

## Evaluating the Predictive Performance for Two Classes

- Important R functions are:
  - Display a cross-tabulation of observed (rows) against predicted (column) class memberships: `gmodels::CrossTable( )`
  - Provide detailed statistics: `caret::confusionMatrix( )`. See online help `help(confusionMatrix)`.
  - Receiver operating characteristic (ROC) curve: `pROC::roc( )`
- For the sake of terminology let us call one outcome as “**positive**”, which is usually an outcome of interest leading to action.

Usually **rare classes** (loan default, insurance fraud, disease outcome in a screening test, spam text messages etc.) are usually labelled “positive”

Positives are usually coded as factor level 1 whereas the negatives are set to 0.
- This convention makes sense under specific test scenarios, but can be arbitrary if both classes are “value-free”.
- For rare positives an intuitive negative prediction will lead to a small error rate equal to the frequency of the rare positives. This provides the motivation for **conditional error rates**.

- For just two classes we get “**confusion matrix**”, which has the observed true classes in its rows and the predicted class in its columns:

	Predicted Negative	Predicted Positive
Observed Negative	<i>True Negative (TN)</i>	<i>False Positive (FP)</i>
Observed Positive	<i>False Negative (FN)</i>	<i>True Positive (TP)</i>

- This allows to calculate several **key statistics**:

Name	Definition	Synonyms
<b>False pos. rate</b>	$FP / (TN + FP)$	Type <i>I</i> error, 1 – specificity (row perspective)
<b>True neg. rate</b>	$TN / (TN + FP)$	<b>specificity</b> (row perspective)
<b>True pos. rate</b>	$TP / (FN + TP)$	<b>sensitivity</b> , 1 – Type <i>II</i> error, power, <b>recall</b> (row perspective)
<b>Pos. pred. value</b>	$TP / (FP + TP)$	<b>Precision</b> , 1 – false discovery proportion (column perspective)
<b>Total Accuracy</b>	$\frac{TP + TN}{TN + FP + FN + TP}$	
<b>Total Error Rate</b>	$\frac{FP + FN}{TN + FP + FN + TP}$	1 – accuracy

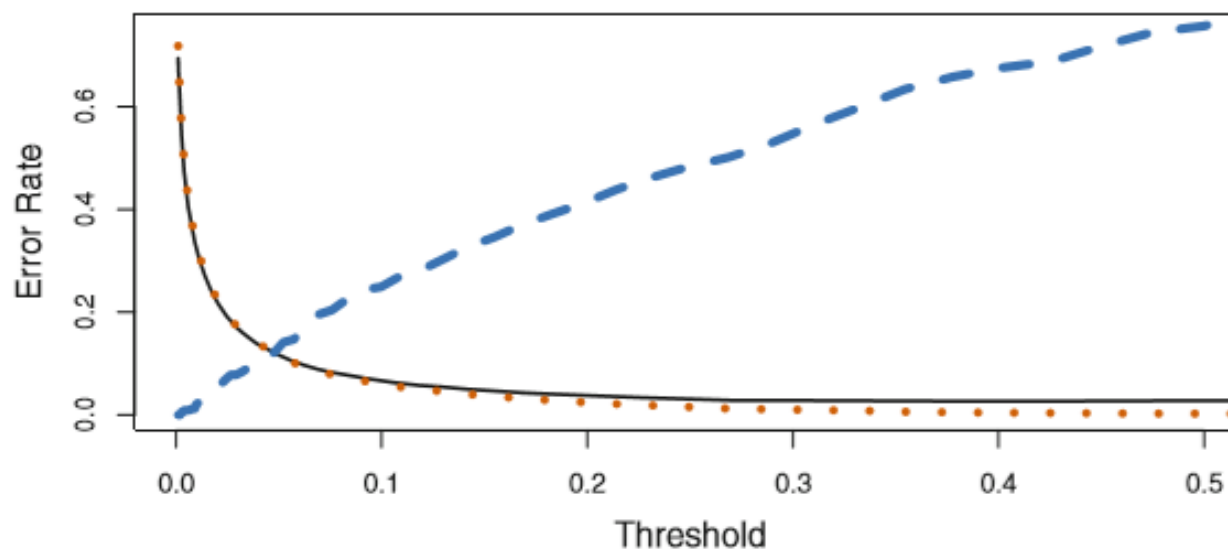
## The Radio Operating Characteristics Curve

- Error rates are affected by the threshold probability  $\delta \in [0,1]$  of assigning an object  $y_i$  either to the positive  $\hat{y}_i = 1$  or negative  $\hat{y}_i = 0$  class in relation to the observed features  $\mathbf{x}_i$ :

$$\hat{y}_i = 0 \text{ if } \Pr(y_i = 1|\mathbf{x}_i) \leq \delta$$

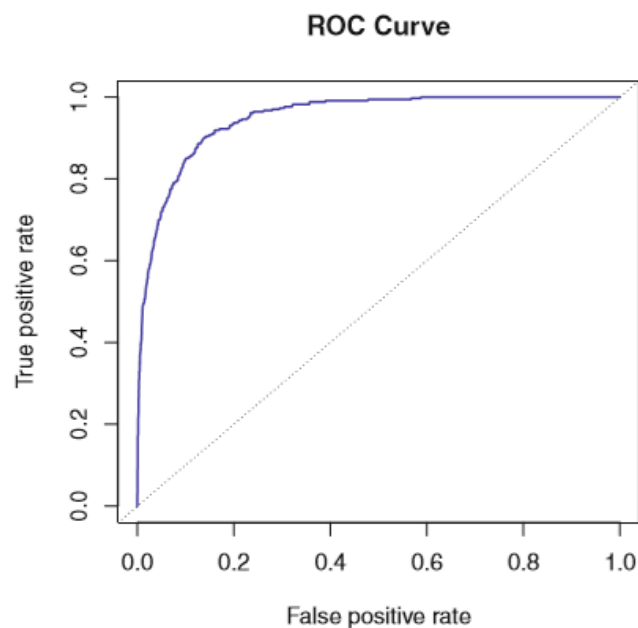
$$\hat{y}_i = 1 \text{ if } \Pr(y_i = 1|\mathbf{x}_i) > \delta$$

- Depending on the threshold  $\delta$  the total error rate, specificity and sensitivity change:
  - For  $\delta = 0$  the false positive rate ( $1 - \text{specificity}$ ) is zero percent and the true positive rate (sensitivity) is zero percent.
  - In contrast, for  $\delta = 1$  the false positive rate is 100 % and the true positive rate is 100%.



