

STA261 - Module 5

Asymptotic Extensions

Rob Zimmerman

University of Toronto

August 3-5, 2021

Limitations of Finite Sample Sizes

- In almost everything we've done so far, we've assumed a sample $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\theta$ of fixed size n
- We've needed to know the distributions of various statistics of X_1, X_2, \dots, X_n
- This requirement has been very limiting, as the distributions of most statistics don't have closed forms (or are unknown entirely)

e.g., $\sum_{i=1}^n \alpha_i \cdot \chi_{(\epsilon_i)}^2$ has no closed form, $\alpha_i \in \mathbb{R}$.

- Even the exact distribution of the sample mean $\frac{1}{n} \sum_{i=1}^n X_i$ is only available for a few parametric families

even though we use \bar{X}_n , like, EVERYWHERE!

On the other hand, $\bar{X}_n \xrightarrow{d} \mathbb{E}(X_i)$ (assuming iid, $E[X] < \infty$, etc)

Driving Up the Sample Size

- On the other hand, we have plenty of *limiting* distributions as $n \rightarrow \infty$
- Example 5.1: $X_1, \dots, X_n \xrightarrow{iid} N(\mu, \sigma^2) \Rightarrow \bar{X}_n \xrightarrow{P} \mu$ by LLN, $\frac{\bar{X}_n - \mu}{\sqrt{\sigma^2/n}} \xrightarrow{D} N(0, 1)$
- Example 5.2: $X_n \sim \text{Bin}(n, p_n)$ with $n \cdot p_n \xrightarrow{n \rightarrow \infty} \lambda$, then $X_n \xrightarrow{d} \text{Poisson}(\lambda)$
- Of course, we never have $n = \infty$ in real life
- But if we have the luxury of a very large sample size, the “difference” between the exact distribution and the limiting distribution should (hopefully) be tolerable
- Since the Normal distribution is particularly nice, we will milk the CLT for all it’s worth

A Review of Standard Limiting Results

$X_n \xrightarrow{d} X$ means that $F_{X_n}(x) \xrightarrow{n \rightarrow \infty} F_X(x)$ at all continuity points of $F_X(\cdot)$.
 $X_n \xrightarrow{p} X$ means that $\forall \epsilon > 0, \lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0$.
 $X_n \xrightarrow{a.s.} X$ means that $P(\lim_{n \rightarrow \infty} X_n = X) = 1$ (not used in our course)

- In the following, let $\{X_n\}_{n \geq 1}$ and $\{Y_n\}_{n \geq 1}$ be sequences of random variables, let X be another random variable, let $x, y \in \mathbb{R}$ be constants, and let $g(\cdot)$ be a continuous function
- Theorem 5.1: If $X_n \xrightarrow{p} X$, then $X_n \xrightarrow{d} X$. If $X_n \xrightarrow{d} x$, then $X_n \xrightarrow{p} x$.
(Converse is not true in general, except when
- Theorem 5.2 (**Slutsky's theorem**): If $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{p} y$, then $Y_n \cdot X_n \xrightarrow{d} y \cdot X$ and $X_n + Y_n \xrightarrow{d} X_n + y$.
- Theorem 5.3 (**Continuous mapping theorem**): If $X_n \xrightarrow{p} X$, then $g(X_n) \xrightarrow{p} g(X)$. If $X_n \xrightarrow{d} X$, then $g(X_n) \xrightarrow{d} g(X)$. *Also a.s. convergence!*

Poll Time!

\bar{X}_n sample mean of $N(3, 1)$

\bar{Y}_n sample mean of $N(5, 1)$

- a) $\bar{X}_n^2 - \bar{Y}_n^2 \rightarrow -16$ in \mathcal{A} but not \mathcal{P}
- b) $\bar{X}_n^2 - \bar{Y}_n^2 \rightarrow -16$ in \mathcal{P} but not \mathcal{A}
- c) $\bar{X}_n^2 - \bar{Y}_n^2 \rightarrow -16$ in \mathcal{A} and \mathcal{P}
- d) None of the above

$$\bar{X}_n \xrightarrow{\text{P}} \mathbb{E}[X_i] = 3 \text{ by WLN}$$

$$\Rightarrow \bar{X}_n^2 \xrightarrow{\text{P}} 9 \text{ by CNT}$$

$$\bar{Y}_n \xrightarrow{\text{P}} \mathbb{E}[Y_i] = 5$$

$$\Rightarrow \bar{Y}_n^2 \xrightarrow{\text{P}} 25$$

$$\Rightarrow \bar{X}_n^2 - \bar{Y}_n^2 \xrightarrow{\mathcal{A}} 9 - 25 = -16 \text{ by Statkey}$$

$$\Rightarrow \bar{X}_n^2 - \bar{Y}_n^2 \xrightarrow{\text{P}} -16 \text{ by Thm S.1}$$

Notation Update

- For the rest of this module, we will accentuate statistics of finite samples with the subscript n (so \mathbf{X} is now \mathbf{X}_n , \bar{X} is now \bar{X}_n , and so on)
- For a generic statistic, we'll write $T_n = T_n(\mathbf{X}_n)$
- If we're talking about a limiting property of a sequence $\{T_n\}_{n \geq 1}$, we'll abuse notation and just write that T_n has that limiting property, when the meaning is clear from context
- Example 5.3: Instead of "the sequence of sample means $\{\bar{X}_n\}_{n \geq 1}$ converges in probability to $E[X]$ ",
we'll say " \bar{X}_n converges in probability to $E[X]$ "
or " $\bar{X}_n \xrightarrow{P} E[X]$ "

Two Big Ones

- Theorem 5.4 (**Weak law of large numbers (WLLN)**): Let X_1, X_2, \dots be a sequence of iid random variables with $\mathbb{E}[X_i] = \mu$. Then

$$\bar{X}_n \xrightarrow{p} \mu.$$

(Also the Strong Law: $\bar{X}_n \xrightarrow{\text{a.s.}} \mu$)

- Theorem 5.5 (**Central limit theorem (CLT)**): Let X_1, X_2, \dots be a sequence of iid random variables with $\mathbb{E}[X_i] = \mu$ and $\text{Var}(X_i) = \sigma^2$. Then

$$\frac{\bar{X}_n - \mu}{\sqrt{\sigma^2/n}} \xrightarrow{d} \mathcal{N}(0, 1).$$

- The CLT is equivalent to $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$, which is the form we'll be using most often

by Slutsky!

