

# Likelihood 2: Wald (and Score) Tests<sup>1</sup>

STA442/2101 Fall 2013

---

<sup>1</sup>See last slide for copyright information.

# Background Reading

Davison Chapter 4, especially Sections 4.3 and 4.4

# Vector of MLEs is Asymptotically Normal

That is, Multivariate Normal

This yields

- ▶ Confidence intervals
- ▶  $Z$ -tests of  $H_0 : \theta_j = \theta_0$
- ▶ Wald tests
- ▶ Score Tests
- ▶ Indirectly, the Likelihood Ratio tests

# Under Regularity Conditions

(Thank you, Mr. Wald)

- ▶  $\hat{\boldsymbol{\theta}}_n \xrightarrow{a.s.} \boldsymbol{\theta}$
- ▶  $\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \xrightarrow{d} \mathbf{T} \sim N_k(\mathbf{0}, \mathcal{I}(\boldsymbol{\theta})^{-1})$
- ▶ So we say that  $\hat{\boldsymbol{\theta}}_n$  is asymptotically  $N_k(\boldsymbol{\theta}, \frac{1}{n}\mathcal{I}(\boldsymbol{\theta})^{-1})$ .
- ▶  $\mathcal{I}(\boldsymbol{\theta})$  is the Fisher Information in one observation.
- ▶ A  $k \times k$  matrix

$$\mathcal{I}(\boldsymbol{\theta}) = \left[ E \left[ - \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(Y; \boldsymbol{\theta}) \right] \right]$$

- ▶ The Fisher Information in the whole sample is  $n\mathcal{I}(\boldsymbol{\theta})$

$$H_0 : \mathbf{L}\boldsymbol{\theta} = \mathbf{h}$$

Suppose  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_7)$ , and the null hypothesis is

- ▶  $\theta_1 = \theta_2$
- ▶  $\theta_6 = \theta_7$
- ▶  $\frac{1}{3}(\theta_1 + \theta_2 + \theta_3) = \frac{1}{3}(\theta_4 + \theta_5 + \theta_6)$

We can write null hypothesis in matrix form as

$$\begin{bmatrix} 1 & -1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & -1 \\ 1 & 1 & 1 & -1 & -1 & -1 & 0 \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \\ \theta_5 \\ \theta_6 \\ \theta_7 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Suppose  $H_0 : \mathbf{L}\boldsymbol{\theta} = \mathbf{h}$  is True, and  $\widehat{\mathcal{I}}(\widehat{\boldsymbol{\theta}})_n \xrightarrow{p} \mathcal{I}(\boldsymbol{\theta})$

By Slutsky 6a (Continuous mapping),

$$\sqrt{n}(\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{L}\boldsymbol{\theta}) = \sqrt{n}(\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h}) \xrightarrow{d} \mathbf{L}\mathbf{T} \sim N_k(\mathbf{0}, \mathbf{L}\mathcal{I}(\boldsymbol{\theta})^{-1}\mathbf{L}')$$

and

$$\widehat{\mathcal{I}}(\widehat{\boldsymbol{\theta}})_n^{-1} \xrightarrow{p} \mathcal{I}(\boldsymbol{\theta})^{-1}.$$

Then by Slutsky's (6c) Stack Theorem,

$$\begin{pmatrix} \sqrt{n}(\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h}) \\ \widehat{\mathcal{I}}(\widehat{\boldsymbol{\theta}})_n^{-1} \end{pmatrix} \xrightarrow{d} \begin{pmatrix} \mathbf{L}\mathbf{T} \\ \mathcal{I}(\boldsymbol{\theta})^{-1} \end{pmatrix}.$$

Finally, by Slutsky 6a again,

$$\begin{aligned} W_n &= n(\mathbf{L}\widehat{\boldsymbol{\theta}} - \mathbf{h})'(\mathbf{L}\widehat{\mathcal{I}}(\widehat{\boldsymbol{\theta}})_n^{-1}\mathbf{L}')^{-1}(\mathbf{L}\widehat{\boldsymbol{\theta}} - \mathbf{h}) \\ &\xrightarrow{d} W = (\mathbf{L}\mathbf{T} - \mathbf{0})'(\mathbf{L}\mathcal{I}(\boldsymbol{\theta})^{-1}\mathbf{L}')^{-1}(\mathbf{L}\mathbf{T} - \mathbf{0}) \sim \chi^2(r) \end{aligned}$$

## The Wald Test Statistic

$$W_n = n(\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h})'(\mathbf{L}\widehat{\mathcal{I}}(\widehat{\boldsymbol{\theta}})_n^{-1}\mathbf{L}')^{-1}(\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h})$$

- ▶ Again, null hypothesis is  $H_0 : \mathbf{L}\boldsymbol{\theta} = \mathbf{h}$
- ▶ Matrix  $\mathbf{L}$  is  $r \times k$ ,  $r \leq k$ , rank  $r$
- ▶ All we need is a consistent estimator of  $\mathcal{I}(\boldsymbol{\theta})$
- ▶  $\mathcal{I}(\widehat{\boldsymbol{\theta}})$  would do
- ▶ But it's inconvenient
- ▶ Need to compute partial derivatives and expected values in

$$\mathcal{I}(\boldsymbol{\theta}) = \left[ E\left[-\frac{\partial^2}{\partial\theta_i\partial\theta_j} \log f(Y; \boldsymbol{\theta})\right] \right]$$

