## Estimation for GLMs

Summary

- Exponential family

- Components of GLM

- IWLS

Reading

- DB Chapter 4

- MN Chapter 2

## Review: Exponential Family

- $Y_i \sim$ Exponential family

  ○ Density;
  $$f_i(y_i|\theta_i, \phi\ = \exp\left\{\frac{y_i\theta_i - b(\theta_i)}{a_i(\phi)} + c_i(y_i, \phi)\right\},$$
  where $\mu_i \equiv E[Y_i] = b'(\theta_i)$ and $V(Y_i) = b''(\theta_i)a_i(\phi)$.

  ○ Link between $\mu_i$ and $\eta_i$ (linear term)

  $$g(\mu_i) = \eta_i = \mathbf{X}_i\beta = \sum_{j=1}^{p} X_{ij}\beta_j$$

  (Very often $X_{i1} = 1$ as intercept).

- Examples of CLMs (with canonical link, i.e. $\theta = \eta$) → Variance function
not Variance

| Model | $\mu$ | $\eta = g(\mu)$ | $V(\mu)$ | $a(\phi)$ |
|-------|-------|-----------------|----------|-----------|
| Linear | mean | $\mu$ | 1 | $\sigma^2$ |
| Logistic | Prob. of success P | $\log\frac{\mu}{1-\mu}$ | $\mu(1-\mu)$ | $\frac{1}{m}$ |
| Poisson | Expected Count $\lambda = \mu$ | $\log(\mu)$ | $\mu$ | 1 |

← Consider $Y/m$ as the response.
$Y \sim Bin(m, p)$

For the Binomial family the distribution of $Y_i/m_i$ is used.

\*HW   Find GLM's implemented in R and SAS

## Estimation for GLMs

- Maximum Likelihood Estimation (MLE) for $\beta$.

  ○ Log-likelihood

$$\log L = \sum_{i=1}^{n} \left\{ \frac{Y_i \theta_i - b(\theta_i)}{a_i(\phi)} + c_i(Y_i, \phi) \right\} = \sum_{i=1}^{p} \log L_i,$$

  ○ Score equation

$$U_j(\beta) = \frac{\partial}{\partial \beta_j} \log L = 0, \quad \text{for } j = 1, \quad , p.$$

  ○ There is no $\beta_j$ in the form of $\log L$. Then, how to estimate $\beta_j$?

    \* Chain rule !

In the Canonical form, $\ell(\beta) = \log L$ is a function of $\theta_i$ orignally,
$\theta_i \longrightarrow \mu_i \longrightarrow \eta_i \longrightarrow \beta$
Since $\theta_i = g(\mu_i)$, $\mu_i = g^{-1}(\eta_i) = h(\eta_i)$, $\eta_i = x_i^T\beta$, then $\ell(\beta)$ is a function of $\beta$

A Question   Since $\theta_i = \eta_i = x_i^T\beta$, why not use $\theta_i \longrightarrow \eta_i \longrightarrow \beta$?

Answer ·  Using $\mu_i = E(Y_i)$, it is convenient to formulate IWLS late

Then, the score functions can be written as:

$$U_j(\beta) = \frac{\partial}{\partial \beta_j} \log L = \sum_{i=1}^{n} \frac{\partial}{\partial \beta_j} \log L_i$$

$$= \sum_i \frac{1}{V_i \cdot (\partial \eta_i / \partial \mu_i)^2} (y_i - \mu_i) \frac{\partial \eta_i}{\partial \mu_i} x_{ij}$$

see proofs below

## Estimation (continued)

- where

Proof: For $i$th obs, let

$\ell_i = \ell_i(y_i; \theta_i; \phi) = \log L_i = \frac{[y_i \theta_i - b(\theta_i)]}{a(\phi)} + c(y_i; \phi)$

where $\theta_i = g(\mu_i) = \eta_i$,

$\mu_i = E(y_i) = h(\theta_i) = h(\eta_i) = g^{-1}(\eta_i)$ (i.e. $h = g^{-1}$),

$\eta_i = x_i^T \beta$,

Recall $\mu_i = b'(\theta_i)$, $V(\mu_i) = b''(\theta_i)$

$\frac{\partial}{\partial \theta_i} \log L_i = \frac{\partial \ell_i}{\partial \theta_i} = \frac{y_i - b'(\theta_i)}{a(\phi)}$

