



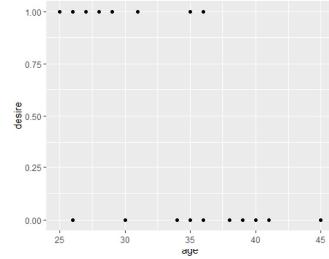
Data 603: Statistical Modelling with Data

Logistic Regression

Part I : Introduction to the Logistic Regression Model

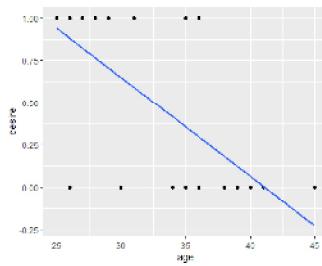
Example: The `desire` data show the distribution of 24 currently married and fecund women interviewed in the Fiji Fertility Survey, according to age, education, desire for more children (wife's perception of husband's desire for additional children). The data are provided in `desire.xlsx` file

```
X1 = age (year)  
X2 = education (0=none, 1=some),  
Y = desire for more children (0=no more, 1=more),  
library("readxl")  
desire <- read_excel("c:/Users/thunida.ngamkham/OneDrive - University of Cal  
gary/dataset603/desire.xlsx")  
library(ggplot2)  
ggplot(data = desire, mapping = aes(x = age, y = desire))+  
geom_point()
```

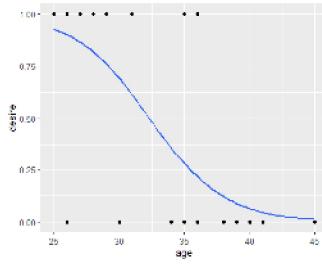


Using linear regression model

```
library("readxl")  
library(ggplot2)  
desire <- read_excel("c:/Users/thunida.ngamkham/OneDrive - University of Cal  
gary/dataset603/desire.xlsx")  
ggplot(data = desire, mapping = aes(x = age, y = desire))+geom_point()  
geom_smooth(method=lm, se=F)
```



```
library("ggplot2")
library("gridExtra")
desire <- read.csv("~/Users/lumtik/ugent/biostatistics - University of Ghent/lectures/ed/les_06/excel")
ggplot(desire, mapping = aes(x = age, y = desire)) + geom_point() +
  stat_smooth(method = "glm", formula = list(family = "binomial"), se = FALSE)
```



The linear regression model discussed in Multiple Regression assumes that the response variable y is **quantitative**. But in many situations, the response variable is instead **qualitative**. For example, eye color is qualitative, taking qualitative values blue, brown, or green. Often qualitative variables are referred to as categorical.

What distinguishes a 'logistic regression model' from the linear regression model is that the outcome variable's in logistic regression is binary or **dichotomous**.

In this topic, we study approaches for predicting qualitative responses, a process that is known as classification. Predicting a qualitative response for an observation can be referred to as classifying that observation, since it involves assigning the observation to a category, or class. On the other hand, often the methods used for classification first predict the probability instead of the categories of a qualitative variable, so the basis for making the classification. In this topic we discuss simple Logistic Regression, and Multiple Logistic Regression for a qualitative binary response.

What is Logistic Regression ?

Logistic Regression seeks to:

1. **Model the probability** of an event occurring depending on the value of the independent variables, which can be categorical or numerical.
2. **Estimate the probability** that an event occurs for a randomly selected observation versus the probability that the event does not occur.
3. **Predict** the effect of a series of variables on a binary response variable.
4. **Classify** observations by estimating the probability that an observation is in a particular category.

Why Not Linear Regression

1. **Linear regression assumptions**: The linear regression model is based on an assumption that the response y is continuous, with errors are normally distributed. If the response variable is binary, this assumption is clearly violated, and so in general, we might expect our inferences to be inaccurate.
 2. **Predicted values may be out of range**: For a binary outcome, the means is the probability of a 1, or success. If we use linear regression to model a binary outcome it is entirely possible to have a linear regression which gives predicted values for some input data which are outside of the [0, 1] range in practice.
- We will also illustrate the concept of classification using the simulated Default data set. We are interested in predicting whether an individual will default on his/her credit card payment on the basis of annual income and monthly credit card balance. In the Default data set, the response `default` has one of two categories, Yes or No. Rather than modelling the response y directly, a logistic regression model's the probability that y belongs to a particular category.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \epsilon$$

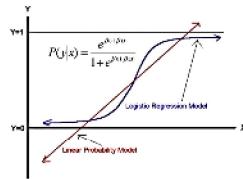


Figure 1 Comparing Graphs between simple linear regression and logistic regression

$$P(Y=1 | X) = \frac{e^{B_0 + B_1 X}}{1 + e^{B_0 + B_1 X}}$$

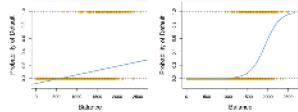


Figure 2 Classification using Default data for simple linear regression and logistic regression
Left figure: Estimated probability of default using linear regression. Since estimated probabilities are negative! The orange ticks indicate the 0/1 values coded for default (No or Yes).

