

# **A Handbook of Statistical Analyses Using R**

---

Brian S. Everitt and Torsten Hothorn



---

## Logistic Regression and Generalised Linear Models: Blood Screening, Women's Role in Society, and Colonic Polyps

---

### 6.1 Introduction

### 6.2 Logistic Regression and Generalised Linear Models

### 6.3 Analysis Using R

#### 6.3.1 ESR and Plasma Proteins

We can now fit a logistic regression model to the data using the `glm` function. We start with a model that includes only a single explanatory variable, `fibrinogen`. The code to fit the model is

```
R> plasma_glm_1 <- glm(ESR ~ fibrinogen, data = plasma,
+                      family = binomial())
```

The formula implicitly defines a parameter for the global mean (the intercept term) as discussed in Chapters ?? and ?. The distribution of the response is defined by the `family` argument, a binomial distribution in our case. (The default link function when the binomial family is requested is the logistic function.)

From the results in Figure 6.2 we see that the regression coefficient for `fibrinogen` is significant at the 5% level. An increase of one unit in this variable increases the log-odds in favour of an ESR value greater than 20 by an estimated 1.83 with 95% confidence interval

```
R> confint(plasma_glm_1, parm = "fibrinogen")
```

```
      2.5 %      97.5 %
0.3387619 3.9984921
```

These values are more helpful if converted to the corresponding values for the odds themselves by exponentiating the estimate

```
R> exp(coef(plasma_glm_1)["fibrinogen"])
```

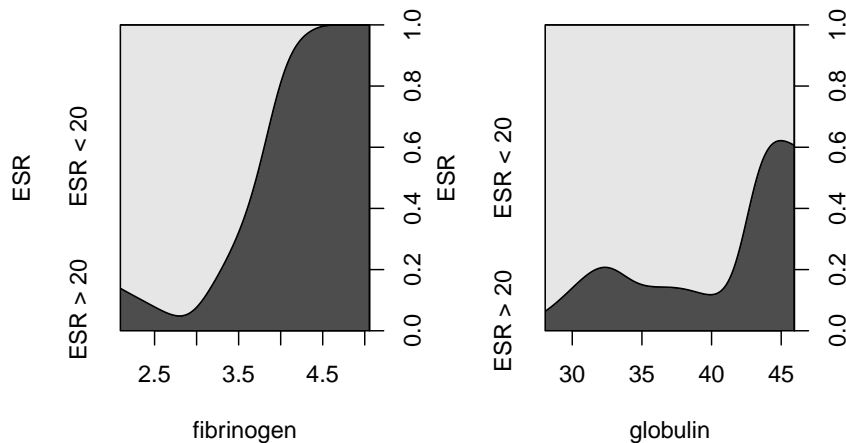
```
fibrinogen
6.215715
```

and the confidence interval

```
R> exp(confint(plasma_glm_1, parm = "fibrinogen"))
```

#### 4 LOGISTIC REGRESSION AND GENERALISED LINEAR MODELS

```
R> data("plasma", package = "HSAUR")
R> layout(matrix(1:2, ncol = 2))
R> cdplot(ESR ~ fibrinogen, data = plasma)
R> cdplot(ESR ~ globulin, data = plasma)
```



**Figure 6.1** Conditional density plots of the erythrocyte sedimentation rate (ESR) given fibrinogen and globulin.

```
2.5 % 97.5 %
1.403209 54.515884
```

The confidence interval is very wide because there are few observations overall and very few where the ESR value is greater than 20. Nevertheless it seems likely that increased values of fibrinogen lead to a greater probability of an ESR value greater than 20.

We can now fit a logistic regression model that includes both explanatory variables using the code

```
R> plasma_glm_2 <- glm(ESR ~ fibrinogen + globulin, data = plasma,
+                       family = binomial())
```

and the output of the `summary` method is shown in Figure 6.3.

