

## GLM Definition

### • 3 components of a GLM

➤ A *random* component for the response:  $Y_{\mathbf{x}} \sim f(y|\theta_{\mathbf{x}}, \phi)$  where

- Canonical parameter  $\theta_{\mathbf{x}}$ : represent the location
- Dispersion parameter  $\phi$ : represent the scale

➤ A *systematic* (linear) component for the predictors:

$$\eta_{\mathbf{x}} = \beta_0 + h_1(\mathbf{x})\beta_1 + \cdots + h_{p-1}(\mathbf{x})\beta_{p-1} \equiv X\beta$$

➤ A link function  $g^*$  such that  $\eta_{\mathbf{x}} = g^*(\theta_{\mathbf{x}})$

### • Random component

➤ Exponential family  $f(y|\theta_{\mathbf{x}}, \phi) = \exp \left[ \frac{y\theta_{\mathbf{x}} - b(\theta_{\mathbf{x}})}{a(\phi)} + c(y, \phi) \right]$

➤ Some distributions in exponential family

- Normal density  $N(\mu_{\mathbf{x}}, \sigma^2)$

$$f(y|\theta_{\mathbf{x}}, \phi) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{(y - \mu_{\mathbf{x}})^2}{2\sigma^2} \right] = \exp \left[ \frac{y\mu_{\mathbf{x}} - \mu_{\mathbf{x}}^2/2}{\sigma^2} - \frac{1}{2} \left( \frac{y^2}{\sigma^2} + \log(2\pi\sigma^2) \right) \right]$$

$$\Rightarrow \theta_{\mathbf{x}} = \mu_{\mathbf{x}}, \phi = \sigma^2, a(\phi) = \phi, b(\theta_{\mathbf{x}}) = \theta_{\mathbf{x}}^2/2, \\ c(y, \phi) = -(y^2/\phi + \log(2\pi\phi))/2$$

- Poisson density  $P(\mu_{\mathbf{x}})$

$\theta_{\mathbf{x}} \rightarrow E(Y_{\mathbf{x}})$

$\phi \rightarrow \text{Var}(Y_{\mathbf{x}})$

$$f(y|\theta_{\mathbf{x}}, \phi) = \frac{e^{-\mu_{\mathbf{x}}} \mu_{\mathbf{x}}^y}{y!} = \exp [y \log \mu_{\mathbf{x}} - \mu_{\mathbf{x}} - \log y!]$$

$$\Rightarrow \theta_{\mathbf{x}} = \log \mu_{\mathbf{x}}, \phi = 1, a(\phi) = 1, b(\theta_{\mathbf{x}}) = \exp(\theta_{\mathbf{x}}), \\ c(y, \phi) = -\log y! \rightarrow \text{if allow } \phi \text{ to vary, it becomes}$$

- Binomial density  $B(n_{\mathbf{x}}, \mu_{\mathbf{x}})$

dispersion parameter. (Recall overdispersion)

$$f(y|\theta_{\mathbf{x}}, \phi) = \binom{n_{\mathbf{x}}}{y} \mu_{\mathbf{x}}^y (1 - \mu_{\mathbf{x}})^{n_{\mathbf{x}} - y}$$

$$\mu_{\mathbf{x}} = \frac{e^{\theta_{\mathbf{x}}}}{1 + e^{\theta_{\mathbf{x}}}} \Rightarrow \theta_{\mathbf{x}} = \log \left( \frac{\mu_{\mathbf{x}}}{1 - \mu_{\mathbf{x}}} \right), \phi = 1, a(\phi) = 1, c(y, \phi) = \log \binom{n_{\mathbf{x}}}{y}$$

$$b(\theta_{\mathbf{x}}) = -n_{\mathbf{x}} \log(1 - \mu_{\mathbf{x}}) = n_{\mathbf{x}} \log(1 + \exp(\theta_{\mathbf{x}}))$$

- The Gamma and inverse Gaussian are other (lesser-used) members of the exponential family

- Some other densities such as the negative binomial and Weibull are not members of exponential family, but are sufficiently close so that GLM can be fit with some simple modifications

related to both  $E(Y_{\mathbf{x}})$  &  $\text{Var}(Y_{\mathbf{x}})$   
 $\rightarrow$  So that we can do  $\chi^2$  goodness-of-fit test

