Generalized Linear Model: Introduction

CF Jeff Lin, MD., PhD.

April 9, 2006, 2004

© Jeff Lin, MD., PhD.

References

Nelder and Wedderburn (1972).
 Generalized linear models.
 J Roy. Statist. Soc. A135:370-384.

McCullagh and Nelder (1989).
 Generalized Linear Models, 2nd ed.
 Chapman and Hall.

Outlines

Statistical Models: systematic and random components

1. Suppose *n* observations with $\underline{\mathbf{X}} = (X_1, \dots, X_n)^T$ and $\underline{\mathbf{Y}} = (Y_1, \dots, Y_n)^T$ of two variables *X* and *Y*.

2.
$$y_i = \alpha + \beta x$$
, $i = 1, ..., n$.

- 3. Approximately linear relationship between y & x
- 4. makes the $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)$, (or fitted values $\hat{\boldsymbol{\mu}} = (\hat{\mu}_1, \dots, \hat{\mu}_n)^T$), close to the observed data.

Modelling: Science and Art

- 1. All models are wrong; some, though, are more useful than others
- 2. Not to fall in love with one model to the excluding of alternatives
- 3. Checks on the fit of a model to the data
- 4. Any statistical model has its own statistical assumptions

General Linear Model (GLM)

ANOVA, ANCOVA, Regression

$$\underline{\mathbf{y}} = \mathbf{X}\underline{\boldsymbol{\beta}} + \underline{\mathbf{e}}$$

$$\underline{\mathbf{e}} = \underline{\mathbf{y}} - \mathbf{X}\underline{\boldsymbol{\beta}}$$

$$y_i = \underline{\mathbf{x}}_i^T\underline{\boldsymbol{\beta}} + e_i$$

$$y_i \sim N(\underline{\mathbf{x}}_i^T\underline{\boldsymbol{\beta}}, \sigma^2), \quad i = 1, 2, \dots, n$$

$$\text{Minimizing} \quad \underline{\mathbf{e}}^T\underline{\mathbf{e}} = \sum_i e_i^2$$

(1)

(2)

(3)

(4)

Generalized Linear Model (GLIM)

1. Response variable $y_i \sim$ an exponential family distribution

- 2. $\mathcal{E}[Y_i] = \mu_i$
- 3. $\underline{\eta} = g(\underline{\mu}) = \mathbf{X}\underline{\beta}$

GLIM Terminology

1. Random component:

 $\underline{\mathbf{Y}}$ are i.i.d. with $\mathcal{E}[\underline{\mathbf{Y}}] = \underline{\mathbf{\mu}}$ and constant variance $\mathbf{Var}[\underline{\mathbf{Y}}] = \sigma^2(\theta, \phi)$

2. Systematic component:

Covariates: $\underline{\mathbf{x}}_1, \underline{\mathbf{x}}_2, \dots, \underline{\mathbf{x}}_p$. Linear predictor: $\underline{\boldsymbol{\eta}}$ Relationship: $\underline{\boldsymbol{\eta}} = \sum_1^p \underline{\mathbf{x}}_j \beta_j$

3. Link function: $g(\underline{\mu}) = \underline{\eta}$

Types of Link Functions

1. **logit**:
$$\eta = \log\left(\frac{\mu}{1-\mu}\right)$$

2. **probit**:
$$\eta = \Phi^{-1}(\mu)$$

3. complementary log-log: $\eta = \log \left[-\log \left(1 - \mu \right) \right]$

4. power family:

$$\eta = \begin{cases} \frac{\mu^{\lambda} - 1}{\lambda}; & \\ \log \mu; & \text{as } \lambda \to 0 \end{cases} \text{ or } \eta = \begin{cases} \mu^{\lambda}; & \lambda \neq 0, \\ \log \mu; & \lambda = 0. \end{cases}$$

1. $\underline{\Upsilon}$ has a distribution in the exponential family,

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

- 2. θ : canonical parameter
- 3. $b(\theta)$: cumulant function

(5)

Normal distribution

$$f_{Y}(y;\theta,\phi) = \frac{1}{\sqrt{2\pi\sigma^{2}}} \exp\left[\frac{-(y-\mu)^{2}}{2\sigma^{2}}\right] \qquad (6)$$
$$= \exp\left[\frac{(y\mu-\mu^{2}/2)}{\sigma^{2}} - \frac{1}{2}\left(\frac{y^{2}}{\sigma^{2}} + \log(2\pi\sigma^{2})\right)\right]$$

so $\theta = \mu$, $\phi = \sigma^2$, and $a(\phi) = \phi$, $b(\theta) = \frac{\theta^2}{2}$, $c(y, \phi) = -\frac{1}{2} \left(\frac{y^2}{\sigma^2} + \log(2\pi\sigma^2) \right)$

