

## 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 (Gaussian) density  $N(\mu_{\mathbf{x}}, \sigma^2)$

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

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

NTHU STAT 5230, 2025, Lecture Notes  
made by S.-W. Cheng (NTHU, Taiwan)

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

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

$$\Rightarrow \theta_{\mathbf{x}} = \log \mu_{\mathbf{x}}, \phi = 1, a(\phi) = 1, b(\theta_{\mathbf{x}}) = \exp(\theta_{\mathbf{x}}), \\ c(y, \phi) = -\log y!$$

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

$$f(y|\theta_{\mathbf{x}}, \phi) = \binom{n_{\mathbf{x}}}{y} \mu_{\mathbf{x}}^y (1 - \mu_{\mathbf{x}})^{n_{\mathbf{x}}-y} \\ = \exp \left[ y \log \left( \frac{\mu_{\mathbf{x}}}{1 - \mu_{\mathbf{x}}} \right) + n_{\mathbf{x}} \log(1 - \mu_{\mathbf{x}}) + \log \binom{n_{\mathbf{x}}}{y} \right]$$

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

- 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

- $\phi$  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

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  $\phi$
- $\text{Var}(Y_x)$  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
- We can introduce weights by setting  $a(\phi) = \phi/w_x$ , where  $w_x$  is a *known* weight that varies between observations

NTHU STAT 5230, 2025, Lecture Notes  
made by S.-W. Cheng (NTHU, Taiwan)

• 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
- There are some convenient and common choice of  $g$

■ Normal (Gaussian) density:

- Standard choice:  $\eta_x = \mu_x$
- Other choice of  $g$  would give

$$Y = g^{-1}(X\beta) + \varepsilon$$

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

$$g^{-1}(Y) = X\beta + \varepsilon$$

■ Poisson density:

- 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  $\mu_x$

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

$$\eta_x = g(\mu_x) = \theta_x$$

which means  $g$  must satisfy  $g(b'(\theta_x)) = \theta_x$

- Examples:

Family	Link	Variance function
Gaussian	$\eta_x = \mu_x$	1
Poisson	$\eta_x = \log(\mu_x)$	$\mu_x$
Binomial	$\eta_x = \log(\mu_x/(1-\mu_x))$	$\mu_x(1-\mu_x)$
Gamma	$\eta_x = \mu_x^{-1}$	$\mu_x^2$
Inverse Gaussian	$\eta_x = \mu_x^{-2}$	$\mu_x^3$

- 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

NTHU STAT 5230, 2025, Lecture Notes  
made by S.-W. Cheng (NTHU, Taiwan)

## Fitting a GLM

- 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

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

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

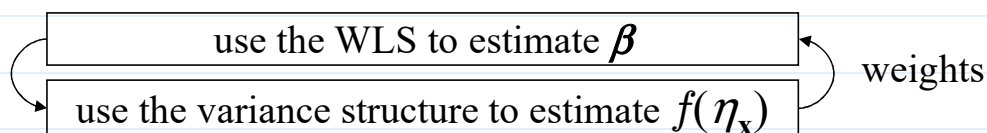
- Likelihood approach: the parameters  $\beta$ , which appear in the link function  $X\beta = \eta_x = g^*(\theta_x)$  of a GLM, can be estimated using maximum likelihood --- maximize the joint log-likelihood as a function of  $\beta$

- 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
- 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  $\hat{\eta}_{x_i} = \mathbf{x}_i^* T \hat{\beta}$



- 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

➤ 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:

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

NTHU STAT 5230, 2025, Lecture Notes  
made by S.-W. Cheng (NTHU, Taiwan)