- $\phi$  is free in normal while fixed at 1 in Poisson and binomial
- Some authors reserve the term *exponential family* distribution for cases where  $\phi$  is not used (such as Poisson and binomial) while using the term *exponential dispersion family* for cases where it is (such as normal)

➤ Some properties of exponential family

Use  $E(\partial \ell / \partial \theta) = 0$  &  
 $E(\partial^2 \ell / \partial \theta^2) + [E(\partial \ell / \partial \theta)]^2 = 0$   
 to find  $E(Y)$  &  $\text{Var}(Y)$ .  
 $\ell$ : log-likelihood function

1. Mean:  $E(Y_x) = \mu_x = b'(\theta_x)$
  2. Variance:  $\text{Var}(Y_x) = b''(\theta_x) a(\phi)$
- $E(Y_x)$  does not depend on  $\phi$  → *depend on  $x$* .
  - $\text{Var}(Y_x)$  is a product of functions of  $\theta_x$  and  $\phi$
  - $b''(\theta_x)$  is called the *variance function* and describes how the variance related to the mean
    - In the normal density case,  $b''(\theta_x) = 1 \Rightarrow$  variance independent of the mean
    - For other distributions such as Poisson and binomial,  $b''(\theta_x)$  is not a constant function  $\Rightarrow$  variance depends on the mean *no parameter.*
  - We can introduce weights by setting  $a(\phi) = \phi / w_x$ , where  $w_x$  is a *known* weight that varies between observations *scale*

## • Link function $\eta_x = g^*(\theta_x)$

➤ We now re-write the link function to describe how the mean response,  $E(Y_x) = \mu_x$ , is linked to  $\eta_x$ , i.e.,  $\eta_x = g(\mu_x)$

➤ In principle,  $g$  can be any monotone, continuous, and differentiable function  *$\mu_x = b'(\theta_x)$ : a function of  $\theta_x$ .*

➤ There are some convenient and common choice of  $g$

- Normal (Gaussian) density:  $Y = X\beta + \varepsilon$

□ Standard choice:  $\eta_x = \mu_x$  *identity link*

□ Other choice of  $g$  would give *known generalization of linear model with non-linear mean*  
 $Y = g^{-1}(X\beta) + \varepsilon$   *$\mu_x = g^{-1}(\eta_x)$*

Notice that this does not correspond directly to a transformation on the response as, for example, in Box-Cox type transformation, i.e.

$$Y' = g(Y) = X\beta + \varepsilon \quad \text{vs} \quad Y = g^{-1}(X\beta) + \varepsilon^*$$

- Poisson density: *like normal after transformation*

□ Identity link  $\eta_x = \mu_x$  would possibly cause negative  $\mu_x$

□ Standard choice:  $\eta_x = \log(\mu_x) \Leftrightarrow \mu_x = \exp(\eta_x) > 0$

□ Log link means that additive effect of  $x$  lead to multiplicative effect on  $\mu_x$   *$\mu_x = e^{\beta_0} \cdot e^{\beta_1 h_1(x)} \dots$*   
 $\eta_x = \beta_0 + \beta_1 h_1(x) + \dots + \beta_{p-1} h_{p-1}(x)$

- Binomial density (treat  $Y_x/n_x$  as the response):
  - Standard choice: logit, probit, complementary log-log
- Canonical link: choose  $g$  such that

$$\eta_x = g(\mu_x) = \theta_x \Rightarrow g^* \text{ is identity function}$$

which means  $g$  must satisfy  $g(b'(\theta_x)) = \theta_x \Rightarrow g = (b')^{-1}$

- Examples:

Family	Link	Variance function
Gaussian	$\eta_x = \mu_x$	(1) $\rightarrow \text{Var}(y)$ is irrelevant to $E(y)$
Poisson	$\eta_x = \log(\mu_x)$	$\mu_x$
Binomial (proportion)	$\eta_x = \log(\mu_x/(1-\mu_x))$	$\mu_x(1-\mu_x)$
Gamma	$\eta_x = \mu_x^{-1}$	$\mu_x^2$
Inverse Gaussian	$\eta_x = \mu_x^{-2}$	$\mu_x^3$

