

Logistic Regression with R: Example One

```
> rm(list=ls()); options(scipen=999) # To avoid scientific notation
> # For Wald tests: Wtest = function(L,Tn,Vn,h=0) # H0: L theta = h
> source("http://www.utstat.utoronto.ca/~brunner/Rfunctions/Wtest.txt")
> # Read the data
> math = read.table("http://www.utstat.utoronto.ca/~brunner/data/legal/mathcat.data.txt")
> head(math); attach(math) # Variable names are now available
  hsgpa hsengl hscalc  course passed outcome
1  78.0    80    Yes Mainstrm    No  Failed
2  66.0    75    Yes Mainstrm    Yes Passed
3  80.2    70    Yes Mainstrm    Yes Passed
4  81.7    67    Yes Mainstrm    Yes Passed
5  86.8    80    Yes Mainstrm    Yes Passed
6  76.7    75    Yes Mainstrm    Yes Passed

> length(hsgpa)
[1] 394
>
> summary(math)
  hsgpa          hsengl          hscalc          course          passed          outcome
Min.   :65.00   Min.   :50.00   No : 21   Catch-up: 35   No :158   Disappeared: 97
1st Qu.:75.70   1st Qu.:71.00   Yes:373  Elite   : 31   Yes:236   Failed      : 61
Median :78.70   Median :77.00                Mainstrm:328   Passed      :236
Mean   :79.74   Mean   :76.29
3rd Qu.:83.00   3rd Qu.:83.00
Max.   :96.20   Max.   :96.00
```

```
> # First, some simple examples to illustrate the methods
> # Two continuous explanatory variables
> modell = glm(passed ~ hsgpa + hsengl, family=binomial)
> summary(modell)
```

```
Call:
glm(formula = passed ~ hsgpa + hsengl, family = binomial)
```

```
Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.5577  -0.9833   0.4340   0.9126   2.2883
```

```
Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) -14.69568    2.00683  -7.323 2.43e-13 ***
hsgpa         0.22982    0.02955   7.776 7.47e-15 ***
hsengl       -0.04020    0.01709  -2.352  0.0187 *
```

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

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

```
Null deviance: 530.66 on 393 degrees of freedom
Residual deviance: 437.69 on 391 degrees of freedom
AIC: 443.69
```

```
Number of Fisher Scoring iterations: 4
```

$$\text{Deviance} = -2[L_M - L_S]$$

Where L_M is the maximum log likelihood of the model, and L_S is the maximum log likelihood of an “ideal” (super) model that fits as well as possible. The greater the deviance, the worse the model fits compared to the “best case.”

Akaike information criterion: $AIC = 2p + \text{Deviance}$,
where p = number of model parameters

```
> betahat1 = modell$coefficients; betahat1
(Intercept)      hsgpa      hsengl
-14.69567812    0.22982332   -0.04020062
>
> # For a constant value of mark in HS English, for every one-point increase
> # in HS GPA, estimated odds of passing are multiplied by ...
> exp(betahat1[2])
hsgpa
1.258378

> # Squares of Z statistics are Wald chi-squares, same null hypothesis
> L0 = rbind(c(0,0,1)) # H0: beta2=0
> Vhat = vcov(modell); Vhat
(Intercept)      hsgpa      hsengl
(Intercept)  4.027354203 -0.0492223614 -0.0021256979
hsgpa       -0.049222361  0.0008734652 -0.0002541750
hsengl      -0.002125698 -0.0002541750  0.0002921532
> Wtest(L0,betahat1,Vhat)
      W      df    p-value
5.53165297 1.00000000 0.01867545
> (-2.352)^2
[1] 5.531904
>
```

$$G^2 = -2 \log \left(\frac{\max_{\theta \in \Theta_0} L(\theta)}{\max_{\theta \in \Theta} L(\theta)} \right)$$

```

> # Deviance = -2LL + c
> # Constant will be discussed later.
> # But recall that the likelihood ratio test statistic is the
> # DIFFERENCE between two -2LL values, so
> # G-squared = Deviance(Reduced)-Deviance(Full)
>
> # Test both explanatory variables at once
> # Null deviance is deviance of a model with just the intercept.
> modell$deviance
[1] 437.6855
> modell$null.deviance
[1] 530.6559
> # G-squared = Deviance(Reduced)-Deviance(Full)
> # df = difference in number of betas
> G2 = modell$null.deviance-modell$deviance; G2
[1] 92.97039
> 1-pchisq(G2,df=2)
[1] 0
>
> # Compare Wald test
> L1 = rbind(c(0,1,0),
+           c(0,0,1))
> Wtest(L1,betahat1,Vhat)
              W              df              p-value
63.7320597548839202773  2.0000000000000000000  0.00000000000000144329
>
> # anova adds explanatory variables in order.
> a1 = anova(modell); a1
Analysis of Deviance Table

Model: binomial, link: logit

Response: passed

Terms added sequentially (first to last)

      Df Deviance Resid. Df Resid. Dev
NULL    393    530.66
hsgpa   1    87.221    392    443.43
hsengl  1     5.749    391    437.69
> # a1 is a matrix
> a1[1,4] - a1[2,4]
[1] 87.22114
> anova(modell,test="Chisq")
Analysis of Deviance Table

Model: binomial, link: logit

Response: passed

Terms added sequentially (first to last)

      Df Deviance Resid. Df Resid. Dev Pr(>Chi)
NULL    393    530.66
hsgpa   1    87.221    392    443.43 <2e-16 ***
hsengl  1     5.749    391    437.69  0.0165 *