$\frac{\partial \theta_i}{\partial \mu_i} = 1 / \frac{\partial \mu_i}{\partial \theta_i} = 1 / b''(\theta_i)$

$\frac{\partial \mu_i}{\partial \eta_i} = 1 / \frac{\partial \eta_i}{\partial \mu_i} = 1 / g'(\mu_i)$

$\frac{\partial \eta_i}{\partial \beta_j} = x_{ij}$, $x_j = (x_{i1}, \dots, x_{ip})^T$

Let $V_i = a(\phi) V(\mu_i) = a(\phi) b''(\theta_i) = V(y_i)$

- Thus, we have

$$U_j(\beta) = \sum_i \frac{\partial \ell_i}{\partial \beta_j} = \sum_i \frac{\partial \ell_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_j}$$

$$= \sum_i \frac{y_i - \mu_i}{V_i} \frac{\partial \mu_i}{\partial \eta_i} x_{ij}$$

* Use $\frac{\partial \mu_i}{\partial \eta_i} = 1 / \frac{\partial \eta_i}{\partial \mu_i} \longrightarrow = \sum_i \frac{1}{V_i (\frac{\partial \eta_i}{\partial \mu_i})^2} (y_i - \mu_i) \cdot \frac{\partial \eta_i}{\partial \mu_i} \cdot x_{ij}$

Here we don't cancel $\frac{\partial \eta_i}{\partial \mu_i}$ to make IWLS form.

## Estimation (continued)

$\longrightarrow$ • i.e., let $W_i = 1/\{V(Y_i)\left(\frac{\partial \mu_i}{\partial \mu_i}\right)^2\}$, with $V(Y_i) = \text{Var}(Y_i) = a_i(\phi)V(\mu_i)$, then

$$U_j(\beta) \equiv \sum_{i=1}^{n} X_{ij}W_i(Z_i - \eta_i),$$

where

$$W_i^{-1} = \left(\frac{\partial \eta_i}{\partial \mu_i}\right)^2 a_i(\phi)V(\mu_i) = (g'(\mu_i))^2 a_i(\phi)V(\mu_i),$$

$$Z_i = \eta_i + (Y_i - \mu_i)\frac{\partial \eta_i}{\partial \mu_i} = \eta_i + (Y_i - \mu_i)g'(\mu_i)$$

$$= \eta_i + \tilde{\varepsilon}_i.$$

we see that $\text{Var}(Z_i) = \text{Var}(Y_i)\left(\frac{\partial \eta_i}{\partial \mu_i}\right)^2 = 1/W_i = \text{Var}(\tilde{\varepsilon}_i)$.

• Note: In CLMs, $\eta_i$ is a linear model (i.e. $= \mathbf{X}_i^T \beta$). So the score function $U_j(\beta)$ above looks like a weighted score function of a linear regression model for $Z_i$ with $W_i = 1/\text{Var}(\tilde{\varepsilon}_i)$. In fact, since

$$Z_i = \eta_i + (Y_i - \mu_i)g'(\mu_i) \triangleq X_i^T\beta + \tilde{\varepsilon}_i$$

$$\text{Var}(\tilde{\varepsilon}_i) = \text{Var}(Y_i)(g'(\mu_i))^2 = 1/W_i,$$

then WLSE $\min_\beta \sum_i W_i(Y_i - X_i^T\beta)^2$

7

## Estimation (continued)

• The score function in matrix notation:

$$U(\beta) = \mathbf{X}^T\mathbf{W}(\mathbf{Z} - \eta),$$

where

$$\mathbf{X}_{n \times p} = (\mathbf{X}_1^T, \ , \mathbf{X}_n^T)^T, \quad \mathbf{W}_{n \times n} = \text{diag}(W_1, \ , W_n),$$

$$\mathbf{Z} = (Z_1, \ , Z_n), \quad \eta = (\eta_1, \ , \eta_n)^T = \mathbf{X}\beta.$$