The coefficient for gamma globulin is not significantly different from zero. Subtracting the residual deviance of the second model from the corresponding value for the first model we get a value of 1.87. Tested using a  $\chi^2$ -distribution with a single degree of freedom this is not significant at the 5% level and so we conclude that gamma globulin is not associated with ESR level. In R, the task of comparing the two nested models can be performed using the `anova` function

---

```
R> summary(plasma_glm_1)

Call:
glm(formula = ESR ~ fibrinogen, family = binomial(), data = plasma)

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept)  -6.8451     2.7703  -2.471  0.0135 *
fibrinogen    1.8271     0.9009   2.028  0.0425 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 30.885  on 31  degrees of freedom
Residual deviance: 24.840  on 30  degrees of freedom
AIC: 28.84

Number of Fisher Scoring iterations: 5
```

---

**Figure 6.2** R output of the `summary` method for the logistic regression model fitted to the `plasma` data.

---

```
R> summary(plasma_glm_2)

Call:
glm(formula = ESR ~ fibrinogen + globulin, family = binomial(),
    data = plasma)

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept) -12.7921     5.7963  -2.207  0.0273 *
fibrinogen   1.9104     0.9710   1.967  0.0491 *
globulin     0.1558     0.1195   1.303  0.1925
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 30.885  on 31  degrees of freedom
Residual deviance: 22.971  on 29  degrees of freedom
AIC: 28.971

Number of Fisher Scoring iterations: 5
```

---

**Figure 6.3** R output of the `summary` method for the logistic regression model fitted to the `plasma` data.

```
R> anova(plasma_glm_1, plasma_glm_2, test = "Chisq")

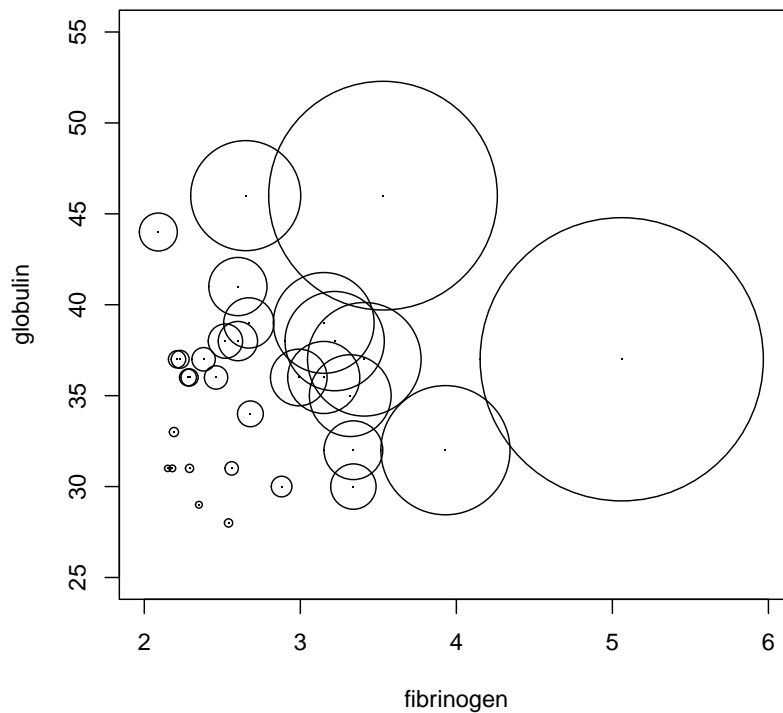
Analysis of Deviance Table

Model 1: ESR ~ fibrinogen
Model 2: ESR ~ fibrinogen + globulin
  Resid. Df Resid. Dev Df Deviance Pr(>Chi)
1         30      24.840
2         29      22.971  1   1.8692  0.1716
```

Nevertheless we shall use the predicted values from the second model and plot

6 LOGISTIC REGRESSION AND GENERALISED LINEAR MODELS

```
R> plot(globulin ~ fibrinogen, data = plasma, xlim = c(2, 6),  
+       ylim = c(25, 55), pch = ".")  
R> symbols(plasma$fibrinogen, plasma$globulin, circles = prob,  
+         add = TRUE)
```



**Figure 6.4** Bubble plot of fitted values for a logistic regression model fitted to the ESR data.

them against the values of *both* explanatory variables using a *bubble plot* to illustrate the use of the `symbols` function. The estimated conditional probability of an ESR value larger 20 for all observations can be computed, following formula (??), by

```
R> prob <- predict(plasma_glm_2, type = "response")
```

and now we can assign a larger circle to observations with larger probability as shown in Figure 6.4. The plot clearly shows the increasing probability of an ESR value above 20 (larger circles) as the values of fibrinogen, and to a lesser extent, gamma globulin, increase.

