred))

print(round(accuracy_score(y_test, logpred),2)*100)
LOOCV = cross_val_score(logmodel, X_train, y_train, cv=k_score, n_jobs=1, scoring = 'accuracy').mean()

[[6909 164]
 [ 598 329]]
90.0
```

```
In [55]: from sklearn import model_selection
from sklearn.neighbors import KNeighborsClassifier

X_trainK, X_testK, y_trainK, y_testK = train_test_split(bank_final, y, test_size = 0.2)

#neighbors
neighbors = np.arange(0,25)

#Cross-entropy loss that will hold our answer
```

```

cv_scores = []

#Perform 10-fold cross validation on training set for odd values of k:
for k in neighbors:
    k_value = k+1
    knn = KNeighborsClassifier(n_neighbors = k_value, weights='uniform', p=2, metric='euclidean')
    kfold = model_selection.KFold(n_splits=10)
    scores = model_selection.cross_val_score(knn, X_train, y_train, cv=kfold, scoring='accuracy')
    cv_scores.append(scores.mean()*100)
    print('k=10.22 cv: 80.25') % (K_value, scores.mean()*100, scores.std()*100)

optimal_k = neighbors[cv_scores.index(max(cv_scores))]
print("The optimal number of neighbors is %d with 40.15%%" % (optimal_k, cv_scores[optimal_k]))

plt.plot(neighbors, cv_scores)
plt.xlabel('Number of Neighbors K')
plt.ylabel('Train Accuracy')
plt.show()

k=1 87.80 (+/- 0.56)
k=2 89.47 (+/- 0.59)
k=3 89.41 (+/- 0.48)
k=4 89.82 (+/- 0.46)
k=5 89.84 (+/- 0.50)
k=6 89.88 (+/- 0.51)
k=7 89.95 (+/- 0.46)
k=8 89.95 (+/- 0.47)
k=9 90.08 (+/- 0.44)
k=10 90.07 (+/- 0.36)
k=11 90.18 (+/- 0.36)
k=12 90.20 (+/- 0.43)
k=13 90.19 (+/- 0.42)
k=14 90.25 (+/- 0.51)
k=15 90.14 (+/- 0.46)
k=16 90.26 (+/- 0.49)
k=17 90.22 (+/- 0.45)
k=18 90.23 (+/- 0.45)
k=19 90.28 (+/- 0.42)
k=20 90.29 (+/- 0.45)
k=21 90.29 (+/- 0.49)
k=22 90.32 (+/- 0.47)
k=23 90.31 (+/- 0.46)
k=24 90.26 (+/- 0.43)
k=25 90.27 (+/- 0.43)

The optimal number of neighbors is 21 with 90.31

Train Accuracy
90.0
89.5
89.0
88.5
88.0
0 5 10 15 20 25
Number of Neighbors K

```

```

In [56]: from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=21)
knn.fit(X_train, y_train)
knnpred = knn.predict(X_test)

print(confusion_matrix(y_test, knnpred))
print(round(accuracy_score(y_test, knnpred),2)*100)
DTREVCV = (cross_val_score(knn, X_train, y_train, cv=k_fold, n_jobs=1, scoring = 'accuracy').mean())

[[ 6951 122]
 [ 671 236]]
90.0

```

```

In [57]: from sklearn.svm import SVC
svc = SVC(kernel = 'rbf')
svc.fit(X_train, y_train)
svcpred = svc.predict(X_test)
print(confusion_matrix(y_test, svcpred))
print(round(accuracy_score(y_test, svcpred),2)*100)
SVCVCV = (cross_val_score(svc, X_train, y_train, cv=k_fold, n_jobs=1, scoring = 'accuracy').mean())

