

22. Additional Topics Related to Likelihood

Information Criteria

Akaike's Information criterion is given by

$$\text{AIC} = -2\ell(\hat{\boldsymbol{\theta}}) + 2k,$$

where $\ell(\hat{\boldsymbol{\theta}})$ is the maximized log likelihood and k is the dimension of the model parameter space.

- $AIC = -2\ell(\hat{\boldsymbol{\theta}}) + 2k$ can be used to determine which of multiple models is “best” for a given data set.
- Small values of AIC are preferred.
- The $+2k$ portion of AIC can be viewed as a penalty for model complexity.

Schwarz's Bayesian Information Criterion is given by

$$\text{BIC} = -2\ell(\hat{\boldsymbol{\theta}}) + k \ln(n)$$

BIC is the same as AIC except the penalty for model complexity is greater for BIC (when $n \geq 8$) and grows with n .

- AIC and BIC can each be used to compare models even if they are not nested (i.e., even if one is not a special case of the other as in our reduced vs. full model comparison discussed previously).
- However, if REML likelihoods are used, compared models must have the same model for the response mean.
- Different models for the mean would yield different error contrasts and different datasets for computation of maximized REML likelihoods.

Large n Theory for MLEs

- Suppose θ is a $k \times 1$ parameter vector.
- Let $\ell(\theta)$ denote the log likelihood function.
- Under regularity conditions discussed in, e.g., Shao, J.(2003) *Mathematical Statistics*, 2nd Ed. Springer, New York; we have the following.

There is an estimator $\hat{\theta}$ that solves the score equations $\frac{\partial \ell(\theta)}{\partial \theta} = \mathbf{0}$ and has the following properties.

1 Consistency of $\hat{\theta}$:

$\hat{\theta}$ is a (weakly) consistent estimator of θ .

This means that $\hat{\theta}$ converges in probability to θ , i.e.,

$$\lim_{n \rightarrow \infty} Pr[||\hat{\theta} - \theta|| > \varepsilon] = 0 \text{ for any } \varepsilon > 0.$$

2 Asymptotic normality of $\hat{\boldsymbol{\theta}}$:

For sufficiently large n , $\hat{\boldsymbol{\theta}} \stackrel{\bullet}{\sim} N(\boldsymbol{\theta}, \mathbf{I}^{-1}(\boldsymbol{\theta}))$, where

$$\begin{aligned}\mathbf{I}(\boldsymbol{\theta}) &= E \left[\left(\frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right) \left(\frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)' \right] \\ &= -E \left[\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \right] \\ &= \left[-E \left\{ \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right\} \right]_{i,j \in \{1, \dots, k\}}\end{aligned}$$

- $I(\theta)$ is known as the *Fisher Information* matrix.
- $I(\theta)$ can be approximated by replacing the unknown θ with $\hat{\theta}$ to obtain $I(\hat{\theta})$.
- An alternative approximation is given by the *observed Fisher Information* matrix:

$$\hat{I}(\hat{\theta}) \equiv \left. \frac{-\partial^2 \ell(\theta)}{\partial \theta \partial \theta'} \right|_{\theta=\hat{\theta}}$$

In practice, when n is sufficiently large, we use the approximation

$$\hat{\boldsymbol{\theta}} \stackrel{\bullet}{\sim} N(\boldsymbol{\theta}, \widehat{\text{Var}}(\hat{\boldsymbol{\theta}})),$$

where $\widehat{\text{Var}}(\hat{\boldsymbol{\theta}})$ can be either $\boldsymbol{I}^{-1}(\hat{\boldsymbol{\theta}})$ or $\hat{\boldsymbol{I}}^{-1}(\hat{\boldsymbol{\theta}})$.

Although such statements do a reasonable job of conveying the idea of approximations we use, they are not mathematically rigorous.

When we say something like

$$\hat{\boldsymbol{\theta}} \stackrel{\bullet}{\sim} N(\boldsymbol{\theta}, \widehat{\text{Var}}(\hat{\boldsymbol{\theta}})) \text{ for sufficiently large } n,$$

we mean that as n grows to infinity,

$$[\widehat{\text{Var}}(\hat{\boldsymbol{\theta}})]^{-1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})$$

converges in distribution to a standard multivariate normal random vector $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I})$:

$$[\widehat{\text{Var}}(\hat{\boldsymbol{\theta}})]^{-1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{d} \mathbf{z}$$

Note that

$$[\widehat{\text{Var}}(\hat{\boldsymbol{\theta}})]^{-1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{d} \mathbf{z}$$

implies

$$\{[\widehat{\text{Var}}(\hat{\boldsymbol{\theta}})]^{-1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})\}' \{[\widehat{\text{Var}}(\hat{\boldsymbol{\theta}})]^{-1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})\} \xrightarrow{d} \mathbf{z}'\mathbf{z}$$

which implies

$$(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})'[\widehat{\text{Var}}(\hat{\boldsymbol{\theta}})]^{-1}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \overset{\bullet}{\sim} \chi_k^2$$

for sufficiently large n .

