

Generalized Linear Models

Robert Edwards

8 March 2019

1 Introduction

In Weeks 3 and 6 we looked at modelling data using linear regression models where we had:

- a **continuous response variable** y and
- one or more **explanatory variables** x_1, x_2, \dots, x_p which were **numerical** and/or **categorical** variables.

Recall that for data $(y_i, x_i), i = 1, \dots, n$, where y is a continuous response variable, we can write a simple linear regression model as follows:

$$y_i = \alpha + \beta x_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

where,

- y_i is the i^{th} observation of the continuous response variable;
- α is the **intercept** of the regression line;
- β is the **slope** of the regression line;
- x_i is the i^{th} observations of the explanatory variable;
- ϵ_i is the i^{th} random component.

Thus, the full probability model for y_i given x_i can be written as

$$y_i | x_i \sim N(\alpha + \beta x_i, \sigma^2)$$

where the mean $\alpha + \beta x_i$ is given by the deterministic part of the model and the variance σ^2 by the random part. Hence we make the assumption that the outcomes y_i are normally distributed with mean $\alpha + \beta x_i$ and variance σ^2 . However, what if our response variable y is not a continuous random variable?

1.1 Generalized Linear Models

The main objective this week is to introduce **Generalised Linear Models (GLMs)**, which extend the linear model framework to response variables that don't follow the normal distribution. GLMs can be used to model non-normal continuous response variables, but they are most frequently used to model binary, categorical or count data. Here we shall focus on binary/categorical response variables. The generalised linear model can be written as:

$$y_i \sim f(g(\mu_i))$$
$$\mu_i = \mathbf{x}_i^T \beta$$

where the response y_i is predicted through the linear combination μ_i of explanatory variables by the link function $g(\cdot)$, assuming some distribution $f(\cdot)$ for y_i , and \mathbf{x}_i^T is the i^{th} row of the design matrix \mathbf{X} . For example, the simple linear regression model above for a continuous response variable has the normal distribution as $f(\cdot)$, with corresponding link function equal to the identity function, that is, $g(\mu_i) = \mu_i$.

What if our response variable y is binary (e.g. yes/no, success/failure, alive/dead)? That is, the independent responses y_i can either be:

- **binary**, taking the value 1 (say success, with probability p_i) or 0 (failure, with probability $1 - p_i$) or
- **binomial**, where y_i is the number of successes in a given number of trial n_i , with the probability of success being p_i and the probability of failure being $1 - p_i$

In both cases the distribution of y_i is assumed to be binomial, but in the first case it is $\text{Bi}(1, p_i)$ and in the second case it is $\text{Bi}(n_i, p_i)$. Hence, a binary response variable y_i has a binomial distribution with corresponding link function $g(\cdot)$ equal to the **logit link** function, that is

$$g(p_i) = \log\left(\frac{p_i}{1 - p_i}\right)$$

which is also referred to as the **log-odds** (since $\frac{p_i}{1-p_i}$ is an odds ratio). Why is such a transformation required when looking at a binary response variable? Well here we are interested in modelling the probability of success p_i , and as we know probabilities must be between 0 and 1 ($p_i \in [0, 1]$). So if we want to model the probability of success using a linear model we need to ensure that the probabilities obtained are between 0 and 1. However, if we just use the identity link function, such that

$$p_i = \mathbf{x}_i^T \beta$$

we would need to ensure that in some way $\mathbf{x}_i^T \beta \in [0, 1]$ that is, the linear combination of the explanatory variables and their corresponding regression coefficients was between 0 and 1. Hence some restrictions of some sort would need to be put in place to ensure this was the case. However, if we use the logit link function, such that

$$\log\left(\frac{p_i}{1 - p_i}\right) = \mathbf{x}_i^T \beta$$

No restrictions need to be in place on our estimates of the parameter vector β , since the inverse of the logit link function will always give us valid probabilities since

$$p_i = \frac{\exp(\mathbf{x}_i^T \beta)}{1 + \exp(\mathbf{x}_i^T \beta)} \in [0, 1]$$

This linear regression model with a binary response variable is referred to as **logistic regression**. As such, when it comes to looking at binary response variables we shall be looking at odds ratios and probabilities of success/failure. The table below is a reminder of the distribution and link function used for the normal model we have previously looked at as well as the logistic regression model we shall be examining for the rest of this week.

