

# 04 - Classification

Logistic Regression, Discriminant Analysis, and Naive Bayes

SYS 6018 | Fall 2020

04-classification.pdf

## Contents

<b>1</b>	<b>Classification Intro</b>	<b>2</b>
1.1	Credit Card Default data (Default) . . . . .	2
<b>2</b>	<b>Classification and Pattern Recognition</b>	<b>5</b>
2.1	Binary Classification . . . . .	5
<b>3</b>	<b>Logistic Regression</b>	<b>9</b>
3.1	Basics . . . . .	9
3.2	Estimation . . . . .	9
3.3	Logistic Regression in Action . . . . .	11
3.4	Logistic Regression Summary . . . . .	12
<b>4</b>	<b>Evaluating Classification Models</b>	<b>13</b>
4.1	Evaluation of Binary Classification Models . . . . .	13
4.2	Common Binary Loss Functions . . . . .	14
4.3	Evaluating Binary Classification Models . . . . .	15
4.4	Performance Metrics . . . . .	16
4.5	Performance over a range of thresholds . . . . .	18
4.6	Summary of Classification Evaluation . . . . .	22
<b>5</b>	<b>Generative Classification Models</b>	<b>23</b>
5.1	Linear/Quadratic Discriminant Analysis (LDA/QDA) . . . . .	23
5.2	Naive Bayes . . . . .	23

---

Some of the figures in this presentation are taken from "An Introduction to Statistical Learning, with applications in R" (Springer, 2013) with permission from the authors: G. James, D. Witten, T. Hastie and R. Tibshirani.

# 1 Classification Intro

## 1.1 Credit Card Default data (Default)

The textbook *An Introduction to Statistical Learning (ISL)* has a description of a simulated credit card default dataset. The interest is on predicting whether an individual will default on their credit card payment.

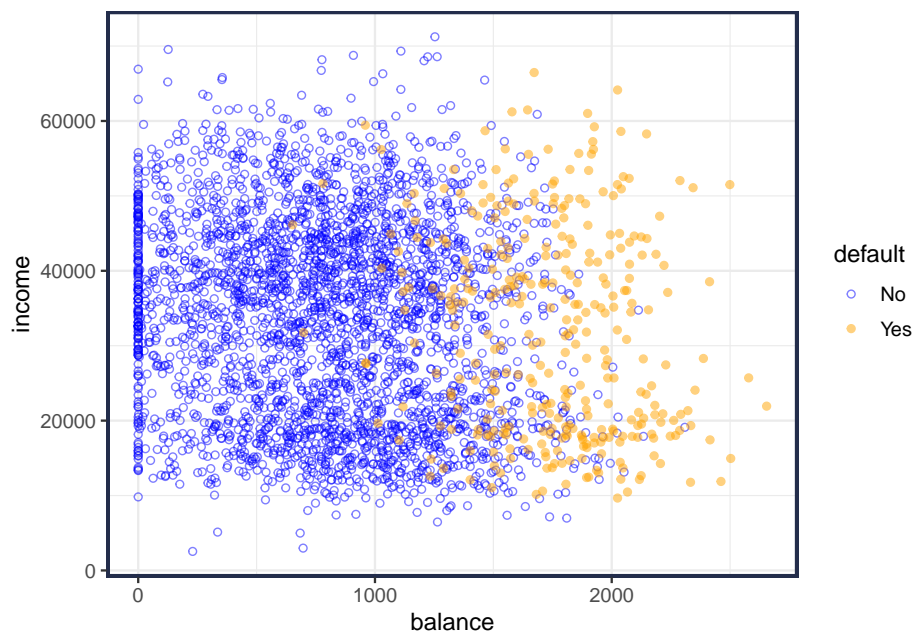
```
1 data(Default, package="ISLR")
```

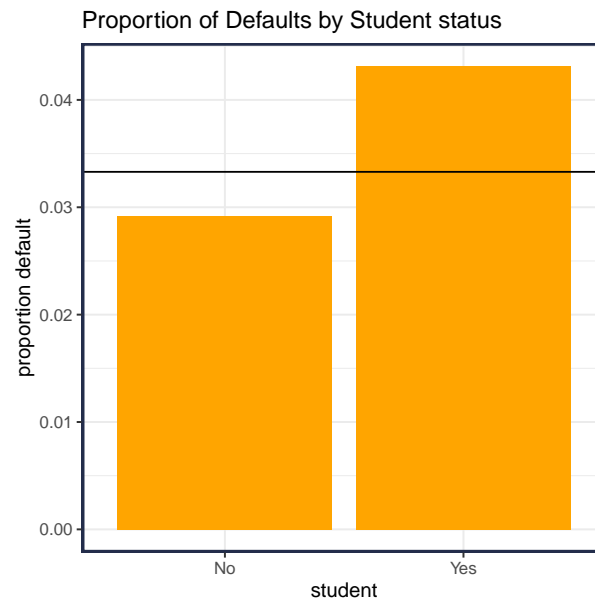
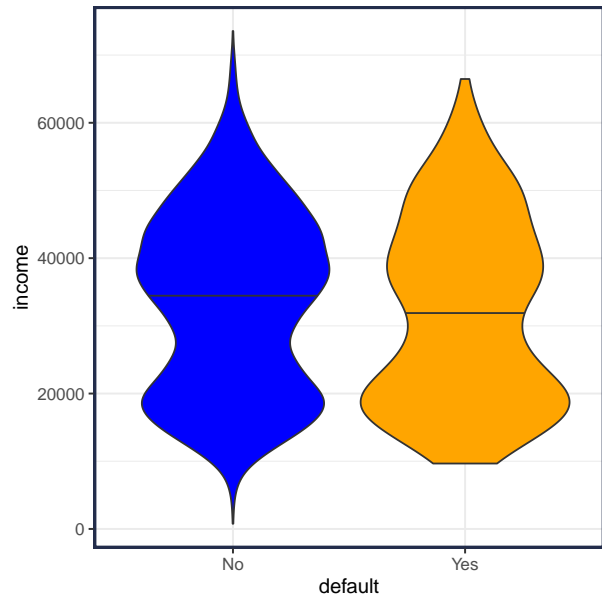
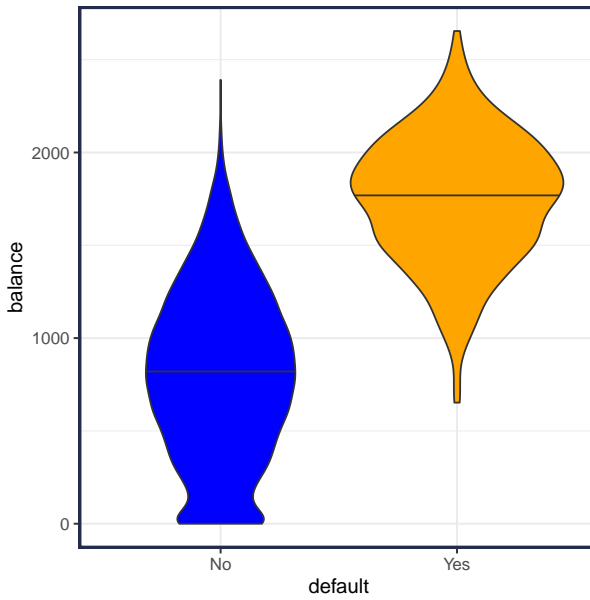
The variables are:

