# Generalised Linear Models (GLMs) I

# Learning goals

- Be able to define Generalised Linear Model (GLM).
- Be able to obtain the cannonical link function for GLM.
- (Challenging) Be able to explain ideas for the iterated weighted least squares (IWLS) algorithm.
- Be able to compute the variance of parameter estimates in GLM using IWLS algorithm.
- Understand (scaled) deviance
  - Be able to define (scaled) deviance.
  - Be able to compute (scaled) deviance.
  - Be able to use it to test model adequacy, perform model selection for nested models and non-nested models.
- Be able to perform diagnostics for GLM using R.
  - Understand what different 'diagnostics measurements' aim to measure.
  - Be able to obtain those diagnostic measurements from the *glm* output.
  - Be able to perform diagnostics for GLM using R.

# Generalised Linear Model

**Definition:** $Y$ is a GLM if it is from an exponential family, and

$$\mu := \mathbb{E}\,Y = g^{-1}(x^T \beta)$$

where

$\qquad V = b'(\theta) \qquad\qquad x^T\beta = g(\mu). \text{ relationship between predictors}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{and } \mu.$

$g$ is a monotonic differentiable function called the *link function*.

$\quad$ $x$ is a vector of independent (predictor) variables, and

$\quad$ $\beta$ is a vector of parameters

$\Rightarrow x^T \beta = g(\hat{\mu})$

$\qquad\qquad \hat{\mu}_i = y_i$

**Remark:** We model *location* using $\eta = x^T \beta$, and let the *scale* sort itself out. That is, we do not model the scale explicitly.

$\mu \Leftrightarrow p.$
for binomial $\mu = mp$

① Interpretation $\quad$ log vs logit.

② estimation of parameter

Recall $Y$ is from an exponential family if

$$f(y; \theta, \phi) = \exp\left[\frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi)\right].$$

**Def:**

If $g(\mu) = g(\mathbb{E}Y) = x^T\beta = \theta$ then $g$ is called the *cannonical* link. Since $\mu = b'(\theta)$, it follows that the canonical link must be $(b')^{-1}$.

$b'(\theta) = \mu$.

$\theta = (b')^{-1}(\mu)$.

$g(\mu) = x^T\beta = \theta$.

$\mu = b'(\theta)$.

$\theta = (b')^{-1}(\mu)$.

$g(\mu) = x^T\beta = \theta$

def. requirement for cannonical link.

# Examples: canonical links

① Poisson

$b(\theta) = e^\theta$

$b'(\theta) = e^\theta = \mu.$

$\theta = \log \mu.$

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

$b'(\theta) = e^\theta = \mu.$

$\theta = \log \mu.$

$g(\mu) = \theta = \log \mu$

① normal

$b(\theta) = \frac{\theta^2}{2}$    $b'(\theta) = \theta$

{[1] $b'(\theta)$ identity function

inverse of identity is still

identity

$g(\mu) = \mu.$

[2] $g(\mu) = \theta$  (by def)

$\mu = \theta$   $g(\theta) = \theta$

$\Rightarrow g(\mu) = \mu$

normal  $\theta = \mu,\ g(\mu) = \mu$

Poisson  $\theta = \log \lambda = \log \mu,\ g(\mu) = \log \mu$

binomial  $\theta = \log \frac{p}{1-p} = \log \frac{\mu}{m-\mu},\ g(\mu) = \log \frac{\mu}{m-\mu}$

✶ ③ binomial

from lab

$\theta = \log \frac{p}{1-p}$

$b(\theta) = m \log(1 + e^\theta)$

$\mu = b'(\theta) = \frac{m e^\theta}{1 + e^\theta}$

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

$\Rightarrow e^\theta = \frac{1}{1 - \frac{\mu}{m}} - 1 = \frac{m - (m - \mu)}{m - \mu} = \frac{\mu}{m - \mu}$

logit link

or $\theta = \log \frac{\frac{\mu}{m}}{1 - \frac{\mu}{m}}$

$\Rightarrow \boxed{\theta = \log \frac{\mu}{m - \mu}}$

$\therefore g(\mu) = \theta = \log \frac{\mu}{m - \mu}$

# Estimation of parameters in GLM

- Iterated Weighted Least Squares (IWLS) algorithm
- Implemented in the *glm* function in R

# Estimation of parameters in GLM using maximum likelihood methods

Suppose we have independent observations $y_i$ from an exponential family, with canonical parameter $\theta_i$ and dispersion parameter $\phi$, for $i = 1, \ldots, n$.