Right figure: Predicted probabilities of default using logistic regression. All probabilities lie between 0 and 1.

For the Default data, logistic regression models the probability of default. For example, the probability of default given balance can be written as

$$\Pr(\text{default} = \text{Yes} | \text{balance})$$

where
 balance = the independent variable
 default = the response variable which is a binary outcome (Yes/No)

The values of $\Pr(\text{default} = \text{Yes} | \text{balance})$ will range between 0 and 1. Then for any given value of balance, a predictor can be made for default .

For example, if one might predict $\text{default} = \text{Yes}$ for any individual for whom $\Pr(\text{default} = \text{Yes} | \text{balance}) = 0.5$. Alternatively, if a company wishes to be conservative in predicting individuals who are at risk for default, then they may choose to use a lower threshold, such as $\Pr(\text{default} = \text{Yes} | \text{balance}) > 0.1$.

The Logistic Model

How should we model the relationship between $p(X) = \Pr(Y = 1|X)$ and X ? For convenience we are using the generic y to denote the response. If we use the approach $p(X) = p_0 + p_1 X$ to predict $\text{default} = \text{Yes}$ using balance, then we obtain the model shown in the left-hand panel of Figure 2.1.

Here we see the problem with this approach: for balances close to zero we predict a negative probability of default. If we were to predict for very large balances, we would get values bigger than 1. These predictions are nonsensical, since of course the true probability of default (regardless of credit card balance) must fall between 0 and 1. This problem is not unique to the credit default data. Any time a straight line is fit to a binary response that's coded as 0 or 1, in principle we can always predict $p(X) < 0$ for some values of X , and $p(X) > 1$ for others (unless the range of X is limited). To avoid this problem, we must model $p(X)$ using a function that gives outputs between 0 and 1 for all values of X . Many functions meet this description.

Simple Logistic Regression Model for a Binary Dependent Variable (Quantitative independent variable)

$$P(y) = P(y = 1|X) = \frac{e^{B_0 + B_1 X}}{1 + e^{B_0 + B_1 X}}$$

where
 $y = 1$ if category A occurs
 $y = 0$ if category B occurs
 $P(y) = P(\text{Category A occurs}) = \pi$

Note that the general logistic model is not a linear function of the β parameters, obtaining the parameter estimate of a nonlinear regression model, such as the logistic model, is a

dataset	y	x
	1	25
	0	50
	1	95
	0	15

$$\log(P(X=x)) = (\beta_0 + \beta_1 x)$$

numerically tedious process and often requires sophisticated computer programs. We use a method called **maximum likelihood estimation** to estimate the β_1 parameter.

The right-hand panel of Figure 2.1 illustrates the fit of the logistic regression model to the Default data. Notice that for low balances we now predict the probability of default as close to 0, but never below zero. Likewise, for high balances we predict a initial probability close to 1, but never above one.

The logistic function will always produce an S-shaped curve of this form, and regardless of the value of X , we will obtain a sensible prediction. We also see that the logistic model is better able to capture the range of probabilities than is the linear regression model.



Figure 2.1

Focusing on the single predictor case, if the parameter $\beta_1 > 0$ then the basic logistic regression model assumes that the probability of success is a monotonically increasing function of X . That is, the probability never decreases as X gets larger, it stays the same or increases. If the parameter $\beta_1 < 0$ the reverse is true.

Figure 5(a) shows a regression line where $\beta_1 = 1$ and $\beta_0 = 0.5$. As is evident, curvature is allowed and predicted probabilities always have a value between 0 and 1. Figure 5(b) shows the regression line when $\beta_1 = -1$ and $\beta_0 = 0.5$. So now, the regression line is monotonically decreasing. The predicted probability never increases.

$$P(\text{success}) = 0.8 \quad \text{the odds} = \frac{P(\text{success})}{P(\text{failure})}$$

$$P(\text{failure}) = 0.2 \quad \frac{P(\text{success})}{P(\text{failure})} = 4$$

The overall odds of winning a lottery prize
are 1 in 13

$$\therefore P(\text{winning}) = \frac{1}{13} \quad \text{odds} = \frac{1}{12}$$

Considering the logistic Regression Model, we find that:

$$\frac{P(y=1|x)}{P(y=0|x)} = \frac{P(y=1|x)}{1-P(y=1|x)}$$

$$= \frac{\pi}{1-\pi}$$

$$= \frac{e^{\beta_0 + \beta_1 x}}{1-e^{\beta_0 + \beta_1 x}}$$

The quantity $\frac{P(y=1|x)}{P(y=0|x)}$ is called the odds, and can take any value between 0 and ∞ .

What is the odds?

To appreciate the logistic model, it's helpful to have an understanding of odds. Most people expect odds to mean the "probabilistic inverse" of success or an event occurring. We automatically think in terms of numbers ranging from 0 to 1, with 0 meaning that the event will certainly not occur and 1 meaning that the event certainly will occur. But there are other ways of representing the chances of events, one of which the odds has a nicely equal claim to being "natural".

For example,

An odds of 4 means we expect 4 times as many occurrences as non-occurrences.

