

①

11 Oct 2001

Back to confidence and testingLikelihood ratio test and extension of anovaProb function $f(y|\theta)$ $\theta \in \Theta$ Data y_1, \dots, y_n Null hypothesis $H_0: \theta \in \Theta_0 \subset \Theta$ nested

Likelihood ratio statistic

$$\lambda_n = \frac{\max_{\theta \in \Theta_0} \prod f(y_i|\theta)}{\max_{\theta \in \Theta} \prod f(y_i|\theta)}$$

$$-2 \log \lambda_n = 2 [\ell(\hat{\theta}) - \ell(\hat{\theta}_0)] \quad \text{difference of deviances}$$

as $n \rightarrow \infty$

$\sim \chi^2_{k-d} \text{ or } \chi^2_{n-r}$, r : number of restrictions
 d : number of free

$$\text{if } \Theta_0 = \{(\theta_1, \dots, \theta_r, \theta_{r+1}, \dots, \theta_k); \theta_j = \theta_{0j}, j=1, \dots, r\}$$

and model in family, ie model correct

nested

There may be a sequence of hypotheses

$$\Theta_1 \subset \Theta_2 \subset \Theta$$

$$\ell(\theta) = \sum_i \log f(y_i|\theta)$$

(1)

11 Oct. 2001

Proof of χ^2 .a) Full mle, k unknowns

$$\text{Set } \hat{\Sigma} = \sum_{i=1}^{n+1} (\hat{\theta} - \theta_0) \sim N(\hat{\theta}, \hat{I}_{\hat{\theta}}^{-1})$$

$$\hat{V} = \frac{1}{n} \sum_{i=1}^{n+1} \left. \frac{\partial \log f(y_i | \theta)}{\partial \theta} \right|_{\hat{\theta}} \quad (*)$$

$$\text{Then } \hat{V} \sim N(\hat{\theta}, \hat{I}_{\hat{\theta}})$$

$$\begin{aligned} \text{where } \hat{I}_{\hat{\theta}} &\sim \text{var} \left\{ \left. \frac{\partial \log f(y | \theta)}{\partial \theta} \right|_{\hat{\theta}} \right\} \\ &= -E \left\{ \left. \frac{\partial^2 \log f(y | \theta)}{\partial \theta \partial \theta^T} \right|_{\hat{\theta}} \right\} \end{aligned}$$

By Taylor expansion

$$\begin{aligned} \hat{\Sigma} &= \sum_{i=1}^{n+1} \left. \frac{\partial \log f(y_i | \theta)}{\partial \theta} \right|_{\hat{\theta}} \sim \sum_{i=1}^{n+1} \left. \frac{\partial \log f(y_i | \theta)}{\partial \theta} \right|_{\theta} \\ &+ \sum_{i=1}^{n+1} \left. \frac{\partial^2 \log f(y_i | \theta)}{\partial \theta \partial \theta^T} \right|_{\theta} (\hat{\theta} - \theta_0)^T \end{aligned}$$

$$\text{So } \hat{V} \sim \hat{I}_{\hat{\theta}} \hat{\Sigma}$$

$$\hat{\Sigma} \sim \hat{I}_{\hat{\theta}}^{-1} \hat{V}$$

②

11 Oct 2001

$$l(\theta) = \sum_i \log f(y_i | \theta)$$

$$l(\theta_0) \approx \sum_i \log f(y_i | \hat{\theta}) + \left(\sum_i \frac{\partial \log f(y_i | \theta)}{\partial \theta} \Big|_{\hat{\theta}} \right) (\theta_0 - \hat{\theta})$$

$$+ \frac{1}{2} (\theta_0 - \hat{\theta})^T \left(\sum_i \frac{\partial^2 \log f(y_i | \theta)}{\partial \theta \partial \theta^T} \Big|_{\hat{\theta}} \right) (\theta_0 - \hat{\theta})$$

$$\text{so } l(\theta_0) \approx l(\hat{\theta}) - \frac{1}{2} (\theta_0 - \hat{\theta})^T \hat{I}_0 (\theta_0 - \hat{\theta})$$

$$\text{or } 2 [l(\hat{\theta}) - l(\theta_0)] \sim \hat{D}^T \hat{I}_0^{-1} \hat{D}$$

$$\sim \chi_k^2$$

$$\sim V^T \hat{I}_0^{-1} V \quad \text{also}$$

This result can be used to:

- i) examine the simple hypothesis $\theta = \theta_0$
- ii) construct a confidence region for θ_0 .
(using \hat{I}_0)
- iii) corresponds to $n = k$ restrictions

(3)

11 Oct. 2001

Now let's suppose that there are $k-r$ free parameters only.