Asymptotic Unbiasedness

- As in Module 2, we're interested in estimators of $\tau(\theta)$
- But now we're concerned with their limiting behaviors as $n \rightarrow \infty$
- For finite n , we insisted that our “best” estimators be unbiased
- In the asymptotic setup, we can relax that slightly
- **Definition 5.1:** Suppose that $\{W_n\}_{n \geq 1}$ is a sequence of estimators for $\tau(\theta)$. If $\text{Bias}_\theta(W_n) \xrightarrow{n \rightarrow \infty} 0$ for all $\theta \in \Theta$, then $\{W_n\}_{n \geq 1}$ is said to be **asymptotically unbiased** for $\tau(\theta)$.

- **Example 5.4:** In the $\text{IN}(\mu, \sigma^2)$ model, $\frac{1}{n+1} \sum_{i=1}^n X_i$ is asymptotically unbiased for μ .
$$\mathbb{E}\left[\frac{1}{n+1} \sum_{i=1}^n X_i\right] = \frac{n}{n+1} \mu, \text{ so } \text{Bias}_\mu\left(\frac{1}{n+1} \sum_{i=1}^n X_i\right) = \mu \cdot \underbrace{\left(\frac{n}{n+1} - 1\right)}_{\xrightarrow{n \rightarrow \infty} 0}$$

Consistency

- $\bar{X}_n \xrightarrow{P} \mu$ is the prototypical example of an estimator converging in probability to the “right thing”
- We have a special name for this
- **Definition 5.2:** A sequence of estimators W_n of $\tau(\theta)$ is said to be **consistent** if $W_n \xrightarrow{P} \tau(\theta)$ for every $\theta \in \Theta$.
for $\tau(\theta)$
- **Example 5.5:**

$N(\mu, \sigma^2)$: $\frac{\bar{X}_n^2}{\bar{X}_n^2}$ is consistent for $\frac{\mu^2}{\mu^2 + \sigma^2}$.

Why? $\bar{X}_n \xrightarrow{P} \mu$ by WLLN

$\Rightarrow \bar{X}_n^2 \xrightarrow{P} \mu^2$ by CMT

$\bar{X}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 \xrightarrow{P} E[X_i^2] = \mu^2 + \sigma^2$ by WLLN

By Slutsky, $\frac{\bar{X}_n^2}{\bar{X}_n^2} \xrightarrow{P} \frac{\mu^2}{\mu^2 + \sigma^2}$. By Theorem S.1, $\frac{\bar{X}_n^2}{\bar{X}_n^2} \xrightarrow{P} \frac{\mu^2}{\mu^2 + \sigma^2}$.

Showing Consistency

- Sometimes it's easy to show consistency directly from the definition
- Example 5.6: Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$, where $\mu \in \mathbb{R}$ and $\sigma^2 > 0$. Is the sample mean \bar{X}_n consistent for μ ?

Let $\varepsilon > 0$. Then $P_\mu(|\bar{X}_n - \mu| < \varepsilon)$

$$\begin{aligned}&= P_\mu(-\varepsilon < \bar{X}_n - \mu < \varepsilon) \\&= P_\mu\left(-\frac{\varepsilon}{\sqrt{\sigma^2/n}} < \frac{\bar{X}_n - \mu}{\sqrt{\sigma^2/n}} < \frac{\varepsilon}{\sqrt{\sigma^2/n}}\right) \\&= P_\mu\left(-\frac{\varepsilon}{\sqrt{\sigma^2/n}} < Z < \frac{\varepsilon}{\sqrt{\sigma^2/n}}\right) \text{ where } Z \sim N(0, 1) \\&= \Phi\left(\frac{\varepsilon\sqrt{n}}{\sigma}\right) - \Phi\left(-\frac{\varepsilon\sqrt{n}}{\sigma}\right) \\&\xrightarrow{n \rightarrow \infty} \Phi(\omega) - \Phi(-\omega) = 1.\end{aligned}$$

$\Rightarrow \bar{X}_n \xrightarrow{P} \mu.$

Showing Consistency

- It's usually easier to use standard limiting results (Slutsky, continuous mapping, etc.) than to go directly from the definition
- Example 5.7: Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$, where $\mu \in \mathbb{R}$ and $\sigma^2 > 0$. Is the sample variance S_n^2 consistent for σ^2 ? Exercise: try doing this from the definition!

$$\begin{aligned} S_n^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \\ &= \frac{n}{n-1} \left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \right] \\ &= \underbrace{\frac{n}{n-1}}_1 \left[\underbrace{\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2}_2 + \underbrace{(\bar{X}_n - \mu)^2}_3 \right] \xrightarrow{\text{SLLN}} 1 \cdot (\sigma^2 + 0) = \sigma^2 \quad \text{by Slutsky (x2)} \\ &\Rightarrow S_n^2 \xrightarrow{P} \sigma^2 \quad \text{by Theorem 5.1.} \end{aligned}$$

$\textcircled{1} \xrightarrow{P} 1$

$\textcircled{2} = \overline{(X-\mu)^2} \xrightarrow{P} \mathbb{E}[(X-\mu)^2] = \sigma^2 \quad \text{by LN}$