MLE and UMVUE

$$\ell(\theta,\phi;y) = \log f_{Y}(y;\theta,\phi) = \left[\frac{y\theta - b(\theta)}{a(\phi)} + c(y,\phi)\right],$$

$$\epsilon\left(\frac{\partial\ell}{\partial\theta}\right) = 0 \qquad (7)$$

$$\epsilon\left(\frac{\partial^{2}\ell}{\partial\theta^{2}}\right) + \epsilon\left(\frac{\partial\ell}{\partial\theta}\right)^{2} = 0 \qquad (8)$$

$$-\epsilon\left(\frac{\partial^{2}\ell}{\partial\theta^{2}}\right) = \epsilon\left(\frac{\partial\ell}{\partial\theta}\right)^{2} \qquad (9)$$

$$\ell(\theta,\phi;y) = \left[\frac{y\theta - b(\theta)}{a(\phi)} + c(y,\phi)\right], \qquad (10)$$

$$\frac{\partial\ell}{\partial\theta} = \frac{y - b'(\theta)}{a(\phi)} \equiv 0 \qquad (11)$$

$$\epsilon\left(\frac{\partial\ell}{\partial\theta}\right) = \frac{\mu - b'(\theta)}{a(\phi)} \equiv 0 \qquad (12)$$

$$\frac{\partial^{2}\ell}{\partial\theta^{2}} = \frac{-b''(\theta)}{a(\phi)} \qquad (13)$$

$$\epsilon\left(\frac{\partial\ell}{\partial\theta}\right)^{2} = \epsilon\left(\frac{y - b'(\theta)}{a(\phi)}\right)^{2} = \frac{\operatorname{Var}[Y]}{a^{2}(\phi)} \qquad (14)$$

$$\mathcal{E}[Y] = \mu = b'(\theta) \tag{15}$$

$$-\frac{b''(\theta)}{a(\phi)} + \frac{\mathbf{Var}[Y]}{a^2(\phi)} \equiv 0 \tag{16}$$

$$\mathbf{Var}[Y] = b''(\theta)a(\phi) = \mathbf{V}(\mu) \tag{17}$$

$$a(\phi): \text{ common form } a(\phi) = \frac{\phi}{w} = \frac{\sigma^2}{w}$$

$$\phi = \sigma^2: \text{ dispersion parameter}$$

$$w: \text{ prior weight}$$

Canonical Links

Sufficient statistic equal in dimension of $\boldsymbol{\beta}$ such that $\theta = \eta$

Normal	$\eta = \mu$
Poisson	$\eta = \log \mu$
Binomial	$\eta = \log[\pi/(1-\pi)]$
Gamma	$\eta = \mu^{-1}$

For the canonical links, the sufficient statistic is $\mathbf{X}^T \underline{\mathbf{Y}}$

Common Univariate Distribution

Table 1: Common Exponential Family Distribution

	Normal	Poisson	Binomial	Gamma
Notation	$N(\mu, \sigma^2)$	$\mathbf{P}(\mu)$	Bin $(n, \pi)/n$	${f G}(\mu,v)$
Dispersion: ϕ	$\phi = \sigma^2$	1	1/n	$\phi = v^{-1}$
Cumulant function: $b(heta)$	$\theta^2/2$	$\exp(\theta)$	$\log(1+e^{\theta})$	$-\log(- heta)$
$c(y; \phi)$	$-\frac{1}{2}\left(\frac{y^2}{\phi} + \log(2\pi\phi)\right)$	$-\log y$	$\log \binom{n}{ny}$	$v\log(vy) - \log y - \log \Gamma(v)$
$\mu(\theta) = \mathcal{E}(Y)$	θ	$\exp(\theta)$	$e^{\theta}/(1+e^{\theta})$	-1/ heta
Canonical link: $\theta(\mu)$	identity	log	logit	reciprocal
Variance function: $\mathbf{V}(\mu)$	1	μ	$\mu(1-\mu)$	μ^2

$$f(y_{i};\theta_{i},\phi) = \exp\left[\frac{y_{i}\theta_{i} - b(\theta_{i})}{a(\phi)} + c(y_{i},\phi)\right]$$
(18)
$$\ell(\underline{\beta}) = \sum_{i} \log f(y_{i};\theta_{i},\phi) = \sum_{i} \ell_{i}(\underline{\beta})$$
(19)

To obtain the likelihood equations, we calculate

$$\frac{\partial \ell_i}{\partial \beta_j} = \frac{\partial \ell_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_j}$$
(20)

Since
$$\mu_i = b'(\theta)$$
 and $\operatorname{Var}(Y_i) = b''(\theta_i)a(\phi)$, then

$$\frac{\partial \ell_i}{\partial \theta_i} = \frac{y_i - b'(\theta_i)}{a(\phi)} = \frac{y_i - \mu_i}{a(\phi)}$$
(21)

$$\frac{\partial \mu_i}{\partial \theta_i} = b''(\theta_i) = \frac{\operatorname{Var}(Y_i)}{a(\phi)}$$
(22)

Also , since $\eta_i = \sum_j \beta_j x_{ij}$,

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

(23)

Finally, $\eta_i = g(\mu_i)$, and $\partial \mu_i / \partial \eta_i$ depends on link function In summary,

$$\frac{\partial \ell_i}{\partial \beta_j} = \frac{(y_i - \mu_i)}{a(\phi)} \frac{a(\phi)}{\mathbf{Var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_i} x_{ij}$$
(24)

Likelihood Equations

$$\sum_{i=1}^{C} \frac{(y_i - \mu_i) x_{ij}}{\operatorname{Var}(Y_i)} \quad \frac{\partial \mu_i}{\partial \eta_i} = 0, \quad j = 1, \dots p.$$
(25)