An odds of 1/5 means that we expect only one-fifth as many occurrences as non-occurrences.

$$\text{In general:} \quad \text{Odds} = \frac{P(y=1|x)}{1-P(y=1|x)} = \frac{\text{probability of event}}{\text{probability of no event}}$$

$$\text{The odds (model)} = e^{\beta_0 + \beta_1 x}$$

$$\frac{P(y=1|x)}{P(y=0|x)} = \frac{e^{\beta_0 + \beta_1 x}}{1-e^{\beta_0 + \beta_1 x}}$$

Considering the Logistic Regression Model, we find that:

$$\frac{P(y=1|x)}{P(y=0|x)} = \frac{P(y=1|x)}{\pi} = \frac{e^{\beta_0 + \beta_1 x}}{1 - e^{\beta_0 + \beta_1 x}}$$

The quantity $\frac{P}{1-P}$ is called the odds, and can take on any value between 0 and ∞ .

What is the odds?

To appreciate the "odds" model, it's helpful to have an understanding of odds. Most people regard probability as the "natural" way to quantify the chances that an event will occur. We automatically think in terms of numbers ranging from 0 to 1, with a 0 meaning that the event will certainly not occur and a 1 meaning that the event certainly will occur. But there are other ways of representing the chances of an event, one of which—the odds—has a nearly equal claim to being "natural".

For example,

An odds of 4 means we expect 4 times as many occurrences as non occurrences.

An odds of 1/5 means that we expect only one-fifth as many occurrences as non occurrences.

$$\text{Odds} = \frac{P(y=1|x)}{1 - P(y=1|x)} = \frac{\text{probability of event}}{\text{probability of no event}}$$

Relationship between Odds and Probability

Odds	Probability
0.1	0.01
0.2	0.05
0.3	0.08
0.4	0.12
0.5	0.17
0.6	0.23
0.7	0.28
0.8	0.33
0.9	0.36

$$\left. \begin{array}{l} \text{odds} < 1 \\ \text{p(success)} < \text{p(failure)} \end{array} \right\} \quad p = 0.5, 1 - p = 0.5$$

$$\left. \begin{array}{l} \text{odds} > 1 \\ \text{p(success)} > \text{p(failure)} \end{array} \right\}$$

$$\text{The odds (model)} = e^{\beta_0 + \beta_1 x}$$

$$\frac{P(Y=1|x)}{P(Y=0|x)} = \frac{e^{\beta_0 + \beta_1 x}}{1 - e^{\beta_0 + \beta_1 x}}$$

$$1 - \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

Values of the odds close to 0 and ∞ indicate very low and very high probabilities of default, respectively.

In general,

If an odds > 1, then the probability of success is higher than failure.

If an odds < 1, then the probability of success is lower than failure.

If an odds = 1, then the probability of success is equal to failure.

For more example about default data, on average 1 in 8 people with an odds of 1/4 will default, since $p(y=1|x)=0.2$ implies an odds of $\frac{0.2}{1-0.2} = 1/4$.

On average 1 in 5 out of every 10 people with an odds of 1 will default, since $p(y=1|x)=0.9$ implies an odds of $\frac{0.9}{1-0.9} = 9$.

By taking the natural logarithm of both sides of the logit model, we arrive at:

$$\begin{aligned} \ln\left(\frac{P(y=1|x)}{1 - P(y=1|x)}\right) &= \ln\left(\frac{P(y=1|x)}{1 - P(y=1|x)}\right) \\ &= \ln\left(\frac{P(y=1|x)}{d}\right) \\ &= \ln\left(\frac{P(y=1|x)}{1-d}\right) \\ &= \ln\left(\frac{P(y=1|x)}{1 - P(y=1|x)}\right) = \beta_0 + \beta_1 x \end{aligned}$$

The \ln is called the log odds or the logit.

This transformation is useful, because it creates a variable with a range from $-\infty$ to $+\infty$. Hence, this transformation solves the problem we encountered in fitting a linear model, to predict the probability of default, because the dependent variable only ranges from 0 to 1. We can get linear predictions and are outside of this range. If we transform linear probabilities to logits, then we do not have this problem because the range of the logit is not restricted. In addition, the interpretation of logits is simple: take the exponential of the logit and you have the odds for the two groups in question.

Note: R function to calculate the analog for a confidence interval for β_1 is provided below.

```
library(TSIR) # for Default data set
summary(Default)
## Default student balance income
## No: 1560 Yes: 1994 Min. : 0.0 Min. : 722
## 1st Qu.: 481.7 1st Qu.: 21342
## Median : 625.0 Median : 34617
## Mean : 858.8 Mean : 34517
## 3rd Qu.: 1164.3 3rd Qu.: 73998
## Max. : 2044.3 Max. : 77504
myLogit <- glm(default ~ balance, data = Default, family = "binomial")
summary(myLogit)
##
```

```
(Intercept) 1.18415e+00 4.58439e-02
## balance 1.002692e-06 1.02951e-06
```

The odds .