$\textcircled{3} = \overline{(\bar{X}-\mu)^2} \xrightarrow{P} \mathbb{E}[\bar{X}-\mu]^2 = 0 \quad \text{by LN + CMT}$

Bringing Back the MSE

- In Module 2, we compared estimators by their MSEs
- To extend that idea to the asymptotic setup, we need a new mode of convergence
- **Definition 5.3:** Suppose that W_n is a sequence of estimators for $\tau(\theta)$. If $\text{MSE}_\theta(W_n) \xrightarrow{n \rightarrow \infty} 0$ for all $\theta \in \Theta$, then W_n is said to **converge in MSE** ~~to~~ $\tau(\theta)$. " $W_n \xrightarrow{\text{MSE}} \tau(\theta)$."
- **Example 5.8:** $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bin}(k_p, p), p \in (0, 1)$. Then $\bar{X}_n \xrightarrow{\text{MSE}} p$, because

$$\begin{aligned}\text{MSE}_p(\bar{X}_n) &= \underbrace{\text{Bias}_p(\bar{X}_n)^2}_{0} + \text{Var}_p(\bar{X}_n) \\ &= \text{Var}_p(\bar{X}_n) \\ &= \frac{1}{n} k_p(1-p) \xrightarrow{n \rightarrow \infty} 0 \quad \forall p \in (0, 1).\end{aligned}$$

So $\bar{X}_n \xrightarrow{\text{MSE}} p$.

Poll Time!

$$\begin{aligned} \text{MSE}_\theta[g(W_n)] &= \mathbb{E}_\theta[(g(W_n) - g(\theta))^2] \\ &= \text{Bias}_\theta(g(W_n))^2 + \text{Var}_\theta(g(W_n)) \end{aligned}$$

$$\text{MSE}_\theta(W_n) = \text{Bias}_\theta(W_n)^2 + \text{Var}_\theta(W_n)$$

$$\begin{array}{ccc} \downarrow^{m \rightarrow \infty} & = \text{Var}_\theta(W_n) \\ 0 & \Leftrightarrow & \downarrow^{n \rightarrow \infty} \\ & & 0 \end{array}.$$

Convergence in MSE is Already Good Enough

- It turns out that convergence in MSE is strong enough to guarantee consistency
- **Theorem 5.6:** If W_n is a sequence of estimators for $\tau(\theta)$ that converges in MSE for all $\theta \in \Theta$, then W_n is consistent for $\tau(\theta)$.

Proof. Let $\varepsilon > 0$, and let $\theta \in \Theta$. By Chebychev inequality,

$$P_\theta(|W_n - \tau(\theta)| \geq \varepsilon) \leq \frac{E_\theta[(W_n - \tau(\theta))^2]}{\varepsilon^2} = \frac{\text{MSE}_\theta(W_n)}{\varepsilon^2} \xrightarrow{n \rightarrow \infty} \frac{0}{\varepsilon^2} = 0.$$

So $W_n \xrightarrow{P} \tau(\theta)$, and so W_n is consistent for $\tau(\theta)$. \square

A Criterion for Consistency

- If we know $\mathbb{E}_\theta [W_n]$ and $\text{Var}_\theta (W_n)$, this next theorem often makes short work out of checking for consistency
- **Theorem 5.7:** If W_n is a sequence of estimators for $\tau(\theta)$ such that $\text{Bias}_\theta (W_n) \xrightarrow{n \rightarrow \infty} 0$ and $\text{Var}_\theta (W_n) \xrightarrow{n \rightarrow \infty} 0$ for all $\theta \in \Theta$, then W_n is consistent for $\tau(\theta)$.

Proof. For any $\theta \in \Theta$, $\text{MSE}_\theta(W_n) = \text{Bias}_\theta(W_n)^2 + \text{Var}_\theta(W_n)$

By assumption, $\text{Bias}_\theta(W_n) \xrightarrow{n \rightarrow \infty} 0$ and $\text{Var}_\theta(W_n) \xrightarrow{n \rightarrow \infty} 0$.

So $\text{MSE}_\theta(W_n) \xrightarrow{n \rightarrow \infty} 0$.

By Theorem 5.6, W_n is consistent for θ . \square

The Sample Mean is Always Consistent

- **Example 5.9:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\theta$, where $\mathbb{E}[X_i] = \mu$. Show that \bar{X}_n is consistent for μ .

By WLLN, $\bar{X}_n \xrightarrow{P} \mathbb{E}[X_i] = \mu$.

So \bar{X}_n is consistent for μ .

The Sample Variance is Always Consistent

- One can (very tediously) show that if X_1, X_2, \dots, X_n are a random sample from a distribution with a finite fourth moment, then

$$\text{Var}(S_n^2) = \frac{\mathbb{E}[(X_i - \mathbb{E}[X_i])^4]}{n} - \frac{\text{Var}(X_i)^2(n-3)}{n(n-1)}$$

and f_θ has a finite fourth moment.

- Example 5.10:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\theta$, where $\mathbb{E}[X_i] = \mu$ and $\text{Var}(X_i) = \sigma^2$. Show that S_n^2 is consistent for σ^2 .

$$\text{Bias}_{\sigma^2}(S_n^2) = 0 \text{ from Assignment 0}$$

$$\text{Var}_{\sigma^2}(S_n^2) = \frac{\mathbb{E}[(X_i - \mu)^4]}{n} - \frac{\sigma^4(n-3)}{n(n-1)} \xrightarrow{n \rightarrow \infty} 0$$

$\downarrow_{n \rightarrow \infty} \quad \downarrow_{n \rightarrow \infty}$

By Theorem 5.7, S_n^2 is consistent for σ^2 .