- *response variable* is categorical (factor) Yes and No, (default)
- the categorical (factor) variable (*student*) is either Yes or No
- the average balance a customer has after making their monthly payment (*balance*)
- the customer's income (*income*)

default	student	balance	income
No	No	729.5	44362
No	Yes	817.2	12106
No	No	1073.5	31767
No	No	529.3	35704
No	No	785.7	38463
No	Yes	919.6	7492

```
1 summary(Default)
2 #> default student balance income
3 #> No :9667 No :7056 Min. : 0 Min. : 772
4 #> Yes: 333 Yes:2944 1st Qu.: 482 1st Qu.:21340
5 #> Median : 824 Median :34553
6 #> Mean : 835 Mean : 33517
7 #> 3rd Qu.:1166 3rd Qu.:43808
8 #> Max. :2654 Max. :73554
```





**Your Turn #1 : Credit Card Default Modeling**

How would you construct a model to predict defaults?

## 2 Classification and Pattern Recognition

- The response variable is categorical and denoted  $G \in \mathcal{G}$ 
  - Default Credit Card Example:  $\mathcal{G} = \{\text{"Yes"}, \text{"No"}\}$
  - Medical Diagnosis Example:  $\mathcal{G} = \{\text{"stroke"}, \text{"heart attack"}, \text{"drug overdose"}, \text{"vertigo"}\}$
- The training data is  $D = \{(X_1, G_1), (X_2, G_2), \dots, (X_n, G_n)\}$
- The optimal decision/classification is often based on the posterior probability  $\Pr(G = g \mid \mathbf{X} = \mathbf{x})$

### 2.1 Binary Classification

- Classification is simplified when there are only 2 classes.
  - Many multi-class problems can be addressed by solving a set of binary classification problems (e.g., [one-vs-rest](#)).
- It is often convenient to *code* the response variable to a binary  $\{0, 1\}$  variable:

$$Y_i = \begin{cases} 1 & G_i = \mathcal{G}_1 \quad (\text{outcome of interest}) \\ 0 & G_i = \mathcal{G}_2 \end{cases}$$

- In the `Default` data, it would be natural to set `default=Yes` to 1 and `default=No` to 0.

#### 2.1.1 Linear Regression

- In this set-up we can run linear regression

$$\hat{y}(\mathbf{x}) = \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j x_j$$

```

1  ## Create binary column (y)
2  Default = Default %>% mutate(y = ifelse(default == "Yes", 1L, 0L))
3
4  ## Fit Linear Regression Model
5  fit.lm = lm(y ~ student + balance + income, data=Default)

```

term	estimate	std.error	statistic	p.value
(Intercept)	-0.0812	0.0084	-9.685	0.0000
studentYes	-0.0103	0.0057	-1.824	0.0682
balance	0.0001	0.0000	37.412	0.0000
income	0.0000	0.0000	1.039	0.2990

### Your Turn #2 : OLS for Binary Responses

1. For the binary  $Y$ , what is linear regression estimating?

2. What is the *loss function* that linear regression is using?
3. How could you create a *hard classification* from the linear model?
4. Does it make sense to use linear regression for binary classification?

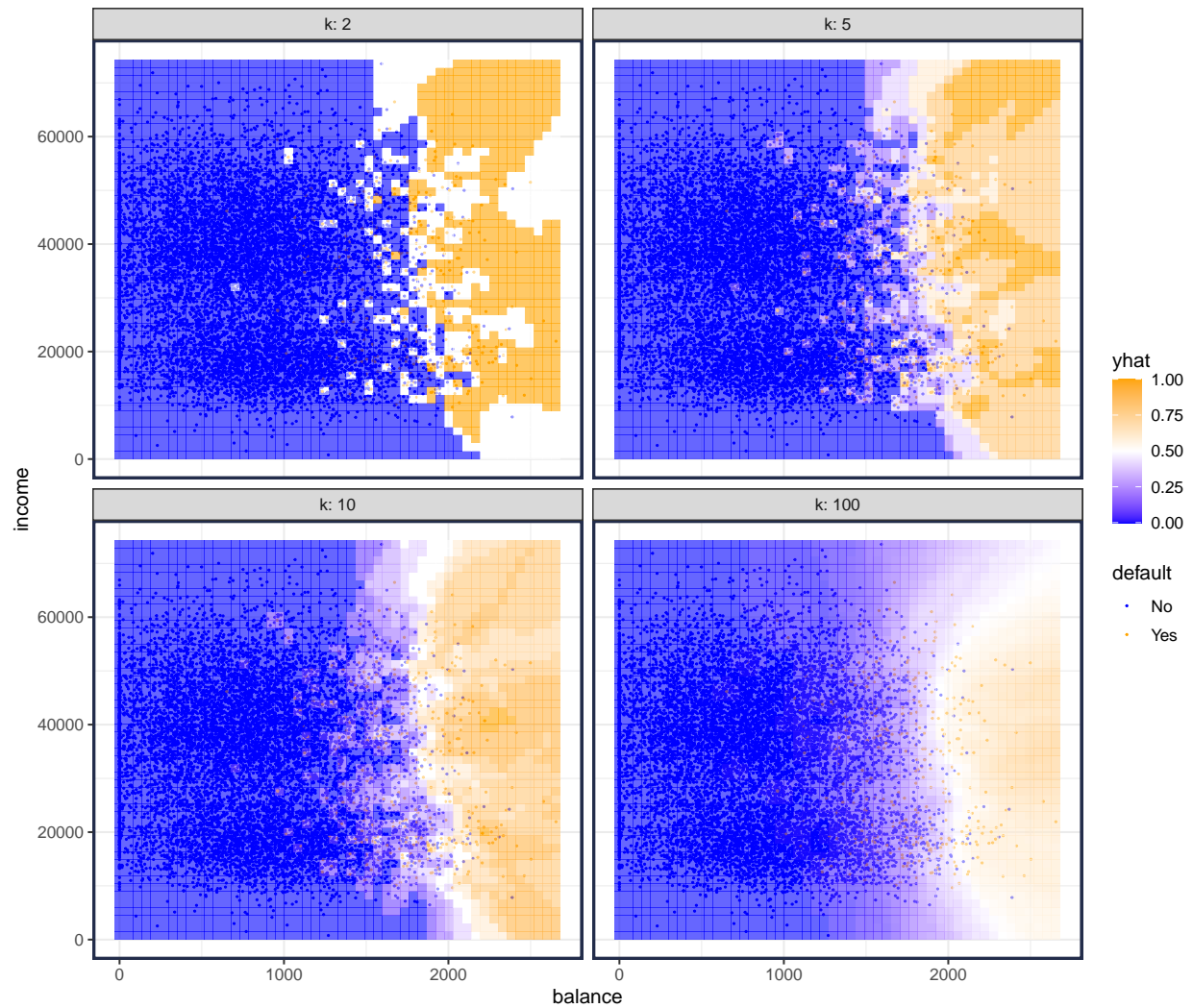
### 2.1.2 *k*-nearest neighbor (kNN)

