

Outline

- 1) Wald test
- 2) Score test
- 3) Generalized likelihood ratio test
- 4) Asymptotic Relative Efficiency

Likelihood-Based Inference

Setting $X_1, \dots, X_n \stackrel{iid}{\sim} p_\theta(x)$, $p_\theta(x)$ "smooth" in θ

Assume $E_\theta \nabla \ell_i(\theta; X_i) = 0$,

$$\text{Var}_\theta [\nabla \ell_i(\theta; X_i)] = -E_\theta \nabla^2 \ell_i(\theta; X_i) = J_i(\theta) > 0,$$

$$\hat{\theta}_{MLE} \xrightarrow{P_\theta} \theta \quad (\text{consistent})$$

Then, if $\theta = \theta_0$:

$$\frac{1}{\sqrt{n}} \nabla \ell_n(\theta_0; X) \Rightarrow N(0, J(\theta_0))$$

$$\frac{1}{n} \nabla^2 J_n(\theta_0; X) \xrightarrow{P} J(\theta_0)$$

Used $O = \nabla \ell_n(\hat{\theta}_n) \approx \nabla \ell_n(\theta_0) + \nabla^2 \ell_n(\theta_0)(\hat{\theta}_n - \theta_0)$

to get $\nabla_n(\hat{\theta}_n - \theta_0) \Rightarrow N(0, J(\theta_0)^{-1})$

Can use this for inference on θ_0 !

Wald - Type Confidence Regions

Assume we have some estimator $\hat{J}_n \succ 0$ s.t.

$\frac{1}{n} \hat{J}_n \xrightarrow{P} J_1(\theta_0) \succ 0$. Then we can plug in:

If $\sqrt{n} (\hat{\theta}_n - \theta_0) \Rightarrow N_d(0, J_1(\theta_0)^{-1})$

then $(J_1(\theta_0))^{1/2} \sqrt{n} (\hat{\theta}_n - \theta_0) \Rightarrow N_d(0, I_d)$

so $\hat{J}_n^{1/2} (\hat{\theta}_n - \theta_0) \Rightarrow N_d(0, I_d)$ (Slutsky)

Leds to test of $H_0: \theta = \theta_0$:

$\|\hat{J}_n^{1/2} (\hat{\theta}_n - \theta_0)\|^2 \Rightarrow \chi_d^2$ (Reject if large)

so, $P_{\theta_0} \left(\hat{J}_n^{1/2} (\hat{\theta}_n - \theta_0) \geq \underbrace{\chi_d^2(\alpha)}_{1-\alpha \text{ quantile}} \right) \rightarrow \alpha$

Note we reject θ_0 iff $\|\hat{J}_n^{1/2} (\hat{\theta}_n - \theta_0)\|^2 > \chi_d^2(\alpha)$

\Leftrightarrow reject θ_0 iff $\theta_0 \notin \hat{\theta}_n + \hat{J}_n^{-1/2} B_d(0)$



More info \Leftrightarrow smaller ellipse (shrinks like $1/\sqrt{n}$)

Options for \hat{J}_n :

1) Most obvious is to "plug in" the MLE:

$$\begin{aligned}\hat{J}_n &= J_n(\hat{\theta}_n) && \text{(MLE for } J_n(\theta)\text{)} \\ &= \text{Var}_{\theta}(\nabla l_n(\theta; X)) \Big|_{\theta=\hat{\theta}_n}\end{aligned}$$

$$(NB) \neq \text{Var}_{\hat{\theta}_n}(\nabla l_n(\hat{\theta}_n(X); X)) = 0$$

$$\text{Or, } \hat{J}_n = -\mathbb{E}_{\theta} \nabla^2 l_n(\theta) \Big|_{\theta=\hat{\theta}_n}$$

2) Observed Fisher info:

$$\hat{J}_n = -\nabla^2 l_n(\hat{\theta}_n; X)$$

Remarks:

- Both have $\frac{1}{n} \hat{J}_n \xrightarrow{P} J_1(\theta_0)$ in "nice" iid sampling setting
- Both make sense outside of iid setting
- Heuristically, plug-in measures info about θ in "typical" data set but obs. info. measures info about θ in "this" data set