$$\text{odds}_1 = 1$$

$$\text{odds}_2 = 1 \times 1.0055 = 1.0055$$

$$\text{odds} = e^{5-2X} = 0.7$$

$$\text{Income} = 10,000$$

$$\text{odds}_1 = 0.7 \quad \frac{1}{10}$$

$$\text{Income} = 10,000 \quad \hat{\beta}_1 \\ \text{odds}_2 = 0.7 \times e^{\hat{\beta}_1} = 0.7 \times 0.8 = 0.56$$

$$\left| \frac{0.56 - 0.7}{0.7} \right| = |-0.2| \text{ decrease} \\ = 0.2 \cdot 100\% = 20\%$$

Inclass Practice Problem

Example: The claims data, showing the distribution of 270 currently married and found women interviewed in the FHI Fertility Survey, according to age, education, desire for more children, the data are provided in desire.xlsx file.

X1= age (year)

X2= education (0=none, 1=some).

Y= desire for more children (0=no more, 1=more).

a) Fit the Logistic Regression Model to predict the probability of desire for more children using age.

b) Construct a 95% confidence interval for the logit model.

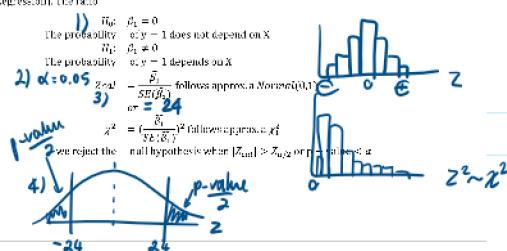
For every 1 year increase in age, we estimate the odds of desire to be multiplied by about 0.703 (and decrease of (1-0.703)=100%-29.5% of the odds).

Testing For The Significance of The Coefficients

After estimating the coefficients, there are several steps involved in assessing the appropriateness, adequacy and usefulness of the model. **Firstly**, the importance of each of the explanatory variables is assessed by carrying out statistical tests of the significance of the coefficients. **Secondly**, the overall goodness of fit of the model is then tested. **Lastly**, the model fit or the ability of the model to discriminate between the two groups defined by the response variable is evaluated.

Significance of The Coefficients

The Wald Z test or the Wald χ^2 test is the test of significance for individual regression coefficients in logistic regression. (Note that we used *t-test* for Simple and Multiple Linear Regression). The ratio



```
library(TSIR) # for Default data set
library(wald) # for wald.test
summary(Default)
```

```
## Default student balance income
## No: 1560 Yes: 1994 Min. : 0.0 Min. : 722
## 1st Qu.: 481.7 1st Qu.: 21342
## Median : 625.0 Median : 34617
## Mean : 858.8 Mean : 34517
## 3rd Qu.: 1164.3 3rd Qu.: 73998
## Max. : 2044.3 Max. : 77504
```

```
myLogit <- glm(default ~ balance, data = Default, family = "binomial")
## The Wald Z test
summary(myLogit)
```

```
## Call:
## glm(formula = default ~ balance, family = "binomial", data = Default)
## Deviance Residuals:
```

```
## Min. 1Q Median 3Q Max
## -2.2087 -0.1492 -0.0589 -0.9221 3.7598
```

```
## Coefficients:
```

```
## Estimate Std. Error z value Pr(>|z|)
## (Intercept) -1.18415 0.04826 -24.95 ***
```

```
## balance 1.00269e-06 1.02951e-06 -22.59 ***
```

```
## Signif. codes: 0 '***' 0.001 '** 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
## (Dispersion parameter for binomial: family taken to be 1)
```

```
## Null deviance: 2948.0 on 9989 degrees of freedom
```

```
## Residual deviance: 2486.9 on 9998 degrees of freedom
```

```
## AIC: 1602.5
```

```
## Number of Fisher Scoring Iterations: 8
```

```
## Wald test statistics for full vs reduced model
```

```
wald.test(y = coefficient(myLogit), sigma = varcov(myLogit), terms = 2) # terms cells
```

```
# which terms in the model are to be tested.
```

```
## Wald test:
```

```
## -----
## ## Chi-squared test:
```

```
## X2 = 622.7, df = 1, p < 2.2e-16
```

```
R function
```

95% CI for the logit model.

$$0.0004 < \beta_1 < 0.00695$$

$$\ominus < \beta_1 < \oplus$$

$$Z_{\text{cal}} = 24.95$$

$$\chi^2_{\text{cal}} = (24.95)^2$$

$$= 622.7$$

$$y = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$y = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$y = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$y = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$y = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$y = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$y = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$y = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$y = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$y = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$y = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

`update(k = conf(...), lSymm = recov(...))`. Terms = ... it performs full or reduced model for the chi-square test.

Deviance: The deviance is negative two times the maximized log-likelihood; the smaller the deviance, the better the fit (see Part 2 for more details).

Note!

Terms tells R which terms in the model are to be tested, in the example, terms 2 represent testing for $\beta_1 = 0$.

From the default data, for the Wald Z test, $Z = -22.95$ with p-value < 0.001, we reject H₀, therefore the probability of default depends on balance. For the Wald χ^2 test, $\chi^2 = 522.7$ with p-value zero. We also conclude that the probability of default depends on balance.