[[6531 542]
 [ 584 343]]
86.0

```

```

In [58]: from sklearn.tree import DecisionTreeClassifier
dtree = DecisionTreeClassifier(criterion='gini') #criterion = entropy, gini
dtree.fit(X_train, y_train)
dtreepred = dtree.predict(X_test)

print(confusion_matrix(y_test, dtreepred))
print(round(accuracy_score(y_test, dtreepred),2)*100)
DTREBVCV = (cross_val_score(dtree, X_train, y_train, cv=k_fold, n_jobs=1, scoring = 'accuracy').mea

[[6608 465]
 [ 483 444]]

```

```

In [59]: from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier(n_estimators = 200) criterion = entropy,gini
rfc.fit(X_train, y_train)
rfcpred = rfc.predict(X_test)

print(confusion_matrix(y_test, rfcpred ))
print(round(accuracy_score(y_test, rfcpred),2)*100)
RFCCV = (cross_val_score(rfc, X_train, y_train, cv=k_fold, n_jobs=1, scoring = "accuracy").mean())

[[6793 280]
 [ 491 436]]
90.0

In [60]: from sklearn.naive_bayes import GaussianNB
gaussiannb = GaussianNB()
gaussiannb.fit(X_train, y_train)
gaussiannbpred = gaussiannb.predict(X_test)
probe = gaussiannb.predict(X_test)

print(confusion_matrix(y_test, gaussiannbpred ))
print(round(accuracy_score(y_test, gaussiannbpred),2)*100)
GAUSIAN = (cross_val_score(gaussiannb, X_train, y_train, cv=k_fold, n_jobs=1, scoring = "accuracy"

[[6272 801]
 [ 417 510]]
83.0

In [62]: from xgboost import XGBClassifier
xgb = XGBClassifier()
xgb.fit(X_train, y_train)
xgbpred = xgb.predict(X_test)

```

```
print(confusion_matrix(y_test, xgbpred))
print(round(accuracy_score(y_test, xgbpred),2)*100)
XGB = (crows_val_score(estimator = xgb, X = X_train, y = y_train, cv = 10).mean())

[16:07:41] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:11
In XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was chan
or to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

[16:07:47]
[ 462 465]
91.0
[16:07:42] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:11
In XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was chan
or to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

[16:07:43] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:11
In XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was chan
or to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

[16:07:44] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:11
In XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was chan
or to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

[16:07:45] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:11
In XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was chan
or to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

[16:07:46] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.5.1/src/learner.cc:11
In XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was chan
or to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
```

```

In XGBoost 1.3.0, the default evaluation metric used with the objective 'b'
or 't' to 'logloss'. Explicitly set eval_metric if you'd like to restore the
[16:07:48] WARNING: C:/Users/Administrator/workspace/xgboost-win64_released
In XGBoost 1.3.0, the default evaluation metric used with the objective 'b'
or 't' to 'logloss'. Explicitly set eval_metric if you'd like to restore the
[16:07:49] WARNING: C:/Users/Administrator/workspace/xgboost-win64_released
In XGBoost 1.3.0, the default evaluation metric used with the objective 'b'
or 't' to 'logloss'. Explicitly set eval_metric if you'd like to restore the
[16:07:50] WARNING: C:/Users/Administrator/workspace/xgboost-win64_released
In XGBoost 1.3.0, the default evaluation metric used with the objective 'b'
or 't' to 'logloss'. Explicitly set eval_metric if you'd like to restore the
[16:07:51] WARNING: C:/Users/Administrator/workspace/xgboost-win64_released
In XGBoost 1.3.0, the default evaluation metric used with the objective 'b'

```

```
In [64]: from sklearn.ensemble import GradientBoostingClassifier
gbk = GradientBoostingClassifier()
gbk.fit(X_train, y_train)
gbkpred = gbk.predict(X_test)
print(confusion_matrix(Y_test, gbkpred))
print(round(accuracy_score(Y_test, gbkpred), 2)*100)
GBRCV = (cross_val_score(gbk, X_train, y_train, cv=k_fold, n_jobs=1, scoring = 'accuracy')).mean()