Choosing Among Consistent Estimators

- Consistency is practically the bare minimum we can ask for from a sequence of estimators
- There are usually plenty of sequences that are consistent for $\tau(\theta)$
- Which one should we use?
- It's tempting to go with whichever has the lowest variance for fixed n , but that would rule out a lot of fine estimators
- Example 5.11: $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Poisson}(\lambda), \lambda > 0$.

Assignment 5: tons & examples to practice with

\bar{X}_n and S_n^2 are both consistent for λ , by our previous examples.

We know that for fixed n , \bar{X}_n is the UMVE (from Module 2), but should we just ignore S_n^2 ?

- $X_1, \dots, X_n \sim N(\mu, \sigma^2)$. S_n^2 and $\hat{\sigma}_{\text{MLE}}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ are both consistent for σ^2 . Which to choose?

Asymptotic Normality

- There's a much more useful criterion, but first we need an important CLT-inspired definition
- **Definition 5.4:** Let T_n be a sequence of estimators for $\tau(\theta)$. If there exists some $\sigma^2 > 0$ such that

$$\sqrt{n}[T_n - \tau(\theta)] \xrightarrow{d} \mathcal{N}(0, \sigma^2),$$

then T_n is said to be **asymptotically normal** with mean $\tau(\theta)$ and **asymptotic variance** σ^2 .

- By virtue of the CLT, most unbiased estimators are asymptotically normal

Why not just talk about the distribution of T_n itself as $n \rightarrow \infty$?

Usually, it's a constant: $\bar{X}_n \xrightarrow{d} \nu$, for example.

Asymptotic Normality: Examples

- **Example 5.12:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bin}(k, p)$. Show that the sample mean \bar{X}_n is asymptotically normal.

By the CLT,

$$\sqrt{n}(\bar{X}_n - E_p(X_i)) \xrightarrow{d} N(0, \text{Var}_p(X_i))$$
$$\Rightarrow \sqrt{n}(\bar{X}_n - kp) \xrightarrow{d} N(0, kp(1-p))$$

So \bar{X}_n is asymptotically normal with mean kp and asymptotic variance $kp(1-p)$.

Asymptotic Normality: Examples

- **Example 5.13:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Exp}(\lambda)$. Show that the second sample moment $\overline{X^2}_n$ is asymptotically normal.

$$\overline{X^2}_n$$

$$\mathbb{E}_\lambda[X_i^2] = \text{Var}_\lambda(X_i) + \mathbb{E}_\lambda[X_i]^2 = \frac{2}{\lambda^2}$$

$$\begin{aligned}\text{Var}_\lambda(X_i^2) &= \mathbb{E}_\lambda[X_i^4] - \mathbb{E}_\lambda[X_i^2]^2 \\ &= \left(\frac{4!}{\lambda^4} \right) - \left(\frac{2}{\lambda^2} \right)^2 = \frac{20}{\lambda^4} \\ &\text{use mgfr!}\end{aligned}$$

By the CLT, $\sqrt{n}(\overline{X^2}_n - \frac{2}{\lambda^2}) \xrightarrow{d} N(0, \frac{20}{\lambda^4})$.

So $\overline{X^2}_n$ is asymptotically normal w/ mean $\frac{2}{\lambda^2}$ and asymptotic variance $\frac{20}{\lambda^4}$.

Asymptotic Distributions

- More generally, we can talk about the limiting distribution of $\sqrt{n}[T_n - \tau(\theta)]$ even when it's not Normal
- Definition 5.5:** Suppose that T_n is a sequence of estimators for $\tau(\theta)$. When it exists, the distribution of $\lim_{n \rightarrow \infty} \sqrt{n}[T_n - \tau(\theta)]$ is called the **asymptotic distribution** (or **limiting distribution**) of T_n .
- So if T_n is an asymptotically normal sequence of estimators for $\tau(\theta)$ with asymptotic variance σ^2 , then its asymptotic distribution is $\mathcal{N}(0, \sigma^2)$
- Example 5.14:** From Example 5.12, \bar{X}_n has asymptotic distribution $N(0, k_p(1-p))$.
- We might prefer to speak of the distribution of T_n itself when n is large
We can say "for large n , the dist. of \bar{X}_n approaches $N(k_p, \frac{k_p(1-p)}{n})$ "
$$\begin{aligned}\sqrt{n}(\bar{X}_n - k_p) &\xrightarrow{\text{dist.}} N(0, k_p(1-p)) \\ \Rightarrow \bar{X}_n &\xrightarrow{\text{dist.}} N(k_p, \frac{k_p(1-p)}{n})\end{aligned}$$

But we can't say "for large n , the dist. of \bar{X}_n is $N(k_p, \frac{k_p(1-p)}{n})$ "

Poll Time!

$$\bar{X}_n \xrightarrow{\text{L}} \underset{\theta}{\mathbb{E}[X_i]} = \theta.$$

The Delta Method

- If some sequence T_n is asymptotically normal for θ and some function $g(\cdot)$ is nice enough, then the next result gives a remarkably easy method of producing an asymptotically normal sequence of estimators of for $g(\theta)$
- Theorem 5.8 (**Delta method**): Suppose that $\theta \in \Theta \subseteq \mathbb{R}$ and $\sqrt{n}(T_n - \theta) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$. If $g : \mathbb{R} \rightarrow \mathbb{R}$ is continuously differentiable with $g'(\theta) \neq 0$, then

$$\sqrt{n}[g(T_n) - g(\theta)] \xrightarrow{d} \mathcal{N}(0, [g'(\theta)]^2 \sigma^2).$$

Assignment 5:
way to handle $g'(\theta) = 0$

Proof. Taylor expand $g(T_n)$ around θ to get $g(T_n) = g(\theta) + g'(\tilde{\theta}_n) \cdot (T_n - \theta)$ for some $\tilde{\theta}_n$ between T_n and θ .

