

➤ Approach 2: (profile likelihood-based method)

- Parameters: $\underline{\theta} = (\underline{\theta}_1, \underline{\theta}_2)$ of current interest
- Likelihood: $\mathcal{L}(\underline{\theta}) = \mathcal{L}(\underline{\theta}_1, \underline{\theta}_2)$

MLE of $\underline{\theta}_2$ when $\underline{\theta}_1$ is fixed. It's a function of $\underline{\theta}_1$

focus only on $\underline{\theta}_1$

Profile likelihood: $\mathcal{L}^*(\underline{\theta}_1) = \max_{\underline{\theta}_2} \mathcal{L}(\underline{\theta}_1, \underline{\theta}_2) = \mathcal{L}(\underline{\theta}_1, \hat{\underline{\theta}}_2(\underline{\theta}_1))$

fixed first

$$\underline{\beta}: \hat{\underline{\beta}}_\lambda = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T [\mathbf{t}_\lambda(\mathbf{Y})] \quad \leftarrow \quad \sigma^2: \hat{\sigma}_\lambda^2 = \text{RSS}_\lambda / (n-p)$$

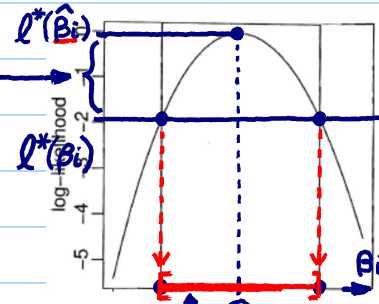
LNp.1-32~33

(recall: the computation of the C.I. for λ in Box-Cox method)

transformation parameter

profile log-likelihood

- For a coefficient β_i , its profile likelihood-based C.I. is



$$2[l^*(\hat{\beta}_i) - l^*(\beta_i^*)] < \chi^2_1(1-\alpha)$$

$$\{\beta_i : l^*(\beta_i) > l^*(\hat{\beta}_i) - (1/2)\chi^2_1(1-\alpha)\}$$

MLE under l^*

compared to Wald test based C.I.

other β_j 's, $j \neq i$, set to the MLE when β_i fixed (i.e., MLE as a function of β_i)

C.I. of β_i $\hat{\beta}_i \leftarrow$ MLE

$H_0: \beta_i = \beta_i^*$ vs. $H_1: \beta_i \neq \beta_i^*$

any constant

e.g.,
• Wald: AB
• profile: some β_i 's

The profile likelihood method is generally preferable for the same Hauck-Donner reason

Acceptance region of LR test under profile likelihood:

Similar method can be generalized to construct

confidence region of several parameters

MLE under profile likelihood

$$-2 \log \left[\frac{l^*(\beta_i^*)}{l^*(\hat{\beta}_i)} \right] < \chi^2_1(1-\alpha)$$

- Some notices about deviance when n_i is small

asymptotic approaches not work

sparse data

$\text{se}(\hat{\beta}_i)$ might be over-estimated (LNp.3-12)

Take the extreme case $n_i = 1$ (binary data)

Each covariate class has only one unit

$$\mathcal{L}(\hat{\underline{\beta}}_{\text{sat}}) = \prod_i \hat{p}_{i,\text{sat}}^{y_i} \times (1 - \hat{p}_{i,\text{sat}})^{1-y_i}$$

Then, $y_i = 0$ or 1 . For a model S (i.e., $\mathbf{X}_S \underline{\beta}_S$) under logit link,

$\sim \text{Bernoulli}(p_i)$

$$D_S = -2 \sum_{i=1}^k \{ \hat{p}_{i,S} \times \text{logit}(\hat{p}_{i,S}) + \log(1 - \hat{p}_{i,S}) \}$$

Check in LNp.3-5

$$\hat{p}_{i,\text{sat}} = y_i / n_i$$

$$-2 \log \left[\frac{\mathcal{L}(\hat{\underline{\beta}}_S)}{\mathcal{L}(\hat{\underline{\beta}}_{\text{sat}})} \right] = -2 \log \left[\frac{\mathcal{L}(\hat{\underline{\beta}}_S)}{\mathcal{L}(\hat{\underline{\beta}}_{\text{sat}})} \right] = -2 l(\hat{\underline{\beta}}_S)$$

When $\mathbf{X}_S^T \mathbf{Y} = \mathbf{X}_S^T \hat{\mathbf{u}}_S$, say $\mathbf{X}_S^T \mathbf{Y} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_k \end{bmatrix}^T \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$, can arbitrarily change \mathbf{Y} , but keep $\mathbf{X}_S^T \mathbf{Y}$ invariant \Rightarrow same D_S for different \mathbf{Y}

This also explains why need to form covariate classes

$$\hat{\underline{\beta}}_S^T \mathbf{Y} = (\mathbf{X}_S \hat{\underline{\beta}}_S)^T \mathbf{Y} = \hat{\underline{\beta}}_S^T \mathbf{X}_S^T \mathbf{Y} = \hat{\underline{\beta}}_S^T \mathbf{X}_S^T \hat{\mathbf{u}}_S = \hat{\underline{\beta}}_S^T \hat{\mathbf{u}}_S = \hat{\underline{\beta}}_S^T (\mathbf{n} \circ \hat{\underline{\mathbf{P}}}_S) = \hat{\underline{\beta}}_S^T \hat{\underline{\mathbf{P}}}_S$$

Why can it overcome the problem in D_S for sparse data?

