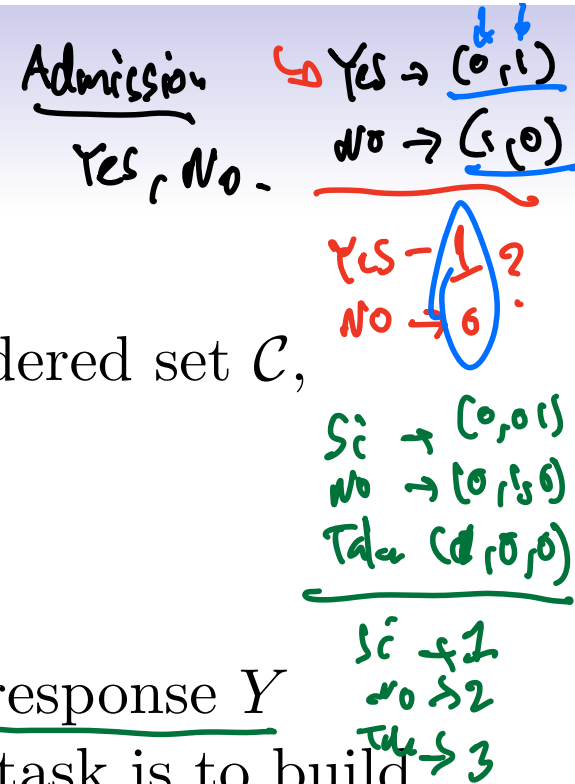


# Classification



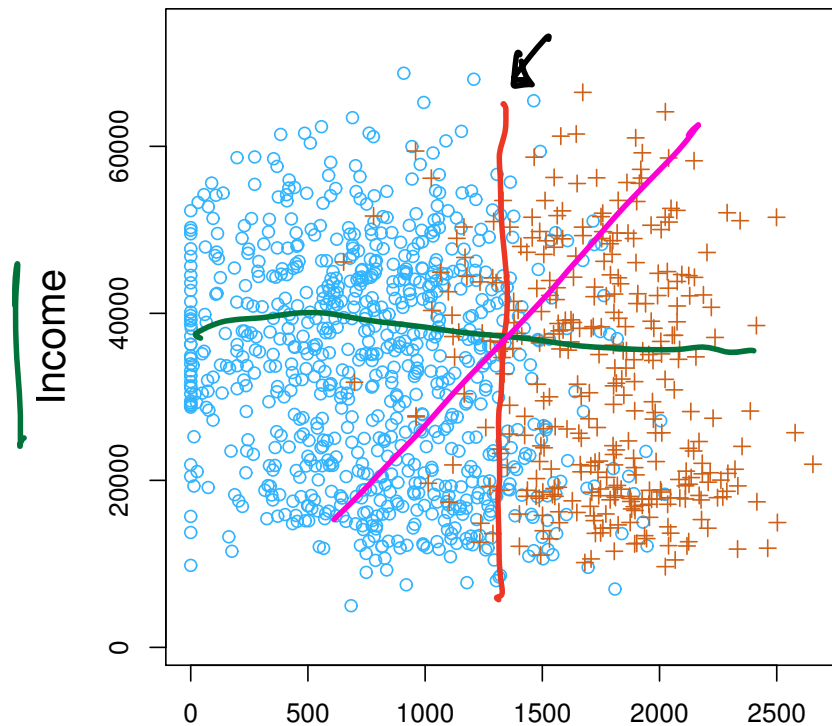
- Qualitative variables take values in an unordered set  $\mathcal{C}$ , such as:  
 $\text{eye color} \in \{\text{brown}, \text{blue}, \text{green}\}$   
 $\text{email} \in \{\text{spam}, \text{ham}\}$ .
- Given a feature vector  $X$  and a qualitative response  $Y$  taking values in the set  $\mathcal{C}$ , the classification task is to build a function  $C(X)$  that takes as input the feature vector  $X$  and predicts its value for  $Y$ ; i.e.  $C(X) \in \mathcal{C}$ .
- Often we are more interested in estimating the probabilities that  $X$  belongs to each category in  $\mathcal{C}$ .

# Classification

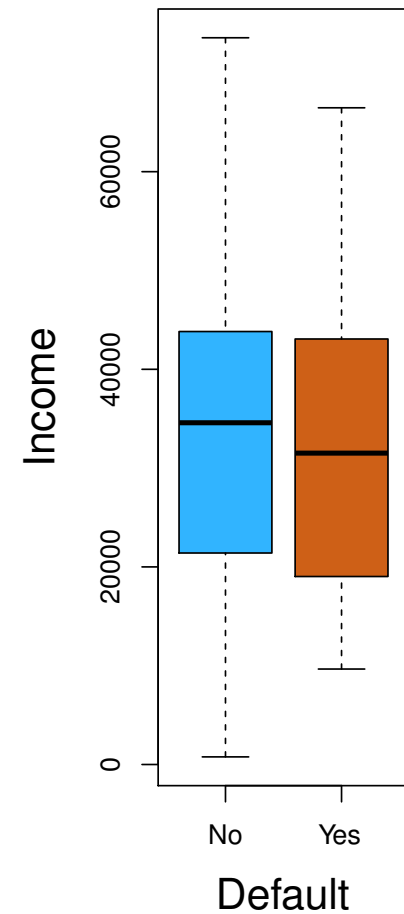
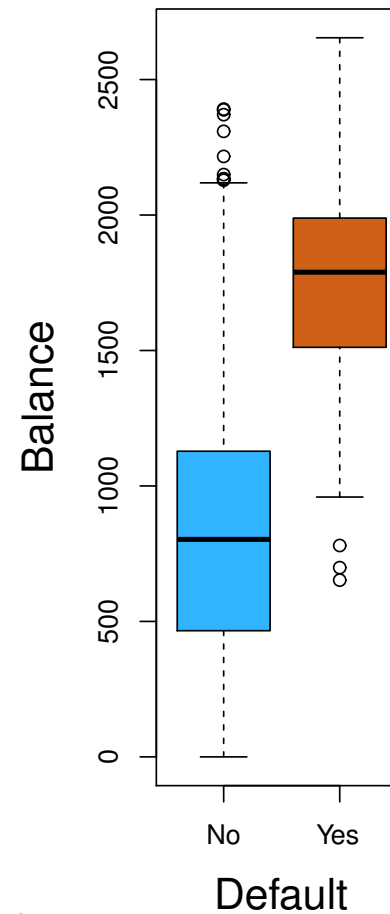
- Qualitative variables take values in an unordered set  $\mathcal{C}$ , such as:  
 $\text{eye color} \in \{\text{brown}, \text{blue}, \text{green}\}$   
 $\text{email} \in \{\text{spam}, \text{ham}\}.$
- Given a feature vector  $X$  and a qualitative response  $Y$  taking values in the set  $\mathcal{C}$ , the classification task is to build a function  $C(X)$  that takes as input the feature vector  $X$  and predicts its value for  $Y$ ; i.e.  $C(X) \in \mathcal{C}$ .
- Often we are more interested in estimating the *probabilities* that  $X$  belongs to each category in  $\mathcal{C}$ .

For example, it is more valuable to have an estimate of the probability that an insurance claim is fraudulent, than a classification fraudulent or not.