Fisher's information matrix:

$$I = -\mathcal{E}\left(\frac{\partial^{2}\ell_{i}}{\partial\beta_{h}\partial\beta_{j}}\right) = -\mathcal{E}\left(\frac{\partial\ell_{i}}{\partial\beta_{h}}\frac{\partial\ell_{i}}{\partial\beta_{j}}\right)$$
(26)
$$= -\mathcal{E}\left[\frac{(Y_{i} - \mu_{i})x_{ih}}{\mathbf{Var}(Y_{i})} \frac{\partial\mu_{i}}{\partial\eta_{i}} \frac{(Y_{i} - \mu_{i})x_{ij}}{\mathbf{Var}(Y_{i})} \frac{\partial\mu_{i}}{\partial\eta_{i}}\right]$$
(27)
$$= \frac{-x_{ih}x_{ij}}{\mathbf{Var}(Y_{i})} \left(\frac{\partial\mu_{i}}{\partial\eta_{i}}\right)^{2} \text{ so that}$$
(28)
$$I = \mathcal{E}\left(\frac{\partial^{2}\ell(\underline{\beta})}{\partial\beta_{h}\partial\beta_{j}}\right) = -\sum_{i=1}^{C} \frac{x_{ih}x_{ij}}{\mathbf{Var}(Y_{i})} \left(\frac{\partial\mu_{i}}{\partial\eta_{i}}\right)^{2}$$
(29)

The Fisher's information matrix I, which has elements

$$\mathbf{I} = \mathcal{E}\left[\frac{-\partial^2 \boldsymbol{\ell}(\underline{\boldsymbol{\beta}})}{\partial \beta_h \partial \beta_j}\right] = \mathbf{X}^T \mathbf{W} \mathbf{X}$$

where \boldsymbol{W} is the diagonal matrix with elements

$$w_i = \frac{\left(\frac{\partial \mu_i}{\partial \eta_i}\right)^2}{\mathbf{Var}(Y_i)}$$

on the main diagonal.

(30)

(31)

Newton-Raphson method,

$$\underline{\boldsymbol{\beta}}^{(m+1)} = \underline{\boldsymbol{\beta}}^{(m)} + (\mathbf{J}^{(m)})^{-1} \mathbf{\mathcal{U}}^{(m)}$$
(32)

 $\boldsymbol{\mathfrak{I}}$ is the sample information matrix having elements

$$\partial^2 \boldsymbol{\ell}(\underline{\boldsymbol{\beta}}) / \partial \beta_h \partial \beta_j, \tag{33}$$

 $\boldsymbol{\mathfrak{U}}$ is the score vector having elements

$$\partial \boldsymbol{\ell}(\underline{\boldsymbol{\beta}}) / \partial \beta_{j},$$
 (34)

 $\mathbf{J}^{(m)}$ and $\mathbf{\mathcal{U}}^{(m)}$ are $\mathbf{\mathcal{J}}$ and $\mathbf{\mathcal{U}}$ evaluated at $\underline{\mathbf{\beta}} = \underline{\mathbf{\beta}}^{(m)}$.

Fisher scoring method

$$\underline{\boldsymbol{\beta}}^{(m+1)} = \underline{\boldsymbol{\beta}}^{(m)} + (\boldsymbol{I}^{(m)})^{-1} \boldsymbol{\mathcal{U}}^{(m)}$$
(35)

or
$$\mathbf{I}^{(m)}\underline{\boldsymbol{\beta}}^{(m+1)} = \mathbf{I}^{(m)}\underline{\boldsymbol{\beta}}^{(m)} + \boldsymbol{\mathfrak{U}}^{(m)}$$
 (36)

where $I^{(m)}$ is the m^{th} approximation for the estimated Fisher's information matrix. $I^{(m)}$ has elements - $\mathcal{E}[\partial^2 \ell(\underline{\beta}) / \partial \beta_h \partial \beta_j]$, evaluated at $\underline{\beta}^{(m)}$.

Recall:
$$I^{(m)}\underline{\beta}^{(m+1)} = I^{(m)}\underline{\beta}^{(m)} + \mathfrak{U}^{(m)}$$
 (37)

The right-hand side of (37) is the vector elements

$$\sum_{j} \left[\sum_{i} \mu_{i} \frac{x_{ih} x_{ij}}{\mathbf{Var}(Y_{i})} \left(\frac{\partial \mu_{i}}{\partial \eta_{i}} \right)^{2} \beta_{j}^{(m)} \right] + \sum_{i} \frac{(y_{i} - \mu_{i}^{(m)}) x_{ih}}{\mathbf{Var}(Y_{i})} \left(\frac{\partial \mu_{i}}{\partial \eta_{i}} \right)$$

where μ_i and $(\partial \mu_i / \partial \eta_i)$ are evaluated at $\underline{\beta}^{(m)}$.

Thus;
$$I^{(m)}\underline{\boldsymbol{\beta}}^{(m)} + \boldsymbol{\mathfrak{U}}^{(m)} = \mathbf{X}^T \mathbf{W}^{(m)} \underline{\mathbf{z}}^{(m)}$$
, (38)

where $\mathbf{W}^{(m)}$ is \mathbf{W} evaluated at $\underline{\boldsymbol{\beta}}^{(m)}$ and $\underline{\mathbf{z}}^{(m)}$ has elements

$$z_{i}^{(m)} = \sum_{j} x_{ij} \beta_{j}^{(m)} + (y_{i} - \mu_{i}^{(m)}) \left(\frac{\partial \eta_{i}^{(m)}}{\partial \mu_{i}^{(m)}}\right)$$
(39)
$$= \eta_{i}^{(m)} + (y_{i} - \mu_{i}^{(m)}) \left(\frac{\partial \eta_{i}^{(m)}}{\partial \mu_{i}^{(m)}}\right)$$
(40)

Fisher scoring have form

$$\left(\mathbf{X}^T \mathbf{W}^{(m)} \mathbf{X} \right) \underline{\boldsymbol{\beta}}^{(m+1)} = \mathbf{X}^T \mathbf{W}^{(m)} \underline{\mathbf{z}}^{(m)}.$$
(41)

The equations have solution

$$\underline{\boldsymbol{\beta}}^{(m+1)} = \left(\mathbf{X}^T \mathbf{W}^{(m)} \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{W}^{(m)} \underline{\mathbf{z}}^{(m)}.$$
(42)

The vector \underline{z} in this formulation is a linearized form of the link function a μ , evaluated at y,

$$g(y_i) \approx g(\mu_i) + (y_i - \mu_i)g'(\mu_i) + \cdots$$

$$\approx \eta_i + (y_i - \mu_i)(\partial \eta_i / \partial \mu_i) = z_i$$
(43)

© Jeff Lin, MD., PhD.