```

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

```
> # Estimate the probability of passing for a student with
> # HSGPA = 80 and HS English = 75
```

$$\pi = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_{p-1} x_{p-1}}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_{p-1} x_{p-1}}}$$

```
>
> x = c(1,80,75); xb = sum(x*modell$coefficients)
> phat = exp(xb)/(1+exp(xb)); phat
[1] 0.6626533

> # An easier way
> gpa80eng75 = data.frame(hsgpa=80,hsengl=75)
> # Default type is estimated logit; type="response" gives estimated probability.
> predict(modell,newdata=gpa80eng75,type="response")
1
0.6626533
>
> # Get standard error too
> predict(modell,newdata=gpa80eng75,type="response",se.fit=T)
$fit
1
0.6626533

$se.fit
1
0.02859302

$residual.scale
[1] 1
```

> # Standard error was calculated using the multivariate delta method.

Let $g : \mathbb{R}^d \rightarrow \mathbb{R}^k$ etc. If $\sqrt{n}(\mathbf{T}_n - \boldsymbol{\theta}) \xrightarrow{d} \mathbf{T}$, then $\sqrt{n}(g(\mathbf{T}_n) - g(\boldsymbol{\theta})) \xrightarrow{d} \dot{g}(\boldsymbol{\theta})\mathbf{T}$,
 where $\dot{g}(\boldsymbol{\theta}) = \left[\frac{\partial g_i}{\partial \theta_j} \right]_{k \times d}$

Here, $g(\boldsymbol{\beta}) = \frac{e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2}}{1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2}}$ and

$$\begin{aligned} \dot{g}(\boldsymbol{\beta}) &= \left(\frac{\partial g}{\partial \beta_0}, \frac{\partial g}{\partial \beta_1}, \frac{\partial g}{\partial \beta_2} \right) \\ &= \left(\frac{e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2}}{(1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2})^2}, \frac{x_1 e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2}}{(1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2})^2}, \frac{x_2 e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2}}{(1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2})^2} \right) \end{aligned}$$

```
> denom = (1+exp(xb))^2
> gdot = x*exp(xb)/denom; gdot
[1] 0.2235439 17.8835124 16.7657928
> gdot = matrix(gdot,nrow=1,ncol=3)
> sqrt(gdot %*% Vhat %*% t(gdot))
[1,] 0.02859302
```

```

>
> ##### Categorical explanatory variables #####
> # Are represented by dummy variables.
> # First look at the data.
>
> coursepassed = table(course,passed); coursepassed
      passed
course   No Yes
Catch-up 27  8
Elite     7 24
Mainstrm 124 204

> prop.table(coursepassed,1) # See proportions of row totals
      passed
course   No      Yes
Catch-up 0.7714286 0.2285714
Elite     0.2258065 0.7741935
Mainstrm 0.3780488 0.6219512

> # Now with logistic regression and dummy variables
>
> # hscal and course are factors with built-in dummy variables
> contrasts(hscal)
      Yes
No      0
Yes     1
> contrasts(course)
      Elite Mainstrm
Catch-up  0         0
Elite     1         0
Mainstrm  0         1
> # Want the reference category for course to be Mainstream
> contrasts(course) = contr.treatment(3,base=3); contrasts(course)
      1 2
Catch-up 1 0
Elite     0 1
Mainstrm 0 0
> # Labels 1 and 2 for the dummy variables are not great
> colnames(contrasts(course)) = c("Catch-up","Elite")
> contrasts(course)
      Catch-up Elite
Catch-up  1     0
Elite     0     1
Mainstrm  0     0

