

Multiple Logistic Regression

Previously...

Last week, we considered logistic regression with a single quantitative explanatory variable:

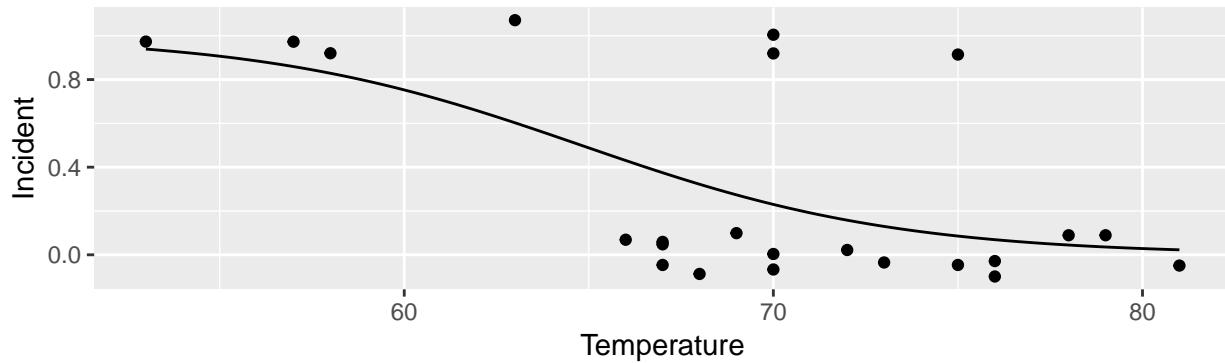
$$Y_i = \begin{cases} 1 & \text{if obs. } i \text{ is in a particular class} \\ 0 & \text{otherwise} \end{cases}$$

X_i = value of quantitative explanatory variable for observation number i

Our model used the following functional form:

$$P(Y_i = 1|X_i) = p(X_i) = \frac{e^{\beta_0 + \beta_1 X_i}}{1 + e^{\beta_0 + \beta_1 X_i}}$$

... and gave us estimated curves like...

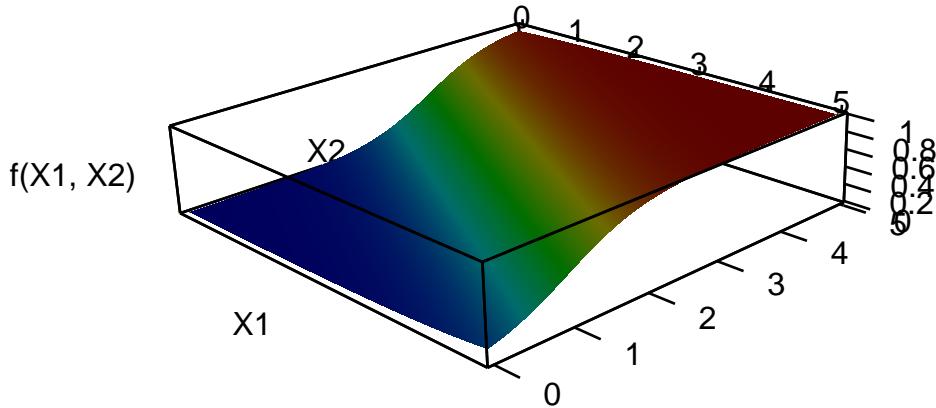


Set Up for Multiple Explanatory Variables

We will now extend this to allow for p explanatory variables which may be either quantitative or categorical.

$$P(Y_i = 1|X_{i1}, \dots, X_{ip}) = p(X_{i1}, \dots, X_{ip}) = \frac{e^{\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}}}{1 + e^{\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}}}$$

Illustration with $p = 2$ explanatory variables:



Running Example

This example is adapted from section 4.3 of ISLR. The ISLR package provides the `Default` data set, which contains information on ten thousand customers; our goal is to predict which customers will default on their credit card debt.

```
library(ggplot2)
library(gridExtra)
library(dplyr)
library(ISLR)

head(Default)

##   default student  balance    income
## 1     No      No 729.5265 44361.625
## 2     No     Yes 817.1804 12106.135
## 3     No      No 1073.5492 31767.139
## 4     No      No  529.2506 35704.494
## 5     No      No  785.6559 38463.496
## 6     No     Yes 919.5885  7491.559
```