## Observed Fisher Information

- ▶ To find  $\widehat{\boldsymbol{\theta}}_n$ , minimize the minus log likelihood.
- ▶ Matrix of mixed partial derivatives of the minus log likelihood is

$$\left[ -\frac{\partial^2}{\partial\theta_i\partial\theta_j} \ell(\boldsymbol{\theta}, \mathbf{Y}) \right] = \left[ -\frac{\partial^2}{\partial\theta_i\partial\theta_j} \sum_{i=1}^n \log f(Y_i; \boldsymbol{\theta}) \right]$$

- ▶ So by the Strong Law of Large Numbers,

$$\begin{aligned} \mathcal{J}_n(\boldsymbol{\theta}) &= \left[ \frac{1}{n} \sum_{i=1}^n -\frac{\partial^2}{\partial\theta_i\partial\theta_j} \log f(Y_i; \boldsymbol{\theta}) \right] \\ &\xrightarrow{a.s.} \left[ E \left( -\frac{\partial^2}{\partial\theta_i\partial\theta_j} \log f(Y; \boldsymbol{\theta}) \right) \right] = \mathbf{I}(\boldsymbol{\theta}) \end{aligned}$$



# A Consistent Estimator of $\mathcal{I}(\boldsymbol{\theta})$

Just substitute  $\widehat{\boldsymbol{\theta}}_n$  for  $\boldsymbol{\theta}$

$$\begin{aligned}\mathcal{J}_n(\widehat{\boldsymbol{\theta}}_n) &= \left[ \frac{1}{n} \sum_{i=1}^n -\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(Y_i; \boldsymbol{\theta}) \right]_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}_n} \\ &\xrightarrow{a.s.} \left[ E \left( -\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(Y; \boldsymbol{\theta}) \right) \right] = \mathcal{I}(\boldsymbol{\theta})\end{aligned}$$

- ▶ Convergence is believable but not trivial to show.
- ▶ Now we have a consistent estimator, more convenient than  $\mathcal{I}(\widehat{\boldsymbol{\theta}}_n)$ : Use  $\widehat{\mathcal{I}}(\boldsymbol{\theta})_n = \mathcal{J}_n(\widehat{\boldsymbol{\theta}}_n)$

# Approximate the Asymptotic Covariance Matrix

- ▶ Asymptotic covariance matrix of  $\hat{\boldsymbol{\theta}}_n$  is  $\frac{1}{n}\boldsymbol{\mathcal{I}}(\boldsymbol{\theta})^{-1}$ .
- ▶ Approximate it with

$$\begin{aligned}\widehat{\mathbf{V}}_n &= \frac{1}{n}\boldsymbol{\mathcal{J}}_n(\widehat{\boldsymbol{\theta}}_n)^{-1} \\ &= \frac{1}{n}\left(\frac{1}{n}\left[-\frac{\partial^2}{\partial\theta_i\partial\theta_j}\ell(\boldsymbol{\theta}, \mathbf{Y})\right]_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}_n}\right)^{-1} \\ &= \left(\left[-\frac{\partial^2}{\partial\theta_i\partial\theta_j}\ell(\boldsymbol{\theta}, \mathbf{Y})\right]_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}_n}\right)^{-1}\end{aligned}$$

# Compare

## Hessian and (Estimated) Asymptotic Covariance Matrix

- ▶  $\widehat{\mathbf{V}}_n = \left( \left[ -\frac{\partial^2}{\partial\theta_i\partial\theta_j} \ell(\boldsymbol{\theta}, \mathbf{Y}) \right]_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}_n} \right)^{-1}$
- ▶ Hessian at MLE is  $\mathbf{H} = \left[ -\frac{\partial^2}{\partial\theta_i\partial\theta_j} \ell(\boldsymbol{\theta}, \mathbf{Y}) \right]_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}_n}$
- ▶ So to estimate the asymptotic covariance matrix of  $\boldsymbol{\theta}$ , just invert the Hessian.
- ▶ The Hessian is usually available as a by-product of numerical search for the MLE.

## Connection to Numerical Optimization

- ▶ Suppose we are minimizing the minus log likelihood by a direct search.
- ▶ We have reached a point where the gradient is close to zero. Is this point a minimum?
- ▶ The Hessian is a matrix of mixed partial derivatives. If all its eigenvalues are positive at a point, the function is concave up there.
- ▶ Its *the* multivariable second derivative test.
- ▶ The Hessian at the MLE is exactly the observed Fisher information matrix.
- ▶ Partial derivatives are often approximated by the slopes of secant lines – no need to calculate them symbolically.