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

➤ IRWLS Procedure

1. Set initial estimates  $\hat{\eta}_{x,0} = g(\hat{\mu}_{x,0})$  and  $\hat{\mu}_{x,0} = g^{-1}(\hat{\eta}_{x,0})$
2. Form the “adjusted dependent variable”
 
$$z_0 = \hat{\eta}_{x,0} + (y_x - \hat{\mu}_{x,0}) \frac{d\eta_x}{d\mu_x} \Big|_{\hat{\eta}_{x,0}}$$
3. Form the weights
 
$$w_0^{-1} = \left(\frac{d\eta_x}{d\mu_x}\right)^2 \Big|_{\hat{\eta}_{x,0}} V(\hat{\mu}_{x,0})$$
4. Use WLS to re-estimate  $\beta$ , then get  $\hat{\eta}_{x,1}$  and  $\hat{\mu}_{x,1}$
5. Iterate steps 2-3-4 until convergence

- 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:

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

$$\widehat{Var}(\beta) = (X^T W X)^{-1} \hat{\phi}$$

where  $\hat{\phi} = X^2 / (k - p)$  and  $X^2$  is the Pearson's  $X^2$  statistic

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

- $\eta_x = \log\left(\frac{\mu_x}{1 - \mu_x}\right)$
- $w_x = n_x \mu_x (1 - \mu_x)$
- $\frac{d\eta_x}{d\mu_x} = \frac{1}{\mu_x(1 - \mu_x)}$
- $\widehat{Var}(\beta) = (X^T W X)^{-1}$
- $V(\mu_x) = \frac{\mu_x(1 - \mu_x)}{n_x}$

➤ 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)

NTHU STAT 5230, 2025, Lecture Notes  
made by S.-W. Cheng (NTHU, Taiwan)

## Hypothesis Tests

• null model and saturated model

➤ null model: the smallest model we will entertain

- Model for no relation between predictors and response
- Usually, it means we fit a common mean  $\mu$  for all  $y_x$
- For some contingency table models, there will be additional parameters that represent row or column totals or other such constraints  $\Rightarrow$  null model has more than one parameter

➤ saturated (full) model: the most complex model

- Model in which data is explained exactly
- Typically,  $k$  parameters for  $k$  data points
- It can be achieved by fitting a sufficiently high-order polynomial or treating quantitative predictors as qualitative predictors or adding enough interactions
- The model tells us no more than the data itself and is usually uninformative

➤ A statistical model  $S$  describes how we partition the data into systematic structure and random variation

- Null model represents one extreme where the data is represented entirely as random variation
- Saturated model represents the data as being entirely systematic
- Model we want usually lie between these two extremes

#### • Deviance

➤ **Q:** how to measure discrepancy between observed and fitted  $y$ ?

➤ Saturated model gives us a measure of how well *any* model could possibly fit  $\Rightarrow$  can consider the difference between the log-likelihood for the saturated and a model  $S$  of interest:

$$2(l(Y, \phi; Y) - l(\hat{\mu}, \phi; Y))$$

(which has a rationale from likelihood-ratio test)

- $l(Y, \phi; Y)$  : the log-likelihood for the saturated model
- $l(\hat{\mu}, \phi; Y)$  : the log-likelihood for the model  $S$

NTHU STAT 5230, 2025, Lecture Notes  
made by S.-W. Cheng (NTHU, Taiwan)

➤ Provided that the observations are independent and for an exponential family distribution with  $a_i(\phi) = \phi/w_i$ ,

$$2(l(Y, \phi; Y) - l(\hat{\mu}, \phi; Y)) = \sum_i 2w_i \left[ y_i \left( \tilde{\theta}_{\mathbf{x}_i} - \hat{\theta}_{\mathbf{x}_i} \right) - b(\tilde{\theta}_{\mathbf{x}_i}) + b(\hat{\theta}_{\mathbf{x}_i}) \right] / \phi \equiv D(Y, \hat{\mu}) / \phi$$

where  $\tilde{\theta}_{\mathbf{x}}$  : the estimates of  $\theta_{\mathbf{x}}$  under the saturated model

$\hat{\theta}_{\mathbf{x}}$  : the estimates of  $\theta_{\mathbf{x}}$  under  $S$

➤  $D(Y, \hat{\mu})$  is called the *deviance* and  $D(Y, \hat{\mu})/\phi$  is called the *scaled deviance*

