

GLM summary

Recall:

Linear Regression Model

→ $y \sim N(X\beta, \Sigma)$ where $\Sigma = \begin{pmatrix} \sigma^2 & & \\ & \ddots & \\ & & \sigma^2 \end{pmatrix}$

↓ response

$y_i = \beta_0 + x_i^T \beta + \varepsilon_i$

↓

$E[y_i] := \mu_i = \beta_0 + x_i^T \beta$

$X = \begin{pmatrix} 1 & x_1^T \\ \vdots & \vdots \\ 1 & x_n^T \end{pmatrix}$

$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$

$$\Rightarrow \log f(y; \beta) = \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(y_i - x_i^T \beta)^2}{2\sigma^2} \right) \right)$$

$$= \sum_{i=1}^N \left[\log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{(y_i - x_i^T \beta)^2}{2\sigma^2} \right]$$

$$\Rightarrow \begin{cases} \hat{\beta} = \frac{(X^T X)^{-1} X^T Y}{\sum_{i=1}^N (y_i - \hat{y}_i)^2} \\ \hat{\sigma}^2 = \frac{RSS}{N} \end{cases} \quad \begin{cases} E[\hat{\beta}] = (X^T X)^{-1} X^T E[Y] = (X^T X)^{-1} X^T X \beta = \beta \\ Var[\hat{\beta}] = (X^T X)^{-1} X^T Var[Y] X (X^T X)^{-1} \\ = \sigma^2 (X^T X)^{-1} \end{cases}$$

Modification

$$E[\hat{\sigma}^2_{MLE}] = \frac{N-p-1}{N} \sigma^2 \Rightarrow \text{Biased estimator}$$
$$\hat{\sigma}^2_{unbias} = \frac{RSS}{N-p-1}$$

Decomposition of $\sum_{i=1}^N (y_i - \bar{y})^2$

only holds for LR without Regularization Term

$$\sum_{i=1}^N (y_i - \bar{y})^2 = \sum_{i=1}^N (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2$$

$$= \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \sum_{i=1}^N (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^N 2(y_i - \hat{y}_i)(\hat{y}_i - \bar{y})$$

$$\Rightarrow TSS = RSS + \text{Regression sum of squares error}$$

Property:

$$\begin{cases} \text{rank}(H) = p+1 \\ H = X(X^T X)^{-1} X^T \in \mathbb{R}^{n \times n} \end{cases}$$

$$HX = X \Rightarrow H \cdot \mathbb{1}_n = \mathbb{1}_n$$

$$\begin{aligned} \hat{\varepsilon}^T (H - \frac{1}{n} \mathbb{I}) y \\ = y^T (I - H) (H - \frac{1}{n} \mathbb{I}) y \end{aligned}$$

① $R^2 = 1 - \frac{RSS}{TSS}$
 $\in [0, 1]$

② when model is larger (more features), R^2 must be bigger !!!

do feature selection

GLM → when response is no longer Gaussian distributed given the params

1) $\eta_i = x_i^T \beta$ → linear predictor [systematic component]

2) $g(\mu_i) = \eta_i$ → link function

3) $y_i \sim \text{exp. family}$

random component
(distributed of response)

$$f(y_i | \theta_i, \phi) = \exp \left\{ \frac{\theta_i y_i - b(\theta_i)}{a_i(\phi)} + c(\phi, y_i) \right\}$$

$$\begin{aligned} \downarrow \\ \left\{ \begin{aligned} \mathbb{E}[Y_i] &= b'(\theta_i) = \mu_i \\ \text{Var}[Y_i] &= a_i(\phi) b''(\theta_i) \end{aligned} \right. \\ \theta_i \xleftarrow{(b')^{-1}(\mu_i)} \mu_i \xleftarrow{g^{-1}(\eta_i)} \eta_i \xleftarrow{x_i^T \beta} \beta \end{aligned}$$