- The *k*-NN method is a non-parametric *local* method, meaning that to make a prediction  $\hat{y}|x$ , it only uses the training data in the *vicinity* of  $x$ .
  - contrast with OLS linear regression, which uses all  $X$ 's to get prediction.
- The model (for regression and binary classification) is simple to describe

$$\begin{aligned} f_{\text{knn}}(x; k) &= \frac{1}{k} \sum_{i: x_i \in N_k(x)} y_i \\ &= \text{Avg}(y_i \mid x_i \in N_k(x)) \end{aligned}$$

- $N_k(x)$  are the set of  $k$  nearest neighbors
  - only the  $k$  closest  $y$ 's are used to generate a prediction
  - it is a *simple mean* of the  $k$  nearest observations
- When  $y$  is binary (i.e.,  $y \in \{0, 1\}$ ), the kNN model estimates

$$f_{\text{knn}}(x; k) \approx p(x) = \Pr(Y = 1 | X = x)$$



### Your Turn #3 : Thoughts about kNN

The above plots show a kNN model using the *continuous* predictors of *balance* and *income*.

- How could you use kNN with the categorical *student* predictor?

- The  $k$ -NN model also has a more general description when the response variables is categorical  $G_i \in \mathcal{G}$

$$\begin{aligned} f_g^{\text{knn}}(x; k) &= \frac{1}{k} \sum_{i: x_i \in N_k(x)} \mathbb{1}(g_i = g) \\ &= \widehat{\text{Pr}}(G_i = g \mid x_i \in N_k(x)) \end{aligned}$$

- $N_k(x)$  are the set of  $k$  nearest neighbors
- only the  $k$  closest  $y$ 's are used to generate a prediction

- it is a *simple proportion* of the  $k$  nearest observations that are of class  $g$



### 3 Logistic Regression

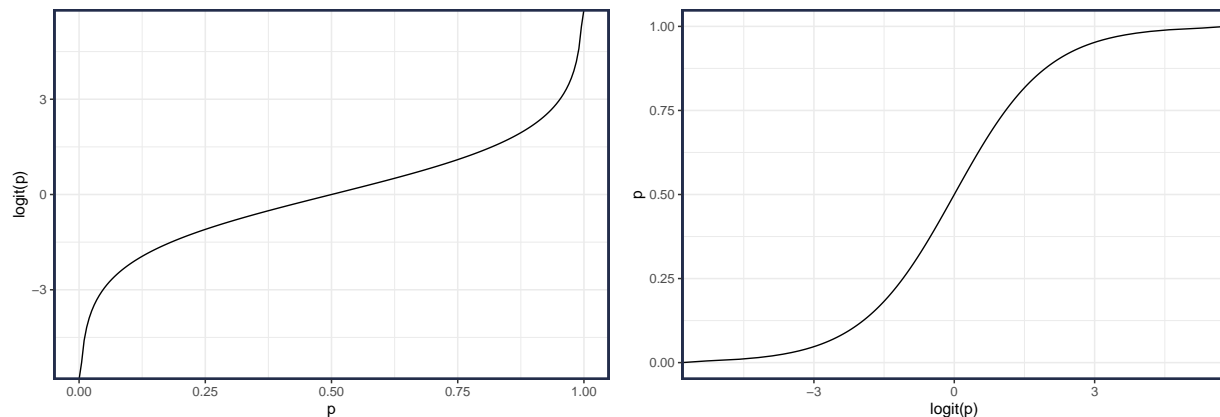
#### 3.1 Basics

- Let  $0 \leq p \leq 1$  be a probability.
- The log-odds of  $p$  is called the *logit*

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right)$$

- The inverse logit is the *logistic function*. Let  $f = \text{logit}(p)$ , then

$$\begin{aligned} p &= \frac{e^f}{1 + e^f} \\ &= \frac{1}{1 + e^{-f}} \end{aligned}$$



- For binary response variables  $Y \in \{0, 1\}$ , linear regression models estimate

$$E[Y \mid X = x] = \Pr(Y = 1 \mid X = x) = \beta^\top x$$

- Logistic Regression models alternatively estimate

$$\log\left(\frac{\Pr(Y = 1 \mid X = x)}{1 - \Pr(Y = 1 \mid X = x)}\right) = \beta^\top x$$

and thus,

$$\begin{aligned} \Pr(Y = 1 \mid X = x) &= \frac{e^{\beta^\top x}}{1 + e^{\beta^\top x}} \\ &= \left(1 + e^{-\beta^\top x}\right)^{-1} \end{aligned}$$

#### 3.2 Estimation

- The data for logistic regression is:  $(\mathbf{x}_i, y_i)_{i=1}^n$  where  $y_i \in \{0, 1\}$ ,  $\mathbf{x}_i = (x_{i0}, x_{i1}, \dots, x_{ip})^\top$ .
- $y_i \mid \mathbf{x}_i \sim \text{Bern}(p_i(\beta))$

$$- p_i(\beta) = \Pr(Y = 1 \mid \mathbf{X} = \mathbf{x}_i; \beta) = \left(1 + e^{-\beta^\top \mathbf{x}_i}\right)^{-1}$$

$$- \beta^\top \mathbf{x}_i = \mathbf{x}_i^\top \beta = \beta_0 + \sum_{j=1}^p x_{ij} \beta_j$$

- Bernoulli Likelihood Function

$$L(\beta) = \prod_{i=1}^n p_i(\beta)^{y_i} (1 - p_i(\beta))^{1-y_i}$$

$$\log L(\beta) = \sum_{i=1}^n \{y_i \ln p_i(\beta) + (1 - y_i) \ln(1 - p_i(\beta))\}$$

- The usual approach to estimating the Logistic Regression coefficients is *maximum likelihood*

$$\begin{aligned} \hat{\beta} &= \arg \max_{\beta} L(\beta) \\ &= \arg \max_{\beta} \log L(\beta) \end{aligned}$$

- We can also view this as the coefficients that minimize the *loss function*, where the loss function is the negative log-likelihood

$$\begin{aligned} \hat{\beta} &= \arg \min_{\beta} \ell(\beta) \\ &= -C \sum_{i=1}^n \{y_i \ln p_i(\beta) + (1 - y_i) \ln(1 - p_i(\beta))\} \end{aligned}$$

- where  $C$  is some constant, e.g.,  $C = 1/n$

- This view facilitates *penalized logistic regression*

$$\hat{\beta} = \arg \min_{\beta} \ell(\beta) + \lambda P(\beta)$$

- Ridge Penalty

$$P(\beta) = \sum_{j=1}^p |\beta_j|^2 = \beta^\top \beta$$

- Lasso Penalty

$$P(\beta) = \sum_{j=1}^p |\beta_j|$$

- Best Subsets

$$P(\beta) = \sum_{j=1}^p |\beta_j|^0 = \sum_{j=1}^p 1_{(\beta_j \neq 0)}$$

- Elastic Net (glmnet form)

$$P(\beta, \alpha) = \sum_{j=1}^p (1 - \alpha) |\beta_j|^2 / 2 + \alpha |\beta_j|$$

### 3.3 Logistic Regression in Action

- In **R**, logistic regression can be implemented with the `glm()` function since it is a type of *Generalized Linear Model*.
- Because logistic regression is a special case of *Binomial* regression, use the `family=binomial()` argument

```
1  #-- Fit logistic regression model
2  fit.lr = glm(y~student + balance + income, data=Default,
3             family="binomial")
```

term	estimate	std.error	statistic	p.value
(Intercept)	-10.8690	0.4923	-22.0801	0.0000
studentYes	-0.6468	0.2363	-2.7376	0.0062
balance	0.0057	0.0002	24.7376	0.0000
income	0.0000	0.0000	0.3698	0.7115

#### Your Turn #4 : Interpreting Logistic Regression

1. What is the estimated probability of default for a Student with a balance of \$1000?
2. What is the estimated probability of default for a *Non-Student* with a balance of \$1000?
3. Why does `student=Yes` appear to lower risk of default, when the plot of student status vs. default appears to increase risk?