Example 1: Two Quantitative Variables

Let's try using `balance` and `income` as explanatory variables.

```
fit <- glm(default ~ balance + income, data = Default, family = binomial)
summary(fit)

##
## Call:
## glm(formula = default ~ balance + income, family = binomial,
##      data = Default)
##
## Deviance Residuals:
##    Min      1Q  Median      3Q     Max
## -2.4725 -0.1444 -0.0574 -0.0211  3.7245
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) -1.154e+01  4.348e-01 -26.545 < 2e-16 ***
## balance      5.647e-03  2.274e-04  24.836 < 2e-16 ***
## income       2.081e-05  4.985e-06   4.174 2.99e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 2920.6 on 9999 degrees of freedom
## Residual deviance: 1579.0 on 9997 degrees of freedom
## AIC: 1585
##
## Number of Fisher Scoring iterations: 8
```

(a) What is the estimated equation for this model?

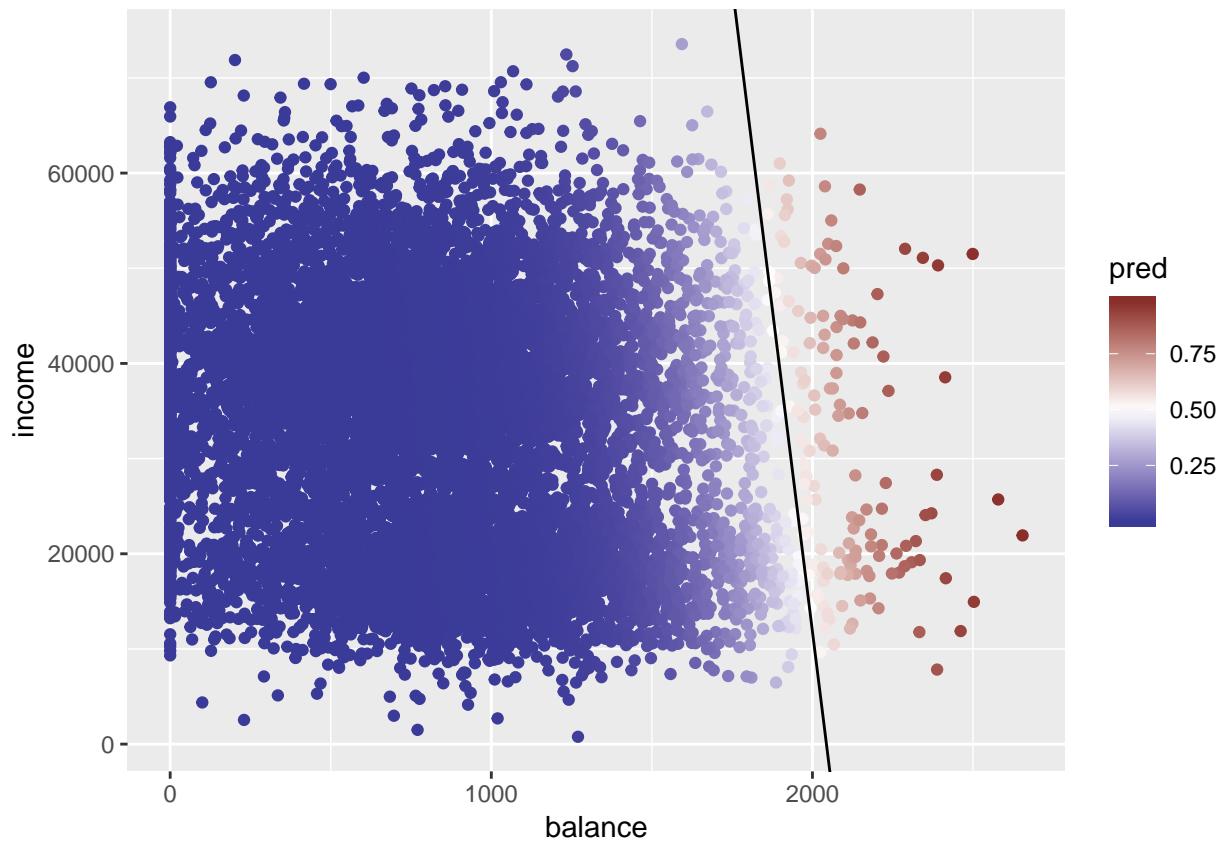
(b) What is the decision boundary?