Making Predictions

Once the coefficients have been estimated, it is a simple matter to compute the probability of default for any given credit card balance. For example, using the coefficient estimates given in the output, we predict that the default probability for an individual with a balance of \$1,000 is

$$\hat{P} = \frac{e^{0.611 + 0.0251 \times 1000}}{1 + e^{0.611 + 0.0251 \times 1000}} = 0.60879$$

which is below 50%. In contrast, the predicted probability of default for an individual with a balance of \$2,000 is much higher, and equals 0.686 or 56.6%.

$$P = \frac{e^{0.611 + 0.0251 \times 2000}}{1 + e^{0.611 + 0.0251 \times 2000}} = 0.686$$

Note! Using R function to calculate \hat{P} when balance = 1000 dollars

```
library(MASS) # for Default data set
mylogit <- glm(defaul ~ balance, data = Default, family = "binomial")
modeldata = data.frame(balance=1000)
predict(mylogit, newdata, type="response")
```

#

0.60879145

A confidence Interval for the odds e^{β_1}

Letting s_{β_1} denote the estimated standard error of $\hat{\beta}_1$, the estimate of β_1 , a 95% confidence interval for the change of the odds is

$e^{\hat{\beta}_1 \pm z_{\alpha/2} s_{\beta_1}}$

```
#>
#> Call:
#> glm(formula = default ~ balance, family = "binomial", data = Default)
#>
#> Deviance Residuals:
#> 0 10 median 20 max
#> -2.2697 -0.1455 -0.2589 -0.9221 3.7595
#>
#> Coefficients:
#> Estimate Std. Error z value Pr(>|z|)
#> (Intercept) -1.665e-01 3.61e-02 -45.93 <2e-16 ***
#> balance 1.099e-01 2.49e-02 44.95 <2e-16 ***
#> ...
#> Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' '
#>
#> Degrees of freedom: 2920 5 on 9489 degrees of freedom
#> Null deviance: 2920.5 on 9489 degrees of freedom
#> Residual deviance: 1396.5 on 9484 degrees of freedom
#> AIC: 1698.5
#>
#> Number of Fisher Scoring Iterations: 8
#>
#> lower.ci upper.ci
#> (Intercept) 2.00994e-02 1.16197e-02 4.88497e-02
#> balance 1.023514e+00 1.025982e+00 1.009594e-02
```

From the output, a 95% confidence interval for β_1 is 1.005 to 1.006 implies that there is a positive relationship between $balance$ and $default$ as this confidence interval doesn't cover 0 and it covers a range greater than 1. This means that an increasing in \$1 in balance increases the odds of $default$ between 0.5% and 0.6%.

odds = 1.0055
95% CI for the odds
1.00508 < odds < 1.005948

less $<$ odds $<$ more
than
1



Model Fit in Logistic Regression Model

In linear regression, R^2 is a very useful quantity, describing the fraction of the variability in the response that the explanatory variables can explain. There are a number of ways one can define an analog to R^2 in the logistic regression case, but none of them are as widely useful as R^2 in linear regression. To evaluate the performance of a logistic regression model, one would select, for example, AIC.



$$100c$$

$$800$$

$$\text{unexplained variation}$$

$$\Pi = \frac{e^{\beta_0}}{1 + e^{\beta_0}}$$

variation in a logistic model

1. Deviance \rightarrow a measure how much UNEXPLAINED variation in a logistic model

Deviance is a measure of goodness of fit of a generalized linear model (GLM). Or rather, R software reports the residual deviance.

- The null deviance and the residual deviance: The null deviance shows how well the response variable is predicted by a model that includes only the intercept.
- The residual deviance indicates how well the response is predicted by the model with independent variables.

For our example, we have a value of 2920.6 points on 9499 degrees of freedom. Including the independent variable (balance) decreased the deviance to 1396.5 points on 9494 degrees of freedom, a significant reduction in deviance.

The Residual Deviance has reduced by 1524.1 points with loss of one degrees of freedom.

2. AIC (Akaike Information Criteria) estimates the relative amount

For our example, we have a value of 2921.6 points in 6999 degrees of freedom, including the independent variable (balance) decreased the Deviance in 1,363.5 points in 998 degrees of freedom, a significant reduction in deviance.

The Residual Deviance has reduced by 1363.5 points with a loss of one degrees of freedom.

2. AIC (Akaike Information Criteria) estimates the relative amount of information lost by a given model.

The Akaike Information Criteria (AIC) provides a method for assessing the quality of your model through comparison of related models. It's based on Deviance. The smaller AIC, the more useful for comparing models, but it's uninterpretable on its own.

Note!

Ridge's scoring algorithm is a derivative of Newton's method for solving maximum likelihood problems numerically.

3. ROC curve

ROC stands for Receiver Operating Characteristic. To explain the ROC curve, we need to understand the important notions of sensitivity and specificity of a test or prediction rule.