Have restrictions such that the original $\theta_1, \dots, \theta_k$ are functions of $s = k-r$ new parameters

$$\text{e.g. } \theta_j = g_j(\beta_1, \dots, \beta_s) \quad j=1, \dots, k$$

$r = k-s$ restrictions

Suppose, for convenience

$$\begin{bmatrix} \beta_1 \\ \vdots \\ \beta_s \end{bmatrix} = \begin{bmatrix} \theta_{s+1} \\ \vdots \\ \theta_k \end{bmatrix}$$

We will need

$$\underset{s}{\overset{k \times s}{M}} = \begin{bmatrix} \frac{\partial g_i}{\partial \beta_j} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

Write underscore when referring to last s coordinates of θ , e.g. $\underline{\theta}_1, \underline{\theta}_s$

Will repeat previous development only using last s coordinates.

(4)

11 Oct. 2001

$$\text{Let } \hat{U} = \sum_i \frac{1}{f_m} \left. \frac{\partial \log f(y_i; \theta)}{\partial \theta} \right|_{\theta_0}$$

(**) (circle)

$$\text{Let } \hat{F} = \sum_m (\hat{\theta} - \theta_0)$$

\hat{J}_0 : information matrix for $\hat{\theta}$

As before

$$\hat{U} \sim \frac{\partial F}{\partial \theta}, \quad \hat{F} \sim \frac{\partial^2}{\partial \theta^2} \hat{U}$$

$$2[\ell(\hat{\theta}_0) - \ell(\theta)] \sim \hat{F}^T \frac{\partial}{\partial \theta} \hat{F}$$

lets connect \hat{U} and V

$$\text{From (4) and (**) } \hat{U} = M^T V \text{ and so}$$

$$\frac{\partial}{\partial \theta} = M^T I_{\hat{U}} M$$

(**) (circle)

Continuing

$$2[\ell(\hat{\theta}_0) - \ell(\theta_0)] \sim \hat{F}^T \frac{\partial}{\partial \theta} \hat{F}$$

$$\sim \hat{U}^T \frac{\partial^2}{\partial \theta^2} \hat{U}$$

$$\sim V^T M \frac{\partial^2}{\partial \theta^2} M^T V$$

(5)

11 Oct. 2021

Subtracting

$$2[\lambda(\hat{0}) - \lambda(\hat{\theta})] \sim V^T \left(I_{n_0}^{-1} - M D_{n_0}^{-1} M^T \right) V$$

$$\text{But } V \sim N_p(0, I_n)$$

$$\text{write } V = I_{n_0}^{k_2} Z \quad \text{with } Z \sim N(0, I)$$

So

$$\textcircled{**} = Z^T \left(I - I_{n_0}^{k_2} M D_{n_0}^{-1} M^T I_{n_0}^{k_2} \right) Z$$

Note that from $\textcircled{***}$

$$I_{n_0}^{k_2} M D_{n_0}^{-1} M^T I_{n_0}^{k_2} I_{n_0}^{k_2} M D_{n_0}^{-1} M^T I_{n_0}^{k_2}$$

$$= I_{n_0}^{k_2} M D_{n_0}^{-1} M^T I_{n_0}^{k_2} \quad \text{i.e. idempotent}$$

So $\textcircled{**}$ is Z^2

$$r = k - n \left(I_{n_0}^{k_2} M D_{n_0}^{-1} M^T I_{n_0}^{k_2} \right)$$

$$= k - \text{tr} \left(D_{n_0}^{-1} M^T I_{n_0}^{k_2} M \right)$$

$$= k - \text{tr} \left(D_{n_0}^{-1} D_{n_0} \right) \quad \text{from } \textcircled{***}$$

$$= k - s = r$$

The nested sequence result follows similarly.

(6)

11 Oct 2001

The deviance of the fitted model $f(y|\hat{\theta})$ is

$$2 [l(\hat{\theta}) - l(\tilde{\theta})] > 0$$

where $\tilde{\theta}$ is the mle for the saturated model,
viz. parameter for each observation

By a big stretch

$$2 [l(\tilde{\theta}) - l(\hat{\theta})] \sim \chi^2_{n-k}$$

People often compare final deviance to

$$E \chi^2_{n-k} = n-k, \text{ but } \exists \text{ problems with approx.}$$

We have seen that

$$2 [l(\hat{\theta}) - l(\tilde{\theta})] \sim \chi^2_{k-s} = \chi^2_p \text{ under null}$$

i.e. the difference in deviances between two nested models has a chi-squared distribution under the null hypothesis.

$$2 [l(\hat{\theta}) - l(\tilde{\theta})] - 2 [l(\hat{\theta}) - l(\hat{\theta}_0)]$$

(H₀)

(H)

(7)

11 Oct 2001

Might have $\Theta_1 \subset \Theta_0 \subset \Theta$

free
params

$$\lambda_1 < \lambda_0 < k$$

constraint/restrictions

$$r_1 < r_0 < 0$$

$$-2[\lambda(\tilde{\theta}) - \lambda(\hat{\theta})] - 2[\lambda(\tilde{\theta}) - \lambda(\hat{\theta}_-)] = D_{\Theta_1} - D_{\Theta_0}$$

$$\sim \chi^2_{D_0 - D_1}$$

Can set up an ANODEV table

Hypothesis	Deviance	Δf	Deviance difference	$\Delta \Delta f$
Θ_1	D_{Θ_1}	$m - D_1$	$D_{\Theta_1} - D_{\Theta_0}$	$D_0 - D_1$
Θ_0	D_{Θ_0}	$m - D_0$	$D_{\Theta_0} - D_{\Theta}$	$k - D_0$
Θ	D_{Θ}	$m - k$		

Deviance gets smaller as bring in more parameters

(8)

11 Oct., 2001

Advantages of divergence.