Furthermore suppose that $y_i$ has mean

*sample specific*

$$\mu_i = b'(\theta_i) = g^{-1}(x_i^T \beta)$$

If $g$ is the canonical link then $\theta_i = x_i^T \beta$.

The log-likelihood is then

$$l(\beta, \phi; y) = \sum_{i=1}^{n} \left( \frac{y_i \theta_i - b(\theta_i)}{a(\phi)} + c(y_i, \phi) \right).$$
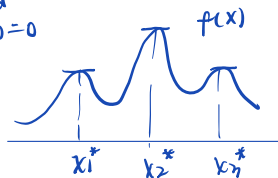
Maximum likelihood methods aim to find the values of parameters that maximise their log likelihood function.

# Aside: Newton-Raphson method

If $f(x)$ has a continuous derivative $f'(x)$, then the problem of finding the value that maximise $f(x)$ is equivalent to 1) finding $x_1^*, \ldots, x_k^*$ that are the roots of $f'(x)$, and then 2) among them, finding the one that maximise $f(x)$.

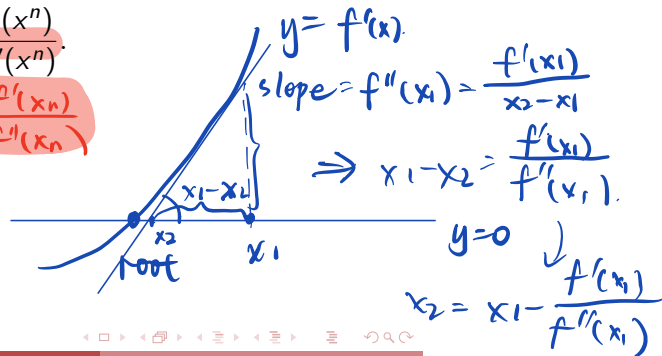We apply the Newton-Raphson method for root-finding to $f'(x)$:

$$x^{n+1} = x^n - \frac{f'(x^n)}{f''(x^n)}.$$

$$X_{n+1} = X_n - \frac{f'(X_n)}{f''(X_n)}$$

(1) find $f'(x) = 0$

$f(x)$

$x_1^* \quad x_2^* \quad x_n^*$

$x$

(2) evaluate $f(x)$. find max

root $x_2$ $x_1$

$x_1 - x_2$

$y = f'(x)$

slope $= f''(x_1) = \frac{f'(x_1)}{x_2 - x_1}$

$\Rightarrow x_1 - x_2 = \frac{f'(x_1)}{f''(x_1)}$

$y = 0$

$x_2 = x_1 - \frac{f'(x_1)}{f''(x_1)}$

# Newton-Raphson method to find MLE

Suppose we wish to find the value of $\boldsymbol{\theta}$ that maximise a log likelihood $l(\boldsymbol{\theta})$ using Newton-Raphson method.

Our update step is

$$\theta^{n+1} = \theta^n - H(\theta^n)^{-1} U(\theta^n)$$

where

$$U(\boldsymbol{\theta}) = \frac{\partial l(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \text{ and } H(\boldsymbol{\theta}) = \frac{\partial^2 l(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} = -\mathcal{J}(\boldsymbol{\theta}).$$

*second derivative.*

If we replace $\mathcal{J}$ by $\mathcal{I}$, the Fisher information, then the algorithm is called Fisher scoring.

The Fisher information is guaranteed to be positive definite (unlike the observed information).

$\theta_{n+1} = \theta_n - \dfrac{f'(\theta_n)}{f''(\theta_n)}$

# Connection with Iterated Weighted Least Squares (IWLS) algorithm?

It turns out that the Fisher scoring applied to GLM can be interpreted as a least squares problem (if you want to know details, read "McCullagh & Nelder on IWLS" posted on LMS).

The *glm* function in R computes MLEs of parameters in GLM using iterated weighted "Least Squares" algorithm.