$$\Rightarrow \sqrt{n}(g(T_n) - g(\theta)) = g'(\tilde{\theta}_n) \cdot \sqrt{n}(T_n - \theta)$$

①: $T_n \xrightarrow{P} \theta$ by Slutsky, we must have $\tilde{\theta}_n \xrightarrow{P} \theta$ too

By CMT, $g'(\tilde{\theta}_n) \xrightarrow{P} g'(\theta)$

②: $\sqrt{n}(T_n - \theta) \xrightarrow{d} N(0, \sigma^2)$ by assumption

By Slutsky,
$$\begin{aligned}\sqrt{n}(g(T_n) - g(\theta)) &\xrightarrow{d} g'(\theta) \cdot N(0, \sigma^2) \\ &= N(0, [g'(\theta)]^2 \sigma^2).\end{aligned}$$

□

The Delta Method: Examples

- **Example 5.15:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$ where $\mu \in \mathbb{R} \setminus \{0\}$ and $\sigma^2 > 0$. Find the limiting distribution of $1/\bar{X}_n$.

Let $g(x) = 1/x$, so $g'(x) = -1/x^2$ for $x \neq 0$.

By CLT, $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2)$.

By the delta method, $\sqrt{n}(g(\bar{X}_n) - g(\mu)) \xrightarrow{d} N(0, [g'(\mu)]^2 \cdot \sigma^2)$
 $\Rightarrow \sqrt{n}\left(\frac{1}{\bar{X}_n} - \frac{1}{\mu}\right) \xrightarrow{d} N(0, \frac{\sigma^2}{\mu^4})$.

So $\frac{1}{\bar{X}_n}$ has limiting distribution $N(0, \frac{\sigma^2}{\mu^4})$.

So for large n , the distribution of $\frac{1}{\bar{X}_n}$ approaches $N\left(\frac{1}{\mu}, \frac{\sigma^2}{\mu^4}\right)$.

The Delta Method: Examples

- **Example 5.16:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$ where $\theta \in (0, 1)$. ~~Find~~
~~the~~ Find the limiting distribution of $\log(1 - \bar{X}_n)$.

let $g(x) = \log(1-x)$, $x \in (0, 1)$.

$$\Rightarrow g'(x) = \frac{-1}{1-x}, \quad x \in (0, 1).$$

By the CLT, $\sqrt{n}(\bar{X}_n - \theta) \xrightarrow{D} N(0, \theta(1-\theta))$.

$$\begin{aligned}\text{By the delta method, } \sqrt{n}(\log(1 - \bar{X}_n) - \log(1 - \theta)) &\xrightarrow{D} N(0, (\frac{-1}{1-\theta})^2 \theta(1-\theta)) \\ &= N(0, \frac{\theta}{1-\theta})\end{aligned}$$

So $\log(1 - \bar{X}_n)$ has asymptotic distribution $N(0, \frac{\theta}{1-\theta})$.

The Delta Method: Examples

- **Example 5.17:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\theta$ where $\mathbb{E}_\theta [X_i] = \mu$ and $\text{Var}_\theta (X_i) = \sigma^2$. If $\tau : \mathbb{R} \rightarrow \mathbb{R}$ is continuously differentiable with $\tau'(\mu) \neq 0$, describe the distribution of $\tau(\bar{X}_n)$ as n becomes large.

By the CLT, $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2)$

By the delta method, $\sqrt{n}(\tau(\bar{X}_n) - \tau(\mu)) \xrightarrow{d} N(0, [\tau'(\mu)]^2 \cdot \sigma^2)$,

which is the asymptotic distribution of $\tau(\bar{X}_n)$.

So as n gets large, the distribution of $\tau(\bar{X}_n)$ approaches $N(\tau(\mu), [\tau'(\mu)]^2 \cdot \sigma^2 / n)$.

Back to Choosing Estimators

- We know that when $T_n = \bar{X}_n$, the CLT says that

$$\frac{T_n - \mathbb{E}_\theta [T_n]}{\sqrt{\text{Var}_\theta (T_n)}} \xrightarrow{d} \mathcal{N}(0, 1)$$

- Recall the Fisher information $I_n(\theta) = \text{Var}_\theta (S(\theta | \mathbf{X}_n)) = \mathbb{E}_\theta [S(\theta | \hat{\mathbf{x}}_n)^2]$
- In Module 2, we said that an unbiased estimator W_n of $\tau(\theta)$ was efficient if its variance attained the Cramér-Rao Lower Bound $[\tau'(\theta)]^2/I_n(\theta)$
- We also noticed that if the X_i 's were iid, then $I_n(\theta) = nI_1(\theta)$ by Theorem 2.10,
under the conditions of the CRLB

Asymptotic Efficiency

- So if we could replace the T_n in the CLT statement with a general unbiased and efficient W_n , it would look like

$$\frac{W_n - \tau(\theta)}{\sqrt{[\tau'(\theta)]^2/n I_1(\theta)}} \xrightarrow{d} \mathcal{N}(0, 1)$$

$\underbrace{/\sqrt{[\tau'(\theta)]^2/n I_1(\theta)}}$ redundant

$= \mathbb{E}_\theta[W_n] \text{ case unbiased}$
 $= \text{Var}_\theta(W_n) \text{ case efficient}$

- Or equivalently

$$\sqrt{n}[W_n - \tau(\theta)] \xrightarrow{d} \mathcal{N}\left(0, \frac{[\tau'(\theta)]^2}{I_1(\theta)}\right)$$

- This is not a *result*, but a *condition* that we can demand of our estimators
- Definition 5.6:** A sequence of estimators W_n is **asymptotically efficient** for $\tau(\theta)$ if

$$\sqrt{n}[W_n - \tau(\theta)] \xrightarrow{d} \mathcal{N}\left(0, \frac{[\tau'(\theta)]^2}{I_1(\theta)}\right)$$

If $\tau(\theta) = \theta$, then $\sqrt{n}(W_n - \theta) \xrightarrow{d} N(0, I_1^{-1}(\theta))$

Asymptotic Efficiency: Examples

- **Example 5.18:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Exp}(\lambda)$, where $\lambda > 0$. Show that $1/\bar{X}_n$ is asymptotically efficient for λ .