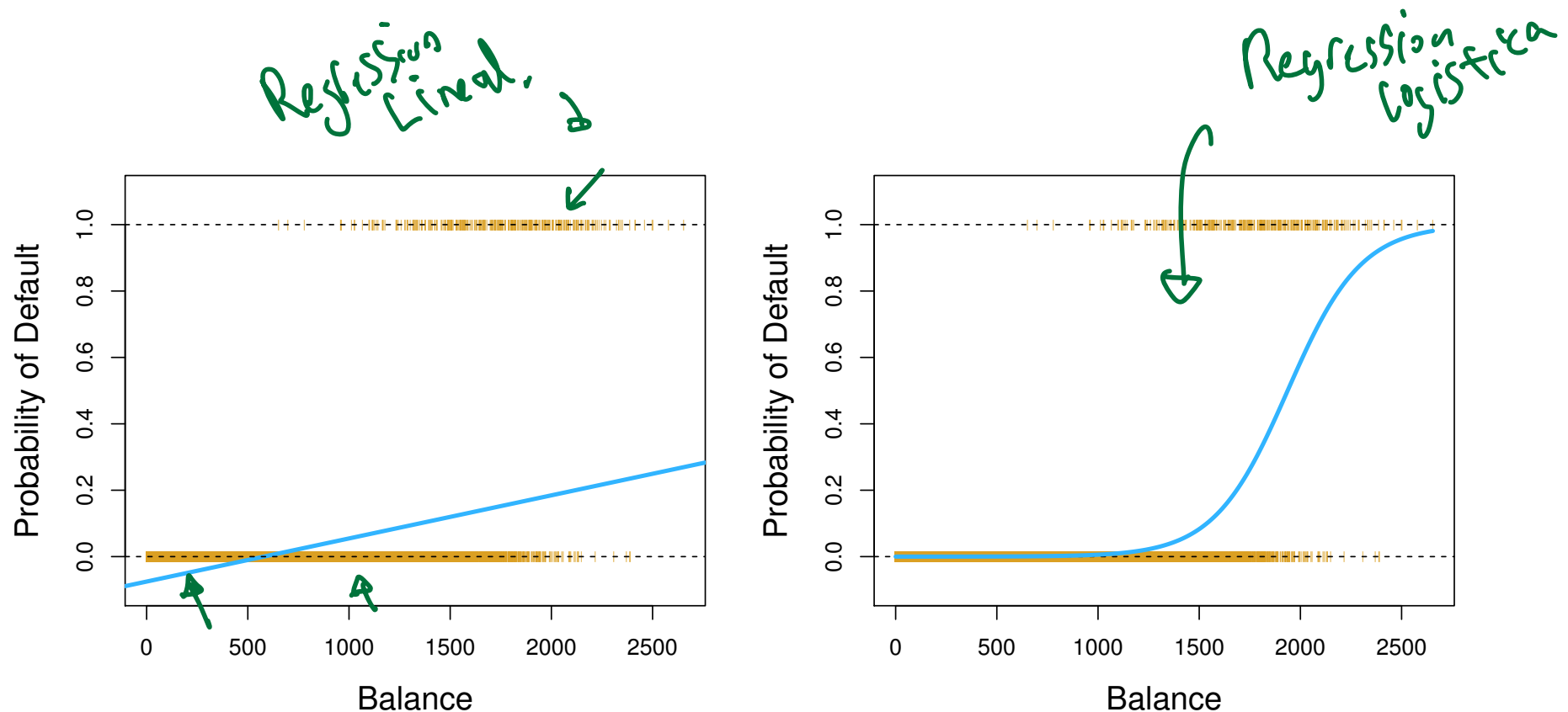
# Example: Credit Card Default



Balance [ deuda en tarjeta de crédito ]




# Linear versus Logistic Regression



The orange marks indicate the response  $Y$ , either 0 or 1. Linear regression does not estimate  $\Pr(Y = 1|X)$  well. Logistic regression seems well suited to the task.

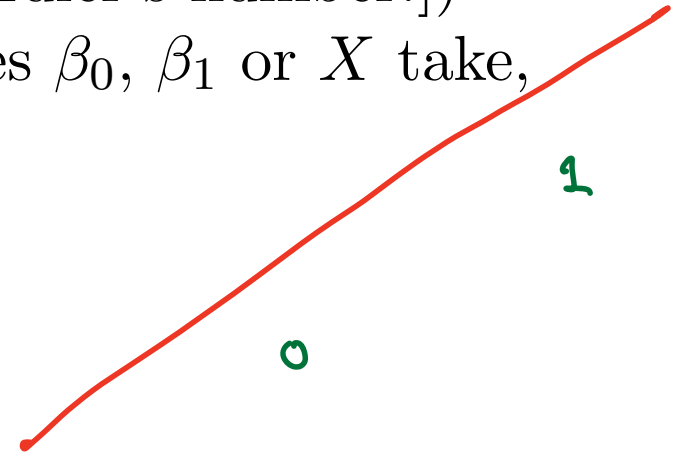
# Logistic Regression

Let's write  $p(X) = \Pr(Y = 1|X)$  for short and consider using **balance** to predict **default**. Logistic regression uses the form

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}.$$


( $e \approx 2.71828$  is a mathematical constant [Euler's number.] )

It is easy to see that no matter what values  $\beta_0$ ,  $\beta_1$  or  $X$  take,  $p(X)$  will have values between 0 and 1.



# Logistic Regression

Let's write  $p(X) = \Pr(Y = 1|X)$  for short and consider using **balance** to predict **default**. Logistic regression uses the form

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}.$$

Limitada a entre 0 y 1.

( $e \approx 2.71828$  is a mathematical constant [Euler's number.] )

It is easy to see that no matter what values  $\beta_0$ ,  $\beta_1$  or  $X$  take,  $p(X)$  will have values between 0 and 1.

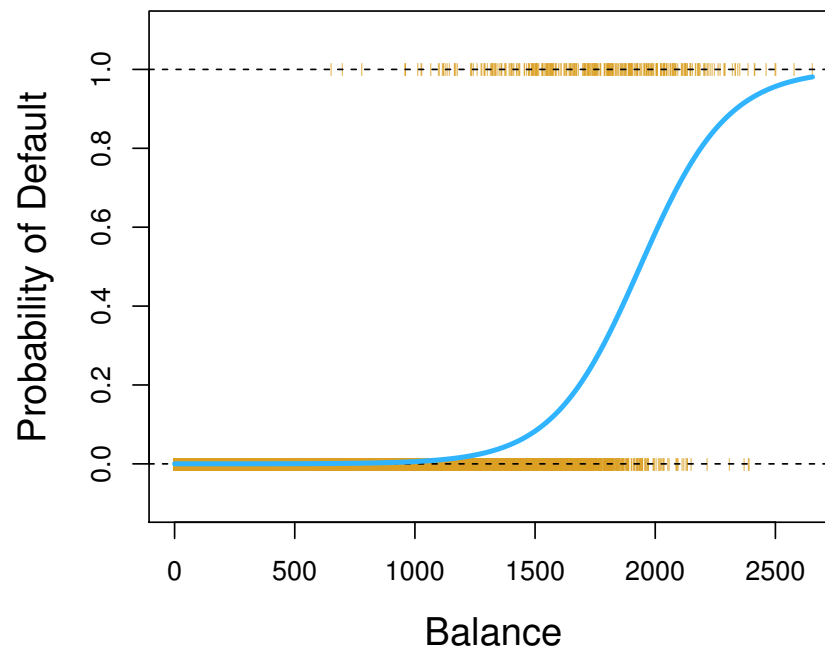
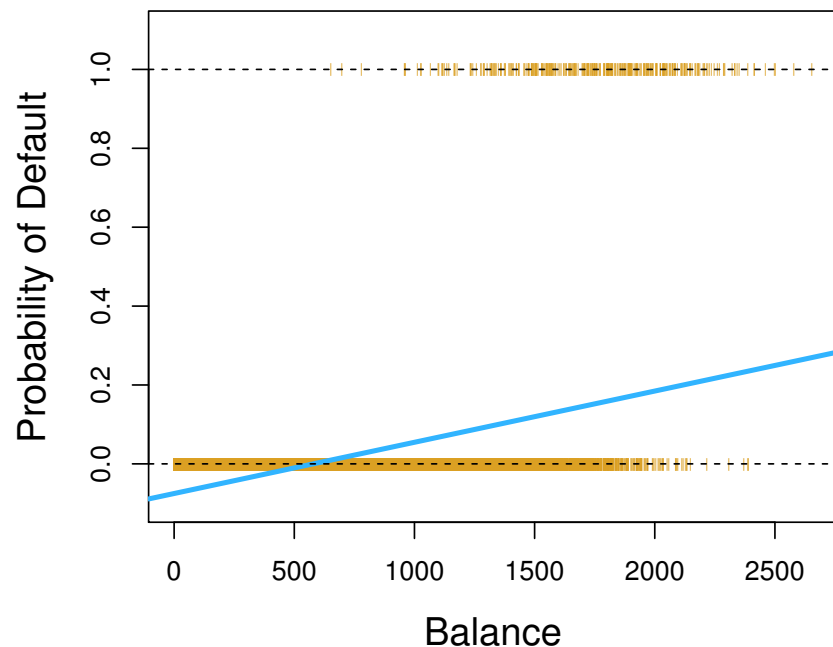
A bit of rearrangement gives