```

```

>
> model2 = glm(passed ~ course, family=binomial); summary(model2)

Call:
glm(formula = passed ~ course, family = binomial)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-1.7251  -1.3948   0.9746   0.9746   1.7181

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)    0.4978    0.1139   4.372 0.0000123 ***
courseCatch-up -1.7142    0.4183  -4.098 0.0000417 ***
courseElite     0.7343    0.4444   1.652  0.0985 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 530.66  on 393  degrees of freedom
Residual deviance: 505.74  on 391  degrees of freedom
AIC: 511.74

Number of Fisher Scoring iterations: 4

> anova(model2, test="Chisq")
> # Both dummy variables are entered at once because course is a factor.

Analysis of Deviance Table

Model: binomial, link: logit
Response: passed

Terms added sequentially (first to last)

      Df Deviance Resid. Df Resid. Dev  Pr(>Chi)
NULL                393     530.66
course  2    24.916    391     505.74 3.887e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

> # Compare a Pearson Chi-squared test of independence.
> chisq.test(coursepassed)

Pearson's Chi-squared test

data:  coursepassed
X-squared = 24.6745, df = 2, p-value = 4.385e-06

```

```

>
> # The estimated odds of passing are __ times as great for students in
> # the catch-up course, compared to students in the mainstream course.
> model2$coefficients
  (Intercept) courseCatch-up   courseElite
    0.4978384   -1.7142338     0.7343053
> exp(model2$coefficients[2])
course1
0.1801017
>
> # Get that number from the contingency table
> addmargins(coursepassed,c(1,2))
      passed
course  No Yes Sum
Catch-up  27  8  35
Elite      7  24  31
Mainstrm 124 204 328
Sum       158 236 394
> pr = prop.table(coursepassed,1); pr # Estimated conditional probabilities
      passed
course  No      Yes
Catch-up 0.7714286 0.2285714
Elite    0.2258065 0.7741935
Mainstrm 0.3780488 0.6219512

> odds1 = pr[1,2]/(1-pr[1,2]); odds1
[1] 0.2962963
> odds3 = pr[3,2]/(1-pr[3,2]); odds3
[1] 1.645161
> odds1/odds3
[1] 0.1801017
> exp(model2$coefficients[2])
course1
0.1801017

```

```
> ##### Now a more realistic analysis #####
>
> model3 = glm(passed ~ hsengl + hsgpa + course, family=binomial)
> summary(model3)
```

```
Call:
glm(formula = passed ~ hsengl + hsgpa + course, family = binomial)
```

```
Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.5404 -0.9852  0.4110  0.8820  2.2109
```

```
Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) -14.18265    2.06382  -6.872 0.00000000000633 ***
hsengl      -0.03534    0.01766  -2.001  0.04539 *
hsgpa        0.21939    0.02988   7.342 0.00000000000021 ***
courseCatch-up -1.29137    0.45190  -2.858  0.00427 **
courseElite   0.75847    0.49308   1.538  0.12399
```

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

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

```
Null deviance: 530.66 on 393 degrees of freedom
Residual deviance: 424.76 on 389 degrees of freedom
AIC: 434.76
```

```
Number of Fisher Scoring iterations: 4
```

```
> anova(model3, test="Chisq")
Analysis of Deviance Table
```

```
Model: binomial, link: logit
```

```
Response: passed
```

```
Terms added sequentially (first to last)
```

	Df	Deviance	Resid. Df	Resid. Dev	Pr(>Chi)
NULL			393	530.66	
hsengl	1	8.286	392	522.37	0.003994 **
hsgpa	1	84.684	391	437.69	< 0.0000000000000022 ***
course	2	12.921	389	424.76	0.001564 **

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

```
>
> # Interpret all the default tests, but watch out!
> summary(glm(passed ~ hsengl, family=binomial))
```

```
Call:
glm(formula = passed ~ hsengl, family = binomial)
```

```
Deviance Residuals:
    Min       1Q   Median       3Q      Max
-1.5895 -1.3039  0.8913  1.0133  1.4060
```

```
Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) -2.29604    0.95182  -2.412  0.01585 *
hsengl       0.03546    0.01247   2.844  0.00446 **
```

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


Repeating a little from earlier ...

(Intercept)	-14.18265	2.06382	-6.872	0.000000000000633	***
hsengl	-0.03534	0.01766	-2.001	0.04539	*
hsgpa	0.21939	0.02988	7.342	0.000000000000021	***
courseCatch-up	-1.29137	0.45190	-2.858	0.00427	**
courseElite	0.75847	0.49308	1.538	0.12399	