Sensitivity and specificity

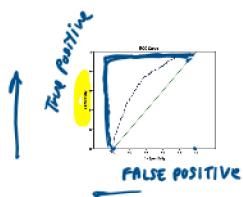
The sensitivity is defined as the probability of the prediction rule or model predicting an observation as 'positive' given that it is truth ($y=1$). In words, the sensitivity is the proportion of truly positive observations which is classified as such by this model or test. Conversely, the specificity is the probability of the model predicting 'negative' given that the observation is negative ($y=0$).

A model needs to not only correctly predict a positive as a positive, but also a negative as a negative.

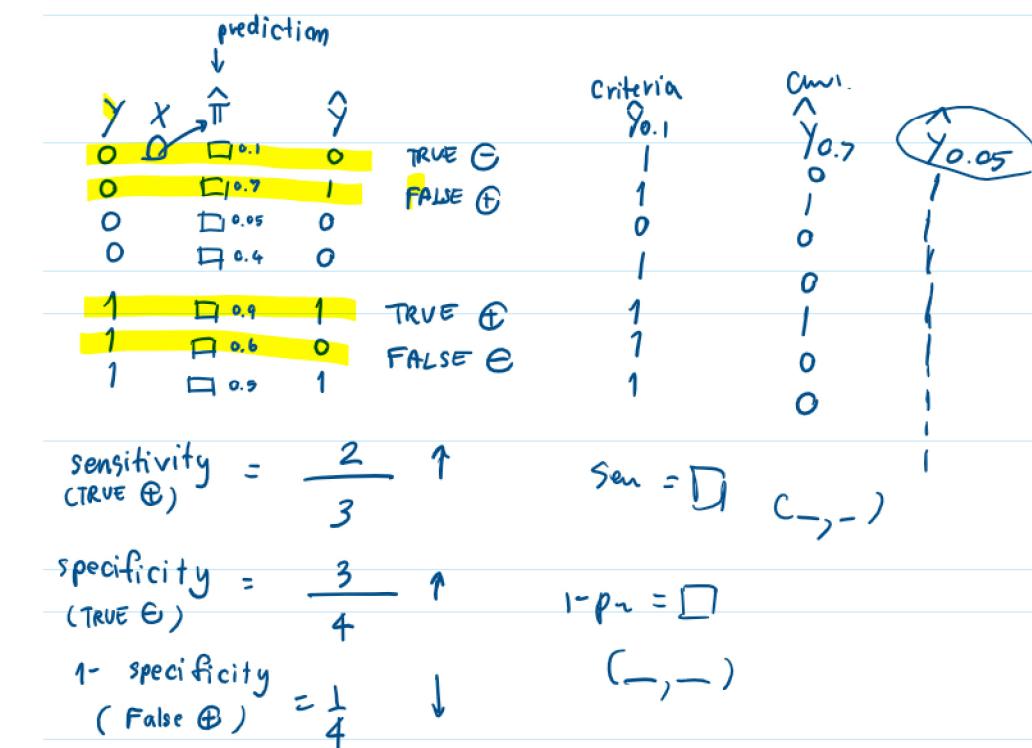
Our model or prediction rule is perfect at classifying observations if it has 100% sensitivity and 100% specificity. In a real study, this is (usually) not attainable. So how can we summarize the discriminative ability of our logistic regression model?

The ROC curve does this by plotting the true positive rate (sensitivity), the probability of predicting a real positive will be a positive, against false positive rate (1-specificity), the probability of predicting a real negative will be a positive. The best decision rule is high on sensitivity and low on 1-specificity. It's a rule that predicts most true positives will be positive and few true negatives will be positive.

How to explore the ROC curve



1. The further the curve is from the diagonal line, the better the model is at discriminating between positives and negatives in general.



2. There are useful statistics that can be calculated from this curve, the Area Under the Curve (AUC). This tells us how well the model predicts the probability of Y.

Inclass Pratice Problem

From the default data,

- write both the logistic regression model of Default on Income and the logit transformation of this logistic regression model.
- Interpret the logistic regression coefficient e^{β_1} in logistic model
- Test if the probability of default depends on Income at $\alpha = 0.05$
- Find a 95% Confidence Interval for the logistic regression coefficient e^{β_1}
- Use the method of Model Fit in Logistic Regression Model to evaluate the performance of a logistic regression model
- Predict the probability of default when Income= 60,000 dollars. Would you consider a person with \$60,000 income defaults on payment?