[[6826  247]
 [ 460 467]]
91.0

In [65]: models = pd.DataFrame({
    'Models': ['Random Forest Classifier', 'Decision Tree Classifier', 'Support Vector Ma
    'K-Nearest Neighbors', 'Logistic Model', 'Gaussian Naïv', 'K-NN', 'Gradient B
    'Score': [RFVCV, DTREVCV, SVCVC, KNNVC, LOGVC, GAUSVC, KGBVC, GBKVC]})
```

```
models.sort_values(by='score', ascending=False)
```

```
Out[45]:
```

	Models	Score
7	Gradient Boosting	0.914308
6	XGBoost	0.910932
4	Logistic Model	0.909666
0	Random Forest Classifier	0.909636
3	K-Neare Neighbors	0.904634
1	Decision Tree Classifier	0.885320
2	Support Vector Machine	0.855550
5	Gaussian NB	0.844432

Accuracy is measured by the area under the ROC curve. An area of 1 represents a perfect test, an area of .5 represents a worthless

A rough guide for classifying the accuracy of a diagnostic test is the traditional academic point system:

```

> .90-1 = excellent (A)
> .80-.90 = good (B)
> .70-.80 = fair (C)
> .60-.70 = poor (D)
> .50-.60 = fail (F)

```

```

In [68]: # XGBBOOST ROC/ AUC , BEST MODEL
from sklearn import metrics
fig, (ax, ax1) = plt.subplots(nrows = 1, ncols = 2, figsize = (15,5))
proba = xgb.predict_proba(X_test)
preds = proba[:,1]
fprxgb, tprxgb, thresholdxgb = metrics.roc_curve(y_test, preds)
roc_aucxgb = metrics.auc(fprxgb, tprxgb)

ax.plot(fprxgb, tprxgb, 'b', label = 'AUC = %0.2f' % roc_aucxgb)
ax.plot([0, 1], [0, 1], 'r--')
ax.set_title('Receiver Operating Characteristic XGBBOOST ', fontsize=10)
ax1.set_xlabel('True Positive Rate', fontsize=7)

```

```
ax.set_xlabel('False Positive Rate', fontsize=15)
ax.legend(loc = 'lower right', prop={'size': 16})

for cat in categories:
    probs = gbk.predict_proba(X_test)
    preds = probs[:,1]
    fprgbk, tprgbk, thresholdgbk = metrics.roc_curve(y_test, preds)
    roc_aucgbk = metrics.auc(fprgbk, tprgbk)
    ax.plot(fprgbk, tprgbk, 'b', label = 'AUC = %.2f' % roc_aucgbk)
    ax.plot([0, 1], [0, 1], 'r--')
    ax.set_title('Receiver Operating Characteristic GRADIENT BOOST', fontsize=10)
    ax.set_xlabel('True Positive Rate', fontsize=10)
    ax.set_ylabel('False Positive Rate', fontsize=15)
```

```
ax1.legend(loc = 'lower right', prop={'size': 16})
plt.subplots_adjust(wspace=1)
```

```
In [69]: #fig, (ax1, ax2, ax3, ax4, ax5) = plt.subplots(nrows = 2, ncols = 3, figsize = (15, 4))
fig, ax_arr = plt.subplots(nrows = 2, ncols = 3, figsize = (20,15))

#LOGMODEL
probs = logmodel.predict_proba(X_test)
preds = probs[:,1]
fprlog, tprlog, thresholdlog = metrics.roc_curve(y_test, preds)
roc_auclog = metrics.auc(fprlog, tprlog)

ax_arr[0,0].plot(fprlog, tprlog, 'b', label = 'AUC = %.2F' % roc_auclog)
ax_arr[0,0].plot([0,1], [0,1], 'g--')
ax_arr[0,0].set_title('Receiver Characteristic Logistic', fontsize=20)
ax_arr[0,0].set_ylabel('True Positive Rate', fontsize=20)
ax_arr[0,0].set_xlabel('False Positive Rate', fontsize=15)
ax_arr[0,0].legend(loc = 'lower right', propset='size': 16)