# Aside: Weighted Least Squares method

*may assume no independence*

Suppose $Y = X\beta + \varepsilon$ where $\varepsilon \sim N(0, \Sigma)$, and $X$ and $\Sigma$ are full rank.

Multiplying by $\Sigma^{-1/2}$ we get $\Sigma^{-1/2}Y = \boxed{\Sigma^{-1/2}X}\beta + \varepsilon'$ where $\varepsilon' \sim N(0, I)$.

What is the least squares estimator of $\beta$?

$$X_*  \qquad Y_* = X_* \beta + \varepsilon'$$

The estimator of $\beta$ that minimises the sum of squares

$$\hat{\beta} = (X_*^T X_*)^{-1} X_*^T Y_*$$

$$(\Sigma^{-1/2}y - \Sigma^{-1/2}X\beta)^T(\Sigma^{-1/2}y - \Sigma^{-1/2}X\beta) = (y - X\beta)^T\Sigma^{-1}(y - X\beta)$$

is $(Y_* - X_* \beta)^T (Y_* - X_* \beta)$

$$(y - X\beta)^T \underbrace{[\Sigma^{-1/2}]^T (\Sigma^{-1/2})}_{\Sigma^{-1}} (y - X\beta)$$

$$
\begin{aligned}
\hat{\beta} &= (X^T\Sigma^{-1}X)^{-1}X^T\Sigma^{-1}y \\
&= (X^T WX)^{-1}X^T Wy,
\end{aligned}
$$

$$[ ? (\Sigma^{-1/2})^T = \Sigma^{-1/2} ]$$

where the $\boxed{\text{weight } W = \Sigma^{-1}}$ → *different weight to different sample.*

# Weighted Least Squares for GLM

$$g(Y_i) \approx \underbrace{g(\mu_i) + (Y_i - \mu_i)g'(\mu_i)}_{\geq 1} \qquad \text{Taylor expansion}$$

Let $Z_i = g(\mu_i) + (Y_i - \mu_i)g'(\mu_i)$ and $\underbrace{\epsilon_i = (Y_i - \mu_i)g'(\mu_i)}$. Then,

$$\underbrace{Z_i} = \underbrace{x_i^T \beta} + \underbrace{\epsilon_i}$$

where

$$\underbrace{\text{Var } \epsilon_i = (g'(\mu_i))^2 \text{Var } Y_i.} \qquad \hat{\beta}$$

If we knew $\text{Var } \epsilon_i$ then the estimator of $\beta$ would be the solution to the weighted least squares problem:

$$\min_{\beta}(z - X\beta)^T \Sigma^{-1}(z - X\beta)$$

where $\Sigma$ is diagonal with $\Sigma_{ii} = \text{Var } \epsilon_i$.

## Weighted Least Squares for GLM

We have $g(Y) \approx Z = X\beta + \varepsilon$ where $\epsilon_i = (Y_i - \mu_i)g'(\mu_i)$ and $\Sigma = \operatorname{Var}\varepsilon$ is diagonal with entries

$$\Sigma_{ii} = (g'(\mu_i))^2 \operatorname{Var} Y_i = (g'(\mu_i))^2 v(\mu_i)a(\phi).$$

This suggests

$$\hat{\beta} = (X^T\Sigma^{-1}X)^{-1}X^T\Sigma^{-1}z.$$

$\hat{\beta_0} \to \hat{\mu_0} \to \hat{z}, \hat{\Sigma}$

$\hat{\beta_1}$

**Problem:** $z_i$ and $\Sigma_{ii}$ depend on $\beta$ because $\mu_i = g^{-1}(x_i^T\beta)$.

**Solution:** iterate! $\Sigma$ depend on $\mu i$, $\mu i$ depend on $\beta$

Note that the $a(\phi)$ factor in the expression for $\hat{\beta}$ cancels out.