By the CLT, $\sqrt{n}(\bar{X}_n - \lambda) \xrightarrow{d} N(0, \lambda^2)$

By the delta method, $\sqrt{n}(1/\bar{X}_n - \lambda) \xrightarrow{d} N(0, \lambda^{-2})$.

$$\begin{aligned}g(x) &= \frac{1}{x}, \quad x \neq 0 \\ \Rightarrow g'(x) &= -\frac{1}{x^2}, \quad x \neq 0 \\ \Rightarrow g'(\lambda) &= -\lambda^{-2} \\ \Rightarrow g'(\lambda) &= \lambda^{-2}\end{aligned}$$

What's $I_2(\lambda)$? $I(\lambda|x) = \log(\lambda) - \lambda x$

$$\Rightarrow S(\lambda|x) = \frac{1}{\lambda} - x$$

$$\Rightarrow -\frac{\partial}{\partial \lambda} S(\lambda|x) = \frac{1}{\lambda^2}$$

$$\Rightarrow I_2(\lambda) = E_x \left[-\frac{\partial}{\partial \lambda} S(\lambda|x) \right] = \frac{1}{\lambda^2}$$

Some thing!
So $1/\bar{X}_n$ is asymptotically efficient for λ .

$$\text{Now } \frac{[g'(\lambda)]^2}{I_2(\lambda)} = \frac{1}{\lambda^{-2}} = \lambda^2$$

Asymptotic Efficiency: Examples

- **Example 5.19:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$, where $\lambda > 0$. Show that \bar{X}_n is asymptotically efficient for λ .

By the CLT, $\sqrt{n}(\bar{X}_n - \lambda) \xrightarrow{d} N(0, \lambda)$

$$\mathbb{E}(\lambda) = \lambda \Rightarrow [\mathbb{E}'(\lambda)]^2 = 1.$$

$$l(\lambda|x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$\Rightarrow l(\lambda|x) = -\lambda + x \cdot \log(\lambda) + c, \quad c \text{ free of } \lambda$$

$$\Rightarrow S(\lambda|x) = -1 + \frac{x}{\lambda}$$

$$\Rightarrow \frac{\partial}{\partial \lambda} S(\lambda|x) = \frac{x}{\lambda^2}$$

$$\Rightarrow I_1(\lambda) = \mathbb{E}_{\lambda} \left[-\frac{\partial}{\partial \lambda} S(\lambda|x) \right] = \frac{1}{\lambda^2} \cdot \mathbb{E}_{\lambda} (x) = \frac{1}{\lambda}$$

So the asymptotic variance of \bar{X}_n is $\lambda = \frac{[\mathbb{E}'(\lambda)]^2}{I_1(\lambda)}$.

So \bar{X}_n is asymptotically efficient for λ .

Large Sample Behaviour for the MLE

- We're ready to see why the MLE is almost always the point estimator of choice when n is large
- To understand this, we need to distinguish between an arbitrary parameter $\theta \in \Theta$ and the true parameter that generated the data, which we will call θ_0
- We'll show that the MLE is asymptotically efficient, under certain regularity conditions
- Under what? REGULARITY CONDITIONS !

Regularity Conditions

- Recall how the Cramér-Rao Lower Bound required some conditions:

$$\text{Var}_{\theta}(\tau(\bar{x})) < \infty \quad \forall \theta \in \mathbb{R} \quad \text{and} \quad \frac{\partial}{\partial \theta} E_{\theta}[\tau(\bar{x})] = \int_{\mathcal{X}} \frac{\partial}{\partial \theta} [\tau(\bar{x}) \cdot f_{\theta}(\bar{x})] d\bar{x}$$

- Such conditions are generically referred to as *regularity conditions*, and they're used to rule out various pathological counterexamples and edge cases
- The exact regularity conditions for our next result are quite technical and not worth getting involved with in this course
- Instead, we will go with four *sufficient* regularity conditions that are relatively easy to check, and which are satisfied by many common parametric models

Poll Time!

The MLE is Often Asymptotically Normal

- **Theorem 5.9:** Let $X_1, X_2, \dots \stackrel{iid}{\sim} f_{\theta_0}$, and let $\hat{\theta}_n(\mathbf{X}_n)$ be the MLE of θ_0 based on a sample of size n . Suppose the following regularity conditions hold:

- ▶ Θ is an open interval (not necessarily finite) in \mathbb{R}
- ▶ The log-likelihood $\ell(\theta | \mathbf{x}_n)$ is three times continuously differentiable in θ
- ▶ The support of f_θ does not depend on θ
- ▶ $I_1(\theta) < \infty$ for all $\theta \in \Theta$

Then

$$\sqrt{n}[\hat{\theta}_n(\mathbf{X}_n) - \theta_0] \xrightarrow{d} \mathcal{N}\left(0, \frac{1}{I_1(\theta_0)}\right).$$

That is, $\hat{\theta}_n(\mathbf{X}_n)$ is a consistent and asymptotically efficient estimator of θ_0 .