$$\log \left( \frac{p(X)}{1 - p(X)} \right) = \beta_0 + \beta_1 X.$$

probabilidad de fraude No pague  
Log Prob.  
probabilidad que no hay fraude. No pague

This monotone transformation is called the log odds or *logit* transformation of  $p(X)$ . (by log we mean *natural log*:  $\ln$ .)

# Linear versus Logistic Regression



Logistic regression ensures that our estimate for  $p(X)$  lies between 0 and 1.

# Maximum Likelihood

verosimilitud

We use maximum likelihood to estimate the parameters.

$$\ell(\beta_0, \beta) = \prod_{i:y_i=1} p(x_i) \prod_{i:y_i=0} (1 - p(x_i)).$$

This *likelihood* gives the probability of the observed zeros and ones in the data. We pick  $\beta_0$  and  $\beta_1$  to maximize the likelihood of the observed data.

3 pacientes  
2 pacientes  
No van a la escuela  
1 paciente sí.

$$p(x; \theta)$$

$$\text{Max } P(x; \theta)$$

$$\text{Max } L(\theta | x)$$

Probabilidad  
que va  
a la  
escuela

$$L(\beta_0, \beta) = \theta (1 - \theta)^2$$

$\uparrow$  Func  
 $\uparrow$  No Func



# Maximum Likelihood

We use maximum likelihood to estimate the parameters.

$$\ell(\beta_0, \beta) = \prod_{i:y_i=1} p(x_i) \prod_{i:y_i=0} (1 - p(x_i)).$$

This *likelihood* gives the probability of the observed zeros and ones in the data. We pick  $\beta_0$  and  $\beta_1$  to maximize the likelihood of the observed data.

Most statistical packages can fit linear logistic regression models by maximum likelihood. In **R** we use the **glm** function.

	Coefficient	Std. Error	Z-statistic	P-value
<b>Intercept</b>	-10.6513	0.3612	-29.5	< 0.0001
<b>balance</b>	0.0055	0.0002	24.9	< 0.0001



# Making Predictions

What is our estimated probability of **default** for someone with a balance of \$1000?

$$\hat{p}(X) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 X}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 X}} = \frac{e^{-10.6513 + 0.0055 \times 1000}}{1 + e^{-10.6513 + 0.0055 \times 1000}} = \underline{0.006}$$

# Making Predictions

What is our estimated probability of **default** for someone with a balance of \$1000?

$$\hat{p}(X) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 X}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 X}} = \frac{e^{-10.6513 + 0.0055 \times 1000}}{1 + e^{-10.6513 + 0.0055 \times 1000}} = \underline{0.006}$$

With a balance of \$2000?

$$\hat{p}(X) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 X}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 X}} = \frac{e^{-10.6513 + 0.0055 \times 2000}}{1 + e^{-10.6513 + 0.0055 \times 2000}} = \underline{0.586}$$

Lets do it again, using **student** as the predictor.

	Coefficient	Std. Error	Z-statistic	P-value
<b>Intercept</b>	-3.5041	0.0707	-49.55	< 0.0001
<b>student [Yes]</b>	0.4049	0.1150	3.52	0.0004

Lets do it again, using **student** as the predictor.

	Coefficient	Std. Error	Z-statistic	P-value
<b>Intercept</b>	-3.5041	0.0707	-49.55	< 0.0001
<b>student [Yes]</b>	0.4049	0.1150	3.52	0.0004

$$\widehat{\Pr}(\text{default}=\text{Yes}|\text{student}=\text{Yes}) = \frac{e^{-3.5041+0.4049 \times 1}}{1 + e^{-3.5041+0.4049 \times 1}} = \underline{0.0431},$$

$$\widehat{\Pr}(\text{default}=\text{Yes}|\text{student}=\text{No}) = \frac{e^{-3.5041+0.4049 \times 0}}{1 + e^{-3.5041+0.4049 \times 0}} = \underline{0.0292}.$$

# Logistic regression with several variables

$$\log \left( \frac{p(X)}{1 - p(X)} \right) = \underbrace{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}$$