$$0.5 = P(Y_i = 1 | X_{i1}, \dots, X_{ip}) = p(X_{i1}, \dots, X_{ip}) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \hat{\beta}_2 X_{i2}}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \hat{\beta}_2 X_{i2}}}$$

Plots

```
df2 <- Default %>%
  mutate(
    pred = predict(fit, type = "response")
  )

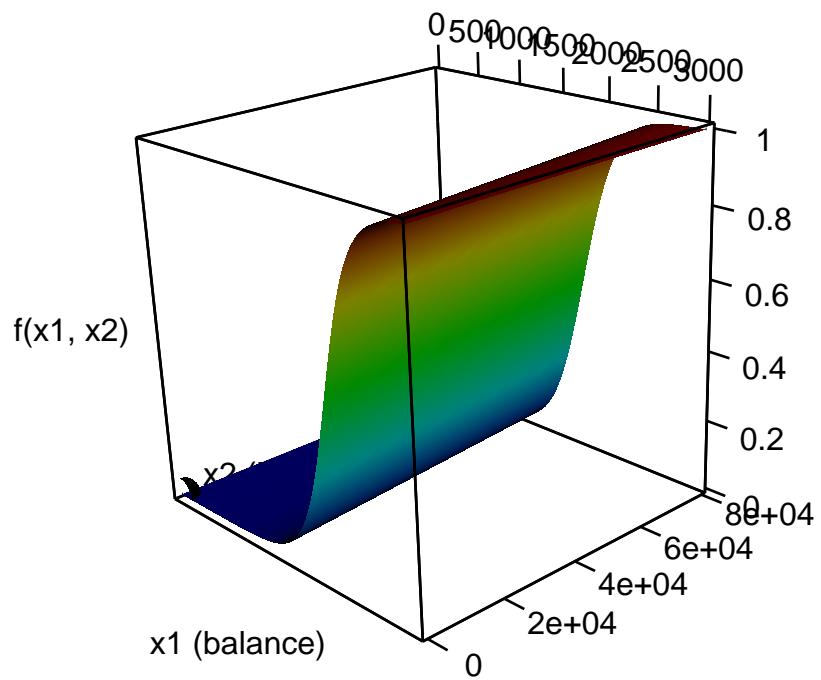
ggplot(data = df2, mapping = aes(x = balance, y = income, color = pred)) +
  geom_point() +
  geom_abline(intercept = 1.154e+01 / 2.081e-05, slope = - 5.647e-03 / 2.081e-05) +
  scale_color_gradient2(low = scales::muted("blue"), high = scales::muted("red"), midpoint = 0.5)
```



```
max_balance <- max(Default$balance)
max_income <- max(Default$income)
background <- expand.grid(
  balance = seq(from = 0, to = max_balance, length = 101),
  income = seq(from = 0, to = max_income, length = 101))
background <- background %>%
  mutate(
    est_prob_default = predict(fit, newdata = background, type = "response"),
    est_default = ifelse(est_prob_default > 0.5, "Yes", "No")
  )

ggplot() +
  geom_point(data = background,
             mapping = aes(x = balance, y = income, color = est_default), size = 0.1, alpha = 0.5) +
  geom_point(data = Default, mapping = aes(x = balance, y = income, color = default)) +
  scale_color_discrete("Default") +
```

```
geom_abline(intercept = 1.154e+01 / 2.081e-05, slope = - 5.647e-03 / 2.081e-05)
```



Example 2: One Categorical Explanatory variable

```
fit <- glm(default ~ student, data = Default, family = binomial)
summary(fit)

##
## Call:
## glm(formula = default ~ student, family = binomial, data = Default)
##
## Deviance Residuals:
##    Min      1Q  Median      3Q     Max 
## -0.2970 -0.2970 -0.2434 -0.2434  2.6585 
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)    
## (Intercept) -3.50413   0.07071 -49.55 < 2e-16 ***
## studentYes   0.40489   0.11502   3.52 0.000431 ***  
## ---        
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 2920.6 on 9999 degrees of freedom
## Residual deviance: 2908.7 on 9998 degrees of freedom
## AIC: 2912.7
##
## Number of Fisher Scoring iterations: 6
```

Similar to the use of categorical explanatory variables in linear models, R has created a new indicator variable for use in the regression:

$$X_{i1} = \text{studentYes}_i = \begin{cases} 1 & \text{if customer } i \text{ is a student} \\ 0 & \text{otherwise} \end{cases}$$

- (a) What is the estimated equation for this model?