Proof (sketch). Take a Taylor series of $\ell'(\hat{\theta}_n)$ around θ_0 . When n is large,

Trust me on this.

$$\begin{aligned}\ell'(\hat{\theta}_n) &\approx \ell'(\theta_0) + (\hat{\theta}_n - \theta_0) \cdot \ell''(\theta_0) \\ \Rightarrow 0 &\approx \ell'(\theta_0) + (\hat{\theta}_n - \theta_0) \cdot \ell''(\theta_0) \\ \Rightarrow (\hat{\theta}_n - \theta_0) &\approx -\frac{\ell'(\theta_0)}{\ell''(\theta_0)}\end{aligned}$$

$$\Rightarrow \sqrt{n}(\hat{\theta}_n - \theta_0) \underset{①}{\approx} \frac{-\frac{1}{\sqrt{n}} l'(\theta_0)}{\frac{1}{n} l''(\theta_0)} \quad ②$$

①: $\frac{1}{\sqrt{n}} l'(\theta_0) = \frac{1}{\sqrt{n}} S(\theta_0 | \vec{x}_n)$

$$= \frac{1}{\sqrt{n}} \sum_{i=1}^n S(\theta_0 | x_i)$$

$$= \sqrt{n} \left(-\frac{1}{n} \sum_{i=1}^n S(\theta_0 | x_i) - 0 \right)$$

$$= \sqrt{n} \left(-\overline{S(\theta_0 | x)}_n - \mathbb{E}_{\theta_0} [S(\theta_0 | x_i)] \right)$$

$\xrightarrow{\text{CLT}}$ $N(0, \text{Var}_{\theta_0}(-S(\theta_0 | x_i)))$ by the CLT

$$= N(0, I_2(\theta_0)).$$

②: $\frac{1}{n} l''(\theta_0) = \frac{1}{n} \frac{\partial^2}{\partial \theta^2} S(\theta | \vec{x}_n) \Big|_{\theta=\theta_0}$

$$= \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} S(\theta | x_i) \Big|_{\theta=\theta_0}$$

$$= \frac{\left(\frac{\partial^2}{\partial \theta^2} S(\theta | x) \Big|_{\theta=\theta_0} \right)_n}{\left(\frac{\partial^2}{\partial \theta^2} S(\theta | x) \Big|_{\theta=\theta_0} \right)_n}$$

$\xrightarrow{\text{WLLN}}$ $\mathbb{E}_{\theta_0} \left[\frac{\partial^2}{\partial \theta^2} S(\theta | x_i) \Big|_{\theta=\theta_0} \right]$ by the WLLN

$$= -I_1(\theta_0) \text{ by Theorem 2.11}$$

By Slutsky's theorem, $\sqrt{n}(\hat{\theta}_n - \theta_0) \underset{\substack{\downarrow \\ -I_1(\theta_0)}}{\approx} \frac{-\frac{1}{\sqrt{n}} l'(\theta_0)}{\frac{1}{n} l''(\theta_0)} \cdot \frac{1}{-I_1(\theta_0)} \cdot N(0, I_2(\theta_0)) = N(0, \frac{1}{I_1(\theta_0)}).$

So $\hat{\theta}_n$ is asymptotically efficient!

Consistency follows immediately by Slutsky. "□"

A Useful Corollary

- Theorem 5.10: Suppose the hypotheses of Theorem 5.9 hold, and that $\tau : \Theta \rightarrow \mathbb{R}$ is continuously differentiable with $\tau'(\theta_0) \neq 0$. Then

$$\sqrt{n}[\tau(\hat{\theta}_n(\mathbf{X}_n)) - \tau(\theta_0)] \xrightarrow{d} \mathcal{N}\left(0, \frac{[\tau'(\theta_0)]^2}{I_1(\theta_0)}\right).$$

That is, $\tau(\hat{\theta}_n(\mathbf{X}_n))$ is a consistent and asymptotically efficient estimator of $\tau(\theta_0)$.

Proof: Exercise! \square

Asymptotically Efficient MLEs: Examples

- **Example 5.20:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$, where $\mu \in \mathbb{R}$ and σ^2 is known. Find the asymptotic distribution of the MLE of μ .

Check conditions & Theorem 5.9:

- ① $\mathbb{R} = \mathbb{R}$ open in \mathbb{R} ✓
- ② $\ell'(\mu|x)$ cts in μ ✓
- ③ Support of $N(\mu, \sigma^2)$ is \mathbb{R} , free of μ ✓
- ④ $I_1(\mu) = \frac{1}{\sigma^2} < \infty \quad \forall \mu \in \mathbb{R}$.

So by Theorem 5.9, $\hat{\mu} = \bar{X}_n$ is asymptotically efficient, with asymptotic distribution $N(0, \sigma^2)$.

OR: Since $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$, $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2)$

$$\hat{\mu} = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

So asymptotically efficient.

Asymptotically Efficient MLEs: Examples

- **Example 5.21:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$, where $p \in (0, 1)$. Find the asymptotic distribution of the MLE of p , and then that of $1/p$.

Check conditions:

- ① $\mathcal{R} = (0, 1)$ open in \mathbb{R} ✓
- ② $l'''(p|x)$ continuous for $p \in (0, 1)$ ✓
- ③ Support of $\text{Bernoulli}(p)$ is $\{0, 1\}$ free of p ✓
- ④ $I(p) = \frac{1}{p(1-p)} < \infty \quad \forall p \in (0, 1)$

$$\begin{aligned}\text{By Theorem 5.9, } \sqrt{n}(\bar{X}_n - p) &\xrightarrow{d} N\left(0, \frac{1}{I(p)}\right) \\ &= N\left(0, \frac{1}{p(1-p)}\right)\end{aligned}$$

Let $\tau(p) = \frac{1}{p}$. Then $\tau'(p) = -\frac{1}{p^2}$,cts on $(0, 1)$.

