# Linear Predictors

The covariates enter the GLM through the linear predictor $\eta$.

The predictor is linear in the coefficients to be estimated, not the covariates.

There are two types of covariates.

1. Variables, which are numbers.

2. Factors, which are categorical and have a finite number of categories.

We will write coefficients as $\beta_1, \beta_2, \dots$

and the covariates as $x_1, x_2, \dots$

# Example

$X_1$ is age

$X_2$ is temperature

} they are both variables

$\eta$ could be

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2$$

$\beta_0$ is ~~the~~ the y-intercept

$\eta$ could be

$$\beta_0 + \beta_1 X_1$$

$\eta$ could be

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2$$

↑ called an interaction term

If $\beta_3 X_1 X_2$ is included, then

so should $\beta_1 X_1$ and $\beta_2 X_2$.

$\eta = \beta_0 + \beta_3 X_1 X_2$ is not allowed.

2

# Notation for $\eta$

| Model | $\eta$ | R Formula |
|---|---|---|
| 1 | $\beta_0$ | $y \sim 1$ |
| age | $\beta_0 + \beta_1 X_1$ | $y \sim X1$ |
| age + age$^2$ ~~age + $\theta$~~ | $\beta_0 + \beta_1 X_1 + \beta_2 X_1^2$ | $y \sim X1 + I(X_1\char`\^2)$ |
| age + temperature | $\beta_0 + \beta_1 X_1 + \beta_2 X_2$ | $y \sim X1 + X2$ |
| age + temperature + age:temperature OR age * temperature | $\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 t_2$ | $y \sim X1 + X2 + X1:X2$ OR $y \sim X1 * X2$ |

3

| Model | $\eta$ | R Formula |
|---|---|---|

age*temperature*weight

$$Y \sim X1 * X2 * X3$$

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$$
$$+ \beta_4 X_1 X_2 + \beta_5 X_1 X_3 + \beta_6 X_2 X_3$$
$$+ \beta_7 X_1 X_2 X_3$$

age+temperature+
weight +
age:temperature

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2$$
$$+ \beta_3 X_3 + \beta_4 X_1 X_2$$

$$Y \sim X1 + X2 + X3$$
$$+ X_1 : X_2$$

4

# Factors

For a factor, we have a coefficient $\alpha_i$ or $\beta_j$ where i or j can take any level (an element of the category). We can also have coefficients for interaction terms.

| Model | $y$ | R Formula |
|---|---|---|
| sex | $\alpha_i$ | $y \sim sex$ |
| smoker/nonsmoker | $\beta_j$ | $y \sim s/n$ |
| sex + s/ns | $\alpha_i + \beta_j$ | $y \sim sex + s/ns$ |
| sex + s/ns + sex:s/ns | $\alpha_i + \beta_j + \gamma_{ij}$ | $y \sim sex + s/ns + sex:s/ns$ |
| OR | | OR |
| sex * s/ns | | $y \sim sex * s/ns$ |

5

We can combine variables and factors.

| Model | $\eta$ | R Formula |
|---|---|---|
| age | $\beta_0 + \beta_1 x_1$ | $Y \sim X1$ |
| sex | $\alpha_i$ | $Y \sim sex$ |
| age + sex | $\alpha_i + \beta x_1$ | $Y \sim age + sex$ |
| age * sex | $\alpha_i + \beta_i x_1$ | $Y \sim age * sex$ |

When we multiply two factors, we can't estimate all $n+m$ parameters, where $i$ has $n$ levels and $j$ has $m$ levels. One parameter must be set to 0 (e.g. $\gamma_{00} = 0$). So are really $n+m-1$ parameters.

The degrees of freedom of $\eta$ is the number of (non-zero) coefficients in it.

# Links

Let $\eta$ denote the linear predictor.

Let $\mu = E(Y)$.

The link connects $\eta$ and $\mu$.

It is a function $g$ for which

$$g(\mu) = \eta.$$

$\mu$ will be our prediction for $Y$.

We would like to take inverses

$$\mu = g^{-1}(\eta).$$

In order to do this, $g$ should be continuous and monotone increasing.