Wald interval for θ_j :

$$\text{If } \hat{\theta}_n \approx N_d(\theta_0, J_n(\theta_0)^{-1})$$

$$\text{then } \hat{\theta}_{n,j} \approx N_d(\theta_{0,j}, \underbrace{(J_n(\theta_0)^{-1})_{jj}}_{\text{s.e.}(\hat{\theta}_{n,j})^2})$$

Leads to univariate interval:

$$C_j = \hat{\theta}_{n,j} \pm \widehat{\text{s.e.}}(\hat{\theta}_{n,j}) \cdot z_{\alpha/2}$$

$$= \hat{\theta}_{n,j} \pm \sqrt{(\hat{J}_n^{-1})_{jj}} \cdot z_{\alpha/2}$$

glm function in R uses these intervals / p-values, with $\hat{J}_n = -\nabla^2 \ell(\hat{\theta}_n)$

Conf. ellipsoid for $\theta_{0,S} = (\theta_{0,j})_{j \in S}$: ($|S| = k$)

$$\hat{\theta}_{n,S} \approx N_k(\theta_{0,S}, (J_n(\theta_0)^{-1})_{SS})$$

$$\Rightarrow C_S = \hat{\theta}_{n,S} + ((\hat{J}_n^{-1})_{SS})^{1/2} B_{\chi_k^{(a)}}(0)$$

More generally, if $\sqrt{n}(\hat{\theta}_n - \theta_0) \Rightarrow N_d(0, \Sigma(\theta_0))$

and $\frac{1}{n} \sum_n \xrightarrow{P_{\theta_0}} \Sigma(\theta_0)$ ($\hat{\theta}_n$ not nec. MLE)

then we can do the same things

Ex Generalized linear model with fixed x

$x_1, \dots, x_n \in \mathbb{R}^d$ fixed

$$y_i \stackrel{\text{ind.}}{\sim} p_{\gamma_i}(y) = e^{\gamma_i y_i - A(\gamma_i)} h(y_i)$$

$$\gamma_i = \beta' x_i \quad (\underline{\text{canonical form}})$$

$$\text{Let } \mu_i(\beta) = \mathbb{E}_\beta y_i \quad (= \mu(\gamma_i(\beta)))$$

(more general: $f(\mu_i) = \beta' x_i$ for link fn f)

Most common examples:

$$\text{Logistic regression: } y_i \stackrel{\text{ind.}}{\sim} \text{Bern}\left(\frac{e^{x_i' \beta}}{1+e^{x_i' \beta}}\right)$$

$$\text{Poisson log-linear model: } y_i \stackrel{\text{ind.}}{\sim} \text{Pois}\left(e^{x_i' \beta}\right)$$

$$l_n(\beta; Y) = \sum_i (x_i' \beta) y_i - A(x_i' \beta) - \log h(y_i)$$

$$\nabla l_n(\beta; Y) = \sum_i y_i x_i - A'(x_i' \beta) \cdot x_i$$

$$= \sum_i (y_i - \mu_i(\beta)) x_i$$

$$-\nabla^2 l_n(\beta; Y) = \sum_i \ddot{A}(x_i' \beta) \cdot x_i x_i'$$

$$= \sum_i \text{Var}_\beta(y_i) \cdot x_i x_i'$$

$$= \text{Var}_\beta(\nabla l_n(\beta; Y))$$

(Not random)

$$(-\nabla^2 \ell_n(\beta))^{\frac{1}{2}} \nabla \ell_n(\beta) \sim (0, I_d) \quad \text{in finite samples}$$

$$\xrightarrow{*} N_d(0, I_d)$$

* Under regularity cond. on $X = \begin{pmatrix} -x_1 & - \\ \vdots & \vdots \\ -x_n & - \end{pmatrix}$

Taylor expansion of ℓ_n leads to

$$\hat{\Sigma}_n^{1/2} (\hat{\beta}_n - \beta) \Rightarrow N_d(0, I_d)$$

Advantages of Wald test:

- 1) Easy to invert, simple conf. regions
- 2) Asymptotically correct

Disadvantages:

- 1) Have to compute MLE
- 2) Depends on parameterization
- 3) Relies on two approximations:
 $\nabla \ell_n \approx \text{Normal}$ and $\ell_n \approx \text{quadratic}$
- 4) Need MLE to be consistent
- 5) Confidence interval / ellipsoid might go outside (H) !