### 6.3.2 Women's Role in Society

Originally the data in Table ?? would have been in a completely equivalent form to the data in Table ?? data, but here the individual observations have been grouped into counts of numbers of agreements and disagreements for the two explanatory variables, `sex` and `education`. To fit a logistic regression model to such grouped data using the `glm` function we need to specify the number of agreements and disagreements as a two-column matrix on the left hand side of the model formula. We first fit a model that includes the two explanatory variables using the code

```
R> data("womensrole", package = "HSAUR")
R> fm1 <- cbind(agree, disagree) ~ sex + education
R> womensrole_glm_1 <- glm(fm1, data = womensrole,
+                           family = binomial())
```

---

```
R> summary(womensrole_glm_1)

Call:
glm(formula = fm1, family = binomial(), data = womensrole)

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept)  2.50937     0.18389  13.646 <2e-16 ***
sexFemale    -0.01145     0.08415  -0.136  0.892
education    -0.27062     0.01541 -17.560 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 451.722  on 40  degrees of freedom
Residual deviance: 64.007  on 38  degrees of freedom
AIC: 208.07

Number of Fisher Scoring iterations: 4
```

---

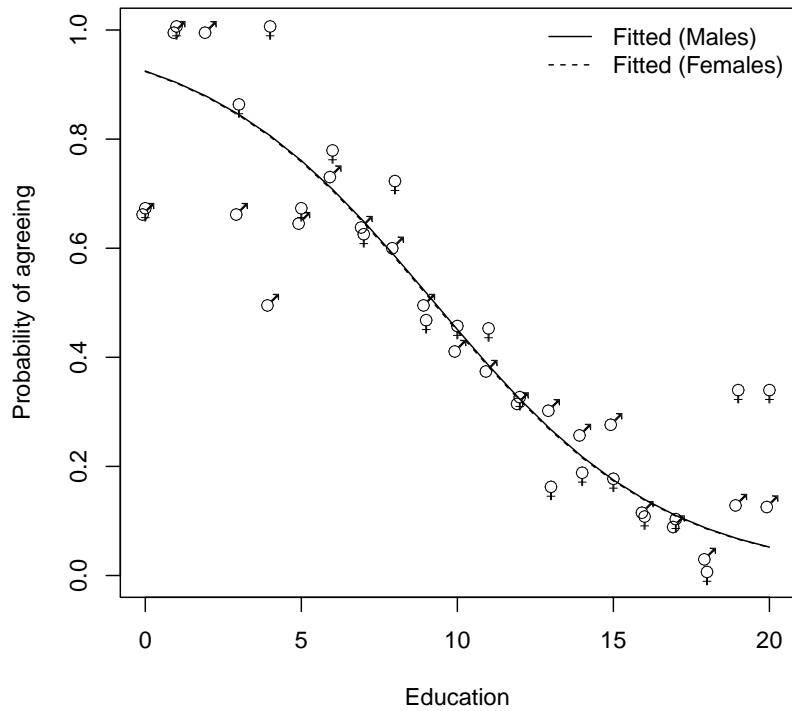
**Figure 6.5** R output of the `summary` method for the logistic regression model fitted to the `womensrole` data.

From the `summary` output in Figure 6.5 it appears that education has a highly significant part to play in predicting whether a respondent will agree with the statement read to them, but the respondent's sex is apparently unimportant. As years of education increase the probability of agreeing with the statement declines. We now are going to construct a plot comparing the observed proportions of agreeing with those fitted by our fitted model. Because we will reuse this plot for another fitted object later on, we define a function which plots years of education against some fitted probabilities, e.g.,

```
R> role.fitted1 <- predict(womensrole_glm_1, type = "response")
and labels each observation with the person's sex:
R> myplot <- function(role.fitted) {
+   f <- womensrole$sex == "Female"
```

8 LOGISTIC REGRESSION AND GENERALISED LINEAR MODELS

R> myplot(role.fitted1)



**Figure 6.6** Fitted (from `womensrole_glm_1`) and observed probabilities of agreeing for the `womensrole` data.

```
+ plot(womensrole$education, role.fitted, type = "n",
+       ylab = "Probability of agreeing",
+       xlab = "Education", ylim = c(0,1))
+ lines(womensrole$education[!f], role.fitted[!f], lty = 1)
+ lines(womensrole$education[f], role.fitted[f], lty = 2)
+ lgtxt <- c("Fitted (Males)", "Fitted (Females)")
+ legend("topright", lgtxt, lty = 1:2, bty = "n")
+ y <- womensrole$agree / (womensrole$agree +
+                          womensrole$disagree)
+ text(womensrole$education, y, ifelse(f, "\\VE", "\\MA"),
+       family = "HersheySerif", cex = 1.25)
+ }
```



The two curves for males and females in Figure 6.6 are almost the same reflecting the non-significant value of the regression coefficient for sex in `womensrole_glm_1`. But the observed values plotted on Figure 6.6 suggest that there might be an interaction of education and sex, a possibility that can be investigated by applying a further logistic regression model using

```
R> fm2 <- cbind(agree,disagree) ~ sex * education
R> womensrole_glm_2 <- glm(fm2, data = womensrole,
+                           family = binomial())
```

The `sex` and `education` interaction term is seen to be highly significant, as can be seen from the `summary` output in Figure 6.7.