	Df	Deviance	Resid. Df	Resid. Dev	Pr(>Chi)
NULL			393	530.66	
hsengl	1	8.286	392	522.37	0.003994 **
hsgpa	1	84.684	391	437.69	< 0.00000000000000022 ***
course	2	12.921	389	424.76	0.001564 **

```
> # Reproduce the Z-test for hsengl
> betahat3 = model3$coefficients; betahat3
(Intercept)      hsengl      hsgpa      course1      course2
-14.18264539 -0.03533871  0.21939002 -1.29136575  0.75846785
>
> V3 = vcov(model3)
> Z = betahat3[2]/sqrt(V3[2,2]) ; Z
      hsengl
-2.001046

> # Do some Wald tests
> # Wald chi-squared for hsengl
> L1 = rbind(c(0,1,0,0,0))
>
> Wtest(L=L1,Tn=betahat3,Vn=V3)
           W           df      p-value
4.00418656 1.00000000 0.04538739
> Z^2
      hsengl
4.004187
> # Test course controlling for hsengl and hsgpa
> # Compare LR G^2 = 12.921, df=2, p=0.001564
> L2 = rbind(c(0,0,0,1,0),
+           c(0,0,0,0,1) )
> Wtest(L=L2,Tn=betahat3,Vn=V3)
           W           df      p-value
11.324864041 2.00000000 0.003474058
```

```
> # How about whether they took HS Calculus?
> model4 = update(model3, ~ . + hscal); summary(model4)
```

```
Call:
glm(formula = passed ~ hsengl + hsgpa + course + hscal, family = binomial)
```

```
Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.5517 -0.9811  0.4059  0.8716  2.2061
```

```
Coefficients:
            Estimate Std. Error z value      Pr(>|z|)
(Intercept) -15.42813    2.20154  -7.008 0.0000000000002419 ***
hsengl      -0.03619    0.01776  -2.038    0.0416 *
hsgpa        0.22036    0.03003   7.337 0.0000000000000219 ***
courseCatch-up -0.88042    0.48834  -1.803    0.0714 .
courseElite   0.79966    0.50023   1.599    0.1099
hscalYes     1.25718    0.67282   1.869    0.0617 .
```

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

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

```
Null deviance: 530.66 on 393 degrees of freedom
Residual deviance: 420.90 on 388 degrees of freedom
AIC: 432.9
```

```
Number of Fisher Scoring iterations: 4
```

```
>
> # Test course controlling for others
> notcourse = glm(passed ~ hsgpa + hsengl + hscal, family = binomial)
> anova(notcourse, model4, test="Chisq")
```

```
Analysis of Deviance Table
```

```
Model 1: passed ~ hsgpa + hsengl + hscal
Model 2: passed ~ hsengl + hsgpa + course + hscal
  Resid. Df Resid. Dev Df Deviance Pr(>Chi)
1      390      427.75
2      388      420.90  2    6.8575  0.03243 *
```

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

```
> # I like Model 3.
```

```

> # I like Model 3. Answer the following questions based on Model 3.
>
> # Controlling for High School english mark and High School GPA,
> # the estimated odds of passing are ___ times as great for students in
> # the Elite course, compared to students in the Catch-up course.
>
> betahat3 = model3$coefficients; betahat3
      (Intercept)          hsengl          hsgpa courseCatch-up   courseElite
-14.18264539   -0.03533871    0.21939002   -1.29136575    0.75846785
> exp(betahat3[5])/exp(betahat3[4])
course2
7.766609
>
> # What is the estimated probability of passing for a student
> # in the mainstream course with 90% in HS English and a HS GPA of 80%?
>
> x = c(1,90,80,0,0); xb = sum(x*model3$coefficients)
> phat = exp(xb)/(1+exp(xb)); phat
[1] 0.54688
>
> # What if the student had 50% in HS English?
> x = c(1,50,80,0,0); xb = sum(x*model3$coefficients)
> phat = exp(xb)/(1+exp(xb)); phat
[1] 0.8322448
>
> # What if the student had -40 in HS English?
> x = c(1,-40,80,0,0); xb = sum(x*model3$coefficients)
> phat = exp(xb)/(1+exp(xb)); phat
[1] 0.9916913

```

```

> # Could do it with predict. A confidence interval would be nice.
> ez = data.frame(hsengl=c(90,50,-40), hsgpa=c(80,80,80),
+               course=c("Mainstrm", "Mainstrm", "Mainstrm"))
> pre = predict(model3,newdata=ez,type="response",se.fit=T); pre

$fit
      1      2      3
0.5468800 0.8322448 0.9916913

$se.fit
      1      2      3
0.06480049 0.06968981 0.01709596

$residual.scale
[1] 1

> # Print CIs nicely
> LowerCL = pre$fit-1.96*pre$se.fit; UpperCL = pre$fit+1.96*pre$se.fit
> Prediction = pre$fit
> results = rbind(Prediction,LowerCL,UpperCL)
> colnames(results) = c("90,80", "50,80", "-40,80")
> results

      90,80      50,80      -40,80
Prediction 0.546880 0.8322448 0.9916913
LowerCL    0.419871 0.6956528 0.9581832
UpperCL    0.673889 0.9688368 1.0251993

```

This handout was prepared by Jerry Brunner, Department of Statistical Sciences, University of Toronto. It is licensed under a Creative Commons Attribution - ShareAlike 3.0 Unported License. Use any part of it as you like and share the result freely. The OpenOffice.org document is available from the course website:

<http://www.utstat.toronto.edu/~brunner/oldclass/2101f19>