# 1. Confusion Matrix

A confusion matrix is an N X N matrix, where N is the number of classes being predicted. For the problem in hand, we have N=2, and hence we get a 2 X 2 matrix. Here are a few definitions, you need to remember for a confusion matrix :

- **Accuracy** : the proportion of the total number of predictions that were correct.
- **Positive Predictive Value or Precision** : the proportion of positive cases that were correctly identified.
- **Negative Predictive Value** : the proportion of negative cases that were correctly identified.
- **Sensitivity or Recall** : the proportion of actual positive cases which are correctly identified.
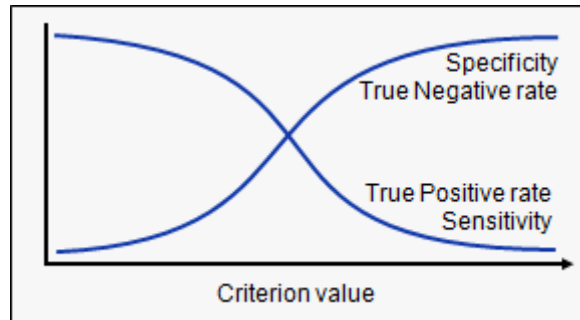- **Specificity** : the proportion of actual negative cases which are correctly identified.

| Confusion Matrix | | Target | | | |
|---|---|---|---|---|---|
| | | Positive | Negative | | |
| **Model** | Positive | a | b | Positive Predictive Value | a/(a+b) |
| | Negative | c | d | Negative Predictive Value | d/(c+d) |
| | | Sensitivity | Specificity | Accuracy = (a+d)/(a+b+c+d) | |
| | | a/(a+c) | d/(b+d) | | |

| Count of ID | Target | | | |
|---|---|---|---|---|
| Model | 1 | 0 | Grand Total | |
| 1 | 3,834 | 639 | 4,473 | 85.7% |
| 0 | 16 | 951 | 967 | 1.7% |
| Grand Total | 3,850 | 1,590 | 5,440 | |
| | 99.6% | 40.19% | | 88.0% |

In general we are concerned with one of the above defined metric. For instance, in a pharmaceutical company, they will be more concerned with minimal wrong positive diagnosis. Hence, they will be more concerned about high Specificity. On the other hand an attrition model will be more concerned with Senstivity.Confusion matrix are generally used only with class output models.
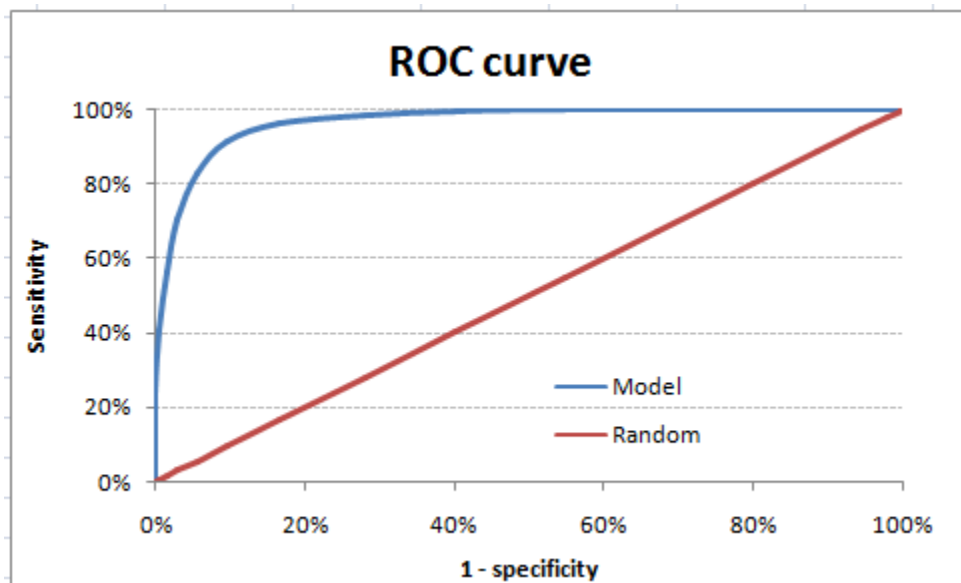
## 2 Area Under the ROC curve (AUC – ROC)

The biggest advantage of using ROC curve is that it is independent of the change in proportion of responders.



The ROC curve is the plot between sensitivity and (1- specificity).

(1- specificity) is also known as false positive rate and sensitivity is also known as True Positive rate. Following is the ROC curve for the case in hand.
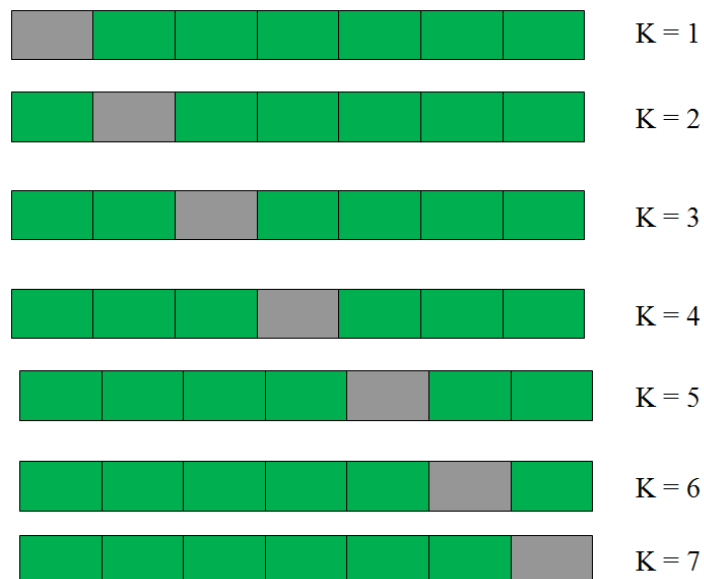
| Count of ID | Target | | | |
|---|---|---|---|---|
| Model | 1 | 0 | Grand Total | |
| 1 | 3,834 | 639 | 4,473 | 85.7% |
| 0 | 16 | 951 | 967 | 1.7% |
| Grand Total | 3,850 | 1,590 | 5,440 | |
| | 99.6% | 40.19% | | 88.0% |

Following are a few thumb rules:

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

## k-fold Cross validation



we divide the entire population into 7 equal samples. Now we train models on 6 samples (Green boxes) and validate on 1 sample (grey box). Then, at the second iteration we train the model with a different sample held as validation. In 7 iterations, we have basically built model on each sample and held each of them as validation. This is a way to reduce the selection bias and reduce the variance in prediction power. Once we have all the 7 models, we take average of the error terms to find which of the models is best.