---

```
R> summary(womensrole_glm_2)

Call:
glm(formula = fm2, family = binomial(), data = womensrole)

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept)    2.09820    0.23550   8.910 < 2e-16 ***
sexFemale       0.90474    0.36007   2.513  0.01198 *
education      -0.23403    0.02019 -11.592 < 2e-16 ***
sexFemale:education -0.08138    0.03109  -2.617  0.00886 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 451.722  on 40  degrees of freedom
Residual deviance:  57.103  on 37  degrees of freedom
AIC: 203.16

Number of Fisher Scoring iterations: 4
```

---

**Figure 6.7** R output of the `summary` method for the logistic regression model fitted to the `womensrole` data.

We can obtain a plot of deviance residuals plotted against fitted values using the following code above Figure 6.9. The residuals fall into a horizontal band between  $-2$  and  $2$ . This pattern does not suggest a poor fit for any particular observation or subset of observations.

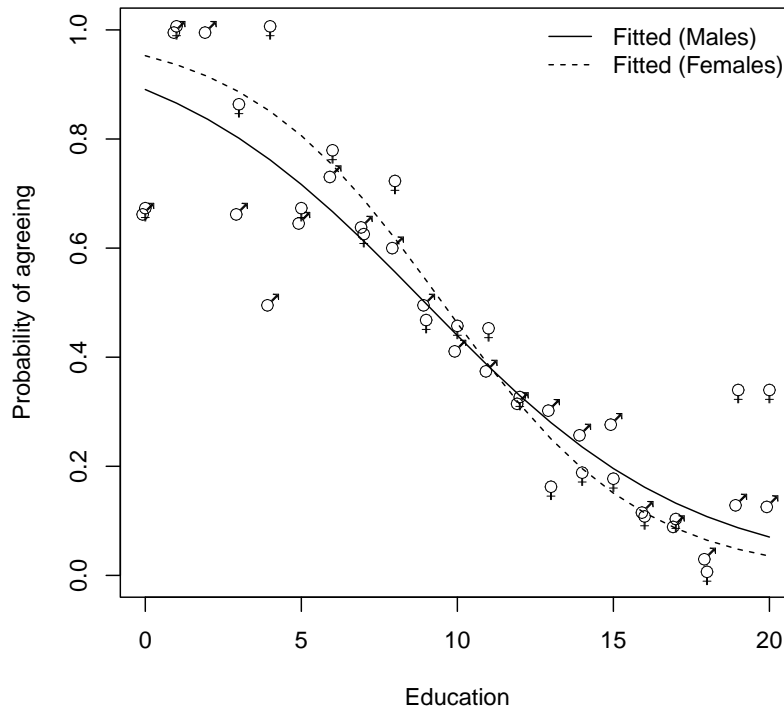
### 6.3.3 Colonic Polyps

The data on colonic polyps in Table ?? involves *count* data. We could try to model this using multiple regression but there are two problems. The first is that a response that is a count can only take positive values, and secondly such a variable is unlikely to have a normal distribution. Instead we will apply a GLM with a log link function, ensuring that fitted values are positive, and a Poisson error distribution, i.e.,

$$P(y) = \frac{e^{-\lambda} \lambda^y}{y!}.$$

10 LOGISTIC REGRESSION AND GENERALISED LINEAR MODELS

```
R> role.fitted2 <- predict(womensrole_glm_2, type = "response")
R> myplot(role.fitted2)
```