A Simple Example

- Suppose $y_1, \dots, y_n \stackrel{i.i.d.}{\sim} \text{Poisson}(\theta)$.
- For $y_i \in \{0, 1, 2, \dots\} \forall i = 1, \dots, n$,

$$L(\theta|\mathbf{y}) = \prod_{i=1}^n \frac{\theta^{y_i} e^{-\theta}}{y_i!}$$

$$\ell(\theta|\mathbf{y}) = \sum_{i=1}^n [y_i \ln(\theta) - \theta - \ln(y_i!)]$$

$$= \ln(\theta) \sum_{i=1}^n y_i - n\theta - \sum_{i=1}^n \ln(y_i!)$$

$$\frac{\partial \ell(\theta|\mathbf{y})}{\partial \theta} = \frac{1}{\theta} \sum_{i=1}^n y_i - n$$

Thus, the score equation is

$$\frac{1}{\theta} \sum_{i=1}^n y_i - n = 0.$$

The only solution to the score equation is $\hat{\theta} = \bar{y}$.

Result (1) on slide 7 implies $\hat{\theta} = \bar{y}$. converges in probability to θ .

In this case, we also know that \bar{y} . converges in probability to θ by the (Weak) Law of Large Numbers (WLLN).

Result (2) on slide 8 implies $\hat{\theta} = \bar{y} \stackrel{\bullet}{\sim} N(\theta, I^{-1}(\theta))$.

$$\begin{aligned} I(\theta) &= -E \left[\frac{\partial^2 \ell(\theta | \mathbf{y})}{\partial \theta \partial \theta} \right] = -E \left[\frac{\partial}{\partial \theta} \left(\frac{1}{\theta} \sum_{i=1}^n y_i - n \right) \right] \\ &= -E \left[-\frac{1}{\theta^2} \sum_{i=1}^n y_i \right] = \frac{1}{\theta^2} \sum_{i=1}^n E(y_i) = \frac{n}{\theta} \end{aligned}$$

Therefore, $I^{-1}(\theta) = \theta/n$ in this case.

Thus, result (2) on slide 8 implies $\hat{\theta} = \bar{y} \stackrel{\bullet}{\sim} N(\theta, \theta/n)$, which is also implied by the Central Limit Theorem (CLT).

To get an estimate of the variance of $\hat{\theta} = \bar{y}_{\cdot}$, we can use

$$I^{-1}(\hat{\theta}) = \hat{\theta}/n = \bar{y}_{\cdot}/n$$

Alternatively, the inverse of the observed Fisher information in this case is

$$\hat{I}^{-1}(\hat{\theta}) = \left[\frac{-\partial^2 \ell(\theta)}{\partial \theta \partial \theta} \Big|_{\theta=\hat{\theta}} \right]^{-1} = \left[\frac{1}{\hat{\theta}^2} \sum_{i=1}^n y_i \right]^{-1} = \left[\frac{n\bar{y}_{\cdot}}{\bar{y}_{\cdot}^2} \right]^{-1} = \bar{y}_{\cdot}/n$$

Thus, $I^{-1}(\hat{\theta}) = \hat{I}^{-1}(\hat{\theta})$ in this case.

Substituting in this consistent estimator for $I^{-1}(\theta)$, we have

$$\hat{\theta} = \bar{y}_{\cdot} \stackrel{\bullet}{\sim} N(\theta, \bar{y}_{\cdot}/n)$$

Wald Tests and Confidence Intervals

Suppose for large n that

$$\hat{\boldsymbol{\theta}} \stackrel{\bullet}{\sim} N(\boldsymbol{\theta}, \widehat{\text{Var}}(\hat{\boldsymbol{\theta}})).$$

Then a confidence interval for $\mathbf{c}'\boldsymbol{\theta}$ that has confidence level approximately equal to $1 - \alpha$ is

$$\mathbf{c}'\hat{\boldsymbol{\theta}} \pm z_{1-\alpha/2} \sqrt{\mathbf{c}'\widehat{\text{Var}}(\hat{\boldsymbol{\theta}})\mathbf{c}},$$

where $z_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile of the $N(0, 1)$ distribution.