#### 3.3.1 Penalized Logistic Regression

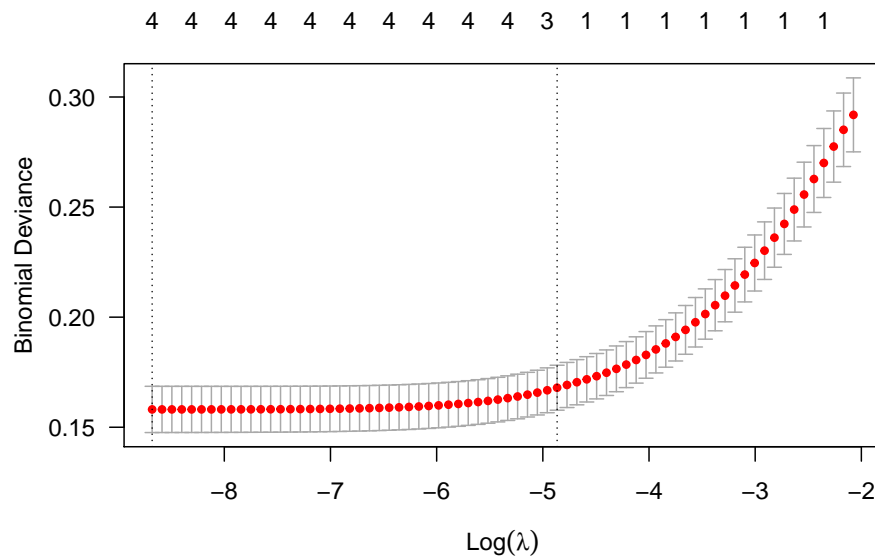
- The `glmnet()` package can estimate logistic regression using an elastic net penalty (e.g., ridge, lasso).

```
1  #-- Fit *penalized* logistic regression model
2  library(glmnet)
3  library(glmnetUtils)
4  set.seed(2020)
5  fit.enet = cv.glmnet(y~student + balance + income, data=Default,
```

```

6         alpha=.5,
7         family="binomial")
8
9  #-- CV performance plot
10 plot(fit.enet, las=1)

```



term	unpenalized	lambda.min	lambda.1se
(Intercept)	-10.869	-11.056	-7.937
studentYes	-0.647	-0.299	-0.041
balance	0.006	0.006	0.004
income	0.000	0.000	0.000
studentNo	NA	0.325	0.044

### 3.4 Logistic Regression Summary

- Logistic Regression (both penalized and unpenalized) estimates a *posterior probability*,  $\hat{p}(x) = \widehat{\Pr}(Y = 1 \mid X = x)$
- This estimate is a function of the estimated coefficients

$$\begin{aligned}
 \hat{p}(x) &= \frac{e^{\hat{\beta}^\top x}}{1 + e^{\hat{\beta}^\top x}} \\
 &= \left(1 + e^{-\hat{\beta}^\top x}\right)^{-1}
 \end{aligned}$$

## 4 Evaluating Classification Models

- Training Data:  $\{X_i, G_i\}$ 
  - $G_i \in \{1, \dots, K\}$  (i.e., there are  $K$  classes)
- Predictor:  $\hat{G}(X)$
- Loss function:  $L(G, \hat{G}(X))$  is the loss incurred by estimating  $G$  with  $\hat{G}$
- Risk is the expected loss (or expected prediction error EPE)
  - Expectation is taken wrt future values of  $(X, G)$

$$\begin{aligned}
 \text{Risk}(\hat{G}) &= \text{EPE} \\
 &= E_{XG} [L(G, \hat{G}(X))] \\
 &= E_X [E_{G|X} [L(G, \hat{G}(X)) | X]] \\
 &= E_X [R_X(\hat{G})]
 \end{aligned}$$

- The Risk at input  $X = x$  is

$$\begin{aligned}
 R_x(\hat{G}) &= E_{G|X=x} [L(G, \hat{G}(x)) | X = x] \\
 &= \sum_{k=1}^K L(G = k, \hat{G}(x)) \Pr(G = k | X = x) \\
 &= \sum_{k=1}^K L(G = k, \hat{G}(x)) p_k(x)
 \end{aligned}$$

- Thus the optimal class label, given  $X = x$ , is

$$\hat{G}(x) = \arg \min_g R_x(g)$$

- The optimal label can only be obtained if you know  $p_k(x) = \Pr(G = k | X = x)$  for all  $k$ .
- Since we won't know  $\{p_k(x)\}_k$ , we have to estimate them.

### 4.1 Evaluation of Binary Classification Models

- We are considering *binary* outcomes, so use the notation  $Y \in \{0, 1\}$
- Let  $p(x) = \Pr(Y = 1 | X = x)$
- The Risk (for a binary outcome) is:

$$\begin{aligned}
 R_x(g) &= L(1, g) \Pr(Y = 1 | X = x) + L(0, g)(1 - \Pr(Y = 1 | X = x)) \\
 &= L(1, g)p(x) + L(0, g)(1 - p(x))
 \end{aligned}$$

- Hard Decision ( $\hat{G}(x) \in \{0, 1\}$ ): choose  $\hat{G}(x) = 1$  if

$$\begin{aligned}
 & R_x(1) < R_x(0) \\
 & L(1, 1)p(x) + L(0, 1)(1 - p(x)) < L(1, 0)p(x) + L(0, 0)(1 - p(x)) \\
 & p(x)(L(1, 1) - L(1, 0)) < (1 - p(x))(L(0, 0) - L(0, 1)) \\
 & p(x)(L(1, 0) - L(1, 1)) \geq (1 - p(x))(L(0, 1) - L(0, 0)) \quad (\text{multiply both sides by } -1) \\
 & \frac{p(x)}{1 - p(x)} \geq \frac{L(0, 1) - L(0, 0)}{L(1, 0) - L(1, 1)}
 \end{aligned}$$