**FIGURE 4.7.** For the **Default** data set, error rates are shown as a function of the threshold value for the posterior probability that is used to perform the assignment. The black solid line displays the overall error rate. The blue dashed line represents the fraction of defaulting customers that are incorrectly classified, and the orange dotted line indicates the fraction of errors among the non-defaulting customers.

- The ROC (radio operating characteristic) plots the false positive rate ( $1 - \text{specificity}$ ) against the sensitivity.



**FIGURE 4.8.** A ROC curve for the LDA classifier on the **Default** data. It traces out two types of error as we vary the threshold value for the posterior probability of default. The actual thresholds are not shown. The true positive rate is the sensitivity: the fraction of defaulters that are correctly identified, using a given threshold value. The false positive rate is 1-specificity: the fraction of non-defaulters that we classify incorrectly as defaulters, using that same threshold value. The ideal ROC curve hugs the top left corner, indicating a high true positive rate and a low false positive rate. The dotted line represents the “no information” classifier; this is what we would expect if student status and credit card balance are not associated with probability of default.

- For an uninformative predictor for which sensitivity = 1 – specificity for any threshold  $\delta \in [0,1]$ .

- For a well-discriminating classifier, the ROC curve takes a rectangular shape:  
The upper left corner denotes that predictor for which all positives are properly classified and no negative is falsely classified.
- The area underneath the ROC curve (AUC) measures the discriminating power of a classifier:
  - A. 0.9 to 1.0: **outstanding**
  - B. 0.8 to 0.9: **good**
  - C. 0.7 to 0.8: **fair**
  - D. 0.6 to 0.7: **poor**
  - E. 0.5 to 0.6: **no discrimination**
- Two competing classifiers may have an identical AUC but different shapes of the ROC curve.

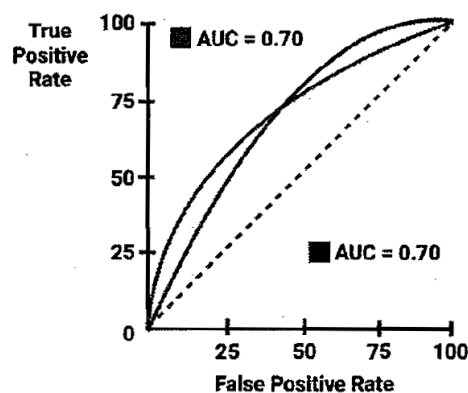


Figure 10.5: ROC curves may have different performance despite having the same AUC

## Appendix: Naïve Bayesian Classifier

- The naïve Bayesian classifier works for nominal scaled features.
- For each feature the conditional probabilities at level given the class membership can be evaluated empirically using a crosstabulation:

	$X = 0$	$X = 1$
$Y = 0$	$\Pr(X = 0 Y = 0)$	$\Pr(X = 1 Y = 0)$
$Y = 1$	$\Pr(X = 0 Y = 1)$	$\Pr(X = 1 Y = 1)$

- Given an observed feature value of  $X$ , the Bayesian probability of a predicted class membership  $Y$  can be calculated:

$$\Pr(\hat{Y} = 0|X = 0) = \frac{\Pr(X = 0|Y = 0) \cdot \Pr(Y = 0)}{\Pr(X = 0)}$$

$$\Pr(\hat{Y} = 0|X = 1) = \frac{\Pr(X = 1|Y = 0) \cdot \Pr(Y = 0)}{\Pr(X = 1)}$$

$$\Pr(\hat{Y} = 1|X = 0) = \frac{\Pr(X = 0|Y = 1) \cdot \Pr(Y = 1)}{\Pr(X = 0)}$$

$$\Pr(\hat{Y} = 1|X = 1) = \frac{\Pr(X = 1|Y = 1) \cdot \Pr(Y = 1)}{\Pr(X = 1)}$$

- Note that  $\Pr(\hat{Y} = 0|X = 0) + \Pr(\hat{Y} = 1|X = 0) = 1$  and  $\Pr(\hat{Y} = 0|X = 1) + \Pr(\hat{Y} = 1|X = 1) = 1$
- For more than one feature  $\mathbf{X} = [X_1, X_2, \dots, X_J]$  and under the assumption that these features are mutually **statistically independent** the product rule of probability calculus can be used to calculate the **joint Bayesian probabilities**:

$$\Pr(\hat{Y}|\mathbf{X}) = \Pr(\hat{Y}|X_1) \cdot \Pr(\hat{Y}|X_2) \cdots \Pr(\hat{Y}|X_J) \cdot$$

- The independence assumption leads to labeling this approach **naïve**.
- A **problem** emerges if any of the conditional probabilities in the crosstabulation is zero due to no observations in the training dataset is observed for this particular combination. This will lead to a joint Bayesian probability  $\Pr(\hat{Y}|\mathbf{X}) = 0$ , because one factor in the product is zero.
- **Solution:** Laplace Estimators. In the crosstabulation  $\gamma$  observations can be added to each cell, which avoids zero conditional probabilities.