Score Test

Test $H_0: \Theta = \Theta_0$ vs. $H_1: \Theta \neq \Theta_0$

We can bypass quadratic approximation entirely by using score as test stat

$$\frac{1}{\sqrt{n}} \nabla \ell_n(\Theta_0; X) \xrightarrow{P_{\Theta_0}} N_d(0, J_1(\Theta_0))$$

$$(\text{or } J_n(\Theta_0)^{-\frac{1}{2}} \nabla \ell_n(\Theta_0; X) \xrightarrow{P_{\Theta_0}} N_d(0, I_d))$$

So, we can reject $H_0: \Theta = \Theta_0$ if

$$\left\| J_n(\Theta_0)^{-\frac{1}{2}} \nabla \ell_n(\Theta_0; X) \right\|_2^2 \geq \chi_d^2(\alpha)$$

$$d=1: \frac{\dot{\ell}_n(\Theta_0)}{\sqrt{J_n(\Theta_0)}} \Rightarrow N(0, 1),$$

can do 1-sided tests

Remarks

- No quadratic approx., no MLE
- No need to estimate Fisher info at Θ_0

Can be generalized to case with nuisance params
Typically estimate via MLE on Θ_0

Score test is invariant to reparameterization:

Assume $d=1$, $\theta = g(s)$, $\dot{g}(s) > 0 \forall s$

$$q_s(x) = P_{g(s)}(x)$$

$$\begin{aligned}\dot{\ell}^{(s)}(s; x) &= \frac{d}{ds} \log P_{g(s)}(x) \\ &= \dot{\ell}^{(\theta)}(g(s); x) \cdot \dot{g}(s)\end{aligned}$$

$$J^{(s)}(s) = J^{(\theta)}(g(s)) \cdot \dot{g}(s)^2$$

$$s_0 \quad \frac{\dot{\ell}^{(s)}(s_0; x)}{\sqrt{J^{(s)}(s_0)}} \stackrel{a.s.}{=} \frac{\dot{\ell}^{(\theta)}(\theta_0; x)}{\sqrt{J^{(\theta)}(\theta_0)}}$$

$$\text{if } \theta_0 = g(s_0)$$

Ex s -parameter exp. fam:

$$X_1, \dots, X_n \stackrel{iid}{\sim} e^{\gamma' \tau(x) - A(\gamma)} h(x)$$

$$\nabla \ell(\gamma; X) = \sum \tau(X_i) - n \mu(\gamma)$$

$$\left\| J_n(\gamma_0)^{-1/2} (\sum \tau(X_i) - n \mu(\gamma_0)) \right\|^2 \Rightarrow \chi_d^2$$

$$\frac{\sum \tau(X_i) - n \mu(\gamma_0)}{\sqrt{n} \operatorname{Var}_{\gamma_0}(\tau(X_i))} \xrightarrow{P_{\gamma_0}} N(0, 1)$$

$$\underline{Ex} \quad X_1, \dots, X_n \stackrel{iid}{\sim} \text{Laplace}(\theta) = \frac{1}{2} e^{-|x-\theta|}$$

Test $H_0: \theta \leq 0$ vs $H_1: \theta > 0$ (right-tailed)

$$\ell_n(\theta; x) = - \sum_{i=1}^n |x_i - \theta| - n \log(2)$$

$$\dot{\ell}_n(\theta; x) = \sum_{i=1}^n \text{sgn}(x_i - \theta) \quad \text{sgn}(z) = \begin{cases} +1 & z > 0 \\ 0 & z = 0 \\ -1 & z < 0 \end{cases}$$

$$\begin{aligned} \dot{\ell}_n(\theta; x) &= \sum_{i=1}^n \text{sgn}(x_i) \\ &= \#\{i : X_i > 0\} - \#\{i : X_i < 0\} \\ &= 2 \underbrace{\#\{i : X_i > 0\}}_{Y \sim \text{Binom}(n, 1/2)} - n \end{aligned} \quad \text{sign test}$$