#RANDOM FOREST
probs = rfc.predict_proba(X_test)
preds = probs[:,1]
fprRFC, tprRFC, thresholdRFC = metrics.roc_curve(y_test, preds)
```

```

roc_curve = metrics.roc_curve('y_test',
                               ax_arr[0,1].plot(tprrfc, tprffc, 'b', label = 'AUC = 0.21*' & roc_aucrfc)
ax_arr[0,1].plot([0,1], [0,1], 'r--')
ax_arr[0,1].set_title('Receiver Operating Characteristic Random Forest', fontsize=20)
ax_arr[0,1].set_ylabel('True Positive Rate', fontsize=20)
ax_arr[0,1].set_xlabel('False Positive Rate', fontsize=15)
ax_arr[0,1].legend(loc = 'lower right', prop={'size': 16})

#NOW
probs = knn.predict_proba(X_test)
preds = probs[:,1]
fpknn, tpknn, thresholdknn = metrics.roc_curve(y_test, preds)
roc_aucknn = metrics.auc(fpknn, tpknn)

ax_arr[0,2].plot(tpknn, tpknn, 'b', label = 'AUC = 0.21*' & roc_aucknn)
ax_arr[0,2].plot([0,1], [0,1], 'r--')
ax_arr[0,2].set_title('Receiver Operating Characteristic KNN', fontsize=20)
ax_arr[0,2].set_ylabel('True Positive Rate', fontsize=20)
ax_arr[0,2].set_xlabel('False Positive Rate', fontsize=15)
ax_arr[0,2].legend(loc = 'lower right', prop={'size': 16})

#DECISION TREE
probs = dtree.predict_proba(X_test)

```

```

probs = probs[:,1]
fpdtrree, tpdtrree, thresholdauc = metrics roc_curve(y_test, probs)
roc_aucdtr = metrics auc(fpdtrree, tpdtrree)

ax_arr[1,0].plot(fpdtrree, tpdtrree, 'b', label = 'AUC = %.2E' % roc_aucdtr)
ax_arr[1,1].plot([0, 1], [0, 1], 'r--')
ax_arr[1,0].set_title('Receiver Operating Characteristic Decision Tree ', fontsize=20)
ax_arr[1,0].set_ylabel('True Positive Rate', fontsize=20)
ax_arr[1,0].set_xlabel('False Positive Rate', fontsize=15)
ax_arr[1,0].legend(loc = 'lower right', prop= {'size': 16})

#GAUSSIAN
probs = gaussiannn_predict_proba(X_test)
probs = probs[:,1]
fpgaug, tpgaug, thresholdgaug = metrics roc_curve(y_test, probs)
roc_aucgaug = metrics auc(fpgaug, tpgaug)

ax_arr[1,1].plot(fpgaug, tpgaug, 'b', label = 'AUC = %.2F' % roc_aucgaug)
ax_arr[1,1].plot([0, 1], [0, 1], 'r--')
ax_arr[1,1].set_title('Receiver Operating Characteristic Gaussian ', fontsize=20)
ax_arr[1,1].set_ylabel('True Positive Rate', fontsize=20)
ax_arr[1,1].set_xlabel('False Positive Rate', fontsize=15)
ax_arr[1,1].legend(loc = 'lower right', prop= {'size': 16})

```

```

# ALL PLOTS -----
ax_arr[1,2].plot(fprgau, tprgau, 'b', label = 'Gaussian', color='black')
ax_arr[1,2].plot(fprdtree, tprdtree, 'b', label = 'Decision Tree', color='blue')
ax_arr[1,2].plot(fprsvm, tprsvm, 'b', label = 'Svm', color='brown')
ax_arr[1,2].plot(fprfc, tprfc, 'b', label = 'Random Forest', color='green')
ax_arr[1,2].plot(fprlog, tprlog, 'b', label = 'Logistic', color='grey')
ax_arr[1,2].set_title('Receiver Operating Characteristic Comparison', fontsize=20)
ax_arr[1,2].set_ylabel('True Positive Rate', fontsize=20)
ax_arr[1,2].set_xlabel('False Positive Rate', fontsize=15)
ax_arr[1,2].legend(loc = 'lower right', prop = 'size': 16)