## So to find the estimated asymptotic covariance matrix

- ▶ Minimize the minus log likelihood numerically.
- ▶ The Hessian at the place where the search stops is exactly the observed Fisher information matrix.
- ▶ Invert it to get  $\widehat{\mathbf{V}}_n$ .
- ▶ This is so handy that sometimes we do it even when a closed-form expression for the MLE is available.

## Estimated Asymptotic Covariance Matrix $\widehat{\mathbf{V}}_n$ is Useful

- ▶ Asymptotic standard error of  $\widehat{\theta}_j$  is the square root of the  $j$ th diagonal element.
- ▶ Denote the asymptotic standard error of  $\widehat{\theta}_j$  by  $S_{\widehat{\theta}_j}$ .
- ▶ Thus

$$Z_j = \frac{\widehat{\theta}_j - \theta_j}{S_{\widehat{\theta}_j}}$$

is approximately standard normal.

## Confidence Intervals and Z-tests

Have  $Z_j = \frac{\hat{\theta}_j - \theta_j}{S_{\hat{\theta}_j}}$  approximately standard normal, yielding

- ▶ Confidence intervals:  $\hat{\theta}_j \pm S_{\hat{\theta}_j} z_{\alpha/2}$
- ▶ Test  $H_0 : \theta_j = \theta_0$  using

$$Z = \frac{\hat{\theta}_j - \theta_0}{S_{\hat{\theta}_j}}$$

## And Wald Tests

Recalling  $\widehat{\mathbf{V}}_n = \frac{1}{n} \mathcal{J}_n(\widehat{\boldsymbol{\theta}}_n)^{-1}$

$$\begin{aligned}W_n &= n(\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h})'(\mathbf{L}\widehat{\mathcal{I}}(\widehat{\boldsymbol{\theta}})_n^{-1}\mathbf{L}')^{-1}(\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h}) \\&= n(\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h})'(\mathbf{L}\mathcal{J}_n(\widehat{\boldsymbol{\theta}}_n)^{-1}\mathbf{L}')^{-1}(\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h}) \\&= n(\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h})' \left( \mathbf{L}(n\widehat{\mathbf{V}}_n)\mathbf{L}' \right)^{-1} (\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h}) \\&= n(\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h})' \frac{1}{n} \left( \mathbf{L}\widehat{\mathbf{V}}_n\mathbf{L}' \right)^{-1} (\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h}) \\&= (\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h})' \left( \mathbf{L}\widehat{\mathbf{V}}_n\mathbf{L}' \right)^{-1} (\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h})\end{aligned}$$



# Score Tests

Thank you Mr. Rao

- ▶  $\hat{\boldsymbol{\theta}}$  is the MLE of  $\boldsymbol{\theta}$ , size  $k \times 1$
- ▶  $\hat{\boldsymbol{\theta}}_0$  is the MLE under  $H_0$ , size  $k \times 1$
- ▶  $\mathbf{u}(\boldsymbol{\theta}) = (\frac{\partial \ell}{\partial \theta_1}, \dots, \frac{\partial \ell}{\partial \theta_k})'$  is the gradient.
- ▶  $\mathbf{u}(\hat{\boldsymbol{\theta}}) = 0$
- ▶ If  $H_0$  is true,  $\mathbf{u}(\hat{\boldsymbol{\theta}}_0)$  should also be close to zero.
- ▶ Under  $H_0$  for large  $N$ ,  $\mathbf{u}(\hat{\boldsymbol{\theta}}_0) \sim N_k(\mathbf{0}, \mathcal{I}(\boldsymbol{\theta}))$ , approximately.
- ▶ And,

$$S = \mathbf{u}(\hat{\boldsymbol{\theta}}_0)' \mathcal{J}_n(\hat{\boldsymbol{\theta}}_0)^{-1} \mathbf{u}(\hat{\boldsymbol{\theta}}_0) \sim \chi^2(r)$$

Where  $r$  is the number of restrictions imposed by  $H_0$

# Three Big Tests

- ▶ Score Tests: Fit just the restricted model
- ▶ Wald Tests: Fit just the unrestricted model
- ▶ Likelihood Ratio Tests: Fit Both

## Comparing Likelihood Ratio and Wald

- ▶ Asymptotically equivalent under  $H_0$ , meaning  $(W_n - G_n) \xrightarrow{P} 0$
- ▶ Under  $H_1$ ,
  - ▶ Both have approximately the same distribution (non-central chi-square)
  - ▶ Both go to infinity as  $n \rightarrow \infty$
  - ▶ But values are not necessarily close
- ▶ Likelihood ratio test tends to get closer to the right Type I error rate for small samples.
- ▶ Wald can be more convenient when testing lots of hypotheses, because you only need to fit the model once.
- ▶ Wald can be more convenient if it's a lot of work to write the restricted likelihood.

## Copyright Information

This slide show was prepared by **Jerry Brunner**, Department of Statistics, University of Toronto. It is licensed under a **Creative Commons Attribution - ShareAlike 3.0 Unported License**. Use any part of it as you like and share the result freely. The  $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$  source code is available from the course website:

<http://www.utstat.toronto.edu/~brunner/oldclass/appliedf13>