#### 4.1.1 Example

- Say we have a goal of estimating if a patient has cancer using medical imaging
  - Let  $G = 1$  for cancer and  $G = 0$  for no cancer
- Suppose we have solicited a loss function with the following values
  - $L(G = 0, \hat{G} = 0) = 0$ : There is no loss for correctly diagnosis a patient without cancer
  - $L(G = 1, \hat{G} = 1) = 0$ : There is no loss (for our model) for correctly diagnosis a patient with cancer
  - $L(G = 0, \hat{G} = 1) = FP$ : There is a cost of  $FP$  units if the model issues a *false positive*, estimating the patient has cancer when they don't
  - $L(G = 1, \hat{G} = 0) = FN$ : There is a cost of  $FN$  units if the model issues a *false negative*, estimating the patient does not have cancer when they really do
  - In these scenarios  $FN$  is often much larger than  $FP$  ( $FN \gg FP$ ) because the effects of not promptly treating (or further testing, etc) a patient is more severe than starting a treatment path for patients that don't actually have cancer
- Our model will decide to issue a positive indication for cancer if  $R_x(1) < R_x(0)$  which occurs when

$$\begin{aligned}
 \frac{p(x)}{1 - p(x)} & \geq \frac{FP}{FN} \\
 p(x) & \geq \frac{FP}{FP + FN}
 \end{aligned}$$

- The ratio of FP to FN is all that matters for the decision. Let's say that  $FP=1$  and  $FN=10$ . Then if  $p(x) \geq 1/11$ , our model will diagnose cancer.
  - Note:  $p(x) = \Pr(Y = 1|X = x)$  is affected by the class prior  $\Pr(Y = 1)$  (e.g., the portion of the population tested with cancer), which is usually going to be small.

## 4.2 Common Binary Loss Functions

- Suppose we are going to estimate a binary response  $Y \in \{0, 1\}$  with a (possibly continuous) predictor  $\hat{y}(x)$
- **0-1 Loss or Misclassification Error**

$$L(y, \hat{y}(x)) = \mathbb{1}(y \neq \hat{y}(x)) = \begin{cases} 0 & y = \hat{y}(x) \\ 1 & y \neq \hat{y}(x) \end{cases}$$

- This assumes  $L(0, 1) = L(1, 0)$  (i.e., false positive costs the same as a false negative)

- Requires that a *hard classification* is made
- The optimal prediction is  $y^*(x) = \mathbb{1}(p(x) > 1 - p(x))$  which is equivalent to  $\mathbb{1}(p(x) > 0.50)$

- **Squared Error**

$$L(y, \hat{y}(x)) = (y - \hat{y}(x))^2$$

- The optimal prediction is  $y^*(x) = E[Y | X = x] = \Pr(Y = 1 | X = x)$

- **Absolute Error**

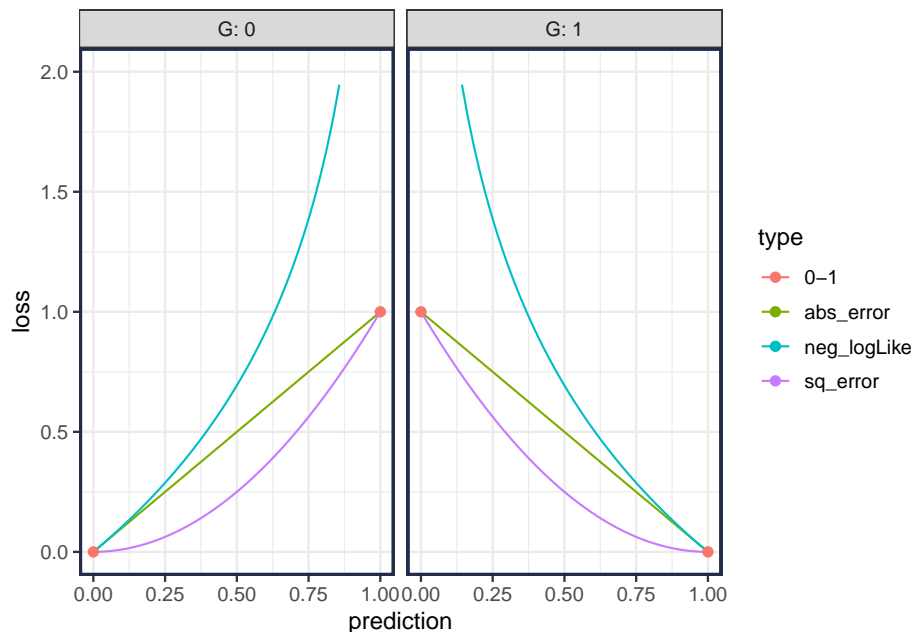
$$L(y, \hat{y}(x)) = |y - \hat{y}(x)|$$

- **Bernoulli negative log-likelihood**

$$L(y, \hat{y}(x)) = -\{y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)\}$$

$$= \begin{cases} -\log \hat{y} & y = 1 \\ -\log(1 - \hat{y}) & y = 0 \end{cases}$$

- Requires  $\hat{y}(x) \in [0, 1]$



### 4.3 Evaluating Binary Classification Models

- Recall, the optimal (*hard classification*) decision is to choose  $\hat{G} = 1$  if:

$$\frac{p(x)}{1 - p(x)} \geq \frac{L(0, 1) - L(0, 0)}{L(1, 0) - L(1, 1)}$$

- Denote  $\gamma(x)$  as the *logit* of  $p(x)$ :

$$\gamma(x) = \log \frac{p(x)}{1 - p(x)} = \log \frac{\Pr(G = 1 | X = x)}{\Pr(G = 0 | X = x)}$$

- Then we get

$$p(x) = \Pr(G = 1 \mid X = x) \\ = \frac{e^{\gamma(x)}}{1 + e^{\gamma(x)}}$$

- And the optimal (*hard classification*) decision can be described in terms of  $\gamma(x)$ :

Choose  $\hat{G}(x) = 1$  if  $\hat{\gamma}(x) > t$ , where  $t$  is a threshold

- If the loss/cost is known, then

$$t^* = \log \left( \frac{L(0, 1) - L(0, 0)}{L(1, 0) - L(1, 1)} \right)$$

- For a given threshold  $t$  and input  $x$ , the hard classification is  $\hat{G}_t(x) = \mathbb{1}(\hat{\gamma}(x) \geq t)$
- Note: we can also use  $\hat{p}(x)$  to make the decision threshold.
  - Just adjust the threshold accordingly.

## 4.4 Performance Metrics

### 4.4.1 Confusion Matrix

- Given a threshold  $t$ , we can make a *confusion matrix* to help analyze our model's performance on data
  - Data =  $\{(X_i, G_i)\}_{i=1}^N$  (ideally this is hold-out/test data)
  - $N_k$  is number of observations from class  $k$  ( $N_0 + N_1 = N$ )

		True Outcome		total
		$G = 1$	$G = 0$	
Model Outcome	$\hat{G}_t = 1$	True Positive (TP)	False Positive (FP)	$\hat{N}_1(t)$
	$\hat{G}_t = 0$	False Negative (FN)	True Negative (TN)	$\hat{N}_0(t)$
total		$N_1$	$N_2$	$N$

Table from: <https://tex.stackexchange.com/questions/20267/how-to-construct-a-confusion-matrix-in-latex>

To illustrate a confusion table in practice let's go back to the `Default` data and see how the basic logistic regression models performs.