plt.subplots_adjust(wspace=0.2)
plt.tight_layout()

```

The image displays three Receiver Operating Characteristic (ROC) curves, each plotting True Positive Rate (TPR) on the y-axis against False Positive Rate (FPR) on the x-axis. The y-axis ranges from 0.0 to 1.0, and the x-axis ranges from 0.0 to 1.0. A red dashed diagonal line represents a random classifier. The left plot shows a blue curve with an AUC of 0.71. The middle plot shows a blue curve with an AUC of 0.82. The right plot shows a blue curve with an AUC of 0.82, with a legend indicating 'Good' (blue), 'Killed' (orange), 'Bad' (green), and 'Lost' (purple) categories.

```
In [70]: from sklearn.metrics import classification_report

In [71]: print('KNN Confusion Matrix\n', confusion_matrix(y_test, knnpred))

KNN Confusion Matrix
[[6951 122]
 [ 671 256]]

In [72]: print('KNN Reports\n', classification_report(y_test, knnpred))

KNN Reports
              precision    recall  f1-score   support

     0               0.91       0.98       0.95       7073
     1               0.68       0.28       0.39        927

 accuracy          0.78       0.63       0.70      8000
 avg prec          0.78       0.63       0.70      8000
```

weighted avg 0.98 0.90 0.88 0.90

Recall

Recall - Specificity $TN / (TN + FP)$ [MATRIX LINE 1]

For all NEGATIVE(0) REAL VALUES how much we predict correct ?

other way to understand, our real test set has 7163+116 = 7279 clients that didn't subscribe(0), and our model predict 98% correct and 116 incorrect

```
In [73]: print(round((7163 / (7163 + 116)),2))
```

0.98

$TP / (TP + FN)$ [MATRIX LINE 2]

For all POSITIVE(1) REAL VALUES how much we predict correct ?

other way to understand, our real test set has 706 + 253 = 959 clients that subscribe(1), and our model predict 26% correct or 253 and 706 incorrect. BUT REMEMBER, Its best we miss by give negative instead of False Positive

```
In [74]: print(round((253 / (253 + 706)),2))
         print(round(metrics.recall_score(y_test, knnpred),2))

0.26
0.28
```

Precision

$TN / (TN + FN)$ (MATRIX COLUMN 1)

For all NEGATIVE(0) PREDICTIONS by our model, how much we predict correct ?

other way to understand, our model pointed 7163 + 706 = 7869 clients that didn't subscribe(0), and our model predict 91% correct or 7163 correct and 706 incorrect

```
In [75]: print(round((7163 / (7163 + 706)),2))
```

0.91

`TN / (TN + FN)` | MATRIX COLUMN 1 |

For all POSITIVE(1) PREDICTIONS by our model, how much we predict correct ?

other way to understand, our model pointed 116 + 253 = 369 clients that subscribe(1), and our model predict 69% correct or 253 and 116 incorrect

```
In [76]: print(round(253 / (253 + 116),2))
print(round(metrics.precision_score(y_test, knnpred),2))
```

0.69
0.68

F1 Score

F1-SCORE F1-Score is a "median" of Recall and Precision, consider this when you want a balance between this metrics $F1 = 2(Precision(0) * Recall(0)) / (Precision(0) + Recall(0))$

```
In [77]: F1_0 = 2*0.91*0.98/(0.91+0.98)
          round(F1_0,2)

Out[77]: 0.94

In [78]: F1_1 = 2*0.69*0.26/(0.69+0.26)
          round(F1_1,2)

Out[78]: 0.38
```

Average Precision

```
In [79]: AVG_precision = (0.91*(7279/8238)) + (0.69*(959/8238))
          round(AVG_precision,2)

Out[79]: 0.88
```