⇒ if canonical link function is used,

then $g(\mu_i) = \theta_i \Rightarrow g = (b')^{-1}$

$$\theta_i = \eta_i = x_i^T \beta$$

$$\Rightarrow f(y_i; \theta_i, \phi) = \exp \left\{ \frac{\beta^T x_i y_i - b(\beta^T x_i)}{a_i(\phi)} + c(y_i, \phi) \right\}$$

Conclusion: (for canonical link)

① $X^T y$ is a Sufficient Statistic for β

② $J(\beta) = - \frac{\partial^2 \log f}{\partial \beta^2}$ & $I(\beta) = \mathbb{E}[J(\beta)]$

observed info matrix

Fisher info matrix

we have $J(\beta) \equiv I(\beta)$

③ Newton-Raphson Method



Fisher-Scoring Algo

Notation $\rightarrow S(\beta) = \frac{\partial \log f}{\partial \beta}$

$$\begin{aligned} l(\beta) &:= \sum_{i=1}^n l_i(\beta) \\ &= \sum_{i=1}^n \log \{ f(y_i | x_i, \beta) \} \end{aligned}$$

\rightarrow Newton-Raphson Method

$$\beta^{k+1} = \beta^k + [J(\beta^k)]^{-1} \cdot S(\beta^k)$$

\rightarrow Fisher-Scoring Method

$$\beta^{k+1} = \beta^k + [I(\beta^k)]^{-1} S(\beta^k)$$

$$S(\beta) = \sum_{i=1}^n s_i(\beta)$$

$$= \sum_{i=1}^n \frac{\partial l_i}{\partial \beta}$$

$$= \sum_{i=1}^n \frac{\partial}{\partial \beta} \log \{ f(y_i | x_i, \beta) \}$$

$$\mathbb{E}_{\beta} [S(\beta)] = 0$$

when have

$$\beta^k$$

$$\downarrow$$

$$\eta^k = x \beta^k$$

$$\downarrow$$

$$g(\mu^k) = \eta^k$$

$$\downarrow$$

$$b'(\theta^k) = \mu^k$$

S(β)

$$\frac{\partial \log f(\beta)}{\partial \beta_j} = \sum_{i=1}^n \frac{y_i - b'(\theta_i)}{a_i(\phi)} \cdot \frac{\partial \theta_i}{\partial \beta_j}$$

$$= \sum_i \frac{y_i - \mu_i}{a_i(\phi)} \frac{\partial \theta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta_j}$$

$$= \sum_i \frac{y_i - \mu_i}{a_i(\phi)} \cdot \frac{1}{b''(\theta_i)} \cdot \frac{1}{g'(\mu_i)} \cdot x_{ij}$$

$$= \frac{1}{\phi} \sum_i (y_i - \mu_i) \cdot w_i \tilde{w}_i(\mu_i) x_{ij}$$

$$\Rightarrow S(\beta) = (x_1, \dots, x_n)^T \tilde{W} (Y - \mu)$$

$$= X^T \tilde{W} (Y - \mu)$$

$$\begin{aligned} \tilde{w}_i(\mu_i) &= \frac{1}{b''(\theta_i) g'(\mu_i)} \\ &= \frac{1}{V(\mu_i) g'(\mu_i)} \end{aligned}$$

① $\frac{\partial^2 \log f}{\partial \beta_j \partial \beta_k} = \frac{1}{\phi} \sum_i w_i x_{ij} \left[(y_i - \mu_i) \frac{\partial \tilde{w}_i}{\partial \beta_k} - \frac{\partial \mu_i}{\partial \beta_k} \cdot \tilde{w}_i(\mu_i) \right]$

$$\Rightarrow \mathbb{E} \left[\frac{\partial^2 \log f}{\partial \beta_j \partial \beta_k} \right] = \frac{1}{\phi} \sum_i -w_i \tilde{w}_i \frac{\partial \mu_i}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta_k} x_{ij}$$

$$= \frac{1}{\phi} \sum_i -w_i \tilde{w}_i(\mu_i) \frac{1}{g'(\mu_i)} \cdot x_{ik} x_{ij}$$

$$= \sum -\hat{w}_i x_{ik} x_{ij}$$

$$\Rightarrow \mathcal{I}(\beta) = \mathbb{E}[J(\beta)] = [x_1, \dots, x_n] \cdot \hat{W} \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

$$= X^T \hat{W} X$$

② When choose canonical link function, then $\tilde{w}_i(\mu_i) = \frac{1}{b''(\theta_i) \cdot g'(\mu_i)}$