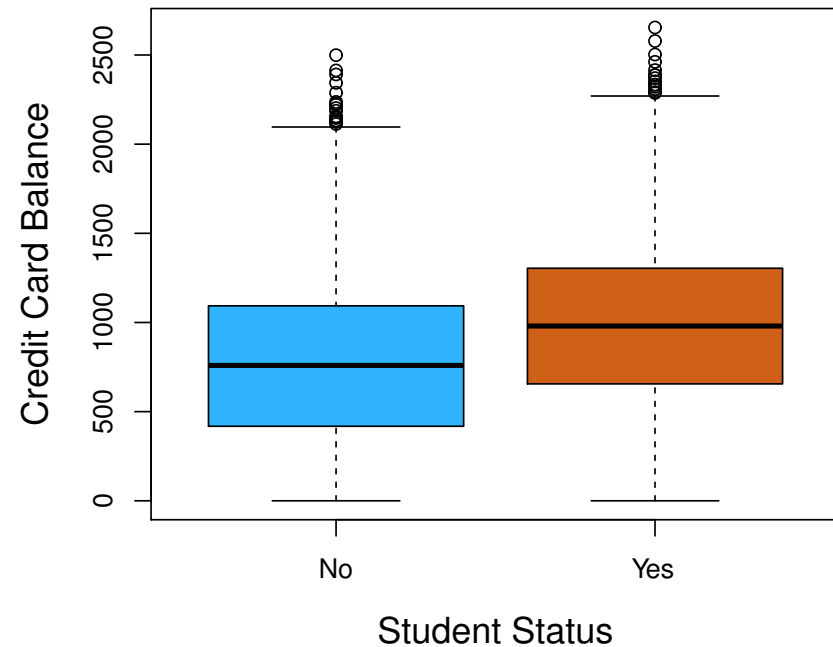
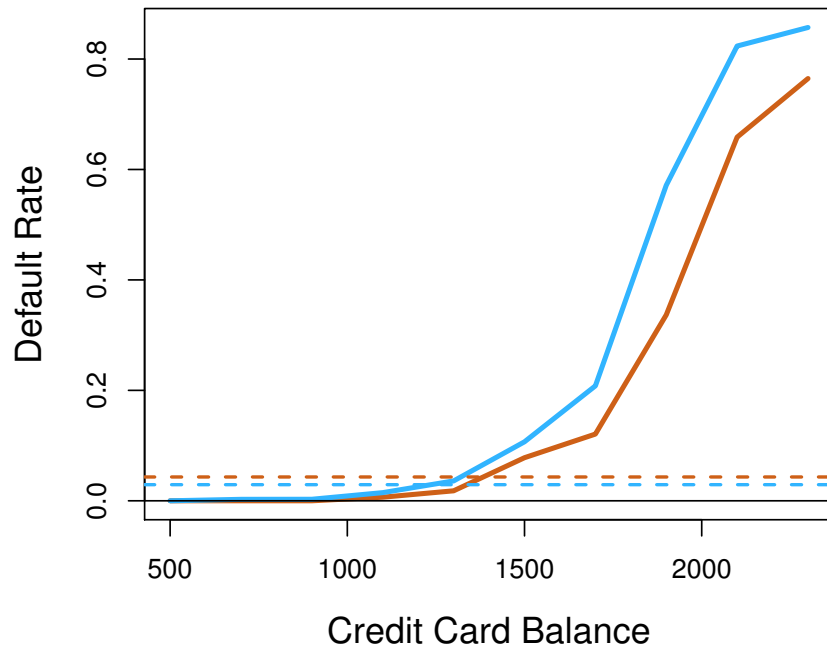
$$\rightarrow p(X) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

	Coefficient	Std. Error	Z-statistic	P-value
Intercept	-10.8690	0.4923	-22.08	< 0.0001
balance	0.0057	0.0002	24.74	< 0.0001
income	0.0030	0.0082	0.37	0.7115
student [Yes]	<u>-0.6468</u>	0.2362	-2.74	0.0062

Why is coefficient for **student** negative, while it was positive before?

1, Heyne 6/24/25

## Confounding



- Students tend to have higher balances than non-students, so their marginal default rate is higher than for non-students.
- But for each level of balance, students default less than non-students.
- Multiple logistic regression can tease this out.

## Example: South African Heart Disease

- 160 cases of MI (myocardial infarction) and 302 controls (all male in age range 15-64), from Western Cape, South Africa in early 80s.
- Overall prevalence very high in this region: 5.1%.
- Measurements on seven predictors (risk factors), shown in scatterplot matrix.
- Goal is to identify relative strengths and directions of risk factors.
- This was part of an intervention study aimed at educating the public on healthier diets.



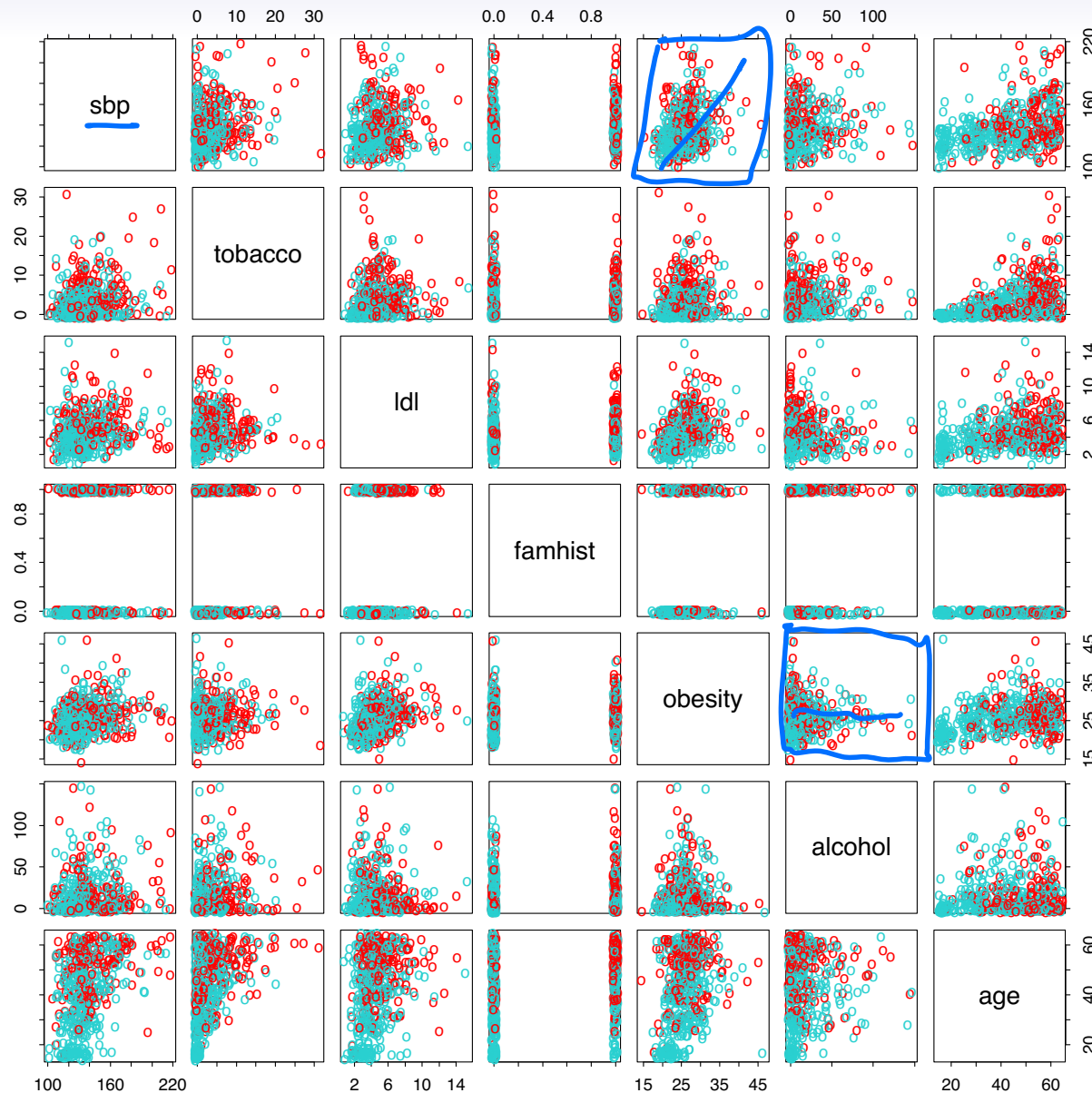


Gráfico de dispersão.