Note: this test is \approx the exact NPM LRT for
 $H_0: \theta = 0$ vs $H_1: \theta = \varepsilon$, for $\varepsilon \downarrow 0$

Intuition: Maximize power for nearby alternatives,
since we'll have power ≈ 1 for $\theta \gg \frac{1}{\sqrt{n}}$

More generally, one-sided score test is "almost"
UMP for nearby alternatives.

$$\log \frac{P_{\theta_0 + \varepsilon}(x)}{P_{\theta_0}(x)} \approx \varepsilon \dot{\ell}_n(\theta_0; x) \quad \text{for small } \varepsilon > 0$$

Ex Pearson's χ^2 test (goodness of fit)

$$N = (N_1, \dots, N_d) \sim \text{Multinom}(n, (\pi_1, \dots, \pi_d))$$

$$= \frac{n! \pi_1^{N_1} \cdots \pi_d^{N_d}}{N_1! \cdots N_d!} \mathbb{1}\{\sum N_i = n\}$$

Note $\sum \pi_j = 1$ so this is a full-rank $(d-1)$ -parameter exp. family, e.g.

$$\pi_j = \begin{cases} \frac{1}{1 + \sum_{k>1} e^{\gamma_k}} & j = 1 \\ \frac{e^{\gamma_j}}{1 + \sum_{k>1} e^{\gamma_k}} & j > 1 \end{cases}$$

$$\nabla \ell_n(\gamma; N) = (N_2, \dots, N_d) - (n\pi_2, \dots, n\pi_d)$$

$$\text{Var}_{\gamma}(\nabla \ell(\gamma)) = \begin{pmatrix} n\pi_2(1-\pi_2) & -n\pi_2\pi_3 & \dots \\ -n\pi_2\pi_3 & n\pi_3(1-\pi_3) & \dots \\ \vdots & \vdots & \ddots & n\pi_d(1-\pi_d) \end{pmatrix}$$

$$= n(\text{diag}(\pi_{2:d}) - \pi_{2:d}^{-1} \pi_{2:d}^{-1})$$

$$\Rightarrow J_n(\gamma)^{-1} = \frac{1}{n} \cdot (\text{diag}(\pi_{2:d})^{-1} - \pi_{2:d}^{-1} 1 1')$$

$$(\text{uses } (A+uv')^{-1} = A^{-1} - \frac{A^{-1}uv'A^{-1}}{1+v'A^{-1}u})$$

Score test of $H_0: \pi = \pi_0$:

(algebra)

$$\nabla \ell_n(\gamma_0) J_n^{-1}(\gamma_0) \nabla \ell_n(\gamma_0) = \sum_{j=1}^d \frac{(N_j - n\pi_{0j})^2}{n\pi_{0j}} \xrightarrow{P_{\pi_0}} \chi^2_{d-1}$$

don't really need
asy. approx

Generalized LRT

Test $H_0: \Theta = \Theta_0$ vs. $H_1: \Theta \neq \Theta_0$

Taylor expand around $\hat{\Theta}_n$:

$$\begin{aligned} l_n(\Theta_0) - l_n(\hat{\Theta}_n) &= \nabla l(\hat{\Theta}_n)^\circ + \frac{1}{2}(\Theta_0 - \hat{\Theta}_n)' \nabla^2 l_n(\tilde{\Theta}_n)(\Theta_0 - \hat{\Theta}_n) \\ &= -\frac{1}{2} \cdot \left\| \underbrace{\left(-\frac{1}{n} \nabla^2 l_n(\tilde{\Theta}_n) \right)^{1/2}}_{\xrightarrow{\rho} J_1(\Theta_0)} \underbrace{(\sqrt{n}(\Theta_0 - \hat{\Theta}_n))}_{\xrightarrow{\rho} N(0, J_1(\Theta_0)')} \right\|_2^2 \\ &\Rightarrow -\frac{1}{2} \chi_d^2 \end{aligned}$$