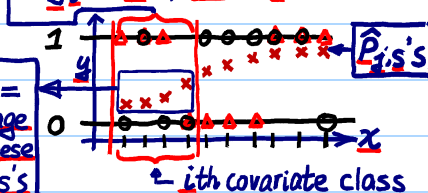
For a D_S to measure goodness-of-fit, it has to compare the fitted values $\hat{y}_{i,S} = \hat{u}_{i,S} = \hat{p}_{i,S}$ to the data y_i (cf., in linear model, we use $y_i - \hat{y}_i$ to compare), but here we have only a function of $\hat{p}_{i,S}$

result in GofE test is suspicious $\because D_S / \hat{\underline{\beta}}_S$ is a degenerate r.v.

Other methods must be used to judge goodness-of-fit for binary data or the case

of small n_i 's [e.g., Hosmer-Lemeshow test, see Hosmer and Lemeshow (2013, 3rd ed.), *Applied Logistic Regression*]

$n_i^* = \#$ of units in i th class
 $y_i^* = \#$ of 1's in i th class



$y_i^* \leftrightarrow \hat{y}_{i,S}^* = n_i^* \hat{p}_{i,S}^*$, i.e., use n_i^* 's (≥ 1) y_i^* 's, $\hat{p}_{i,S}^*$'s to measure goodness-of-fit

Reading: Faraway (2006, 1st ed.), 2.1, 2.2, 2.3, Appendix A

- a comparison

p. 3-15

Q: Why no estimation/inference related to $\text{Var}(Y)$?
 $\text{Var}(y_i) = \frac{\mu_i(n_i - \mu_i)}{n_i}, \mu_i = E(y_i)$

	Linear model	Binomial
	$Y = X\beta + \varepsilon$	$Y \sim N(X\beta, \sigma^2 I)$
		GLM
estimation of β	LS = MLE	MLE
Goodness of fit	RSS, $\hat{\sigma}^2$ & $R^2 = 1 - \text{RSS}/\text{TSS}$	Deviance
or Lack of fit	$y_i \leftrightarrow \hat{y}_i (y_i - \hat{y}_i), \text{RSS}_S \leftrightarrow \text{RSS}_{\text{sat}}$	$y_i \leftrightarrow \hat{y}_i \left(\frac{y_i}{\hat{y}_i} \right), 2 \log \left(\frac{\mathcal{L}_{\text{sat}}}{\mathcal{L}_S} \right)$
$H_0: S \text{ vs. } H_1: L \setminus S$	$(*) - \frac{(\text{RSS}_S - \text{RSS}_L) / (df_S - df_L)}{\text{RSS}_L / df_L} \sim F(H_0)$ exact	$D_S - D_L \xrightarrow{a} \chi^2_{df_S - df_L}(H_0)$ asymptotic (n_i 's large)
$H_0: \beta_i = 0$	from $(*)$ ← equivalent $(\square) - \hat{\beta}_i / \text{se}(\hat{\beta}_i) \sim t(H_0)$ exact	from (Δ) ← different $\hat{\beta}_i / \text{se}(\hat{\beta}_i) \xrightarrow{a} N(0,1)(H_0)$ asymptotic
Confidence interval or region	from $(*)$ ← equivalent from (\square) ← estimate $\pm (c.v.) \times \text{se}(\text{estimate})$	profile likelihood method (related to (Δ)) from (\circ) ← different

Tolerance distribution

- Consider the following example:

latent continuous r.v.
 $F = \text{cdf of } N(0,1), T \sim N(\mu, \sigma^2)$

$\text{eg. } T_j$'s $\xrightarrow{i.i.d.} F\left(\frac{t - \mu}{\sigma}\right)$
 a (known) cdf \uparrow location parameter μ , scale parameter σ

observed data
 assume x_i 's are different

	X	$Z(1 \times X, 0 \times V)$	Y
1st covariate class	x_1	$z_{11} z_{12} \dots z_{1n}$	y_1
\vdots	\vdots	\vdots	\vdots
k th covariate class	x_k	$z_{k1} z_{k2} \dots z_{kn}$	y_k

y_i : # of students got x_i
 z_{ij} : 1 if student j got x_i , 0 otherwise
 from different students
 1st student (T_1), 2nd student (T_2), ...
 not observe

➤ Students answers k questions on a test

➤ The j th student has an aptitude T_j

➤ The i th question has a fixed difficulty $x_i, i=1, \dots, R$

If changed to (x_{i1}, \dots, x_{im}) wrong

➤ The j th student will get the i th answer correct only if $T_j \leq x_i$

➤ The probability that a randomly selected student will get the i th answer wrong is:

GLM for binomial response
 y_i 's and covariate x_i 's

response: $y_i \sim B(n_i, p_i)$
 $p_i = P(T_j \leq x_i) = F((x_i - \mu)/\sigma)$
 $\Rightarrow F^{-1}(p_i) = (-\mu/\sigma) + (1/\sigma)x_i \equiv \beta_0 + \beta_1 \times x_i \equiv \eta_i$
 link function $g \Leftrightarrow F = g^{-1}$
 linear structure covariate

If changed to $T_j \leq \sum_k \beta_k R_k(x_i)$
 Become $\eta_i = \sum_k \beta_k R_k(x_i)$

➤ The distribution of T_j is called tolerance distribution, which arose from toxicity studies where the aptitude would be replaced with the tolerance of the insects. $T_j \leq x_i$ [die survive]

• $F = \Phi(\eta) \Rightarrow$ probit link

• $F = e^\eta / (1 + e^\eta) \Rightarrow$ logit link

• $F = 1 - e^{-e^\eta}$ (heavier tail)

\Rightarrow complementary log-log link

useful in interpreting

β under logit link

\rightarrow Odds

Recall.

linear $\rightarrow \eta_i = \beta_0 + \beta_1 \eta_1(x_i) + \dots + \beta_p \eta_p(x_i) = \eta_i^T \beta$

nonlinear $\rightarrow \tilde{\mu}_i / n_i = p_i = g^{-1}(\eta_i) = g^{-1}(\eta_i^T \beta)$

• Definition of odds o_x (賠率)

➤ Sometimes, a better scale than probability to present chance

➤ Express the payoffs for bets

$$O(A) = \frac{P(A)}{P(A^c)}$$

1:1

If it's a fair bet

■ even bet: pay \$1 for every \$1 bet, odds=1 $\Rightarrow p=0.5$

■ 3-1 against bet: pay \$3 for every \$1 bet, odds=1/3 $\Rightarrow p=0.25$

■ 1-3 on bet: pay \$1 for every \$3 bet, odds=3 $\Rightarrow p=0.75$

$$o_x = p_x / (1 - p_x) \Rightarrow p_x = o_x / (1 + o_x)$$

$\frac{1}{0} \Leftrightarrow \frac{\infty}{0}$

3:1

$$0 \leq p_x \leq 1$$

$$0 \leq o_x < \infty$$

$$-\infty < \log(o_x) < \infty$$

➤ A mathematical advantage of odds: unbounded above and unbounded below after taking log.

Recall. problem with upper/lower bounds (LNp 3-12)

• odds and logit link (logistic regression)

➤ Consider the model

an advantage of using logit link

its value can be interpreted as "log odds" at \underline{x}

$$\log(o_{\underline{x}}) = \log\left(\frac{p_{\underline{x}}}{1 - p_{\underline{x}}}\right) = \text{logit}(p_{\underline{x}}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 = \eta_{\underline{x}}$$

$$\log(o_{\underline{x}'}) = \text{logit}(p_{\underline{x}'}) = \beta_0 + \beta_1(x_1 + 1) + \beta_2 x_2 = \eta_{\underline{x}'} = \eta_{\underline{x}} + \beta_1$$

When changing from \underline{x} to \underline{x}' .

$$\text{taking exp} \left[\begin{aligned} \log(o_{\underline{x}'}) &= \log(o_{\underline{x}}) + \beta_1 \\ \Rightarrow o_{\underline{x}'} &= o_{\underline{x}} \times e^{\beta_1} \end{aligned} \right] \Rightarrow \left[\begin{aligned} \beta_1 &= \log(o_{\underline{x}'} / o_{\underline{x}}) = \log\left(\frac{o_{\underline{x}'}}{o_{\underline{x}}}\right) \leftarrow \text{log odds ratio of } \underline{x}' \text{ relative to } \underline{x} \\ e^{\beta_1} &= \frac{o_{\underline{x}'}}{o_{\underline{x}}} \leftarrow \text{odd ratio} \end{aligned} \right]$$

at two covariate values $\underline{x} = (x_1, x_2)$ and $\underline{x}' = (x_1 + 1, x_2)$.

they are already ratios

➤ The coefficient β_1 can be interpreted as: compare 2 odds in terms of ratio

A unit increase in x_1 with x_2 held fixed

Probit:

$$\eta = \Phi^{-1}(p)$$

CLL:

$$\eta = \log[-\log(1 - p)]$$

■ increases the log-odds of successes by β_1 , or meaning in terms of p ?

■ increases the odds of successes by a factor of $\exp(\beta_1)$ especially in interpreting β

➤ The usually interpretational difficulties regarding causation apply as in standard regression

check LNp.1-19~21

➤ No such simple interpretation exists for other links, e.g.,

check in LNp. 3-5