(b) What is the predicted probability of default for a non-student?

```
predict(fit, newdata = data.frame(student = "No"), type = "response")

##           1
## 0.02919501
# compare to...
exp(-3.50413) / (1 + exp(-3.50413))

## [1] 0.02919495
Default %>%
  filter(student == "No") %>%
  summarize(prop_default = mean(default == "Yes"))

##   prop_default
## 1 0.02919501
```

(c) What is the predicted probability of default for a student?

```
predict(fit, newdata = data.frame(student = "Yes"), type = "response")

##           1
## 0.04313859
# compare to...
exp(-3.50413 + 0.40489) / (1 + exp(-3.50413 + 0.40489))

## [1] 0.04313862
Default %>%
  filter(student == "Yes") %>%
  summarize(prop_default = mean(default == "Yes"))

##   prop_default
## 1 0.04313859
```

(d) Does someone's student status have a statistically significant association with whether or not they default?

Note about decision boundaries

- In this example, our predicted class is 0 for all values of x_i !
- In general, a decision boundary need not exist; this often happens with categorical explanatory variables.

Example 3: All 3 Explanatory Variables

```
fit <- glm(default ~ student + balance + income, data = Default, family = binomial)
summary(fit)

##
## Call:
## glm(formula = default ~ student + balance + income, family = binomial,
##      data = Default)
##
## Deviance Residuals:
##    Min      1Q  Median      3Q     Max
## -2.4691 -0.1418 -0.0557 -0.0203  3.7383
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) -1.087e+01  4.923e-01 -22.080 < 2e-16 ***
## studentYes   -6.468e-01  2.363e-01  -2.738  0.00619 **
## balance      5.737e-03  2.319e-04  24.738 < 2e-16 ***
## income       3.033e-06  8.203e-06   0.370  0.71152
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 2920.6 on 9999 degrees of freedom
## Residual deviance: 1571.5 on 9996 degrees of freedom
## AIC: 1579.5
##
## Number of Fisher Scoring iterations: 8
```

(a) What is the estimated equation for this model?

(b) Does an individual's student status have a statistically significant association with whether or not they default? Compare your estimate to that from example 2.

(c) Does an individual's income have a statistically significant association with whether or not they default? Compare to your result from example 1.

(d) What is the estimated equation for non-students?

(e) What is the estimated equation for students?