Scatterplot matrix of the *South African Heart Disease* data. The response is color coded — The cases (MI) are red, the controls turquoise. **famhist** is a binary variable, with 1 indicating family history of MI.

```
> heartfit <- glm(chd ~ ., data = heart, family = binomial)
> summary(heartfit)
```

Call:

```
glm(formula = chd ~ ., family = binomial, data = heart)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-4.1295997	0.9641558	-4.283	1.84e-05	***
sbp	0.0057607	0.0056326	1.023	0.30643	
tobacco	0.0795256	0.0262150	3.034	0.00242	**
ldl	0.1847793	0.0574115	3.219	0.00129	**
famhistPresent	0.9391855	0.2248691	4.177	2.96e-05	***
obesity	-0.0345434	0.0291053	-1.187	0.23529	
alcohol	0.0006065	0.0044550	0.136	0.89171	
age	0.0425412	0.0101749	4.181	2.90e-05	***

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 596.11 on 461 degrees of freedom  
 Residual deviance: 483.17 on 454 degrees of freedom  
 AIC: 499.17

Cual variable afecta el ataque al corazón mas

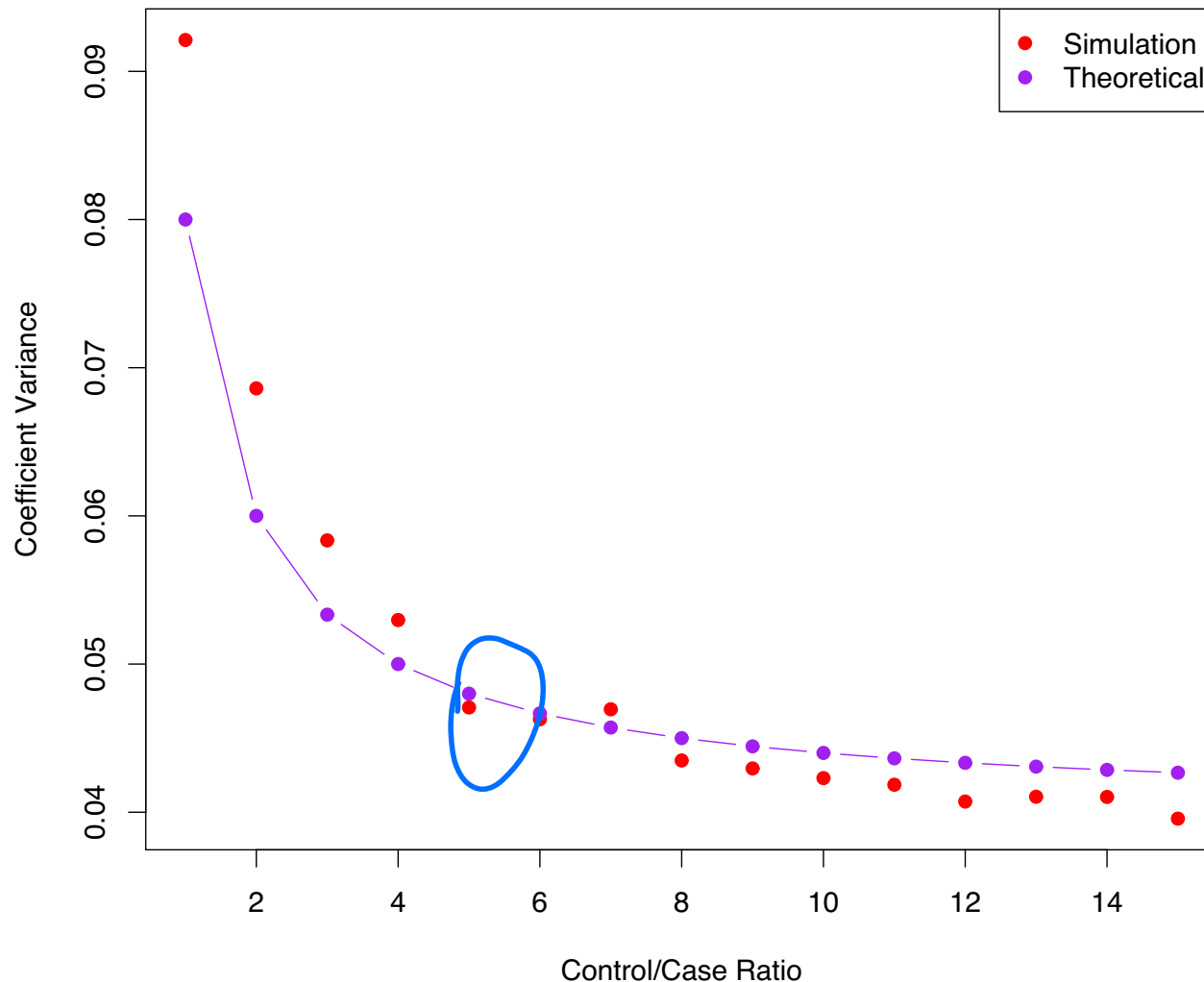
# Case-control sampling and logistic regression

- In South African data, there are 160 cases, 302 controls —  $\tilde{\pi} = 0.35$  are cases. Yet the prevalence of MI in this region is  $\pi = 0.05$ .
- With case-control samples, we can estimate the regression parameters  $\beta_j$  accurately (if our model is correct); the constant term  $\beta_0$  is incorrect.
- We can correct the estimated intercept by a simple transformation

$$\hat{\beta}_0^* = \hat{\beta}_0 + \log \frac{\pi}{1 - \pi} - \log \frac{\tilde{\pi}}{1 - \tilde{\pi}}$$

- Often cases are rare and we take them all; up to five times that number of controls is sufficient. See next frame

# Diminishing returns in unbalanced binary data




Sampling more controls than cases reduces the variance of the parameter estimates. But after a ratio of about 5 to 1 the variance reduction flattens out.

# Logistic regression with more than two classes

So far we have discussed logistic regression with two classes. It is easily generalized to more than two classes. One version (used in the R package **glmnet**) has the symmetric form

$$\Pr(Y = k|X) = \frac{e^{\beta_{0k} + \beta_{1k}X_1 + \dots + \beta_{pk}X_p}}{\sum_{\ell=1}^K e^{\beta_{0\ell} + \beta_{1\ell}X_1 + \dots + \beta_{p\ell}X_p}}$$

Here there is a linear function for *each* class.



# Logistic regression with more than two classes

So far we have discussed logistic regression with two classes. It is easily generalized to more than two classes. One version (used in the R package **glmnet**) has the symmetric form

$$\Pr(Y = k|X) = \frac{e^{\beta_{0k} + \beta_{1k}X_1 + \dots + \beta_{pk}X_p}}{\sum_{\ell=1}^K e^{\beta_{0\ell} + \beta_{1\ell}X_1 + \dots + \beta_{p\ell}X_p}}$$

Here there is a linear function for *each* class.

(The *mathier* students will recognize that some cancellation is possible, and only  $K - 1$  linear functions are needed as in 2-class logistic regression.)

# Logistic regression with more than two classes

So far we have discussed logistic regression with two classes. It is easily generalized to more than two classes. One version (used in the R package **glmnet**) has the symmetric form