ML Estimation:

Iterative Re-Weighted Least Squares (IRWLS)

This z_i "adjusted" or "working" dependent variable \underline{z} has i^{th} element approximated by $z_i^{(m)}$ for the m^{th} cycle of the iterative scheme.

$$\underline{\boldsymbol{\beta}}^{(m+1)} = \left(\mathbf{X}^T \mathbf{W}^{(m)} \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{W}^{(m)} \underline{\mathbf{z}}^{(m)}.$$
(44)

We regress $\underline{\mathbf{z}}^{(m)}$ on the X with weight $\mathbf{W}^{(m)}$ to obtained a new estimate $\underline{\boldsymbol{\beta}}^{(m+1)}$. This estimate yields a new linear predictor value $\underline{\boldsymbol{\eta}}^{(m+1)} = \mathbf{X}\underline{\boldsymbol{\beta}}^{(m+1)}$ and a new-adjusted-dependent-variable value $\underline{\boldsymbol{z}}^{(m+1)}$ for the next cycle.

ML Estimation

The asymptotic covariance matrix of $\hat{\underline{\beta}}$ is the inverse of the information matrix, estimated by

$$\widehat{\mathbf{Cov}}(\widehat{\underline{\beta}}) = \left(\mathbf{X}^T \,\widehat{\mathbf{W}} \, \mathbf{X} \right)^{-1} \tag{45}$$

where $\widehat{\mathbf{W}}$ is evaluated at $\underline{\hat{\boldsymbol{\beta}}}$. The form of \mathbf{W} depends on the link chosen for the model.

Simplification for Canonical Links

$$\eta_i = \theta_i = \sum_j x_{ij} \beta_j \tag{46}$$

when $a(\phi)$ in the likelihood function is identical for all observations, the **kernel** of the log likelihood is $\sum y_i \theta_i$, which simplifies to

$$\sum_{i} y_i \left(\sum_{j} x_{ij} \beta_j \right) = \sum_{j} \beta_j \left(\sum_{i} y_i x_{ij} \right)$$
(47)

Sufficient statistics for estimating $\boldsymbol{\beta}$ in the **GLIM** are then

$$\sum_{i} y_i x_{ij}, \quad j = 1, \dots, p \tag{48}$$

Simplification for Canonical Links

For the canonical link,

$$\frac{\partial \mu_i}{\partial \eta_i} = \frac{\partial \mu_i}{\partial \theta_i} = \frac{\partial b'(\theta_i)}{\partial \theta_i} = b''(\theta_i)$$
(49)

so we simplify the likelihood equation as

$$\frac{\partial \ell_i}{\partial \beta_j} = \frac{(y_i - \mu_i)}{a(\phi)} \frac{a(\phi)}{\operatorname{Var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_i} x_{ij}$$

$$= \frac{(y_i - \mu_i)}{\operatorname{Var}(Y_i)} b''(\theta_i) x_{ij}$$

$$= \frac{(y_i - \mu_i)x_{ij}}{a(\phi)}$$
(50)
(51)

Simplification for Canonical Links

The second derivatives of the log likelihood

$$\frac{\partial^2 \ell_i}{\partial \beta_h \partial \beta_j} = -\frac{x_{ij}}{a(\phi)} \left(\frac{\partial \mu_i}{\partial \beta_h}\right)$$

These do not depend on the observations y_i , so

$$\frac{\partial^2 \boldsymbol{\ell}(\underline{\boldsymbol{\beta}})}{\partial \beta_h \partial \beta_j} = \mathcal{E}\left[\frac{\partial^2 \boldsymbol{\ell}(\underline{\boldsymbol{\beta}})}{\partial \beta_h \partial \beta_j}\right] \implies \boldsymbol{\Im} = \boldsymbol{I}$$
(54)

The Newton-Raphson and Fisher scoring algorithm are identical.

(53)

Goodness-of-Fit Tests

A saturated GLIM has many parameters as observations, giving a perfect fit. (vs. nested unsaturated (reduced) model)

- Let $\tilde{\theta}$ denote the estimate of θ for all the **saturated model**
- Let $\hat{\theta}$ denote the estimate of θ for the **reduced model**
- Then the ratio

 $-2\log \left(\frac{\text{maximum likelihood under reduced model}}{\text{maximum likelihood under saturated model}}\right)$

describes lack of fit.

Goodness-of-Fit

When the random component has $a(\phi) = \phi/w$, this measure equals

$$-2 \log(\text{likelihood ratio}) = 2 \sum w_i \left[\frac{y_i(\tilde{\theta}_i - \hat{\theta}) - (b(\tilde{\theta}) - b(\hat{\theta}))}{\phi} \right]$$
(55)
$$= \frac{D(\underline{\mathbf{y}}; \hat{\boldsymbol{\mu}})}{\phi}$$
(56)

 $\left(\frac{D(\underline{\mathbf{y}};\underline{\boldsymbol{\mu}})}{\phi}\right)$ is called the **scaled deviance**.

