# Likelihood 2: Wald Tests[1]
## STA442/2101 Fall 2014

---

# 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)

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

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

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

# $\widehat{\boldsymbol{\theta}}_n$ is asymptotically $N_k\left(\boldsymbol{\theta}, \frac{1}{n}\boldsymbol{\mathcal{I}}(\boldsymbol{\theta})^{-1}\right)$

- Asymptotic covariance matrix of $\widehat{\boldsymbol{\theta}}_n$ is $\frac{1}{n}\boldsymbol{\mathcal{I}}(\boldsymbol{\theta})^{-1}$, and of course we don't know $\boldsymbol{\theta}$.

- For tests and confidence intervals, we need a good *approximate* asymptotic covariance matrix,

- Based on a consistent estimate of the Fisher information matrix.

- $\boldsymbol{\mathcal{I}}(\widehat{\boldsymbol{\theta}}_n)$ would do.

- But it's inconvenient: Need to compute partial derivatives and expected values in

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

and then substitute $\widehat{\boldsymbol{\theta}}_n$ for $\boldsymbol{\theta}$.

# Another approximation of the asymptotic covariance matrix

Approximate

$$\frac{1}{n}\boldsymbol{\mathcal{I}}(\boldsymbol{\theta})^{-1} = \left[n\, E[-\frac{\partial^2}{\partial\theta_i\partial\theta_j}\log f(Y;\boldsymbol{\theta})]\right]^{-1}$$

with

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

Details of why it's a good approximation are omitted.

## 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.

▶ Partial derivatives are often approximated by the slopes of secant lines – no need to calculate them symbolically.

▶ It's *the* multivariable second derivative test.

# So to find the estimated asymptotic covariance matrix

- ▶ Minimize the minus log likelihood numerically.
- ▶ The Hessian at the place where the search stops is usually available.
- ▶ 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 $Z_j = \frac{\widehat{\theta}_j - \theta_j}{S_{\widehat{\theta}_j}}$ approximately standard normal, yielding

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

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

## And Wald Tests

$$W_n = (\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h})^\top \left(\mathbf{L}\widehat{\mathbf{V}}_n\mathbf{L}^\top\right)^{-1} (\mathbf{L}\widehat{\boldsymbol{\theta}}_n - \mathbf{h})$$

A very important special case of the earlier

$$
\begin{aligned}
W_n &= n\,(\mathbf{L}\mathbf{T}_n - \mathbf{h})^\top \left(\mathbf{L}\widehat{\boldsymbol{\Sigma}}_n\mathbf{L}^\top\right)^{-1} (\mathbf{L}\mathbf{T}_n - \mathbf{h}) \\
&= (\mathbf{L}\mathbf{T}_n - \mathbf{h})^\top \left(\mathbf{L}\frac{1}{n}\widehat{\boldsymbol{\Sigma}}_n\mathbf{L}^\top\right)^{-1} (\mathbf{L}\mathbf{T}_n - \mathbf{h})
\end{aligned}
$$

# Comparing Likelihood Ratio and Wald tests

- Asymptotically equivalent under $H_0$, meaning $(W_n - G_n^2) \xrightarrow{p} 0$
- Under $H_1$,
  - Both have the same approximate distribution (non-central chi-square).
  - Both go to infinity as $n \to \infty$.
  - 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