➤ Deviance for the common GLM

Family	deviance
Gaussian	$\sum_i (y_i - \hat{\mu}_{\mathbf{x}_i})^2$
Poisson	$2 \sum_i [y_i \log(y_i / \hat{\mu}_{\mathbf{x}_i}) - (y_i - \hat{\mu}_{\mathbf{x}_i})]$
Binomial	$2 \sum_i [y_i \log(y_i / \hat{\mu}_{\mathbf{x}_i}) + (n - y_i) \log((n - y_i) / (n - \hat{\mu}_{\mathbf{x}_i}))]$
Gamma	$2 \sum_i [-\log(y_i / \hat{\mu}_{\mathbf{x}_i}) + (y_i - \hat{\mu}_{\mathbf{x}_i}) / \hat{\mu}_{\mathbf{x}_i}]$
Inverse Gaussian	$\sum_i (y_i - \hat{\mu}_{\mathbf{x}_i})^2 / (\mu_{\mathbf{x}_i}^2 y_i)$



➤ Pearson's  $X^2$  statistic 
$$X^2 = \sum_i \frac{(y_i - \hat{\mu}_{\mathbf{x}_i})^2}{V(\hat{\mu}_{\mathbf{x}_i})}$$

it is an alternative measure of discrepancy that is sometimes used in replace of the deviance

- Goodness-of-fit test: whether the current model  $S$  fit the data

➤ Given the model  $S$  is correct,

- $D(Y, \hat{\mu}_S)/\phi \stackrel{a}{\sim} \chi_{df_S}^2$
- $X^2 \stackrel{a}{\sim} \chi_{df_S}^2$

➤ For Gaussian GLM, cannot use the test because do not know the value of the dispersion parameter  $\phi$

(**Q**: why not replace  $\phi$  by an estimate  $\hat{\phi}$  ?)

➤ For the binomial and the Poisson,  $\phi=1$ , so the test is practical

- Difference-in-deviance test: compare two nested models  $S \subset L$

➤ Given the model  $S$  is correct,

$$(D(Y, \hat{\mu}_S) - D(Y, \hat{\mu}_L))/\phi \stackrel{a}{\sim} \chi_{df_S - df_L}^2$$

NTHU STAT 5230, 2025, Lecture Notes  
made by S.-W. Cheng (NTHU, Taiwan)

➤ For the Gaussian model and other models where the dispersion  $\phi$  is *not known*, this chi-square test cannot be directly used

- We can insert an estimate of  $\phi$  and

$$\frac{D(Y, \hat{\mu}_S) - D(Y, \hat{\mu}_L)}{\hat{\phi}_L} \stackrel{a}{\sim} F_{df_S - df_L, df_L}$$

where  $\hat{\phi}_L = X_L^2/df_L$  and  $X_L^2$  is the Pearson's  $X^2$  statistic under the model  $L$

- For the Gaussian model,  $\hat{\phi}_L = RSS_L/df_L$  and the resulting  $F$ -statistic has an exact  $F$  distribution under the model  $S$

➤ Goodness-of-fit test:  $L$ =saturated model

- Notes:

➤ The null distribution in the goodness-of-fit and difference-in-deviance test is only asymptotically correct

➤ The approximation is better when comparing models than for the goodness of fit statistic

- Wald test for individual  $\beta_j$ :

- $\hat{\beta}_j / se(\hat{\beta}_j) \stackrel{a}{\sim} N(0, 1)$
- $\hat{\beta}_j / (se(\hat{\beta}_j)\hat{\phi}) \stackrel{a}{\sim} t_{df_L}$

- variable selection

- stepwise methods

- Can sequentially (forward or backward or a mix of both) apply difference-in-deviance test to compare nested models (in much the same manner as in standard regression models)
- The usual concerns about the validity of multiple testing and missing good model carry over

- criterion-based methods

$$\text{AIC}_S = \text{Deviance}_S + 2p$$

$$\text{BIC}_S = \text{Deviance}_S + p \log k$$

where  $p$  is the number of parameters in the model  $S$  and  $k$  is the number of covariate classes

- Choose the model with the smallest AIC or BIC
- AIC will tend to pick a larger model than the BIC
- AIC and BIC can be used to compare non-nested model