Inclass Practice Problem

Experience in hiring: Suppose you are investigating years of experience the hiring practices of a particular firm. Is there any sufficient evidence to indicate that the years of experience is an important predictor of hiring status? If yes, interpret the logistic regression coefficient e^{β_1} in logistic model. The data are provided in DISCRIM.csv file

- Write both the logistic regression model of Hire on Experience and the logit transformation of this logistic regression model.
- Interpret the logistic regression coefficient e^{β_1} in logistic model
- Test if the probability of hiring depend on experience at $\alpha = 0.05$
- Find a 95% Confidence interval for the odds ratio for the logistic regression coefficient e^{β_1}
- Using the method of Model Fit in Logistic Regression Model to evaluate the performance of a logistic regression model
- Predict the probability of hiring when experience= 1 years. Would you consider a person with a 7 years job experience will be hired?

$$a) \hat{P} = P(Y=1|Income)$$

$$= \frac{-3.094 - 0.000008959X}{1 + e^{-3.094 - 0.000008959X}}$$

$$\text{logit} = -3.094 - 0.000008959X$$

$$\text{the odds} = \frac{e^{\beta_1}}{e^{0}} = e^{-3.094 - 0.000008959X}$$

$$b) e^{\beta_1} = e^{-0.000008959} = 0.99999 / 65$$

$$d) 95\% \text{ CI for } e^{\beta_1}$$

$$0.99998340 < e^{\beta_1} < 0.99999989$$

e)

ROC

AIC

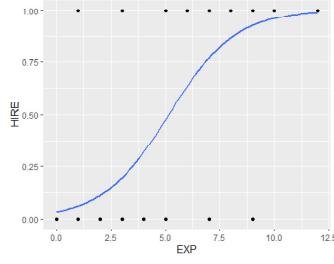
AUC = 0.5327

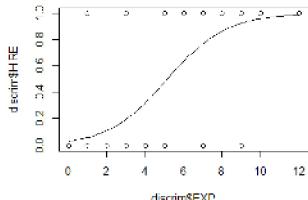
Visualizing the data and logistic regression model

The data and logistic regression model can be plotted with ggplot2 or base graphics:

```
library(ggplot2)
discrim=read.csv("c:/Users/thuntida.ngamkham/OneDrive - University of Calgary
/dataset603/DISCRIM.csv", header = TRUE)
mylogit<-glm(HIRE ~EXP, data = discrim, family = "binomial")
```

```
#option1 using ggplot function
ggplot(discrim, aes(x=EXP, y=HIRE)) + geom_point() +
stat_smooth(method="glm", method.args=list(family="binomial"), se=FALSE)
```





Simple Logistic Regression Model with a Qualitative Independent variable.

We can use qualitative predictors with the logistic regression model using the dummy variable approach as well. For example, the Default data set contains the qualitative variable student. To fit the model we simply create a dummy variable that takes on a value of 1 for students and 0 for non-students. The logistic regression model that results from predicting probability of default from student status can be seen in the output:

```
library(DAAG)
# for Default data set
summary(Default)
## #> #>   default    student    balance     income
## #> #>   No. 10662  Yes.10596  Min. : 0.0  Min. : 772
## #> #>   Min. : 0.0000  1st Qu.: 481.7  1st Qu.: 2124.8
## #> #>   Median : 025.0  Median : 343.0
## #> #>   Mean   : 035.4  Mean   : 331.27
## #> #>   3rd Qu.: 1168.3  3rd Qu.: 4388.0
## #> #>   Max. : 3265.3  Max. : 77554
## #>
## #> mylogit <- glm(default ~ student, data = Default, family = "binomial")
## #> summary(mylogit)
## #>
## #> Call:
## #> glm(formula = default ~ student, family = "binomial", data = Default)
```

```
## #> Deviance Residuals:
## #>   Min. 1Q Median 3Q Max
## #> -2.397 -0.2978 -0.2454  0.2424  2.6595
## #>
## #> Coefficients:
## #>   Estimate Std. Error z value Pr(>z)
## #> (Intercept) -0.0821  0.0781 -1.03 59 < 2e-16 ***
## #> studentYes  0.40498  0.11502  3.52 0.000421 ***
## #> ...
## #> Signif. codes: 0 '***' 0.001 '** 0.01 '*' 0.1 ' '
## #>
## #> Dispersion parameter for binomial: family taken to be 1
## #>
## #> Null deviance: 2928.9 on 9998 degrees of freedom
## #> Residual deviance: 2928.7 on 9998 degrees of freedom
## #> AIC: 2912.7
## #>
## #> Number of Fisher Scoring iterations: 6
```

From the output, the logit is

$$\hat{y} = -0.0821 + 0.40498x$$

and the estimated logistic regression model is

$$\hat{p} = \frac{e^{\hat{y}}}{1 + e^{\hat{y}}} = \frac{e^{-0.0821 + 0.40498x}}{1 + e^{-0.0821 + 0.40498x}}$$

Interpretations of Logistic Regression Coefficients in the Logistic Model for Qualitative independent variable

By computing e^{β_1} (awful of the coefficient), so we interpret in terms of the odds ratio.

Odd ratio if x_1 is the probability of default for students and x_2 is for non-students, the odds ratio for student nonstudent defined as $\frac{x_1}{x_2}(3-x_2)$

$$\frac{x_1}{x_2} = \frac{0.40498}{0.0781} = 5.1665$$

Note! R function to calculate the odds ratio:

```
library(DAAG)
# for Default data set
mylogit <- glm(default ~ factor(student), data = Default, family = "binomial")
sum.coef<-summary(mylogit)$coef
datc.exp(sum.coef[,1])
print(datc)
```

```
#> (Intercept) factor(student)Yes
#> 0.08007259 1.4991382
```

The odds ratio of $e^{1.4991382}$ tells us that the predicted odds of a default for students are 1.499128 times the odds for non-students. In other words, the odds of a default for students are 0.913848 higher than the odds for non-students.