Test stat: $2(l_n(\hat{\Theta}_n; x) - l_n(\Theta_0; x)) \xrightarrow{\rho_{\Theta_0}} \chi_d^2$

Composite vs. Composite:

$$H_0: \Theta \in \Theta_0 \quad \text{vs} \quad H_1: \Theta \in \Theta \setminus \Theta_0 ,$$

Assume : $\mathbb{U} = \mathbb{R}^d$, Θ_0 d_0 -dim manifold

- $\Theta_0 \subset \text{relint}(\Theta)$
- $\hat{\Theta}_n \xrightarrow{P_{\Theta_0}} \Theta_0$
- Likelihood "smooth"

$$\text{Then } 2(l_n(\hat{\Theta}_n) - l_n(\hat{\Theta}_0)) \Rightarrow \chi^2_{d-d_0}$$

$$\text{where } \hat{\Theta}_0 = \arg \min_{\Theta \in \Theta_0} l_n(\Theta; x)$$

Why? Assume wlog $\Theta_0 = O$, $J_1(O) = I_d$ (reparam.)

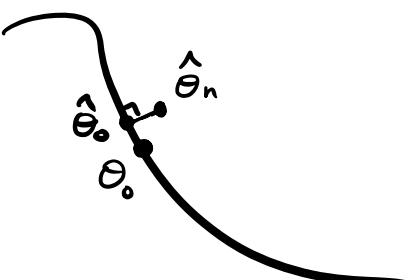
$$\text{Then } \hat{\Theta}_n \approx N_d(\Theta_0, \frac{1}{n} I_d)$$

And locally, $\nabla^2 l_n(\Theta) \approx n I_d$ near Θ_0

$$l_n(\Theta) - l_n(\hat{\Theta}_n) \approx \frac{n}{2} \|\Theta - \hat{\Theta}_n\|^2$$

$$\hat{\Theta}_0 \approx \arg \min_{\Theta \in \Theta_0} \|\Theta - \hat{\Theta}_n\| = \text{Proj}_{\Theta_0}(\hat{\Theta}_n)$$

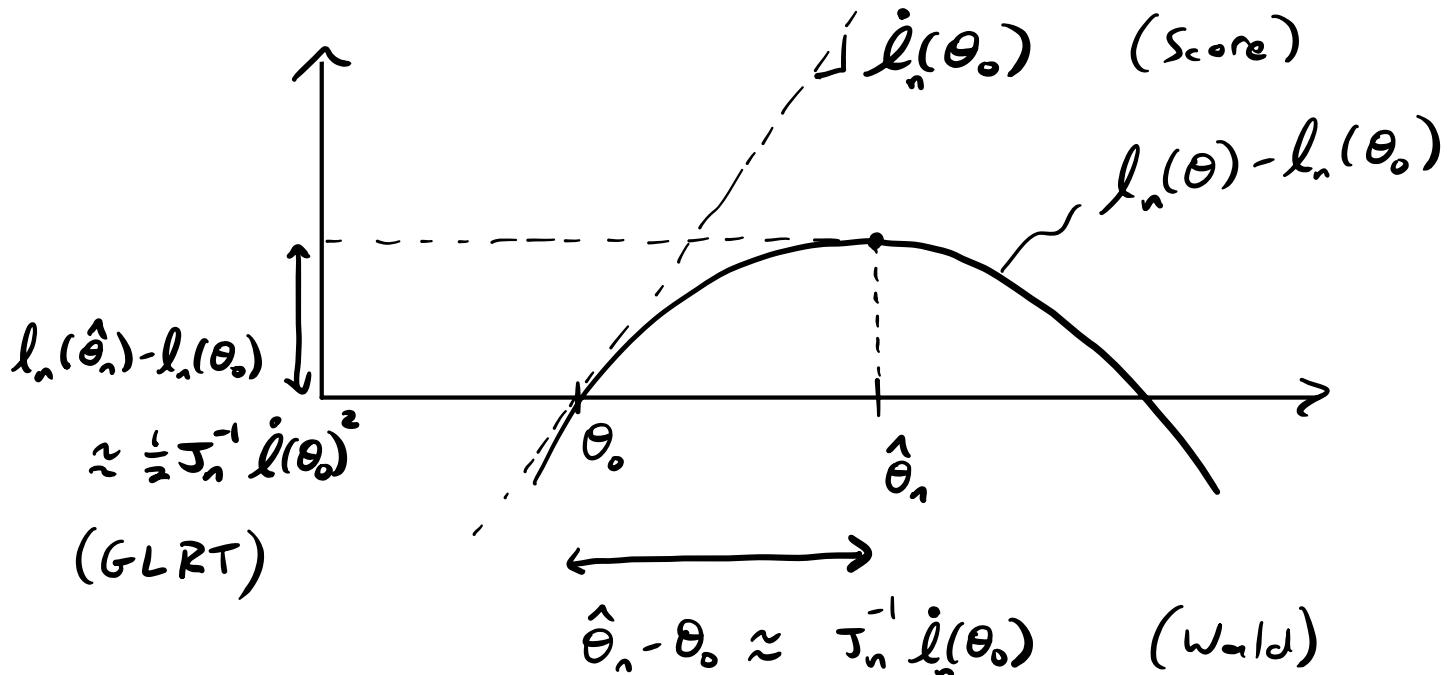
$$2(l_n(\hat{\Theta}_0) - l_n(\hat{\Theta}_n)) \approx n \|\hat{\Theta}_n - \text{Proj}_{\Theta_0}(\hat{\Theta}_n)\|^2$$