The usual link, called the canonical link

is $\theta(\mu)$

Unless, otherwise specified, the canonical

link will be used.

| Distribution | Canonical Link | R name |
|---|---|---|
| normal | $g(\mu) = \mu$ | identity |
| Poisson | $g(\mu) = \log \mu$ | log |
| Binomial | $g(\mu) = \ln\left(\frac{\mu}{1-\mu}\right)$ | logit |
| Gamma | $g(\mu) = \frac{1}{\mu}$ | inverse |

Recall! $\mu = g^{-1}(\eta)$

| Link | Inverse |
|---|---|
| identity | $g^{-1}(\eta) = \eta$ |
| log | $g^{-1}(\eta) = e^{\eta}$ |

(this makes sense because
for Poisson $\mu = \lambda > 0$)

| | |
|---|---|
| logit | set $g(\mu) = \ln\left(\frac{\mu}{1-\mu}\right) = \eta$ |

$$\Rightarrow \frac{\mu}{1-\mu} = e^{\eta} \Rightarrow \mu = \frac{e^{\eta}}{1+e^{\eta}}$$

(this makes sense because
for Binomial $\mu = p \in [0, 1]$)

~~Inverse~~
Link
inverse

# Using glm in R

To fit a GLM in R we use the command

$$\text{model} \leftarrow \text{glm}\left( Y \sim \cdots, \text{family} = \cdots , \right)$$

linear
predictor $(\eta)$

normal
poisson
binomial
gamma

link = ··· )
    identity
    log
    logit
    inverse

data = ···
    data.frame

if link is omitted
the canonical link
is used

The estimates of the parameters in $\eta$

are obtained by

$$summary(model)$$

The estimates are obtained by using maximum likelihood and numerical methods.

Call the estimates of parameters $\alpha_i$, call them $\hat{\alpha}_i$.

Using these $\hat{\alpha}_i$ in $\eta$ gives an estimate of $\eta$, called $\hat{\eta}$.

The prediction for $Y$, called fitted value, is

$$\hat{\mu} = g^{-1}(\hat{\eta}).$$

# Significance of the Parameters

If a parameter is not significant, then it should not be in the model.

It is a fact that

$$\hat{\beta} \sim N(\beta, \; CRLB(\beta))$$

We use this fact to test

$$H_0 : \beta = 0 \quad vs \quad H_1 : \beta \neq 0$$

$H_0$ implies $\beta$ is not significant.

Under $H_0$, $\quad \tilde{\beta} \sim N(0, \; CRLB(\beta))$

Let $\widehat{CRLB(\beta)}$ be $CRLB(\beta)$ with $\beta$ replaced by $\hat{\beta}$.

Under $H_0$

$$\frac{\hat{\beta} - 0}{\sqrt{\widehat{CRLB(\beta)}}} \sim N(0,1)$$

We reject $H_0$ at the 5% level if

$$\left| \frac{\hat{\beta}}{\sqrt{\widehat{CRLB}(\beta)}} \right| > 1.96$$

$$\Rightarrow \quad |\hat{\beta}| > 1.96 \sqrt{\widehat{CRLB}(\beta)}$$

or roughly $\quad |\hat{\beta}| > 2 \sqrt{\widehat{CRLB}(\beta)}$

In R the p-values for all coefficients

are obtained by

$$summary \ (model)$$

We want measures of models

saying how good they are.

12

## The Saturated Model

The saturated model is a benchmark against which we can compare any other model. If a model has as many parameters as there are data points, it will fit the data perfectly.

I.e. $\hat{\mu}_i = y_i$ for all $i$,

where $y_i$ are the observed data.

Each model has a log-likelihood

$$\ell(\underline{y}; \theta, \phi) = \ln f_{\underline{y}}(\underline{y}; \theta, \phi)$$

# Example

Claim amounts per year are exponential and independent.

$$f_{Y_i}(y_i) = \frac{1}{\mu_i} \exp\left(-\frac{y_i}{\mu_i}\right)$$

$$\ln \prod_{i=1}^{n} f_{Y_i}(y_i) = -\sum_{i=1}^{n} \frac{y_i}{\mu_i} - \sum_{i=1}^{n} \ln(\mu_i)$$

For the saturated model $\hat{\mu}_i = y_i$

We obtain estimated log likelihood by substituting the estimated parameters $\hat{\alpha}_i$ for the parameters $\alpha_i$

In this example, we obtain the likelihood

$$l_s = -\sum_{i=1}^{n} 1 - \sum_{i=1}^{n} \ln(y_i)$$

$$= -n - \sum_{i=1}^{n} \ln(y_i)$$

14