Plots

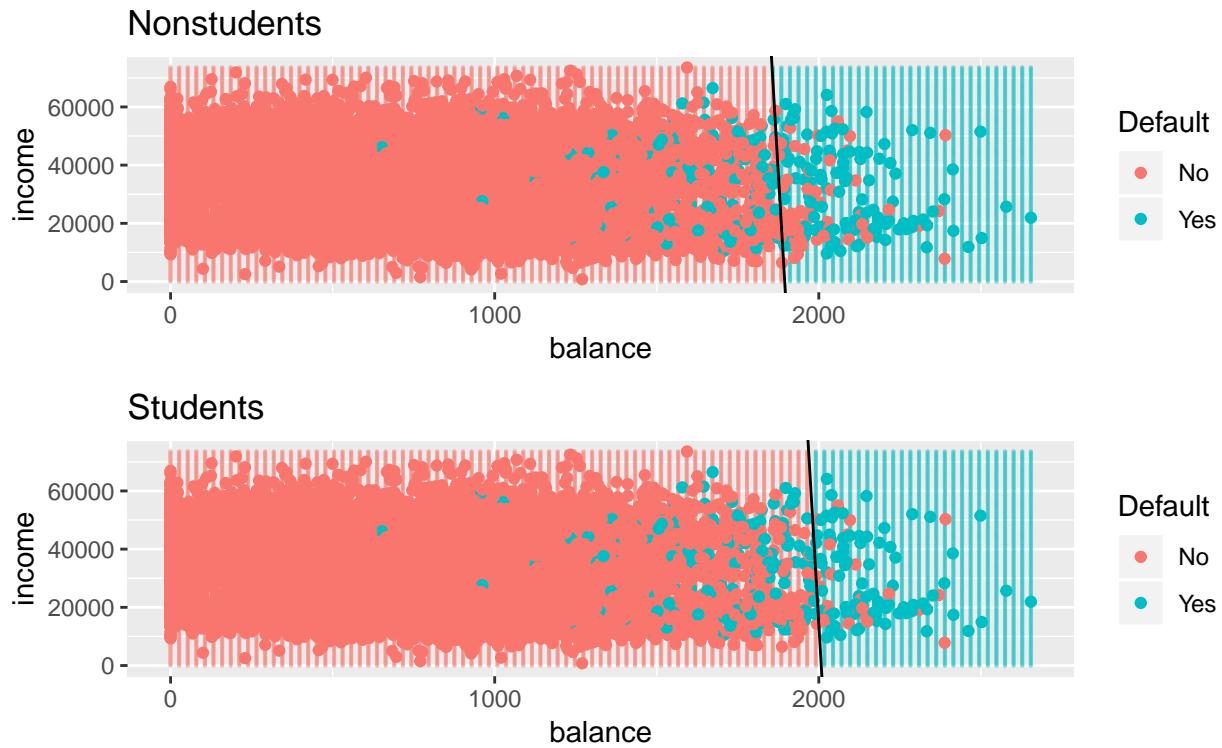
```
background <- expand.grid(
  balance = seq(from = 0, to = max_balance, length = 101),
  income = seq(from = 0, to = max_income, length = 101))
background_nonstudents <- background %>%
  mutate(student = "No")
background_students <- background %>%
  mutate(student = "Yes")

background_nonstudents <- background_nonstudents %>%
  mutate(
    est_prob_default = predict(fit, newdata = background_nonstudents, type = "response"),
    est_default = ifelse(est_prob_default > 0.5, "Yes", "No")
  )
background_students <- background_students %>%
  mutate(
    est_prob_default = predict(fit, newdata = background_students, type = "response"),
    est_default = ifelse(est_prob_default > 0.5, "Yes", "No")
  )

p_nonstudents <- ggplot() +
  geom_point(data = background_nonstudents,
    mapping = aes(x = balance, y = income, color = est_default), size = 0.1, alpha = 0.5) +
  geom_point(data = Default, mapping = aes(x = balance, y = income, color = default)) +
  scale_color_discrete("Default") +
  geom_abline(intercept = 1.087e+01 / 3.033e-06, slope = - 5.737e-03 / 3.033e-06) +
  ggtitle("Nonstudents")

p_students <- ggplot() +
  geom_point(data = background_students,
    mapping = aes(x = balance, y = income, color = est_default), size = 0.1, alpha = 0.5) +
  geom_point(data = Default, mapping = aes(x = balance, y = income, color = default)) +
  scale_color_discrete("Default") +
  geom_abline(intercept = (1.087e+01 + 6.468e-01) / 3.033e-06, slope = - 5.737e-03 / 3.033e-06) +
  ggtitle("Students")

grid.arrange(p_nonstudents, p_students, ncol = 1)
```



Ethical Considerations

In the U.S., there is a history of discrimination against demographic groups in granting loans. The Equal Credit Opportunity Act was passed in 1974, and “makes it illegal to consider makes it unlawful for any creditor to discriminate against any applicant, with respect to any aspect of a credit transaction, on the basis of race, color, religion, national origin, sex, marital status, or age (provided the applicant has the capacity to contract)” (https://en.wikipedia.org/wiki/Equal_Credit_Opportunity_Act).

Our model uses the covariates `balance`, `income`, and `student` to predict probability of loan default, which are allowed by the law. However, it's a fact that some of the covariates in our model (like `balance` and `income`) are correlated with protected characteristics like race, sex, or marital status. At a population level, the model we've developed would deem women and people of color creditworthy at lower rates than other groups.