- In order to evaluation on hold-out data, split the data into train/test (used about 90% training, 10% testing), fit a logistic regression model on training data, and make predictions on the test data
- Note that only 3.3% of the data is default.



- Using a threshold of  $\hat{p}(x) \geq 0.10$  to make a hard classification.

```

1  #-- train/test split
2  set.seed(2019)
3  test = sample(nrow(Default), size=1000)
4  train = -test
5
6  #-- fit model on training data
7  fit.lm = glm(y~student + balance + income, family='binomial',
8              data=Default[train, ])
9
10 #-- Get predictions (of p(x)) on test data
11 p.hat = predict(fit.lm, newdata=Default[test, ], type='response')
12
13
14 #-- Make Hard classification (use .10 as cut-off)
15 G.hat = ifelse(p.hat >= .10, 1, 0)
16
17 #-- Make Confusion Table
18 G.test = Default$y[test] # true values
19
20 table(predicted=G.hat, truth = G.test) %>% addmargins()
21 #>           truth
22 #> predicted    0    1  Sum
23 #>           0   896    7  903
24 #>           1    68   29   97
25 #>           Sum  964   36 1000

```

#### 4.4.2 Metrics

Metric	Definition	Estimate
Risk/Exp.Cost	$\sum_{i=0}^1 \sum_{j=0}^1 L(i, j) P_X(G(X) = i, \hat{G}_t(X) = j)$	$\frac{1}{N} \sum_{i=1}^N L(G_i, \hat{G}_t(x_i))$
Mis-classification Rate	$P_{XG}(\hat{G}_t(X) \neq G(X)) =$ $P_X(\hat{G}_t(X) = 0, G(X) = 1) +$ $P_X(\hat{G}_t(X) = 1, G(X) = 0)$	$\frac{1}{N} \sum_{i=1}^N \mathbb{1}(\hat{G}_t(x_i) \neq G_i)$
False Positive Rate (FPR) {1-Specificity}	$P_X(\hat{G}_t(X) = 1 \mid G(X) = 0)$	$\frac{1}{N_0} \sum_{i: G_i=0} \mathbb{1}(\hat{G}_t(x_i) = 1)$
True Positive Rate (TPR) {Hit Rate, Recall, Sensitivity}	$P_X(\hat{G}_t(X) = 1 \mid G(X) = 1)$	$\frac{1}{N_1} \sum_{i: G_i=1} \mathbb{1}(\hat{G}_t(x_i) = 1)$
Precision TP/(TP + FP)	$P_X(G(X) = 1 \mid \hat{G}_t(X) = 1)$	$\frac{1}{\hat{N}_1(t)} \sum_{i: \hat{G}_t(x_i)=1} \mathbb{1}(G_i = 1)$

- Note: Performance estimates are best carried out on *hold-out* data!
- See [Wikipedia Page: Confusion Matrix](#) for more metrics

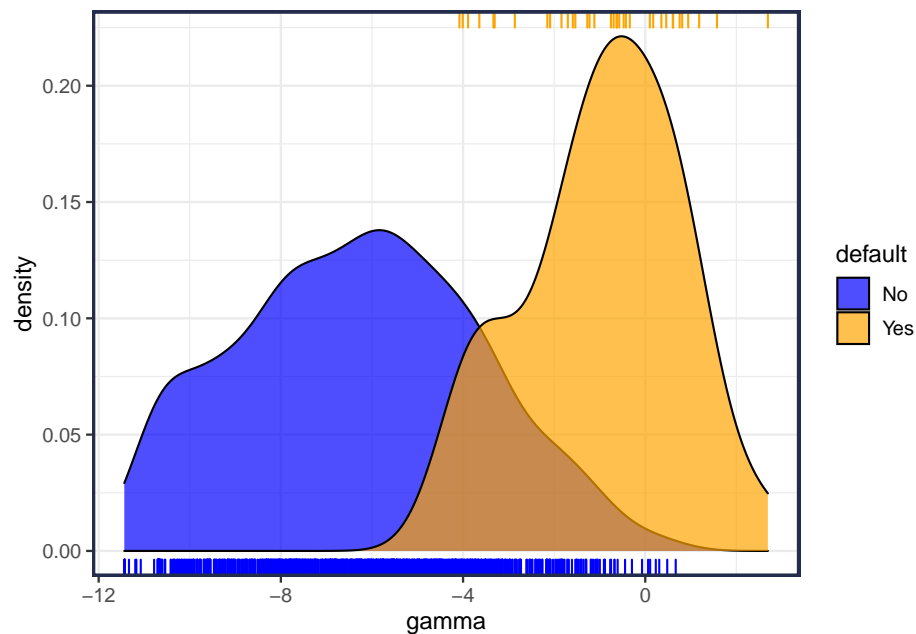
## 4.5 Performance over a range of thresholds

In the previous example, a hard classification was made using a threshold of  $\hat{p}(x) \geq 0.10$ . But performance varies as we adjust the threshold. Let's explore!

I'll use  $\hat{\gamma}(x)$  instead of  $\hat{p}(x)$  for this illustration.

```
1 #-- Get predictions (of gamma(x)) on test data
2 gamma = predict(fit.lm, newdata=Default[test,], type='link')
```

- The model is unable to perfectly discriminate between groups, but the *defaults* do get scored higher in general:
  - As a reference point, note that  $\gamma(x) = 0 \rightarrow \Pr(Y = 1 | X = x) = 1/2$

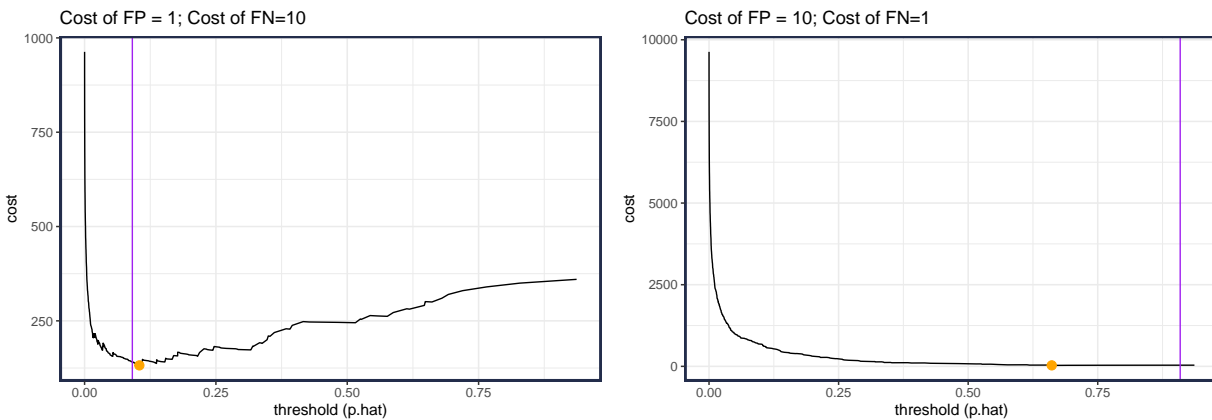


- We can calculate performance over a range of thresholds.
  - Unless the test data is too large, use all unique values of the training data as the thresholds. If too large, manually create threshold sequence.

```
1 #-- Get performance data (by threshold)
2 perf = tibble(truth = G.test, gamma, p.hat) %>%
3   #- group_by() + summarize() in case of ties
4   group_by(gamma, p.hat) %>%
5   summarize(n=n(), n.1=sum(truth), n.0=n-sum(truth)) %>% ungroup() %>%
6   #- calculate metrics
7   arrange(gamma) %>%
8   mutate(FN = cumsum(n.1),      # false negatives
9          TN = cumsum(n.0),      # true negatives
10         TP = sum(n.1) - FN,     # true positives
11         FP = sum(n.0) - TN,     # false positives
12         N = cumsum(n),          # number of cases predicted to be 1
13         TPR = TP/sum(n.1), FPR = FP/sum(n.0)) %>%
14   #- only keep relevant metrics
15   select(-n, -n.1, -n.0, gamma, p.hat)
```