1. Invariant under 1-1 parametrizations of the model
- "
2. Unifies a variety of model selection problems: linear model (anova), glm, graphical models ...
3. Has interpretation in terms of K-L divergence
4. lot known about large sample distribution

(9)

11 Oct. 2001

Distribution of the deviance / likelihood ratio statistic when the model is not true

We saw that

MLE $\hat{\theta} \xrightarrow{P} \theta_*$ where θ_* maximizes

$$E_p \{ \log f(Y|\theta) \}$$

$$\underset{n}{\text{L}} \ell(\theta) = \underset{n}{\text{L}} \sum \log f(Y_i|\theta)$$

$$\xrightarrow{P} E_p \{ \log f(Y|\theta) \}$$

so can anticipate

$$\underset{n}{\text{L}} \ell(\hat{\theta}) \xrightarrow{P} E_p \{ \log f(Y|\theta_*) \}$$

Likewise

$$\underset{n}{\text{L}} \ell(\underline{\theta}) \rightarrow E_p \{ \log f(Y|\underline{\theta}_*) \}$$

$$\text{so } -\frac{2}{n} \log \lambda_n \rightarrow 2 E_p \{ \log f(Y|\theta_*) \} - 2 E_p \{ \log f(Y|\underline{\theta}_*) \}$$

Going to ∞ when $\theta_* \notin \Theta_0$

$\underline{\theta}_*$ gives $\max_{\theta \in \Theta_0} E_p \{ \log f(Y|\theta) \}$

(10)

11 Oct. 2001

When $\theta_x \in \Theta_0$

$$-2 \log \hat{\gamma}_n \sim c_1 z_1^2 + \dots + c_r z_r^2 \quad n=k$$

"Foutz & Srivastava (1978) Canadian J. Stat 6, 273-9.

(1)

14 Oct. 2001

Example. Gauss-Markov + normal model

$$y_i = \alpha + (\tilde{x}_i - \bar{x})' \beta + \epsilon_i \quad \epsilon_i \sim N(0, \sigma^2)$$

$$H_0: \beta = \beta_0$$

$$\Theta = \begin{bmatrix} \alpha \\ \beta \\ \sigma \end{bmatrix}$$

$$L(\Theta) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2\sigma^2} \|y_i - \alpha - (\tilde{x}_i - \bar{x})' \beta\|^2 \right\}$$

$$l(\Theta) = -\frac{m}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^m \|y_i - \alpha - (\tilde{x}_i - \bar{x})' \beta\|^2 + C$$

$$\text{"Full" mle: } \hat{\alpha} = \bar{y}$$

$$\hat{\tilde{X}}' \hat{\tilde{X}} \hat{\beta} = \hat{\tilde{X}}' \hat{y}$$

$$\hat{\tilde{X}} = \tilde{X} - \bar{\tilde{X}}$$

$$\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m \|y_i - \bar{y} - (\tilde{x}_i - \bar{\tilde{x}})' \hat{\beta}\|^2$$

$$l(\hat{\Theta}) = -\frac{m}{2} \log \hat{\sigma}^2 + C'$$

(2)

14 Oct. 01

Mle under H_0 : $\hat{\alpha} = \bar{y}$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \|y_i - \bar{y}\|^2$$

$$n \ell(\hat{\theta}_0) = -\frac{n}{2} \log \hat{\sigma}^2$$

$$2[\ell(\hat{\theta}) - \ell(\hat{\theta}_0)] = n \log \hat{\sigma}^2 \sum_{i=1}^n \|y_i - \bar{y}\|^2 - \sum_{i=1}^n \|y_i - \bar{y} - (x_i - \bar{x})' \hat{\beta}\|^2$$

Now

$$\sum_{i=1}^n \|y_i - \bar{y}\|^2 = \sum_{i=1}^n \|(x_i - \bar{x})' \hat{\beta}\|^2 + \sum_{i=1}^n \|y_i - \bar{y} - (x_i - \bar{x})' \hat{\beta}\|^2$$

$$SST = SSR + SSE$$

Test statistic is based on SSR/SSE

We saw its exact distribution under H_0 .

$$\frac{SSR/p}{SSE/(n-p-1)} \sim F_{p, n-p-1}$$

(3)

14 Oct. 01

Relation to previous approximation:

It was based on $n \rightarrow \infty$

$$\hat{\sigma}^2 \xrightarrow{\text{Pois}} \sigma^2$$

$$n \log \frac{SSR + SSE}{SSE} = n \log \left(1 + \frac{SSR}{SSE} \right)$$

$$\sim \frac{SSR}{SSE/m}$$

$$\sim \frac{SSR}{\sigma^2}$$

$$\sim \chi_p^2$$