$$\Pr(Y = k|X) = \frac{e^{\beta_{0k} + \beta_{1k}X_1 + \dots + \beta_{pk}X_p}}{\sum_{\ell=1}^K e^{\beta_{0\ell} + \beta_{1\ell}X_1 + \dots + \beta_{p\ell}X_p}}$$

Here there is a linear function for *each* class.

(The *mathier* students will recognize that some cancellation is possible, and only  $K - 1$  linear functions are needed as in 2-class logistic regression.)

Multiclass logistic regression is also referred to as *multinomial regression*.

# Discriminant Analysis

Here the approach is to model the distribution of  $X$  in each of the classes separately, and then use *Bayes theorem* to flip things around and obtain  $\Pr(Y|X)$ .

When we use normal (Gaussian) distributions for each class, this leads to linear or quadratic discriminant analysis.

However, this approach is quite general, and other distributions can be used as well. We will focus on normal distributions.



# Bayes theorem for classification

Thomas Bayes was a famous mathematician whose name represents a big subfield of statistical and probabilistic modeling. Here we focus on a simple result, known as Bayes theorem:

$$\Pr(Y = k | X = x) = \frac{\Pr(X = x | Y = k) \cdot \Pr(Y = k)}{\Pr(X = x)}$$

↕  
clase  
Emergencia  
o No

↕  
Datos

↕ Prob. de los variables

Prob. de los variables  
↕  
Si va a la emergencia

↕  
Prevalencia de las clases.

Como evaluamos  
el LDA.

## [LDA on Credit Data]

		<i>True Default Status</i>		
		No	Yes	Total
<i>Predicted Default Status</i>	No	9644	252	9896
	Yes	23	81	104
Total		9667	333	10000

Matrix de confusión.

$(23 + 252)/10000$  errors — a 2.75% misclassification rate!

Some caveats:

- This is *training* error, and we may be overfitting.

# LDA on Credit Data

		<i>True Default Status</i>		
		No	Yes	Total
<i>Predicted Default Status</i>	No	9644	252	9896
	Yes	23	81	104
Total		9667	333	10000

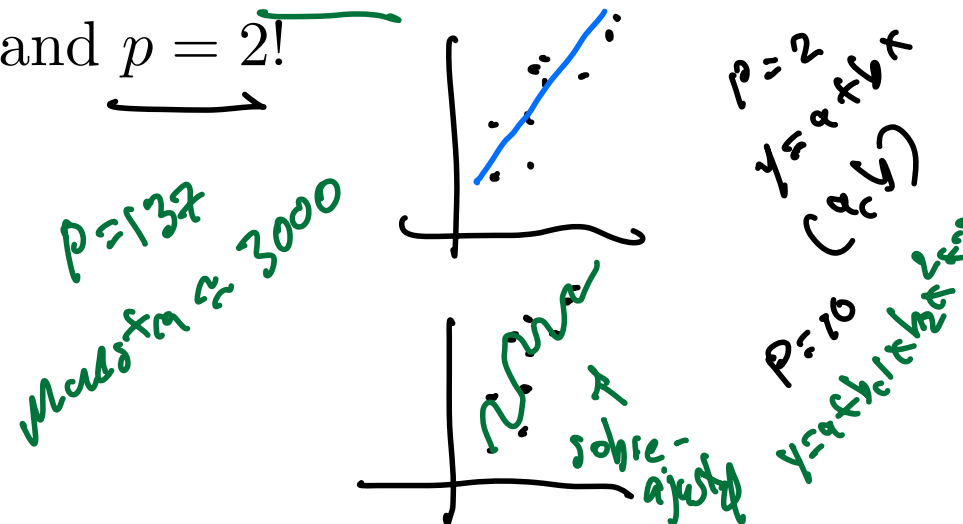
$(23 + 252)/10000$  errors — a 2.75% misclassification rate!

Some caveats:

- This is *training* error, and we may be overfitting. Not a big concern here since  $n = 10000$  and  $p = 2$ !

Numero de variables      pequeno

Tamaño de set de entrenamiento      grande



## LDA on Credit Data

		<i>True Default Status</i>		
		No	Yes	Total
<i>Predicted Default Status</i>	No	9644	252	9896
	Yes	23	81	104
Total		9667	333	10000

$(23 + 252)/10000$  errors — a 2.75% misclassification rate!

Some caveats:

- This is *training* error, and we may be overfitting. Not a big concern here since  $n = 10000$  and  $p = 2$ !
- If we classified to the prior — always to class **No** in this case — we would make  $333/10000$  errors, or only 3.33%.

## LDA on Credit Data

		<i>True Default Status</i>		
		No	Yes	Total
<i>Predicted Default Status</i>	No	9644	252	9896
	Yes	23	81	104
Total		9667	333	10000

$(23 + 252)/10000$  errors — a 2.75% misclassification rate!

Some caveats:

- This is *training* error, and we may be overfitting. Not a big concern here since  $n = 10000$  and  $p = 2$ !
- If we classified to the prior — always to class **No** in this case — we would make  $333/10000$  errors, or only 3.33%.
- Of the true **No**'s, we make  $23/9667 = 0.2\%$  errors; of the true **Yes**'s, we make  $252/333 = 75.7\%$  errors!