### 4.5.1 Cost Curves

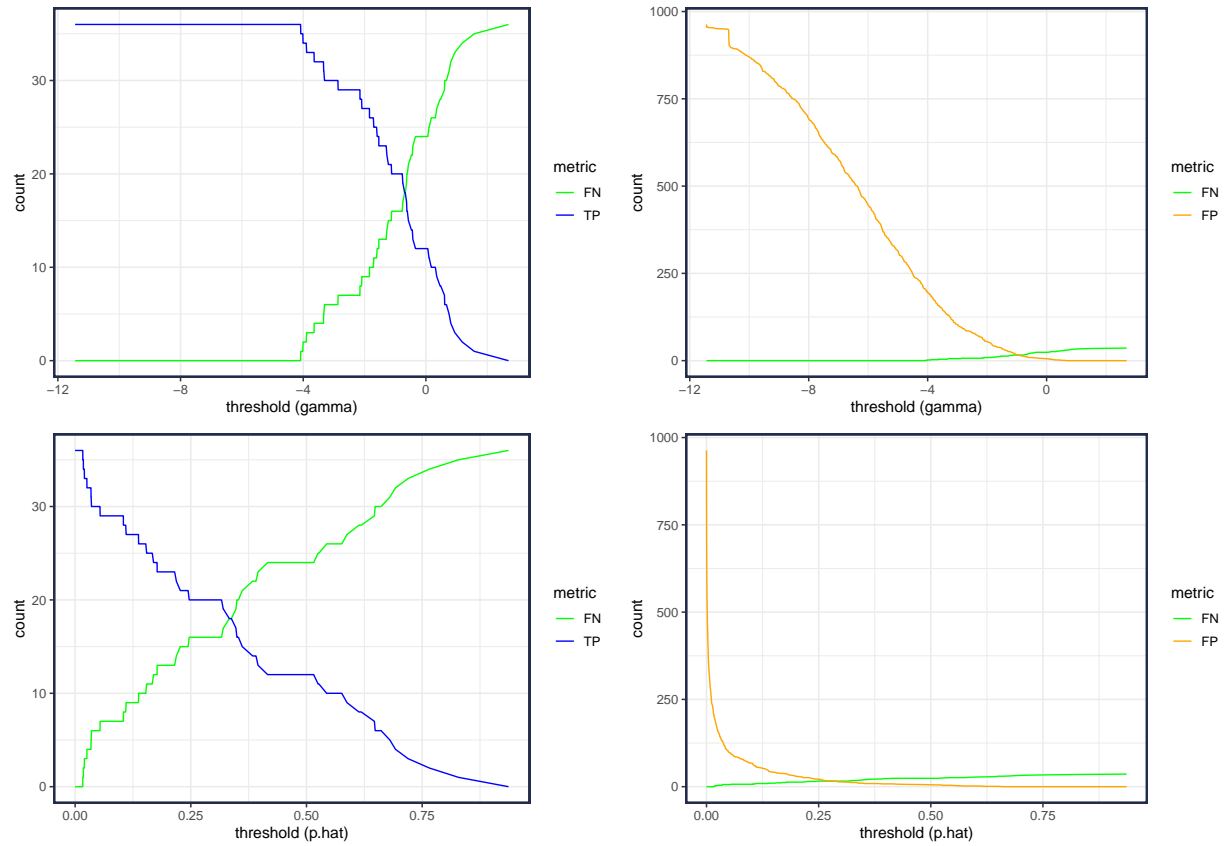
- note: the **purple** is the *theoretical* optimal threshold (using  $t^* = \log FP/FN$ ) and the **orange** point is at the optimal value using the model