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

$$= \frac{\partial \theta_i}{\partial \mu_i} \cdot \frac{\partial \mu_i}{\partial \eta_i}$$

$$= \frac{\partial \theta_i}{\partial \eta_i} = 1$$

observed info matrix

$$\Rightarrow \frac{\partial^2 \log f}{\partial \beta_j \partial \beta_k} = \frac{1}{\phi} \sum_i -\frac{\partial \mu_i}{\partial \beta_k} \cdot \tilde{w}_i w_i x_{ij}$$

Coincidence

$$= \mathbb{E} \left[\frac{\partial^2 \log f}{\partial \beta_j \partial \beta_k} \right]$$

Fisher info matrix

when picking canonical link function

\Rightarrow Newton-Raphson Algo \Leftrightarrow Fisher-scoring algo

Details of Fisher-Scoring Algo

$$\beta^{k+1} = \beta^k + [\mathcal{I}(\beta^k)]^{-1} S(\beta^k)$$

$$= \beta^k + (X^T \hat{W} X)^{-1} X^T \tilde{W} (y - \mu^k)$$

$$= (X^T \hat{W} X)^{-1} (X^T \hat{W} X \beta^k + X^T \tilde{W} (y - \mu^k))$$

$$= (X^T \hat{W} X)^{-1} X^T \hat{W} (\eta^k + \hat{W}^{-1} \tilde{W} (y - \mu^k))$$

where

$$\tilde{W} = \text{diag} \left(\frac{w_i}{\phi b''(\theta_i) g'(\mu_i)} \right)$$

$$\hat{W} = \text{diag} \left(\frac{w_i}{\phi b''(\theta_i) g'(\mu_i)^2} \right)$$

$$\eta^k = X \beta^k$$

$$:= (X^T \hat{W} X)^{-1} X^T \hat{W} z^k$$

where $z^k = \eta^k + \hat{W}^{-1} \tilde{W} (y - \mu)$

linear predictor of (i-th) sample
based on k-iter β

$$\eta_i^k = x_i^T \beta^k$$

$$g(\mu_i^k) = \eta_i^k = x_i^T \beta^k$$

can be viewed as:

$$\mathcal{D} = \{(x_i, z_i^k)\}_{i=1}^n$$

$$\& z^k \sim N(X\beta, \hat{W})$$

$$\hat{\beta}_{MLE} \rightarrow \beta^{k+1}$$

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

$$\approx g(y_i)$$

$g(y_i)$ 在 μ_i^k 处的 Taylor approximation $\rightarrow z_i^k$

adjusted dependent variable

Fisher-scoring Algo

(Iterative WLS)

① initial β^0 , then μ^0 comes from $g(\mu^0) = X^T \beta^0$

② given β^k & μ^k , do the following calculation:

a) $z_i^k = \eta_i^k + g'(\mu_i^k) \cdot (y_i - \mu_i^k)$

where $\eta_i^k = x_i^T \beta^k$

b) $(\hat{W}^k)_i = \frac{w_i}{\phi V(\mu_i^k) g'(\mu_i^k)^2}$

c) $\beta^{k+1} = \hat{\beta}_{MLE}$ for $z^k \sim N(X\beta, W^k)$

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

Weighted Least Square Solution

adjusted dependent variable!

Remark: for each iteration, (k-th)

z^k (response) is updated

W^k (variance) is updated

(IWLS) \Rightarrow iterative reweighted LS

Calculation is based on β^k (and $\mu^k = g^{-1}(X\beta^k)$)