$$\begin{aligned} &= n \|\text{Proj}_{\Theta_0}^\perp(\hat{\Theta}_n)\|^2 \\ &\Rightarrow \chi^2_{d-d_0} \end{aligned}$$


Asymptotic Equivalence

Recall quadratic approx. picture ($d=1$):

$$\ell_n(\theta) - \ell_n(\theta_0) \approx \dot{\ell}_n(\theta_0)(\theta - \theta_0) + \frac{1}{2} J_n(\theta_0) (\theta - \theta_0)^2$$



For large n ,

$$\begin{aligned} \ell_n(\hat{\theta}_n) - \ell_n(\theta_0) &\approx \| J_n(\theta_0)^{-1/2} (\hat{\theta}_n - \theta_0) \|^2 \\ (\text{GLRT}) &\qquad \approx \qquad \approx \\ \| J_n^{1/2} (\hat{\theta}_n - \theta_0) \|^2 &\qquad \| J_n(\theta_0)^{-1/2} \nabla \ell_n(\theta_0) \|^2 \\ (\text{wald}) &\qquad \qquad \qquad (\text{score}) \end{aligned}$$

Asymptotic Relative Efficiency (ARE)

Suppose $\hat{\Theta}_n^{(i)}$, $i=1,2$ are two asy. Normal estimators of $\theta \in \mathbb{R}$, with

$$\sqrt{n}(\hat{\Theta}_n^{(i)} - \theta) \Rightarrow N(0, \sigma_i^2)$$

The ARE of $\hat{\Theta}^{(2)}$ wrt $\hat{\Theta}^{(1)}$ is σ_1^2 / σ_2^2
 e.g. if $\sigma_2^2 = 2\sigma_1^2$ then $\hat{\Theta}^{(2)}$ is 50% as efficient

Interpretation: Suppose $\sigma_1^2 / \sigma_2^2 = \gamma \in (0, 1)$

Then for large n ,

$$\hat{\Theta}_{[gn]}^{(1)}(x_1, \dots, x_{[gn]}) \xrightarrow{D} \hat{\Theta}_n^{(2)}(x_1, \dots, x_n) \approx N(\theta, \frac{\sigma_2^2}{n})$$

Using $\hat{\Theta}^{(2)}$ is like throwing away $100(1-\gamma)\%$ of the data and then using $\hat{\Theta}^{(1)}$