 $D(\mathbf{y}; \hat{\boldsymbol{\mu}})$ is called the **deviance**.

The greater the scaled deviance, the poorer the fit.

Goodness-of-Fit Tests

For two models, when the second is a special case of the first, the difference

$$D(\underline{\mathbf{y}}; \underline{\hat{\boldsymbol{\mu}}}_{2}) - D(\underline{\mathbf{y}}; \underline{\hat{\boldsymbol{\mu}}}_{1})$$

$$= 2\sum w_{i} \Big[y_{i}(\hat{\theta}_{1i} - \hat{\theta}_{2i}) - \Big(b(\hat{\theta}_{1i} - b(\hat{\theta}_{2i}) \Big) \Big] \sim \boldsymbol{\chi}^{2}$$
(57)

also has the form of the deviance.

Under regular conditions, **the difference in scaled deviances has approximately a chi-squared distribution**, with degrees of freedom equal to the difference between the number of parameters in the two models.

Testing Hypothesis

 $\underline{\boldsymbol{\beta}} = (\underline{\boldsymbol{\beta}}_1^T, \underline{\boldsymbol{\beta}}_2^T)^T$ where $\underline{\boldsymbol{\beta}}_1$ is the $q \times 1$ vector of coefficients of interest and $\underline{\boldsymbol{\beta}}_2$ is the $(p-q) \times 1$ vector of other covariate coefficients. Also we partition the information matrix into q and (p-q) rows and columns as follows:

$$\begin{aligned}
\mathbf{J}(\underline{\boldsymbol{\beta}})_{p \times p} &= \begin{pmatrix} \mathbf{J}_{11}(\underline{\boldsymbol{\beta}}) & \mathbf{J}_{12}(\underline{\boldsymbol{\beta}}) \\
\mathbf{J}_{21}(\underline{\boldsymbol{\beta}}) & \mathbf{J}_{22}(\underline{\boldsymbol{\beta}}) \end{pmatrix} \\
\text{where} \quad [\mathbf{J}_{rs}(\underline{\boldsymbol{\beta}})] &= \begin{pmatrix} -\frac{\partial^2 \boldsymbol{\ell}(\underline{\boldsymbol{\beta}})}{\partial \beta_r \partial \beta_s} \end{pmatrix}, \quad r, s = 1, \dots p.
\end{aligned}$$
(58)

Testing Hypothesis

Let

$$\mathbf{J}^{-1}(\underline{\boldsymbol{\beta}}) = \begin{pmatrix} \mathbf{J}^{11}(\underline{\boldsymbol{\beta}}) & \mathbf{J}^{12}(\underline{\boldsymbol{\beta}}) \\ \mathbf{J}^{21}(\underline{\boldsymbol{\beta}}) & \mathbf{J}^{22}(\underline{\boldsymbol{\beta}}) \end{pmatrix}$$

be the partition of its inverse.

Let $\hat{\boldsymbol{\beta}}_{(p \times 1)} = (\hat{\boldsymbol{\beta}}_1^T, \hat{\boldsymbol{\beta}}_2^T)^T$ be the partitioned maximum likelihood estimate for $\boldsymbol{\beta}$.

(59)

Testing Hypothesis

 $H_0: \underline{\beta}_1 = \underline{\beta}_{01}$ vs. $H_A: \underline{\beta}_1 \neq \underline{\beta}_{01}$

1. The Wald test statistic has the form

$$X_W^2 = (\underline{\hat{\beta}}_1 - \underline{\beta}_{01})^T [\mathbf{J}^{11}(\underline{\hat{\beta}})]^{-1} (\underline{\hat{\beta}}_1 - \underline{\beta}_{01}) \sim \boldsymbol{\chi}_q^2$$
(60)

Note that this statistic depends upon the entire vector $\underline{\hat{\beta}}$ in the inverse information calculation.
Testing Hypothesis

 $H_0: \underline{\beta}_1 = \underline{\beta}_{01}$ vs. $H_A: \underline{\beta}_1 \neq \underline{\beta}_{01}$

2. The Likelihood ratio test statistic is

$$X_{LR}^{2} = 2 \left\{ \ell(\underline{\hat{\beta}}) - \ell(\underline{\beta}_{01}, \underline{\hat{\beta}}_{2}(\underline{\beta}_{01})) \right\} \sim \chi_{q}^{2}$$
(61)

where $\underline{\hat{\beta}}_{2}(\underline{\beta}_{01})$ be the maximum partial likelihood estimate of $\underline{\beta}_{2}$ with $\underline{\beta}_{1}$ fixed at $\underline{\beta}_{01}$.