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

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

$$\theta = \text{logit}(\mu) = \log\left(\frac{\mu}{1-\mu}\right)$$

$$\text{logit}\left(\frac{e^\theta}{1+e^\theta}\right) = \log\left(\frac{e^\theta/(1+e^\theta)}{(1+e^\theta)/(1+e^\theta)}\right)$$

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

- If a canonical link is used

- $X^T Y$  is sufficient for  $\beta$
- Mathematically and computationally convenient
- Often physically justified
- However, it is not required to always use the canonical link and sometimes context may compel another choice

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

$$b''(\theta) = -e^\theta(1+e^\theta)^{-2}e^\theta + e^\theta \cdot (1+e^\theta)^{-1}$$

$$= -\mu^2 + \mu$$

$$= \mu(1-\mu)$$

## Fitting a GLM ← estimating $\beta$

- The log-likelihood for a single observation  $(\mathbf{x}_i, y_i)$  is:

$$\log L(\theta_{\mathbf{x}_i}, \phi; y_i) = w_i \left[ \frac{y_i \theta_{\mathbf{x}_i} - b(\theta_{\mathbf{x}_i})}{\phi} \right] + c(y_i, \phi)$$

where  $a_i(\phi) = \phi/w_i$  and  $w_i$  is a known weight that varies between observations

*Wi not contain any parameters.*

*a(φ)*

- For independent observations  $(\mathbf{x}_i, y_i)$ ,  $i=1, 2, \dots, k$ , the joint log-likelihood is:

$$\sum_{i=1}^k \log L(\theta_{\mathbf{x}_i}, \phi; y_i) = \sum_{i=1}^k w_i \left[ \frac{y_i \theta_{\mathbf{x}_i} - b(\theta_{\mathbf{x}_i})}{\phi} \right] + \sum_{i=1}^k c(y_i, \phi) \equiv l(\underline{\mu}_{\mathbf{x}}, \phi; Y)$$

$$g(\eta_x) = g^*(\theta_x)$$

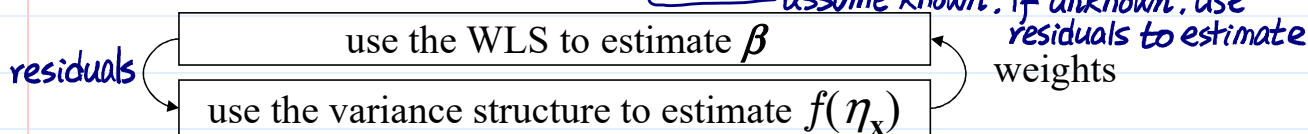
- Likelihood approach: the parameters  $\beta$ , which appear in the link function  $X\beta = \eta_x = g^*(\theta_x)$  of a GLM, can be estimated using maximum likelihood --- maximize the joint log-likelihood as a function of  $\beta$

- We can maximize the joint log-likelihood *analytically* and find an exact solution for the MLE of  $\beta$  --- Gaussian GLM is the only common case where this is possible  $\Rightarrow \beta_{MLE} = \text{least square estor}$
- Typically, we must use numerical optimization - IRWLS

# Iteratively Re-Weighted Least Squares (IRWLS)

➤ Consider the Gaussian linear model  $Y = X\beta + \varepsilon$

- Suppose  $\text{Var}(y_x) = \text{Var}(\varepsilon) \propto f(\eta_x)$  where  $\hat{\eta}_{x_i} = \mathbf{x}_i^* \hat{\beta}$



- Can use weights  $w_i$  where  $w_i^{-1} = f(\hat{\eta}_{x_i})$
- Since the weights are a function of  $\hat{\beta}$ , an iterative fitting procedure would be needed
  - set the weights all equal to one  $\Rightarrow$  estimate  $\hat{\beta} \Rightarrow$  use  $\hat{\beta}$  to re-compute the weights  $\Rightarrow$  re-estimate  $\hat{\beta} \Rightarrow$  repeat until converge

$$g(\mu_x) = \eta_x = X\beta$$

➤ similar idea can be applied to fit a GLM

- Roughly speaking, want to regress  $g(y_x)$  on  $X$  with weights inversely proportional to  $\text{Var}(g(y_x))$
- However,  $g(y_x)$  might not make sense in some cases, e.g., binomial GLM  $\Rightarrow$  linearize  $g(y_x)$  as follows: regress  $z_x$  on  $X$

$$g(y_x) \approx g(\mu_x) + (y_x - \mu_x)g'(\mu_x) = \eta_x + (y_x - \mu_x) \frac{d\eta_x}{d\mu_x} \equiv z_x$$