# Types of errors

**False positive rate:** The fraction of negative examples that are classified as positive — 0.2% in example.

**False negative rate:** The fraction of positive examples that are classified as negative — 75.7% in example.

We produced this table by classifying to class **Yes** if

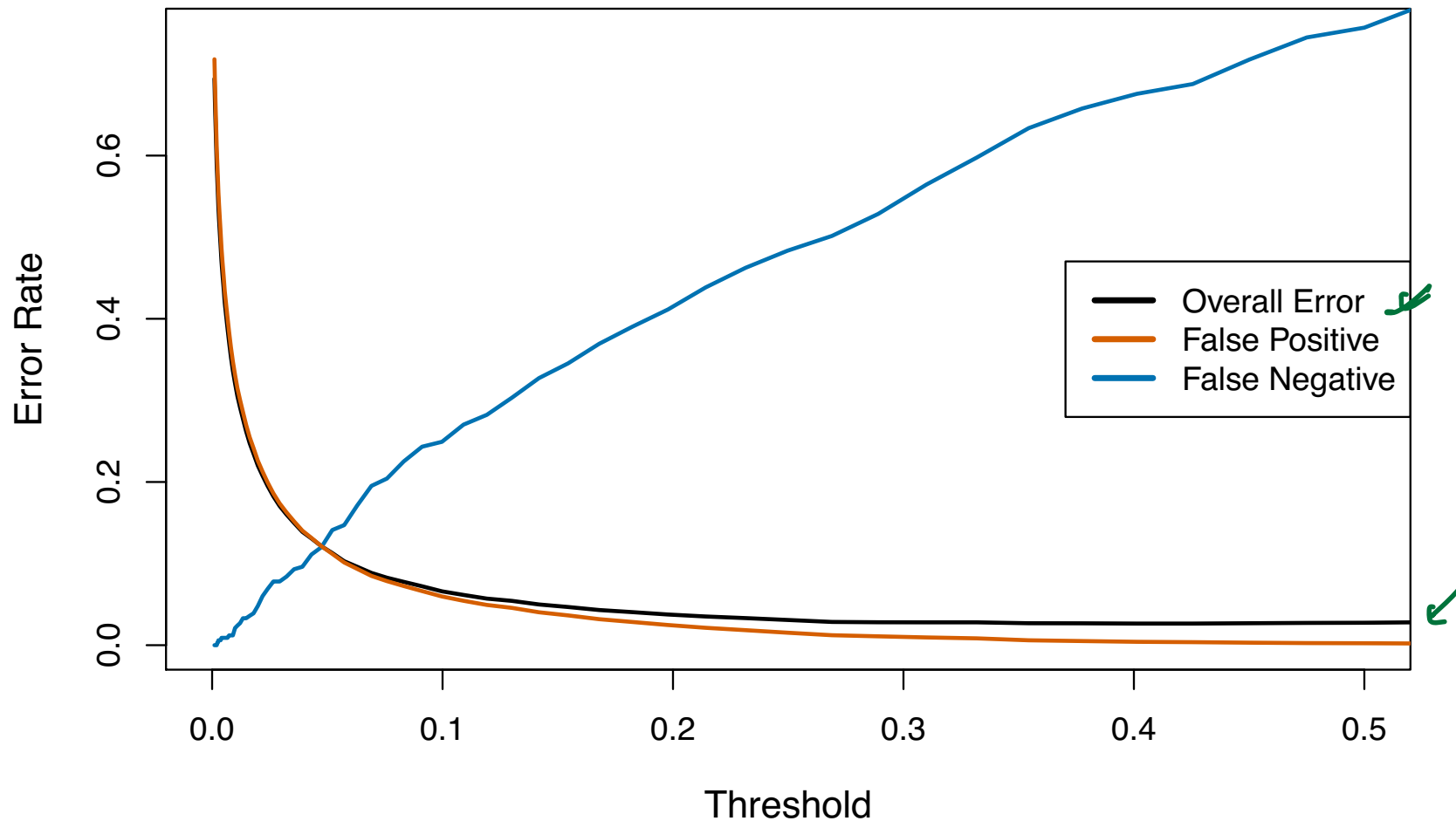
$$\widehat{\Pr}(\text{Default} = \text{Yes} | \text{Balance}, \text{Student}) \geq 0.5$$

We can change the two error rates by changing the threshold from 0.5 to some other value in  $[0, 1]$ :

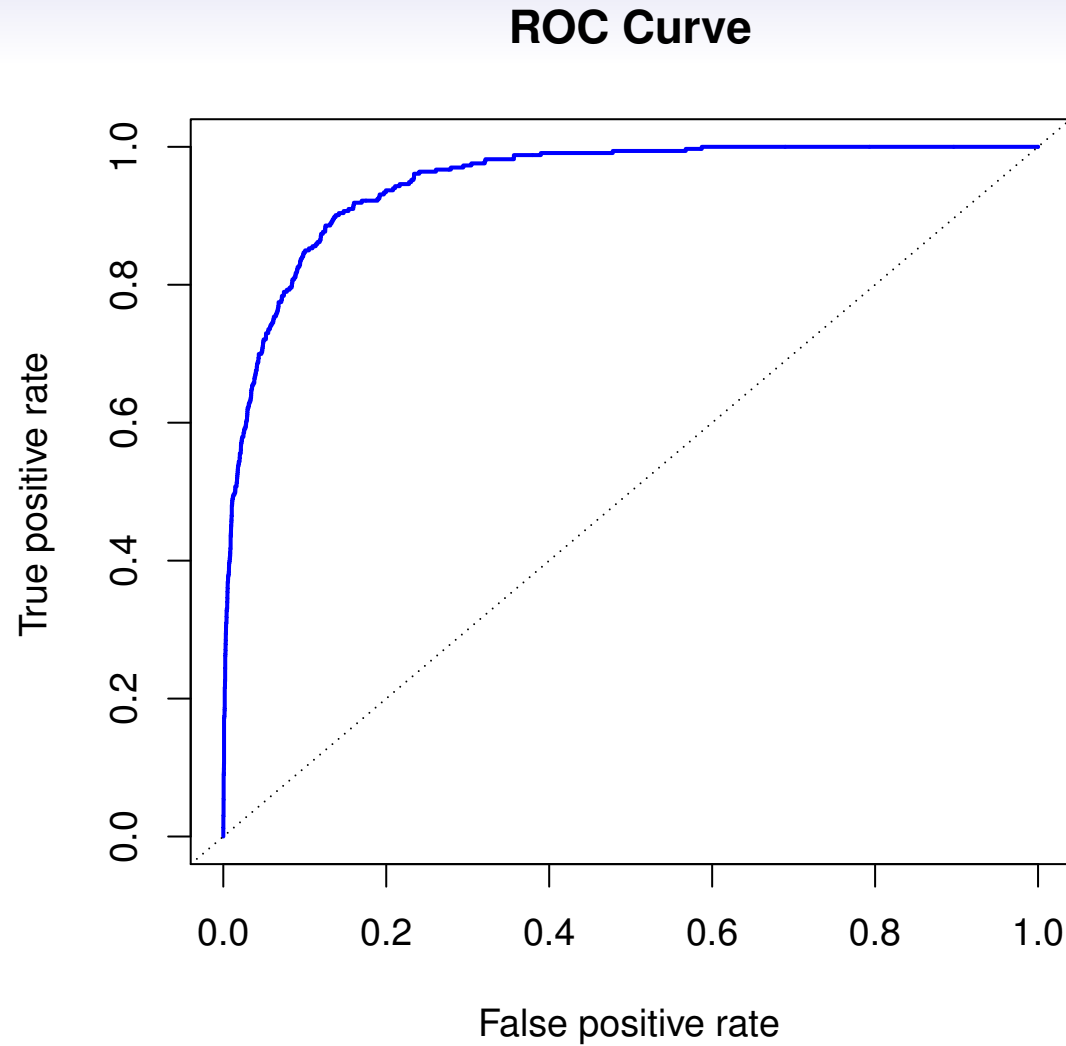

$$\widehat{\Pr}(\text{Default} = \text{Yes} | \text{Balance}, \text{Student}) \geq \textit{threshold},$$

and vary *threshold*.

## Varying the *threshold*

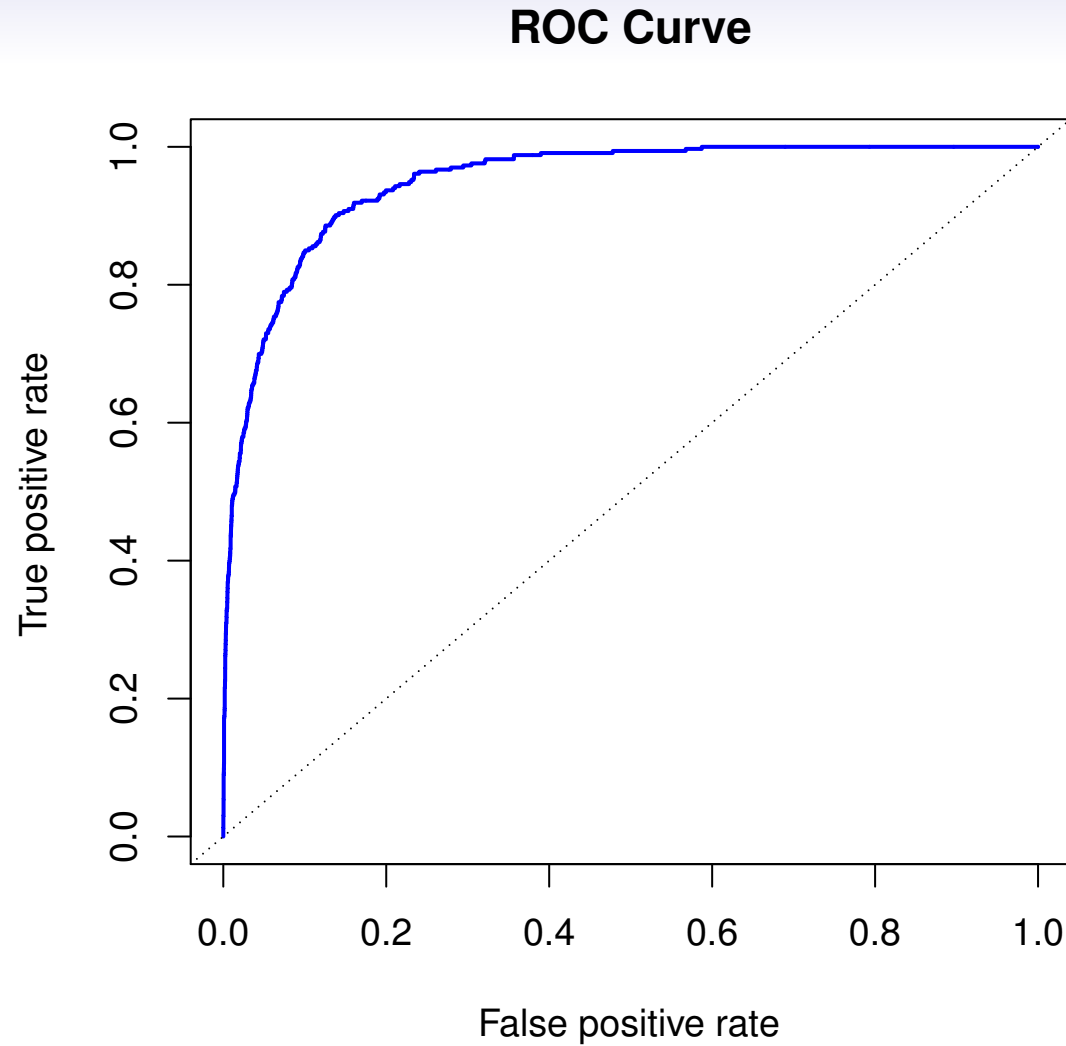


In order to reduce the false negative rate, we may want to reduce the threshold to 0.1 or less.



The *ROC plot* displays both simultaneously.





The *ROC plot* displays both simultaneously.

Sometimes we use the *AUC* or *area under the curve* to summarize the overall performance. Higher *AUC* is good.

## Logistic Regression versus LDA

For a two-class problem, one can show that for LDA

$$\log \left( \frac{p_1(x)}{1 - p_1(x)} \right) = \log \left( \frac{p_1(x)}{p_2(x)} \right) = \underbrace{c_0 + c_1 x_1 + \dots + c_p x_p}$$

So it has the same form as logistic regression.

The difference is in how the parameters are estimated.

# Logistic Regression versus LDA

For a two-class problem, one can show that for LDA

$$\log \left( \frac{p_1(x)}{1 - p_1(x)} \right) = \log \left( \frac{p_1(x)}{p_2(x)} \right) = c_0 + c_1 x_1 + \dots + c_p x_p$$

So it has the same form as logistic regression.

The difference is in how the parameters are estimated.

- Logistic regression uses the conditional likelihood based on  $\Pr(Y|X)$  (known as *discriminative learning*).

# Logistic Regression versus LDA

For a two-class problem, one can show that for LDA

$$\log \left( \frac{p_1(x)}{1 - p_1(x)} \right) = \log \left( \frac{p_1(x)}{p_2(x)} \right) = c_0 + c_1 x_1 + \dots + c_p x_p$$

So it has the same form as logistic regression.

The difference is in how the parameters are estimated.

- Logistic regression uses the conditional likelihood based on  $\Pr(Y|X)$  (known as *discriminative learning*).
- LDA uses the full likelihood based on  $\Pr(X, Y)$  (known as *generative learning*).

# Logistic Regression versus LDA

For a two-class problem, one can show that for LDA

$$\log \left( \frac{p_1(x)}{1 - p_1(x)} \right) = \log \left( \frac{p_1(x)}{p_2(x)} \right) = c_0 + c_1 x_1 + \dots + c_p x_p$$

So it has the same form as logistic regression.

The difference is in how the parameters are estimated.

- Logistic regression uses the conditional likelihood based on  $\Pr(Y|X)$  (known as *discriminative learning*).
- LDA uses the full likelihood based on  $\Pr(X, Y)$  (known as *generative learning*).
- Despite these differences, in practice the results are often very similar.

# Logistic Regression versus LDA

For a two-class problem, one can show that for LDA

$$\log \left( \frac{p_1(x)}{1 - p_1(x)} \right) = \log \left( \frac{p_1(x)}{p_2(x)} \right) = c_0 + c_1 x_1 + \dots + c_p x_p$$

So it has the same form as logistic regression.

The difference is in how the parameters are estimated.

- Logistic regression uses the conditional likelihood based on  $\Pr(Y|X)$  (known as *discriminative learning*).