The score equation is

• Assuming full rank of $\mathbf{X}$, $\quad U(\beta) = X^T W(Z - \eta) = X^T W(Z - X\beta) = 0$

$$U(\hat{\beta}) = 0 \quad \rightarrow \quad \hat{\beta} = (X^T W X)^{-1}(X^T W)Z$$

○ Problem: $Z$ and $W$ depend on $\beta$.

8

## Estimation: Iterative WLS

- Iterative Weighted Least Squares (IWLS) algorithm for GLMs

  ○ Step 1. Initialization

$$\eta = g(\mathbf{Y}), \quad (\mathbf{Y} \text{ initializes } \mu) \text{ or } \eta = \mathbf{X}\beta^{(0)},$$

←  *In Binomial and Poisson distributions, to avoid log(0) for small values of $\psi$ in log($\psi$), use a small $\varepsilon > 0$, let $\eta = log(\varepsilon)$*

where $\beta^{(0)}$ is an initial value of $\beta$.

Note: you may need some modification to avoid $\log(0)$.

  ○ Repeat the followings until changes in $\beta$ are small (e.g., $\|\hat{\beta}^{(t+1)} - \hat{\beta}^t\| < 10^{-5}$ in SAS):

*Note. $\eta = log(\mathbf{Y})$ is a vector*

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

(ii) $\mathbf{Z} = \eta + (\mathbf{Y} - \mu)g'(\mu)$ (element-wise multiplication)

(iii) $\mathbf{W} = \text{diag}(\{[g'(\mu)]^2 a(\phi)V(\mu)\}^{-1}$ (element-wise multiplication)

(iv) $\beta = (\mathbf{X}^T\mathbf{W}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{W}\mathbf{Z}$.

(v) $\eta = \mathbf{X}\beta$.

(vi) Go back to (i).

---

## IWLS (example)

- Independent r.v. $Y_1, \quad , Y_n \sim \text{Poisson}(\lambda_i)$ with $p$ covariates $X_i^T$

$$f(y_i|\theta, \phi) = \frac{\lambda_i^{y_i} e^{\lambda_i}}{y_i!},$$

with $\mu_i = \lambda_i$.

  ○ $\theta = \log(\lambda)$, $b(\theta) = e^\theta$  $a(\phi) = 1$.

  ○ Canonical link: $\eta = \theta = \log(\lambda) \rightarrow$ Log-linear regression

$$\log(\mu_i) = \mathbf{X}_i^T \beta,$$

○ MLE of $\beta$   For the Poisson model, $\theta_i = \log(\mu_i)$, $b(\theta_i) = e^{\theta_i}$

$a(\phi) = 1$,    $\mu_i = b'(\theta_i)$

$$\frac{\partial}{\partial \theta_i} \log L_i = \frac{y_i - b'(\theta_i)}{a(\phi)} = y_i - e^{\theta_i}$$

$$\frac{\partial \theta_i}{\partial \mu_i} = \frac{1}{\frac{\partial \mu_i}{\partial \theta_i}} = \frac{1}{e^{\theta_i}}$$

$$\frac{\partial \mu_i}{\partial \eta_i} = \frac{1}{\frac{\partial \eta_i}{\partial \mu_i}} = \frac{1}{g'(\mu_i)} = \frac{1}{1/\mu_i} = \mu_i = e^{\theta_i}$$