Testing Hypothesis

 $H_0: \underline{\beta}_1 = \underline{\beta}_{01}$ vs. $H_A: \underline{\beta}_1 \neq \underline{\beta}_{01}$

3. The score test statistic is

$$X_{SC}^{2} = \boldsymbol{\mathcal{U}}_{1}[\underline{\boldsymbol{\beta}}_{01}, \underline{\hat{\boldsymbol{\beta}}}_{2}(\underline{\boldsymbol{\beta}}_{01})]^{T} [\boldsymbol{\mathcal{I}}^{11}[\underline{\boldsymbol{\beta}}, \underline{\hat{\boldsymbol{\beta}}}_{2}(\underline{\boldsymbol{\beta}}_{01})] \boldsymbol{\mathcal{U}}_{1}[\underline{\boldsymbol{\beta}}_{01}, \underline{\hat{\boldsymbol{\beta}}}_{2}(\underline{\boldsymbol{\beta}}_{01})] \\ \sim \boldsymbol{\chi}_{q}^{2}$$

$$(62)$$

where $\mathcal{U}_1[\underline{\beta}_{01}, \underline{\hat{\beta}}_2(\underline{\beta}_{01})]$ is the $(q \times 1)$ of scores for $\underline{\beta}_1$, evaluated at the hypothesized value of $\underline{\beta}_{01}$ and at the restricted partial maximum likelihood estimator for β_2 .

Testing Hypothesis

Asymptotic, as $n \to \infty$ and $m \to \infty$, all three of these statistics have an approximate chi-squared distribution with q degree of freedom, when the null hypothesis is true.

Processes in Model Fitting of GLIM

- 1. Model selection
- 2. Parameter estimation
- 3. Prediction of future values

We must anticipate that, cluster around the "best" model will be a set of alternatives almost as good and not statistically distinguishable.

Model Selection

- 1. Assume independent (or at least uncorrelated) observations
- 2. Error structure: a single error term in the model
- 3. Scale depends on the purpose
- 4. Additivity effect
- 5. Choice of independent (X) variables (or covariates)

Estimation in Model Fitting

- 1. Maximizing the likelihood or log likelihood
- 2. Minimize the goodness-of-fit criterion

$$D^{\star}(\underline{\mathbf{y}};\underline{\boldsymbol{\mu}}) = 2[\boldsymbol{\ell}(\underline{\mathbf{y}};\underline{\mathbf{y}}) - \boldsymbol{\ell}(\underline{\boldsymbol{\mu}},\underline{\mathbf{y}})]$$
(63)

D^{*}(<u>y</u>; <u>μ</u>) the scaled deviance
 For normal-theory linear regression models with known variance σ², th deviance is identical to the residual of sum of squares and minimum deviance is synonymous with least squares.

Prediction in Model Fitting

- 1. Prediction
- 2. Calibration

Snoring and Heart Disease Example

- 1. Table 2 is from an epidemiological survey of 2484 subjects to investigate snoring as a risk factor for heart disease.
- 2. Those surveyed were classified according to their spouses' report how much they snored.
- 3. The model states that the probability of heart disease $\pi(x)$ is how much they snored related to the level of snoring x.

Snoring and Heart Disease Example

- 1. We treat the rows of the table as independent binomial samples with that probability as the parameter.
- 2. We use score (0, 2, 4, 5) for the snoring categories, treating the last two levels as closer than the other adjacent.

Snoring and Heart Disease Example

Table 2:RelationshipbetweenSnoring and Heart Disease

	Heart Disease	
Snoring	yes	No
Never	24	1355
Occasionally	35	603
nearly every night	21	192
Every night	30	224

Linear Probability Model

 One approach to modelling the effect of X uses the form of ordinary regression, by which the expected value of Y is a linear function of X. The model

$$\pi(x) = \alpha + \beta x \tag{64}$$

is called a **linear probability model**, because the probability of success changes linearly in x.

- 2. The parameter β represents the change in the probability per unit change in x.
- 3. This is a GLIM with binomial random component and identity link function.

Linear Probability Model: Disadvantages

- 1. Unfortunately, this model has a major structural defect.
- 2. Probabilities fall between 0 and 1, whereas linear functions take values over the entire real line.
- 3. This model can be valid over a finite range of x values.

Linear Probability Model: Disadvantages

- 1. It looks like an ordinary regression, least squares estimators of the model parameters are not optimal.
- The variance of the binary outcome for each subject,
 Var(Y) = π(x)[1 π(x)], is not constant for all x, but rather depends on x through its influence on π(x).
- 3. Because of the non-constant variance, maximum likelihood (ML) estimators for this model, like most GLIM, can have smaller standard errors than least squares estimators.

Linear Probability Model: Snoring Data

1. For the snoring data the linear probability model is

 $\hat{\pi}(x) = 0.0172 + 0.0198x.$

- 2. The estimated probability of heart disease is about 0.02 for nonsnorers, it increases 2(0.0198) = 0.04 for occasional snores.
- 3. The standard error of the slope estimate of 0.0198equals 0.0028.

(65)

Logistic Regression Model

1. In practice, nonlinear relationships between $\pi(x)$ and x are often monotonic, with $\pi(x)$ increasing continuously as x increases, or $\pi(x)$ decreasing continuously as x increases.

2. The S-shaped curves displayed are often realistic shapes for the relationship.

© Jeff Lin, MD., PhD.



Figure 1: Logistic Regression Function

Logistic Regression Model

1. The most important function having this shape has the model form

$$\log\left(\frac{\pi(x)}{1-\pi(x)}\right) = \alpha + \beta x.$$
(66)

Logistic Regression Model

- 1. This is called the **logistic regression function**, and is often called **logit model**.
- 2. In GLIM ,the random component for the (success, failure) determinations is binomial.
- 3. The link function is the logit transformation $\log[\pi/(1-\pi)]$, symbolized by logit (π) .
- 4. The logit is the natural parameter of the binomial distribution, so the logit link is its canonical link.

