22. Additional Topics Related to Likelihood

Information Criteria

Akaike's Information criterion is given by

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

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

- AIC = −2ℓ(θ̂) + 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

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

BIC is the same as AIC except the penalty for model complexity is greater for BIC (when $n \ge 8$) 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 $\frac{\partial \ell(\theta)}{\partial \theta} = 0$ and has the following properties.

• 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.,

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

Symptotic normality of $\hat{\boldsymbol{\theta}}$:
For sufficiently large $n, \hat{\boldsymbol{\theta}} \stackrel{\bullet}{\sim} N(\boldsymbol{\theta}, \boldsymbol{I}^{-1}(\boldsymbol{\theta})), \text{ where }$

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

Ι

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

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

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

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

where $\widehat{\operatorname{Var}}(\hat{\boldsymbol{\theta}})$ 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

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

we mean that as *n* grows to infinity,

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

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

$$[\widehat{\operatorname{Var}}(\hat{\boldsymbol{\theta}})]^{-1/2}(\hat{\boldsymbol{\theta}}-\boldsymbol{\theta}) \stackrel{d}{\to} \boldsymbol{z}$$

Note that

$$[\widehat{\operatorname{Var}}(\hat{\boldsymbol{\theta}})]^{-1/2}(\hat{\boldsymbol{\theta}}-\boldsymbol{\theta}) \stackrel{d}{\rightarrow} \boldsymbol{z}$$

implies

$$\{ [\widehat{\operatorname{Var}}(\hat{\boldsymbol{\theta}})]^{-1/2} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \}' \{ [\widehat{\operatorname{Var}}(\hat{\boldsymbol{\theta}})]^{-1/2} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \} \stackrel{d}{\to} \boldsymbol{z}' \boldsymbol{z}$$

which implies

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

for sufficiently large n.

A Simple Example

• Suppose
$$y_1, \ldots, y_n \stackrel{i.i.d.}{\sim} \text{Poisson}(\theta)$$
.

• For $y_i \in \{0, 1, 2, ...\} \ \forall i = 1, ..., n$, $L(\theta|\mathbf{y}) = \prod_{i=1}^{n} \frac{\theta^{y_i} e^{-\theta}}{y_i!}$ $\ell(\theta|\mathbf{y}) = \sum [y_i \ln(\theta) - \theta - \ln(y_i!)]$ $= \ln(\theta) \sum y_i - n\theta - \sum^{n} \ln(y_i!)$ $\frac{\partial \ell(\theta|\mathbf{y})}{\partial \theta} = \frac{1}{\theta} \sum_{i=1}^{n} y_i - n$

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 $\hat{\theta} = \bar{y} \cdot \stackrel{\bullet}{\sim} N(\theta, I^{-1}(\theta))$.

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

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 $\hat{\theta} = \bar{y}$, we can use

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

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

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

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

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

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

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

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

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

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

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

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

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.,

for
$$\boldsymbol{ heta} \in \mathbb{R}^k, \boldsymbol{g}(\boldsymbol{ heta}) = egin{bmatrix} g_1(\boldsymbol{ heta}) \\ g_2(\boldsymbol{ heta}) \\ \vdots \\ g_m(\boldsymbol{ heta}) \end{bmatrix}$$

for some functions g_1, \ldots, g_m .

Suppose g is differentiable with derivative matrix

$$oldsymbol{D} \equiv \left[egin{array}{cccc} rac{\partial g_1(oldsymbol{ heta})}{\partial heta_1} & \cdots & rac{\partial g_m(oldsymbol{ heta})}{\partial heta_1} \ dots & \ddots & dots \ rac{\partial g_1(oldsymbol{ heta})}{\partial heta_k} & \cdots & rac{\partial g_m(oldsymbol{ heta})}{\partial heta_k} \end{array}
ight]$$

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

$$\boldsymbol{g}(\hat{\boldsymbol{ heta}}) pprox \boldsymbol{g}(\boldsymbol{ heta}) + \boldsymbol{D}'(\hat{\boldsymbol{ heta}} - \boldsymbol{ heta})$$

which implies

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

and

$$\operatorname{Var}[\boldsymbol{g}(\hat{\boldsymbol{ heta}})] pprox \operatorname{Var}[\boldsymbol{g}(\boldsymbol{ heta}) + \boldsymbol{D}'(\hat{\boldsymbol{ heta}} - \boldsymbol{ heta})] = \boldsymbol{D}' \operatorname{Var}(\hat{\boldsymbol{ heta}}) \boldsymbol{D}$$

• If $\hat{\boldsymbol{\theta}} \stackrel{\bullet}{\sim} N(\boldsymbol{\theta}, \operatorname{Var}(\hat{\boldsymbol{\theta}}))$, it follows that $\boldsymbol{g}(\hat{\boldsymbol{\theta}}) \stackrel{\bullet}{\sim} N(\boldsymbol{g}(\boldsymbol{\theta}), \boldsymbol{D}' \operatorname{Var}(\hat{\boldsymbol{\theta}}) \boldsymbol{D}).$

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

$$\boldsymbol{g}(\hat{\boldsymbol{\theta}}) \stackrel{\bullet}{\sim} N(\boldsymbol{g}(\boldsymbol{\theta}), \hat{\boldsymbol{D}}' \widehat{\operatorname{Var}}(\hat{\boldsymbol{\theta}}) \hat{\boldsymbol{D}})$$

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

Reduced Model Maximized Likelihood 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 $\theta_1 k_1 \times 1$ and $\theta_2 k_2 \times 1$, where $k = k_1 + k_2$ and $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$.
- Consider a test of $H_0: \boldsymbol{\theta}_1 = \boldsymbol{d}_1$.

• 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 $\chi^2_{k_1}$ under the null hypothesis by our previous

result when *n* is "sufficiently large."

Also,

$$Pr\left\{2\left[\ell(\hat{\boldsymbol{\theta}})-\ell\left(\left[\begin{array}{c}\boldsymbol{\theta}_{1}\\\hat{\boldsymbol{\theta}}_{2}(\boldsymbol{\theta}_{1})\end{array}\right]\right)\right] \leq \chi^{2}_{k_{1},1-\alpha}\right\}\approx 1-\alpha$$

which implies

$$Pr\left\{\ell\left(\left[\begin{array}{c}\boldsymbol{\theta}_{1}\\\hat{\boldsymbol{\theta}}_{2}(\boldsymbol{\theta}_{1})\end{array}\right]\right)\geq\ell(\hat{\boldsymbol{\theta}})-\frac{1}{2}\chi_{k_{1},1-\alpha}^{2}\right\}\approx1-\alpha.$$

- Thus, the set of values of θ_1 that, when maximizing over θ_2 , yield a maximized likelihood within $\frac{1}{2}\chi^2_{k_1,1-\alpha}$ 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

$$\ell\left(\left[egin{array}{c} oldsymbol{ heta}_1 \ \hat{oldsymbol{ heta}}_2(oldsymbol{ heta}_1) \end{array}
ight)
ight)$$

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.