Likewise, a test of $H_0 : \mathbf{c}'\boldsymbol{\theta} = d$ can be based on the test statistic

$$\frac{\mathbf{c}'\hat{\boldsymbol{\theta}} - d}{\sqrt{\mathbf{c}'\widehat{\text{Var}}(\hat{\boldsymbol{\theta}})\mathbf{c}}},$$

which has a distribution that is approximately $N(0, 1)$ under H_0 .

Likewise, if C is a $q \times k$ matrix of rank q , a test of $H_0 : C\theta = d$ can be based on the test statistic

$$(C\hat{\theta} - d)'[C\widehat{\text{Var}}(\hat{\theta})C']^{-1}(C\hat{\theta} - d),$$

which has a distribution that is approximately χ_q^2 under H_0 .

Multivariate Delta Method

- Suppose g is a function from \mathbb{R}^k to \mathbb{R}^m , i.e.,

$$\text{for } \boldsymbol{\theta} \in \mathbb{R}^k, \mathbf{g}(\boldsymbol{\theta}) = \begin{bmatrix} g_1(\boldsymbol{\theta}) \\ g_2(\boldsymbol{\theta}) \\ \vdots \\ g_m(\boldsymbol{\theta}) \end{bmatrix}$$

for some functions g_1, \dots, g_m .

- Suppose g is differentiable with derivative matrix

$$\mathbf{D} \equiv \begin{bmatrix} \frac{\partial g_1(\boldsymbol{\theta})}{\partial \theta_1} & \cdots & \frac{\partial g_m(\boldsymbol{\theta})}{\partial \theta_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_1(\boldsymbol{\theta})}{\partial \theta_k} & \cdots & \frac{\partial g_m(\boldsymbol{\theta})}{\partial \theta_k} \end{bmatrix}.$$

Now suppose $\hat{\theta}$ has mean θ and variance $\text{Var}(\hat{\theta})$.
Then Taylor's Theorem implies

$$\mathbf{g}(\hat{\theta}) \approx \mathbf{g}(\theta) + \mathbf{D}'(\hat{\theta} - \theta)$$

which implies

$$E[\mathbf{g}(\hat{\theta})] \approx \mathbf{g}(\theta) + \mathbf{D}'E(\hat{\theta} - \theta) = \mathbf{g}(\theta)$$

and

$$\text{Var}[\mathbf{g}(\hat{\theta})] \approx \text{Var}[\mathbf{g}(\theta) + \mathbf{D}'(\hat{\theta} - \theta)] = \mathbf{D}'\text{Var}(\hat{\theta})\mathbf{D}.$$

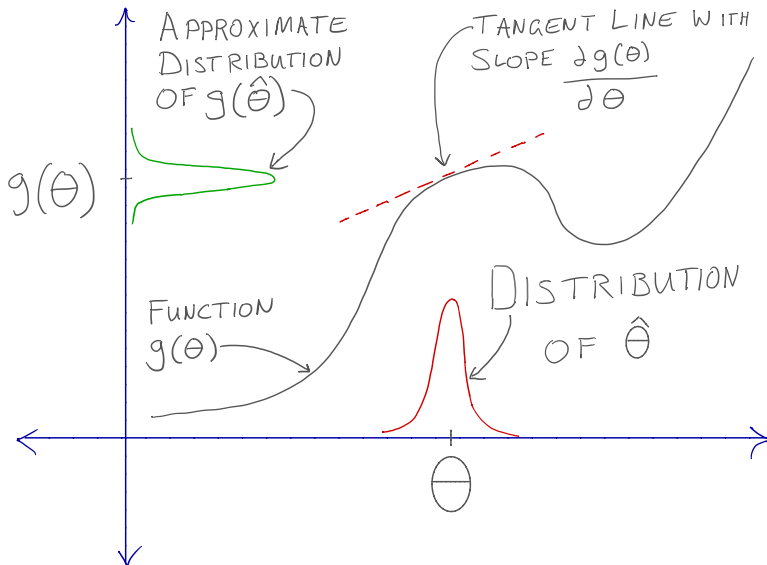
- If $\hat{\theta} \overset{\bullet}{\sim} N(\theta, \text{Var}(\hat{\theta}))$, it follows that

$$g(\hat{\theta}) \overset{\bullet}{\sim} N(g(\theta), D' \text{Var}(\hat{\theta}) D).$$