$$\frac{\partial \eta_i}{\partial \beta_j} = x_{ij}$$

○ The score function is

$$U_j(\beta) = \frac{\partial \ell_i}{\partial \theta_i} \cdot \frac{\partial \theta_i}{\partial \mu_i} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta_j}$$

$$= (y_i - \mu_i) x_{ij}$$

In the tform of IWLS, $U_j(\beta) = W_i(Z_i - \eta_i) x_{ij}$,

where $W_i = \frac{1}{Var(Y_i)\left(\frac{\partial \eta_i}{\partial \mu_i}\right)^2} = \frac{1}{\mu_i \cdot (1/\mu_i)^2} = \mu_i = e^{\theta_i}$,

$Z_i = x_i^T \beta + (y_i - \mu_i)\left(\frac{\partial \eta_i}{\partial \mu_i}\right) = \eta_i + (y_i - \mu_i) \cdot \frac{1}{e^{\theta_i}}$

One can see that $Var(Z_i) = Var(Y_i)\left(\frac{\partial \eta_i}{\partial \mu_i}\right)^2 = 1/W_i$

Using proc/IML to implement IWLS

See Table4-3.sas
for the plot
in Figure 4.5,
and Table4-4.sas
for the results
from IWLS in
Table 4.4, P.68

[In addition I
have Table4-3 R
to do the same
(calculation]

**IWLS vs. Other Algorithms**

- IWLS solves the score equations.
- There are other commonly used methods for solving score equations:
  ○ Fisher scoring:
  $$\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} + (\mathbf{I}^{(t)})^{-1} U^{(t)},$$
  where
  $$\mathbf{I}^{(t)} = -E\left[\frac{\partial^2 \log L}{\partial \beta \partial \beta^T}\right]_{\beta = \hat{\beta}^{(t)}},$$
  the total information matrix for the sample.
  ○ Newton-Raphson:
  $$\hat{\beta}^{(t+1)} = \hat{\beta}^t + (\hat{\mathbf{I}}^{(t)})^{-1} U^{(t)},$$
  where
  $$\hat{\mathbf{I}}^{(t)} = -\frac{\partial^2 \log L}{\partial \beta \partial \beta^T}\Big|_{\beta = \hat{\beta}^{(t)}}$$
  where the information matrix I.

Ex See data in Table 4.3
and Figure 4.5 in page 67
3rd edt
Poisson Regression with
the identity link

In this example,
$g(\mu_i) = \eta_i \Rightarrow \mu_i = \eta_i$,
then $\frac{\partial \eta_i}{\partial \mu_i} = 1$ and

$Z_i = \eta_i + (y_i - \mu_i)\frac{\partial \eta_i}{\partial \mu_i}$
$= \eta_i + (y_i - \mu_i) \cdot 1$
$= y_i$

$W_i = \frac{1}{Var(Y_i)\left(\frac{\partial \eta_i}{\partial \mu_i}\right)^2}$
$= \frac{1}{Var(Y_i)} = \frac{1}{\mu_i} = \frac{1}{\beta_1 + \beta_2 x_i}$

The IWLS is
$J^{(m-1)} b^{(m)} = X^T W Z^{(m-1)}$

$$J = X^T W X = \begin{bmatrix} \sum_{i=1}^{N} \frac{1}{b_1 + b_2 x_i} & \sum_{i=1}^{n} \frac{x_i}{b_1 + b_2 x_i} \\ \sum_{i=1}^{N} \frac{x_i}{b_1 + b_2 x_i} & \sum_{i=1}^{N} \frac{x_i^2}{b_1 + b_2 x_i} \end{bmatrix}, \quad X^T W Z = \begin{bmatrix} \sum_{i=1}^{N} \frac{y_i \cdot 12}{b_1 + b_2 x_i} \\ \sum_{i=1}^{N} \frac{x_i y_i}{b_1 + b_2 x_i} \end{bmatrix}$$