Logistic Regression Model: Snoring Data

1. For the snoring data the logistic regression model is

$$logit [\hat{\pi}(x)] = -3.87 + 0.397x.$$

2. This gives a model of

$$\operatorname{logit}[\pi(x)] = \log\left(\frac{\pi(x)}{1 - \pi(x)}\right) = \alpha + \beta x,\tag{68}$$

where x = snore with *ML* estimates $\alpha = .3.87$ and $\beta = 0.397$.

(67

Logistic Regression Model: Snoring Data

- 1. The positive value of $\hat{\beta} = 0.40$ reflects the increased chance of heart disease at higher levels of snoring.
- 2. Since $\pi(x) = \frac{\exp(\alpha + \beta x)}{1 \exp(\alpha + \beta x)}$, the estimated probability of heart disease is about 0.02 for nonsnorers, it increases 2(0.0198) = 0.04 for occasional snores.

- 1. Let X denote random variable, and let x denote a potential value for X.
- 2. The cumulative distribution function (*cdf*) F(x) for X is defined as

$$F(x) = P(X \le x), \quad -\infty < x < \infty.$$
(69)

3. Such a function, plotted as a function of *x*, has appearance, like S-shaped.

- 1. As x increases, F(x) increases gradually from 0 to 1, since $P(X \le x)$ increases as x increases.
- 2. This subsets a class of models for binary responses whereby the dependence of $\pi(x)$ on x has form

$$\pi(x) = F(x),\tag{70}$$

where F is a cdf for some distribution.

- 1. The logistic regression curve has this form.
- 2. When $\beta > 0$, F(x) is the *cdf* of a two-parameter logistic distribution.
- 3. When $\beta < 0$, the formula for $1 \pi(x)$ has the logistic *cdf* appearance.
- 4. Each choice of α and of $\beta > 0$ corresponds distribution with a symmetric, bell shape.
- In fact, it looks similar to a normal distribution but with slightly thicker tail. Model from (70) occurs naturally when a tolerance distribution applies to subjects' response.

- 1. For instance, in a toxicology study, suppose that researchers spray an insectide at various dosage levels on batches of mosquitoes.
- 2. For each mosquito, the response is whether it dies.
- 3. Each mosquito may have a certain tolerance to the insecticide, such that it dies if the dosage level exceeds its tolerance and survives if the dosage level is less than its tolerance

- 1. Tolerance would vary among mosquitoes.
- If a *cdf F* describes the distribution of tolerances, then the model for the probability π(x) of death at dosage level x necessarily has form (70).

- 1. When F is the cdf of a normal distribution, model type (70) is called the **probit model**.
- 2. The link function for the model is then called the **probit link**.
- 3. The probit model has alternative expression

$$\operatorname{probit}\left[\pi(x)\right] = \alpha + \beta x. \tag{71}$$

4. The probit link applied to a probability $\pi(x)$ transforms it to the standard normal *z*-score at which the left-tail probability equal to $\pi(x)$.

- 1. For instance, probit (0.05) = -1.645, probit (0.5) = 0, probit (0.95) = 1.645, and probit (0.975) = 1.96.
- 2. The probit model is a GLM with binomial random component and probit link.

1. For the snoring data the probit regression model is

probit $[\hat{\pi}(x)] = -2.061 + 0.188x.$

- 2. At snoring level x = 0, the fitted probit equals -2.061 + 0188(0) = -2.061.
- 3. The fitted probability $\hat{\pi}(x)$ is the left-tail probability for the standard normal distribution at -2.061, which equals 0.02.
- 4. At snoring level x = 5, the fitted probability equals -2.061 + 0188(5) = -1.12, which corresponding to a fitted probability 0.131.

(72)

- 1. For practical purposes, probit and logistic regression curves Look the same.
- 2. It is rare, and requires enormous sample sizes, to find data for which a logistic regression model fits well but the probit model fits poorly, or conversely.
- 3. When both model fit well, slope estimates in logistic regression models are roughly about 1.6 2.0 times those in probit model.

- 1. The probit transform maps $\pi(x)$ so that the regression curve for $\pi(x)$ (or $1 \pi(x)$), when $\beta < 0$ has appearance of the normal *cdf* with mean $-\frac{\alpha}{\beta}$ and standard deviation $\sigma = \frac{1}{|\beta|}$.
- 2. For the snoring data, the probit fit corresponds to a normal *cdf* having mean of $-\frac{\hat{\alpha}}{\hat{\beta}} = \frac{2.061}{0.188} = 11.0$ and standard deviation of $\frac{1}{|\beta|} = \frac{1}{0.188} = 5.3.$

- 1. The predicted probability of heart disease equals $\frac{1}{2}$ at snoring level x = 11.0; that is x = 11.0 has a fitted probit -2.061 + 0.188(11) = 0, which is the z-score corresponding to a left-tail probability of $\frac{1}{2}$.
- 2. The fitted probit value of -2.061 at x = 0 means that 0 is 2.06 standard deviations below the mean of a normal distribution with mean 11.0 and standard deviation 5.3.
- 3. Since snoring level is restricted to the range 0 5, for these data, well below 11.0, then fitted probability over this range are quite small.



Figure 2: Predicted probabilities for logistic, probit and linear regression models.

© Jeff Lin, MD., PhD.

Poisson Regression Model

- 1. The Poission distribution has a ppositive mean, it is more common to model the log of the mean.
- 2. Like the linear predictor $\alpha + \beta x$, the log of the mean can take any real value.
- 3. The log mean is the natural parameter for the poisson distribution, and the log link is the canonical link with a Poisson random component.
- A Poisson loglinear model is a GLIM that assume a Poisson distribution for Y and uses the log link.