- In practice, we often need to estimate D by replacing θ in D with $\hat{\theta}$ to obtain \hat{D} .
- Similarly, we often need to replace $\text{Var}(\hat{\theta})$ with an estimate $\widehat{\text{Var}}(\hat{\theta})$.

$$g(\hat{\theta}) \overset{\bullet}{\sim} N(g(\theta), \hat{D}' \widehat{\text{Var}}(\hat{\theta}) \hat{D})$$

Delta Method Example with $k = 1$



Likelihood Ratio Based Inference

Suppose we wish to test the null hypothesis that a reduced model provides an adequate fit to a dataset relative to a more general full model that includes the reduced model as a special case.

- Define Λ as

$$\frac{\text{Reduced Model Maximized Likelihood}}{\text{Full Model Maximized Likelihood}}.$$

- Λ is known as the *likelihood ratio*.
- $-2 \ln(\Lambda)$ is known as the *likelihood ratio test statistic*.
- Tests based on $-2 \ln(\Lambda)$ are called *likelihood ratio tests*.

- Under the regularity conditions in Shao (2003) mentioned previously, the likelihood ratio test statistic $-2 \ln(\Lambda)$ is approximately distributed as central $\chi^2_{k_f - k_r}$ under the null hypothesis, where k_f and k_r are the dimensions of the parameter space under the full and reduced models, respectively.
- This approximation can be reasonable if n is “sufficiently large.”

Likelihood Ratio Tests and Confidence Regions for a Subvector of the Full Model Parameter Vector θ

- Suppose θ is $k \times 1$ vector and is partitioned into vectors θ_1 $k_1 \times 1$ and θ_2 $k_2 \times 1$, where $k = k_1 + k_2$ and $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$.
- Consider a test of $H_0 : \theta_1 = d_1$.

- Suppose $\hat{\theta}$ is the MLE of θ and $\hat{\theta}_2(\theta_1)$ maximizes $\ell\left(\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}\right)$ over θ_2 for any fixed value of θ_1 .
- Then $2\left[\ell(\hat{\theta}) - \ell\left(\begin{bmatrix} d_1 \\ \hat{\theta}_2(d_1) \end{bmatrix}\right)\right]$ is approximately $\chi_{k_1}^2$ under the null hypothesis by our previous result when n is “sufficiently large.”

Also,

$$Pr \left\{ 2 \left[\ell(\hat{\boldsymbol{\theta}}) - \ell \left(\begin{bmatrix} \boldsymbol{\theta}_1 \\ \hat{\boldsymbol{\theta}}_2(\boldsymbol{\theta}_1) \end{bmatrix} \right) \right] \leq \chi_{k_1, 1-\alpha}^2 \right\} \approx 1 - \alpha$$

which implies

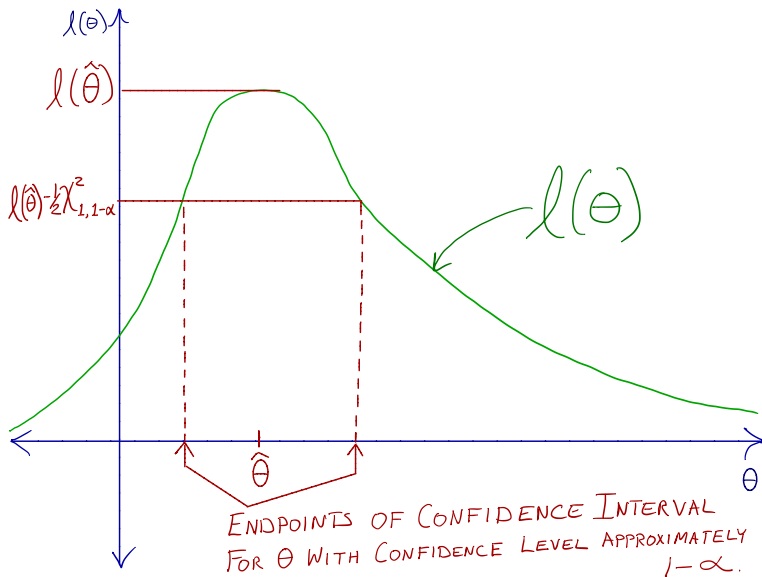
$$Pr \left\{ \ell \left(\begin{bmatrix} \boldsymbol{\theta}_1 \\ \hat{\boldsymbol{\theta}}_2(\boldsymbol{\theta}_1) \end{bmatrix} \right) \geq \ell(\hat{\boldsymbol{\theta}}) - \frac{1}{2} \chi_{k_1, 1-\alpha}^2 \right\} \approx 1 - \alpha.$$