Gaussian Linear Model in LNp.7-7

$$\begin{aligned} \square E(z_x) &= \eta_x \equiv X\beta \\ \square \text{Var}(z_x) &= \left( \frac{d\eta_x}{d\mu_x} \right)^2 \text{Var}(y_x) \propto \left( \frac{d\eta_x}{d\mu_x} \right)^2 V(\mu_x) \\ \square \widehat{\text{Var}}(z_x) &= \left( \frac{d\eta_x}{d\mu_x} \right)^2 V(\mu_x) \Big|_{\mu_x = \hat{\mu}_x} \equiv \frac{1}{w_x} \end{aligned}$$

$z_x = g(\mu_x)$  (link function)  
 $V(\mu_x)$  (variance function) (LNp 7-3)  
 depends on  $\mu_x$  (mean of  $y_x$ ), i.e., depends on  $z_x$   
 e.g.  $\hat{\mu}_{x,0} = y_x$

## IRWLS Procedure

- estimation
- Set initial estimates  $\hat{\eta}_{x,0} = g(\hat{\mu}_{x,0})$  and  $\hat{\mu}_{x,0} = g^{-1}(\hat{\eta}_{x,0})$
  - Form the "adjusted dependent variable"

$$z_0 = \hat{\eta}_{x,0} + (y_x - \hat{\mu}_{x,0}) \frac{d\eta_x}{d\mu_x} \Big|_{\hat{\eta}_{x,0}}$$

- Form the weights

$$w_0^{-1} = \left( \frac{d\eta_x}{d\mu_x} \right)^2 \Big|_{\hat{\eta}_{x,0}} V(\hat{\mu}_{x,0})$$

- Use WLS to re-estimate  $\beta$ , then get  $\hat{\eta}_{x,1}$  and  $\hat{\mu}_{x,1}$

- Iterate steps 2-3-4 until convergence

Q: how much information about the distribution of  $y$  been used in the IRWLS?

- Note: the fitting procedure use only  $\eta_x = g(\mu_x)$  and  $V(\mu_x)$ 
  - $\Rightarrow$  It requires no further knowledge of the distribution of  $y_x$
  - $\Rightarrow$  Quasi-likelihood approach

➤ Estimates of variance:  $\text{WLS} \rightarrow \widehat{\text{Var}}(\hat{\beta}) = (X^T W X)^{-1} \cdot \hat{\sigma}^2$

- Comparable to the form used in weighted least squares
- The weights are now a function of the response for a GLM

- Gaussian
- Binomial with dispersion para.
- Poisson with dispersion para.

$$\widehat{\text{Var}}(\hat{\beta}) = (X^T W X)^{-1} \hat{\phi} \begin{bmatrix} \hat{w}_1 & \hat{w}_2 & 0 \\ 0 & \dots & \hat{w}_k \end{bmatrix}$$

where  $\hat{\phi} = X^2 / (k - p)$  and  $X^2$  is the Pearson's  $X^2$  statistic  $\hat{\phi} = 1$  in binomial & Poisson (Lnp.7-2)

➤ Example: binomial response (treat the proportion (=count/ $n_x$ ) as the response, not the count)

- $\eta_x = \log\left(\frac{\mu_x}{1 - \mu_x}\right) \Leftrightarrow \eta = g(\mu)$
- $\frac{d\eta_x}{d\mu_x} = \frac{1}{\mu_x(1 - \mu_x)}$
- $V(\mu_x) = \frac{\mu_x(1 - \mu_x)}{n_x}$
- $w_x = \frac{1}{\left[\left(\frac{d\eta}{d\mu}\right)^2 V(\mu)\right]^{-1}} = n_x \mu_x (1 - \mu_x)$
- $\widehat{\text{Var}}(\hat{\beta}) = (X^T W X)^{-1}$

➤ Some notes about IRWLS:

- In most cases, the convergence is usually fast
- If there is a failure to converge, this is often a sign of some problem with the model specification or unusual feature of the data (e.g., data is linearly separable)