© Jeff Lin, MD., PhD.

Poisson Regression Model

- 1. Let μ denote the expected value for a Poisson variate, Y, and let X denote an explanatory variable.
- 2. The Poisson loglinear model has form

$$\log(\mu) = \alpha + \beta x. \tag{73}$$

3. For this model, the man satisfies the exponential relationship

$$\mu = \exp(\alpha + \beta x) = e^{\alpha} (e^{\beta})^x.$$
(74)

Poisson Regression Model

1. A one-unit increase in X has a multiplicative impact of e^{β} on μ .

- 2. The mean of Y at x = 1 equals the mean of Y at x multiplied by e^{β} .
- 3. If $\beta = 0$, then $e^{\beta} = e^0 = 1$ and the multiplicative factor is 1; that is, the mean of Y does not changes as X changes.

4. If $\beta > o$, then $e^{\beta} > 1$, and the mean of Y increases as X increases.

5. If $\beta < o$, then $e^{\beta} < 1$, and the mean of Y decreases as X increases.

Horseshoe Crab Data

- 1. Table 3 (in file crab.txt) is a study of nesting horseshoe crabs is that each female horseshoe crab had a male crab resident in her nest.
- 2. The study investigated factors affecting whether the female crab had any other males, called **satellites**, residing nearby.
- 3. The response outcome for each female crab is her number of satellites.
- 4. Explanatory variables are the female crab's color, spine condition, carapace width, and weight.
- 5. This data set comes from a study on 173 female horseshoe crabs.
Horseshoe Crab Data

Table 3: Variables descriptions of Crab Data

Variable	Description
С	= color (light-medium, medium, dark-medium)
S	= spine condition (both good, one worn or broken, both broken)
W	= width of carapace in cm
Wt	= weight in kg
Sa	= number of satellites (male residing nearby)

Horseshoe Crab Data

- 1. Figure 3 plots the response counts of satellites against width, with numbered symbols indicating the number of observations at each point.
- 2. The substantial variability makes it difficult to discern a clear trend.

Crab Data: satellites by width of female crab



Figure 3: Number of satellites by width of female crab.

Horseshoe Crab Data

- To get a clear picture, we grouped the female crabs into width categories (≤ 23.25, 23.25 24.25, 24.25 25.25, 25.25 26.25, 26.25 27.25, 27.25 28.25, 28.25 29.25, > 29.25) and calculated the sample mean number of satellites for female crabs in each category.
- 2. Figure 4 plots these sample mean against the sample mean width in each category.

Smoothings of crab counts



Figure 4: Number of satellites by width of female crab.

- 1. The Goal of this study was to find a model for the number of satellites.
- 2. Let μ denote the expected number of satellites for a female crab, and let x denote her width.
- 3. Assuming that the number of satellites has a Poisson distribution, we get a model (73) (using the log link) of

$$\log(\mu) = \alpha + \beta x,\tag{75}$$

where $\alpha = -3.305$ and $\beta = 0.164$.

- 1. The effect $\hat{\beta} = 0.164$ of width has an asymptotic (large-sample) standard error [s.e. = 0.020.
- 2. Since $\hat{\beta} > 0$, width has a positive estimated effect on the number of satellites.

- 1. The model fitted value at any width level is an estimated mean number of satellites, $\hat{\mu}$.
- 2. The **fitted value**, $\hat{\mu}$ from $\hat{\mu} = \exp(\hat{\alpha} + \hat{\beta}x)$, at the mean width of x = 26.3 is

$$\hat{\mu} = \exp(\hat{\alpha} + \hat{\beta}x) = \exp[-3.305 + 0.164(26.3)] = 2.74.$$
 (76)

- 1. For this model, $\exp[\hat{\beta}] = \exp[0.164] = 1.18$ represents the multiplicative effect on the fitted value for each 1-unit increases in x.
- 2. For instance, the fitted value at x = 27.3 = 26.3 + 1 is exp[-3.305 + 0.164(27.3)] = 3.23, which equals $1.18 \times 2.74 = 3.23$.
- 3. A 1-cm increase in width yields an 18% increase in the estimated mean number of satellites.

- 1. Figure 4 shows that $\mathcal{E}(Y)$ may grow approximately linearly with width.
- 2. This suggests the Poisson GLIM identity link.
- 3. It has ML fit

$$\mu = \alpha + \beta x, \tag{77}$$

where $\alpha = -11.525$ and $\beta = 0.55$ (s.e. = 0.058).

- 1. The effect of X on μ in this model is additive, rather than multiplicative.
- 2. A 1-cm increase in width has a predicted increase of $\hat{\beta} = 0.55$ in the expected number of satellites.
- 3. For instance, the fitted value at the mean width of x = 26.3 is $\hat{\mu} = -11.53 + 0.55(26.3) = 2.93$; at x = 27.3, it is 2.93 + 0.55 = 3.48.

- 1. The fitted value are positive at all widths observed in the sample, and the model provides a simple description of the effect of width on the number of satellites:
- 2. Approximately 2 cm increase in width give an additional satellite male.

- 1. A comparison of the two models' predictions, Figure 5 plots $\hat{\mu}$ against width for the models with log link and identity link.
- 2. Although they diverge somewhat for relatively small and large widths, they provide similar predictions over the width range in which most observations occur.

Compare log and identity link for Crab Data



Figure 5: Estimated mean number of satellites for log and identity links for Crab Data.

© Jeff Lin, MD., PhD.