**Figure 6.8** Fitted (from `womensrole_glm_2`) and observed probabilities of agreeing for the `womensrole` data.

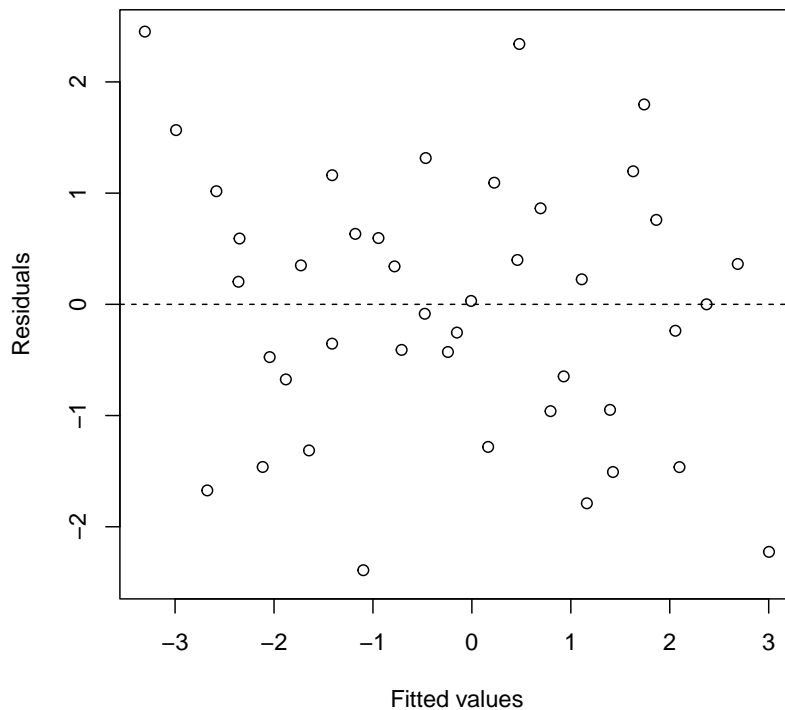
This type of GLM is often known as *Poisson regression*. We can apply the model using

```
R> data("polyps", package = "HSAUR")
R> polyps_glm_1 <- glm(number ~ treat + age, data = polyps,
+                       family = poisson())
```

(The default link function when the Poisson family is requested is the log function.)

We can deal with overdispersion by using a procedure known as *quasi-likelihood*, which allows the estimation of model parameters without fully knowing the error distribution of the response variable. [McCullagh and Nelder \(1989\)](#) give full details of the quasi-likelihood approach. In many respects it

```
R> res <- residuals(womensrole_glm_2, type = "deviance")
R> plot(predict(womensrole_glm_2), res,
+       xlab="Fitted values", ylab = "Residuals",
+       ylim = max(abs(res)) * c(-1,1))
R> abline(h = 0, lty = 2)
```



**Figure 6.9** Plot of deviance residuals from logistic regression model fitted to the `womensrole` data.

simply allows for the estimation of  $\phi$  from the data rather than defining it to be unity for the binomial and Poisson distributions. We can apply quasi-likelihood estimation to the `colonic polyps` data using the following R code

```
R> polyyps_glm_2 <- glm(number ~ treat + age, data = polyyps,
+                      family = quasipoisson())
R> summary(polyyps_glm_2)
```

*Call:*

```
glm(formula = number ~ treat + age, family = quasipoisson()),
```

## 12 LOGISTIC REGRESSION AND GENERALISED LINEAR MODELS

```
R> summary(polyps_glm_1)

Call:
glm(formula = number ~ treat + age, family = poisson(), data = polyps)

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept)  4.529024   0.146872  30.84 < 2e-16 ***
treatdrug   -1.359083   0.117643  -11.55 < 2e-16 ***
age         -0.038830   0.005955   -6.52 7.02e-11 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for poisson family taken to be 1)

Null deviance: 378.66 on 19 degrees of freedom
Residual deviance: 179.54 on 17 degrees of freedom
AIC: 273.88

Number of Fisher Scoring iterations: 5
```

**Figure 6.10** R output of the `summary` method for the Poisson regression model fitted to the `polyps` data.

```
data = polyps)

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  4.52902   0.48106   9.415 3.72e-08 ***
treatdrug   -1.35908   0.38533  -3.527 0.00259 **
age         -0.03883   0.01951  -1.991 0.06284 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for quasipoisson family taken to be 10.72805)

Null deviance: 378.66 on 19 degrees of freedom
Residual deviance: 179.54 on 17 degrees of freedom
AIC: NA

Number of Fisher Scoring iterations: 5
```

The regression coefficients for both explanatory variables remain significant but their estimated standard errors are now much greater than the values given in Figure 6.10. A possible reason for overdispersion in these data is that polyps do not occur independently of one another, but instead may ‘cluster’ together.

---

## Bibliography

---

McCullagh, P. and Nelder, J. A. (1989), *Generalized Linear Models*, London, UK: Chapman & Hall/CRC.