Making Predictions

Once the coefficients have been estimated, it is a simple matter to compute the probability of default for each dummy variable (0 or 1). For example, using the coefficient estimates given in the output, we predict that the default probability for a student and non-student:

$$\begin{aligned} p &= \frac{e^{1.4991382}}{1 + e^{1.4991382}} \text{ which } \\ p(\text{default|student} = \text{Yes}) &= \frac{e^{1.4991382}}{1 + e^{1.4991382}} = 0.3431 \\ p(\text{default|student} = \text{No}) &= \frac{e^{-0.08007259}}{1 + e^{-0.08007259}} = 0.0292 \end{aligned}$$

The predicted probability of default for a student with a balance of 4,375 is while 2.9% for non-student. This indicates that students tend to have higher default probabilities than non-students.

Inclass Practice Problem

Gender discrimination in hiring: Suppose we are investigating allegations of gender discrimination in the hiring practices of a particular firm. An equal-rights group claims that females are just like men in the firm's hiring rules. Is there any sufficient evidence to indicate the gender bias in hiring in the firm's hiring rules?

$$\hat{p} = \frac{\hat{\beta}_0 + \hat{\beta}_1 x}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 x}}$$

The odds = $c^{\hat{\beta}_0 + \hat{\beta}_1 x}$

$$\text{The odds of student} = e^{\hat{\beta}_0 + \hat{\beta}_1} = \frac{\hat{p}_{\text{st}}}{1 - \hat{p}_{\text{st}}}$$

$$\text{The odds of nonstudent} = e^{\hat{\beta}_0} = \frac{\hat{p}_{\text{nonst}}}{1 - \hat{p}_{\text{nonst}}}$$

$$\text{The odds of st} = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1}}{e^{\hat{\beta}_0}} = \frac{e^{\hat{\beta}_1}}{1 - e^{\hat{\beta}_1}}$$

$$\text{The odds of nonst} = \frac{e^{\hat{\beta}_0}}{e^{\hat{\beta}_0 + \hat{\beta}_1}} = \frac{1}{e^{\hat{\beta}_1}}$$

$$\text{The odds ratio} = e^{\hat{\beta}_1}$$

not students.

Inclass Practice Problem

Gender discrimination in hiring. Suppose you are investigating allegations of gender discrimination in the hiring practices of a particular firm. An equal rights group claims that females are less likely to be hired than males. Is there any sufficient evidence to indicate that gender is an important predictor of hiring status?

DISCRIM.CSV

$$\text{Male} = 1 \quad \hat{P}_1 = e^{\beta_1} / (e^{\beta_1} + e^0) = e^{1.7918} / (e^{1.7918} + e^0)$$
$$\text{Female} = 0 \quad \hat{P}_0 = 1 - \hat{P}_1 = 1 - e^{1.7918} / (e^{1.7918} + e^0)$$
$$\text{The odds ratio} = \frac{\hat{P}_{\text{Male}}}{1 - \hat{P}_{\text{Male}}} = \frac{e^{1.7918}}{1 - e^{1.7918}} = \frac{e^{1.7918}}{e^{1.7918} - 1} = 6$$

Example: The desire data, showing the distribution of 24 currently married and second women, interviewed in the 2002 Fertility Survey, according to age, education, desire for more children, the data are provided in `desire.csv` file

X1= age (year)

X2= education (0=none, 1=some).

Y= desire for more children (0=no, more, 1=more).

Is there any sufficient evidence to indicate that education is an important predictor of desire for children?

```
# Fit an logistic regression model
library("rms")
desire <- read_excel("c:/Users/thanithi.aganthan/Desktop/University of California/Statistical Methods/desire.xlsx")
mylogit<-glm(desire~factor(eduation),data=desire,family = "binomial")
summary(mylogit)

##
## Call:
## glm(formula = desire ~ factor(eduation), family = "binomial",
##     data = desire)
## 
## Deviance Residuals:
##    Min      Q1      Median      Q3      Max 
## -2.73440 -2.03652 -0.09578  2.73808  1.48258 
## 
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)    
## (Intercept) -2.5221   0.3407 -7.366  0.2062 ***
## factor(eduation) 1.8549   0.3710  5.004  0.0001 *** 
## 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Dispersion parameter for binomial family taken to be 1:
## 
## Null deviance: 23.271  on 23  degrees of freedom
## Residual deviance: 28.650  on 22  degrees of freedom
## AIC: 32.67
## 
## Number of Fisher Scoring Iterations: 4
sum.coef<-summary(mylogit)$coef
coef<-exp(sum.coef[,1])
print(ct)
##
## (Intercept) Factor(eduation)'
```

The odds ratio of $e^1 = 7$ tells us that the predicted odds of desire for more children for educated person are 7 times the odds for uneducated person. In other words, the odds of desire for more children for educated person are 60% higher than the odds for uneducated person.