# Iterated Weighted Least Squares (IWLS) algorithm

Notice:

$$\hat{\beta} = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} z$$

$$= (X^T W X)^{-1} X^T W z$$

$$z_i = g(\mu_i) + (y_i - \mu_i) g'(\mu_i)$$

$$W_{ii} = \frac{1}{\Sigma_{ii}} = \frac{1}{(g'(\mu_i))^2 v(\mu_i)}$$

$$\hat{\mu}_i = g^{-1}(x_i^T \beta)$$

---

**IWLS algorithm (for finding MLE of $\beta$)**

1. Start with $\hat{\mu}^0 = y$.
2. Given $\hat{\mu}^n$ calculate, for each $i$,
   $z_i^n = g(\hat{\mu}_i^n) + (y_i - \hat{\mu}_i^n) g'(\hat{\mu}_i^n)$ and
   $W_{ii}^n = \frac{1}{\Sigma_{ii}} = \frac{1}{g'(\hat{\mu}_i^n)^2 v(\hat{\mu}_i^n)}$.
3. Put $\hat{\beta}^{n+1} = (X^T W^n X)^{-1} X^T W^n z^n$ and
   for each $i$, $\hat{\mu}_i^{n+1} = g^{-1}(x_i^T \hat{\beta}^{n+1})$.
4. If $\hat{\beta}^{n+1}$ is sufficiently close to $\hat{\beta}^n$ then stop, otherwise
   return to (2).

---

For GLM, IWLS algorithm is equivalent to Fisher scoring.
Example: see the section "IWLS" in Bliss.pdf.

# Variance of $\hat{\boldsymbol{\beta}}$ from IWLS algorithm

Suppose that the IWLS algorithm converges to the estimate $\hat{\boldsymbol{\beta}}$, then

$$\hat{\boldsymbol{\beta}} = \underbrace{(X^T\hat{\Sigma}^{-1}X)^{-1}X^T\hat{\Sigma}^{-1}}_{A}\underbrace{z}_{z}$$

$\mathrm{Var}(\hat{\beta}) = A\,\mathrm{var}(z)\,A^T$

where, elementwise, we have

$$
\begin{aligned}
\hat{\mu}_i &= g^{-1}(x_i^T\hat{\beta}) \qquad \Sigma_i \\
z_i &= g(\hat{\mu}_i) + \boxed{(y_i - \hat{\mu}_i)g'(\hat{\mu}_i)} \\
\hat{\Sigma}_{ii} &= (g'(\hat{\mu}_i))^2 v(\hat{\mu}_i).
\end{aligned}
$$

Since $\mathrm{Var}\,z = \hat{\Sigma}$ we have

$$
\begin{aligned}
\mathrm{Var}\,\hat{\boldsymbol{\beta}} &= [(X^T\hat{\Sigma}^{-1}X)^{-1}X^T\hat{\Sigma}^{-1}]\hat{\Sigma}[(X^T\hat{\Sigma}^{-1}X)^{-1}X^T\hat{\Sigma}^{-1}]^T \\
&= (X^T\hat{\Sigma}^{-1}X)^{-1}
\end{aligned}
$$

Note that the $a(\phi)$ term in $\hat{\Sigma}$ does not cancel here, as it did in the IWLS algorithm, so we need to estimate it.

## Estimator for $a(\phi)$

$Y_i$ has mean $\mu_i$ and variance $v(\mu_i)a(\phi)$. So

$$\sum_i \frac{(Y_i - \hat{\mu}_i)^2}{v(\hat{\mu}_i)a(\phi)} \approx \chi^2_{n-p}$$

where $p$ is the number of parameters used to estimate $\mu$.

Let

$$X^2 := \sum_i \frac{(Y_i - \hat{\mu}_i)^2}{v(\hat{\mu}_i)}.$$

Then, $X^2/(n-p)$ will be an estimator for $a(\phi)$.

$X^2$ is called Pearson's $\chi^2$ statistic, and it can be shown that $X^2/(n-p)$ is a consistent estimator for $a(\phi)$.

# Variance of $\hat{\boldsymbol{\beta}}$ from IWLS algorithm

$$\mathrm{Var}\,\hat{\boldsymbol{\beta}} \;=\; (X^T\hat{\Sigma}^{-1}X)^{-1},$$

where

$$
\begin{aligned}
\hat{\Sigma}_{ii} &= (g'(\hat{\mu}_i))^2 v(\hat{\mu}_i)a(\hat{\phi}) \\
&= (g'(\hat{\mu}_i))^2 v(\hat{\mu}_i)X^2/(n-p).
\end{aligned}
$$