Model	Random Component	Systematic Component	Link Function
Normal	$y_i \stackrel{\text{indep}}{\sim} N(\mu_i, \sigma^2)$	$\mathbf{x}_i^T \beta = \beta_0 + \beta_1 x_i + \dots$	$g(\mu_i) = \mu_i$
Logistic	$y_i \stackrel{\text{indep}}{\sim} \text{Bi}(1, p_i)$	$\mathbf{x}_i^T \beta = \beta_0 + \beta_1 x_i + \dots$	$g(\mu_i) = \log\left(\frac{\mu_i}{1 - \mu_i}\right) = \log\left(\frac{p_i}{1 - p_i}\right)$

2 Binary Logistic Regression with One Numerical Explanatory Variable

Here we shall begin by fitting a logistic regression model with one numerical explanatory variable. Let's return to the `evals` data from the `moderndive` package that we examined in Week 3.

2.1 Teaching Evaluation Scores

Recall from previous weeks that student feedback in higher education is extremely important when it comes to the evaluation of teaching techniques, materials, and improvements in teaching methods and technologies. However, there have been studies into potential bias factors when feedback is provided, such as the physical appearance of the teacher; see Economics of Education Review for details. Here, we shall return to the study of student evaluations of $n = 463$ professors from The University of Texas at Austin.

Previously, we looked at **teaching score** as our continuous response variable and **beauty score** as our explanatory variable. Now we shall consider **gender** as our response variable, and hence shall have a binary response variable (female/male). We will examine if there is any difference in **gender** by **age** of the teaching instructors within the `evals` data set.

First, let's start by selecting the variables of interest from the `evals` data set:

```
evals.gender <- evals %>%  
  select(gender, age)
```

Now, let's look at a boxplot of **age** by **gender** to get an initial impression of the data:

```
ggplot(data = evals.gender, aes(x = gender, y = age, fill = gender)) +  
  geom_boxplot() +  
  labs(x = "Gender", y = "Age") +  
  theme(legend.position = "none")
```

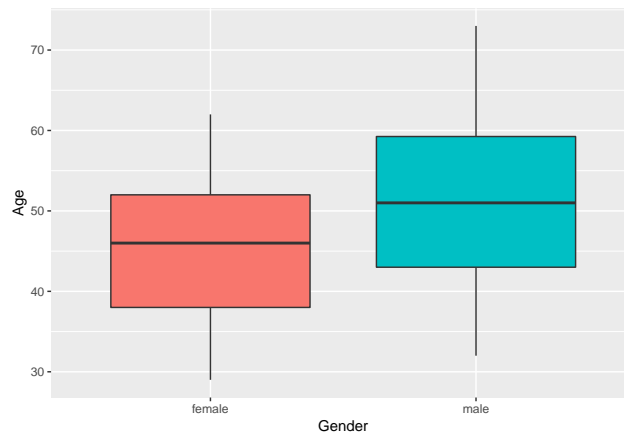


Figure 1: Teaching instructor age by gender

Here we can see that the male teaching instructors tend to be older than that of their female colleagues. Now, let's fit a logistic regression model to see whether age is a significant predictor of the odds of a teaching instructor being male or female.

2.2 Log-odds

To fit a logistic regression model we will use the generalised linear model function `glm`, which acts in a very similar manner to the `lm` function we have used previously. We only have to deal with an additional argument. The logistic regression model with **gender** as the response and **age** as the explanatory variable is given by:

```
model <- glm(gender ~ age, data = evals.gender, family = binomial(link = "logit"))
```

Here we include the additional `family` argument, which states the distribution and link function we would like to use. Hence `family = binomial(link = "logit")` states we have a binary response variable, and thus have a binomial distribution, with its corresponding **logit link** function. Now, let's take a look at the summary produced from our logistic regression model:

```
model %>%  
  summary()
```

Call:

```
glm(formula = gender ~ age, family = binomial(link = "logit"),  
    data = evals.gender)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.7134	-1.1815	0.7238	1.0180	1.4778

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-2.69795	0.51194	-5.270	1.36e-07 ***
age	0.06296	0.01059	5.948	2.71e-09 ***

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

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 630.30 on 462 degrees of freedom
Residual deviance: 591.41 on 461 degrees of freedom
AIC: 595.41

Number of Fisher Scoring iterations: 4

Firstly, the baseline category for our binary response is **female**. This is due to the default baseline in R being taken as the one which comes first alphabetically, which can be seen from the `levels` function:

```
levels(evals.gender$gender)
```

```
[1] "female" "male"
```

This means that estimates from the logistic regression model are for a change on the **log-odds** scale for **males** in comparison to the response baseline **females**. That is

$$\ln\left(\frac{p}{1-p}\right) = \alpha + \beta \cdot \text{age} = -2.7 + 0.06 \cdot \text{age}$$

where $p = \text{Prob}(\text{Male})$ and $1 - p = \text{Prob}(\text{Female})$. Hence, the **log-odds** of the instructor being male increase by 0.06 for every one unit increase in **age**. This provides us with a point estimate of how the log-odds changes with age, however, we are also interested in producing a 95% confidence interval for these log-odds. This can be done using the `confit` function in the `MASS` package:

```
confint(model) %>%
  kable(booktabs = TRUE) %>%
  kable_styling(latex_options = "hold_position",
    format = "latex")
```

```
confint(model)
```

To understand how these endpoints are calculated, consider the following code:

```
mod.coef.logodds <- model %>%
  summary() %>%
  coef()
```

```
age.logodds.lower <- mod.coef.logodds["age", "Estimate"] - 1.96 * mod.coef.logodds["age", "Std. Error"]
age.logodds.upper <- mod.coef.logodds["age", "Estimate"] + 1.96 * mod.coef.logodds["age", "Std. Error"]
```

Hence the point estimate for the log-odds is 0.06, which has a corresponding 95% confidence interval of (0.0422178, 0.0837117). This can be displayed graphically using the `plot_model` function from the `sjPlot` package simply passing our `model` as an argument:

```
plot_model(model, show.values = TRUE, transform = NULL, title = "Log-Odds (Male instructor)", show.p = TRUE)
```



Figure 2: The log-Odds of age for Male instructors

Some of the interesting arguments that can be passed to the `plot_model` function are:

- `show.values = TRUE/FALSE`: Whether the log-odds/odds values should be displayed;
- `show.p = TRUE/FALSE`: Adds asterisks that indicate the significance level of estimates to the value labels;
- `transform`: A character vector naming the function that will be applied to the estimates. The default transformation uses `exp` to display the odds ratios, while `transform = NULL` displays the log-odds; and
- `vline.color`: colour of the vertical “zero effect” line.

Further details on using `plot_model` can be found [here](#) and [here](#)

Now, let’s add the estimates of the log-odds to our data set:

```
evals.gender <- evals.gender %>%
  mutate(logodds.male = predict(model))
```

What is the log-odds of a 63 year old instructor being male?

- The log-odds of a 62 year-old instructor being male is 1.21

What is the log-odds of a 29 year old instructor being male?

- The log-odds of a 29 year-old instructor being male is -0.87

2.3 Odds

Typically we would like to work on the **odds** scale as it is easier to interpret an odds-ratio as opposed to the log-odds-ratio. To obtain the odds we simply exponentiate the log-odds, that is

$$\frac{p}{1-p} = \exp(\alpha + \beta \cdot \text{age})$$

```
model %>%  
  coef() %>%  
  exp()
```

```
(Intercept)      age  
0.06734369  1.06498927
```

On the odds scale, the value of the intercept (-2.7) gives the odds of a teaching instructor being male given their `age = 0`, which is obviously not a viable age for a teaching instructor, and hence why this value is very close to zero. For `age` we have an odds of 0.06, which indicates that for every 1 unit increase in age, the odds of the teaching instructor being male increase by a factor of 0.06. So how is this calculated? Let's look at the odds-ratio obtained from instructors aged 51 to 52 years old, that is, a 1 units difference:

$$\frac{\text{Odds}_{\text{age}=52}}{\text{Odds}_{\text{age}=51}} = \left(\frac{\frac{p_{\text{age}=52}}{1-p_{\text{age}=52}}}{\frac{p_{\text{age}=51}}{1-p_{\text{age}=51}}} \right) = \frac{\exp(\alpha + \beta \cdot 52)}{\exp(\alpha + \beta \cdot 51)} = \exp(\beta \cdot (52 - 51))$$

```
evals.gender <- evals.gender %>%  
  mutate(probs.male = fitted(model))  
  
ggplot(data = evals.gender, aes(x = age, y = probs.male)) +  
  geom_dotplot(dotsize = 0.6, alpha = 0.2, aes(fill = gender), binwidth = 1, stackgroups = TRUE) +  
  geom_smooth(method = "glm",  
             method.args = list(family = "binomial"),  
             se = FALSE) +  
  labs(x = "Age", y = "Probability of instructor being male")
```