- LDA uses the full likelihood based on  $\Pr(X, Y)$  (known as *generative learning*).
- Despite these differences, in practice the results are often very similar.

Footnote: logistic regression can also fit quadratic boundaries like QDA, by explicitly including quadratic terms in the model.

# Multinomial Logistic Regression

Logistic regression is frequently used when the response is binary, or  $K = 2$  classes. We need a modification when there are  $K > 2$  classes. E.g. **stroke**, **drug overdose** and **epileptic seizure** for the emergency room example.

The simplest representation uses different linear functions for each class, combined with the *softmax* function to form probabilities:

$$\Pr(Y = k | X = x) = \frac{e^{\beta_{k0} + \beta_{k1}x_1 + \dots + \beta_{kp}x_p}}{\sum_{l=1}^K e^{\beta_{l0} + \beta_{l1}x_1 + \dots + \beta_{lp}x_p}}.$$

# Multinomial Logistic Regression

Logistic regression is frequently used when the response is binary, or  $K = 2$  classes. We need a modification when there are  $K > 2$  classes. E.g. **stroke**, **drug overdose** and **epileptic seizure** for the emergency room example.

The simplest representation uses different linear functions for each class, combined with the *softmax* function to form probabilities:

$$\Pr(Y = k | X = x) = \frac{e^{\beta_{k0} + \beta_{k1}x_1 + \dots + \beta_{kp}x_p}}{\sum_{l=1}^K e^{\beta_{l0} + \beta_{l1}x_1 + \dots + \beta_{lp}x_p}}. \quad \downarrow$$

- There is a redundancy here; we really only need  $K - 1$  functions (see the book for details).



# Multinomial Logistic Regression

Logistic regression is frequently used when the response is binary, or  $K = 2$  classes. We need a modification when there are  $K > 2$  classes. E.g. **stroke**, **drug overdose** and **epileptic seizure** for the emergency room example.

The simplest representation uses different linear functions for each class, combined with the *softmax* function to form probabilities:

$$\Pr(Y = k | X = x) = \frac{e^{\beta_{k0} + \beta_{k1}x_1 + \dots + \beta_{kp}x_p}}{\sum_{l=1}^K e^{\beta_{l0} + \beta_{l1}x_1 + \dots + \beta_{lp}x_p}}.$$

- There is a redundancy here; we really only need  $K - 1$  functions (see the book for details).
- We fit by maximizing the *multinomial* log likelihood (cross-entropy) — a generalization of the binomial.

# Multinomial Logistic Regression

Logistic regression is frequently used when the response is binary, or  $K = 2$  classes. We need a modification when there are  $K > 2$  classes. E.g. **stroke**, **drug overdose** and **epileptic seizure** for the emergency room example.

The simplest representation uses different linear functions for each class, combined with the *softmax* function to form probabilities:

$$\Pr(Y = k | X = x) = \frac{e^{\beta_{k0} + \beta_{k1}x_1 + \dots + \beta_{kp}x_p}}{\sum_{l=1}^K e^{\beta_{l0} + \beta_{l1}x_1 + \dots + \beta_{lp}x_p}}.$$

- There is a redundancy here; we really only need  $K - 1$  functions (see the book for details).
- We fit by maximizing the *multinomial* log likelihood (cross-entropy) — a generalization of the binomial.
- An example is given in Chapter 10 where we fit the 10-class model to the **MNIST** digit dataset.