- Thus, the set of values of θ_1 that, when maximizing over θ_2 , yield a maximized likelihood within $\frac{1}{2}\chi_{k_1, 1-\alpha}^2$ of the likelihood maximized over all θ , form a $100(1 - \alpha)\%$ confidence region for θ_1 .
- Such a confidence region is known as a *profile likelihood confidence region* because

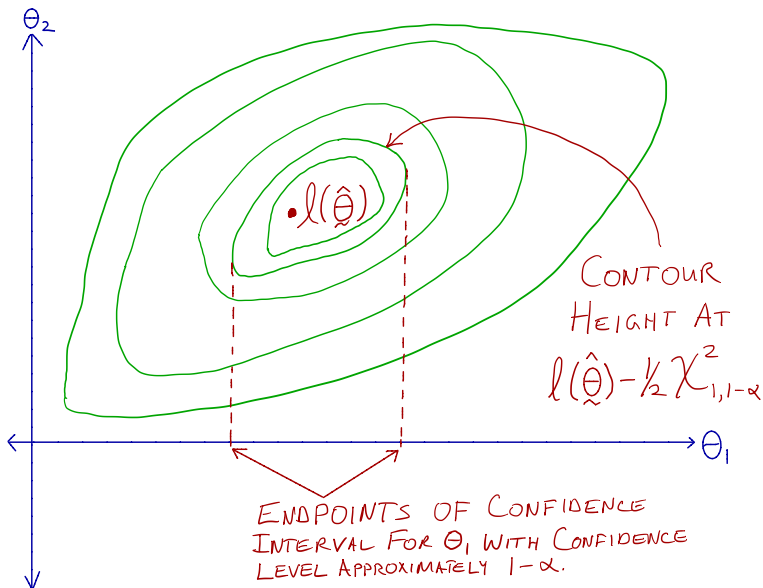
$$\ell \left(\begin{bmatrix} \theta_1 \\ \hat{\theta}_2(\theta_1) \end{bmatrix} \right)$$

is the *profile log likelihood* for θ_1 .

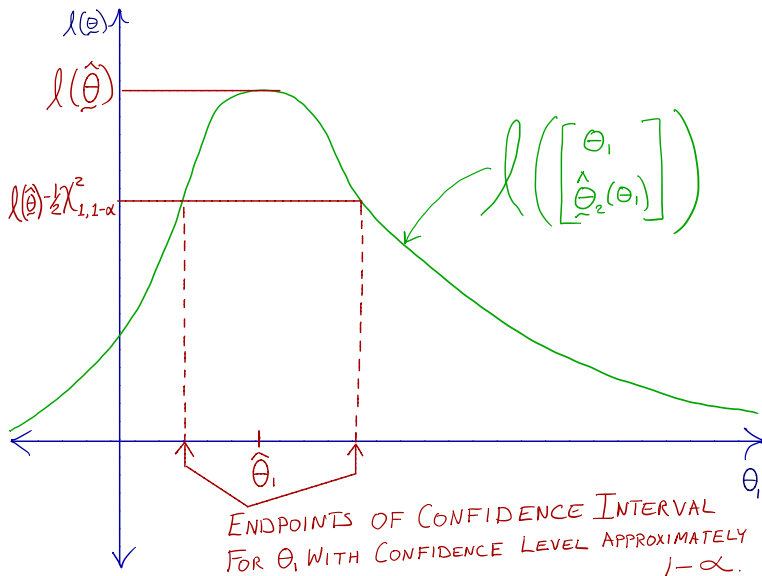
Sketch for the Case of $k = 1$



Sketch for the Case of $k = 2$



Sketch for the Case of $k_1 = 1$ and k_2 Arbitrary



Warnings

- The normal and χ^2 approximations mentioned in these notes may be crude if sample sizes are not sufficiently large.
- The regularity conditions mentioned in these notes do not hold if the true parameter falls on the boundary of the parameter space. Thus, as an example, testing $H_0 : \sigma_u^2 = 0$ is not covered by the methods presented here.