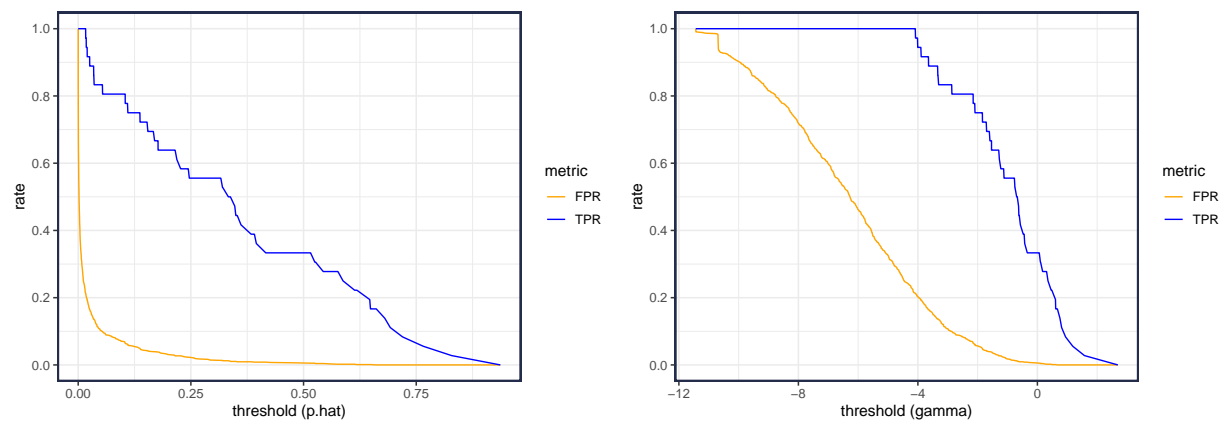
#### Optimal Threshold

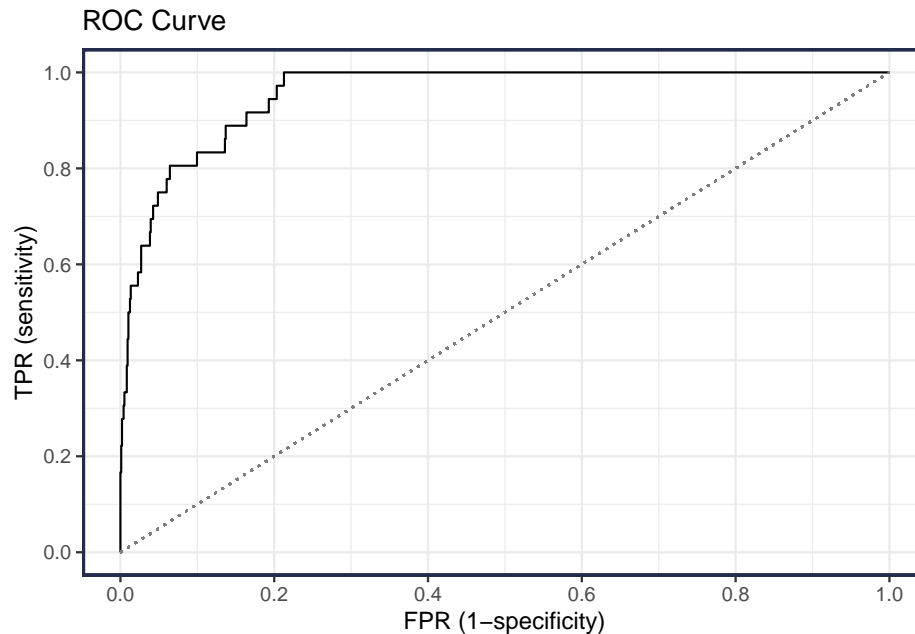
- The *theoretically* optimal threshold is based on the *true*  $\gamma(x) = \log \frac{p(x)}{1-p(x)}$  (for a given cost ratio of FP to FN)
- The observed optimal threshold will differ when the model's estimate  $\hat{\gamma}(x) \neq \gamma(x)$ 
  - Hopefully, they are close and it won't make much difference which one you use. But I'd take the estimated threshold if I had sufficient data.
- Note that the estimated values depend on the prior class probabilities. If you suspect these may differ in the future, then you should adjust the threshold.

### 4.5.2 General Performance as function of threshold (select metrics)



### 4.5.3 ROC Curves (Receiver Operating Characteristic)



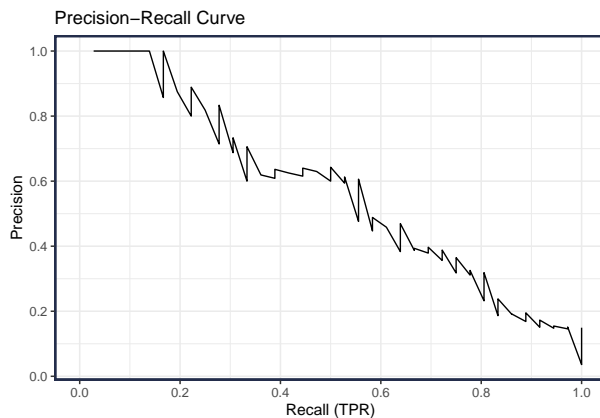
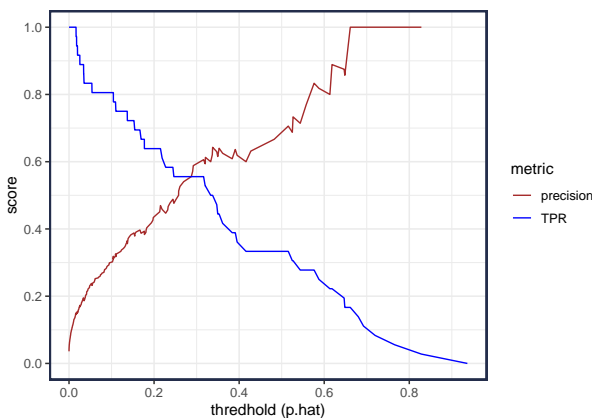


### AUROC

- The *area under the ROC curve* (AUROC) is a popular performance metric
- I don't think it is a great way to compare classifiers for several reasons
  - The main reason is that in a real application you can almost always come up with an estimated cost/loss for the different decisions
  - To say it another way, comparisons should be made at a single point on the curve; the entire FPR region should not factor into the comparison.
- The AUROC is equal to the probability that a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one.
  - AUROC is proportional to the [Mann-Whitney U statistic](#)

#### 4.5.4 Precision Recall Curves

- Popular for information retrieval/ranking
- The *precision* metric is not monotonic wrt threshold, hence the sawteeth pattern.

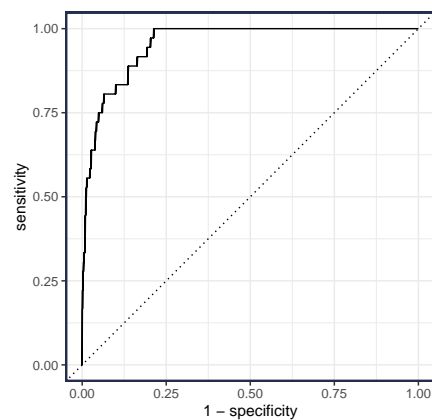
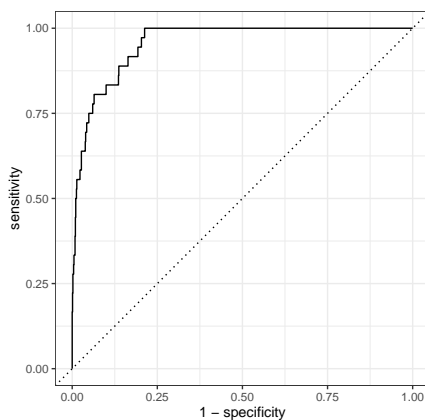


### 4.5.5 R Code

Once we have the FP, TP, TN, FN values for a set of thresholds (like what is in the `perf` object), then you have everything you need to calculate any metric (e.g., gain, lift, F1, ...).

- But I will mention the `yardstick` R package which offers some functionality you may find convenient
- Here is a list of the [metrics included in the yardstick package](#)

```
1 library(yardstick) # for evaluation functions
2
3 #-- ROC plots
4 ROC = tibble(truth = factor(G.test, levels=c(1,0)), gamma) %>%
5   yardstick::roc_curve(truth, gamma)
6
7 autoplot(ROC) # autoplot() method
8
9 ROC %>% # same as autoplot()
10   ggplot(aes(1-specificity, sensitivity)) + geom_line() +
11   geom_abline(lty=3) +
12   coord_equal()
13
14 #-- Area under ROC (AUROC)
15 tibble(truth = factor(G.test, levels=c(1,0)), gamma) %>%
16   roc_auc(truth, gamma)
17 #> # A tibble: 1 x 3
18 #>   .metric .estimator .estimate
19 #>   <chr>   <chr>      <dbl>
20 #> 1 roc_auc binary      0.955
21
22 roc_auc_vec(factor(G.test, 1:0), gamma)
23 #> [1] 0.9552
```



## 4.6 Summary of Classification Evaluation

- Use cost! The other metrics are probably not going to give you what you really want.

- For Binary Classification Problems, the optimal decision is to choose  $\hat{G}(x) = 1$  if

$$\frac{p(x)}{1 - p(x)} \geq \frac{L(0, 1) - L(0, 0)}{L(1, 0) - L(1, 1)} \\ = \frac{\text{FP} - \text{TN}}{\text{FN} - \text{TP}}$$

- Consider the connection to Decision Theory, make the decision that maximizes *expected utility*. The losses define the utility.

## 5 Generative Classification Models

STILL IN PRODUCTION 2020-09-17

### 5.1 Linear/Quadratic Discriminant Analysis (LDA/QDA)

### 5.2 Naive Bayes