$$\begin{aligned}\text{So by Theorem 5.10, } \sqrt{n}\left(\frac{1}{\bar{X}_n} - \frac{1}{p}\right) &\xrightarrow{d} N\left(0, \frac{1}{p^2} + p(1-p)\right) \\ &= N\left(0, \frac{1-p}{p^3}\right)\end{aligned}$$

$$\begin{aligned}L(p|x) &= p^x(1-p)^{1-x} \\ l(p|x) &= x \cdot \log(p) + (1-x) \cdot \log(1-p) \\ l''(p|x) &= \frac{2x}{p^3} - \frac{2(1-x)}{(1-p)^3} \\ S(p|x) &= \frac{x}{p} - \frac{1-x}{1-p} \\ S'(p|x) &= -\frac{x}{p^2} + \frac{(1-x)}{(1-p)^2} \\ I_{\tau}(p) &= -E_p[S'(\tau(p|x))] = \frac{1}{p(1-p)}\end{aligned}$$

The MLE Isn't Always Asymptotically Normal

- **Example 5.22:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Unif}(0, \theta)$, where $\theta > 0$. Show that the MLE of θ is not asymptotically normal.

$\hat{\theta}_{\text{MLE}} = X_{(n)}$. If $\sqrt{n}(X_{(n)} - \theta) \xrightarrow{d} N(0, ?)$, then $Y_n := \sqrt{n}(\theta - X_{(n)}) \xrightarrow{d} N(0, ?)$ too.

But $P_\theta(Y_n \leq y)$

$$= P_\theta(\theta - X_{(n)} \leq y/\sqrt{n})$$

$$= P_\theta(-X_{(n)} \leq y/\sqrt{n} - \theta)$$

$$= P_\theta(X_{(n)} \geq \theta - y/\sqrt{n})$$

$$= 1 - P_\theta(X_{(n)} \leq \theta - y/\sqrt{n})$$

$$= 1 - \left(\frac{\theta - y/\sqrt{n}}{\theta}\right)^n$$

$$= 1 - (1 - y/\sqrt{n})^n$$

$$\xrightarrow{n \rightarrow \infty} \begin{cases} 1, & y \geq 0 \\ 0, & y < 0 \end{cases} \quad (\text{degenerate at } 0)$$

If n instead of \sqrt{n} ,

$$1 - (1 - y/n)^n \xrightarrow{n \rightarrow \infty} 1 - e^{-y} \sim \text{Exp}(1).$$

Approximate Tests and Intervals

- We've seen that a lot of statistics are asymptotically normal
- What about test statistics?
- If we're willing to approximate a test statistic (whose exact distribution we might not know for fixed n) by one with a Normal distribution, we can perform tests and create intervals that we couldn't have before
- As in Modules 3 and 4, we'll start off with tests and then use the test statistics from those to construct confidence intervals

Wilks' Theorem

- Recall the LRT statistic for testing $H_0 : \theta = \theta_0$ versus $H_A : \theta \neq \theta_0$ was given by $\lambda(\mathbf{X}_n) = \frac{L(\theta_0 | \mathbf{X}_n)}{L(\hat{\theta} | \mathbf{X}_n)}$, where $\hat{\theta} = \hat{\theta}(\mathbf{X}_n)$ is the unrestricted MLE of θ based on \mathbf{X}_n
- Amazingly, the LRT statistic always converges in distribution to a known distribution, regardless of the statistical model (assuming it's nice enough)
- **Theorem 5.11 (Wilks' theorem):** Let $X_1, X_2, \dots \stackrel{iid}{\sim} f_\theta$, where the model satisfies the same regularity conditions as in Theorem 5.9. If we test $H_0 : \theta = \theta_0$ versus $H_A : \theta \neq \theta_0$ using $\lambda(\mathbf{X}_n)$, then

$$-2 \log (\lambda(\mathbf{X}_n)) \xrightarrow{d} \chi^2_{(1)}. \quad \text{under } H_0.$$

Poll Time!

$$\lambda(\vec{x}) = \frac{\sup_{\theta \in \Theta_0} L(\theta | \vec{x}_0)}{\sup_{\theta \in \Theta} L(\theta | \vec{x})} \in (0, 1)$$

$$\Rightarrow \log(\lambda(\vec{x})) \in (-\infty, 0)$$

$$\Rightarrow -2 \log(\lambda(\vec{x})) \in (0, \infty) \text{ always positive!}$$

Approximate LRTs: Examples

- **Example 5.23:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$, where $p \in (0, 1)$. Construct an approximate size- α LRT of $H_0 : p = p_0$ versus $H_A : p \neq p_0$.

Example 3.23: $\lambda(\vec{x}_n) = \left(\frac{p_0}{\bar{x}_n}\right)^{\bar{n}X_n} \left(\frac{1-p_0}{1-\bar{x}_n}\right)^{n-\bar{n}X_n}$

$$\Rightarrow \log(\lambda(\vec{x}_n)) = n \left[\bar{x}_n \log\left(\frac{p_0}{\bar{x}_n}\right) + (1-\bar{x}_n) \cdot \log\left(\frac{1-p_0}{1-\bar{x}_n}\right) \right]$$

$$\Rightarrow -2 \cdot \log(\lambda(\vec{x}_n)) = -2n \left[\bar{x}_n \log\left(\frac{p_0}{\bar{x}_n}\right) + (1-\bar{x}_n) \cdot \log\left(\frac{1-p_0}{1-\bar{x}_n}\right) \right]$$

By Wilks' theorem, $R = \{ \vec{x} \in \mathcal{X}^n : -2n(\bar{x} \cdot \log(p_0/\bar{x}) + (1-\bar{x}) \cdot \log(\frac{1-p_0}{1-\bar{x}})) \geq \chi_{c\alpha, 1-\alpha}^2 \}$

is the rejection region of an approximate size- α test of $H_0 : p = p_0$ vs $H_A : p \neq p_0$.

Approximate LRTs: Examples

- **Example 5.24:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$, where $\mu \in \mathbb{R}$. Construct an approximate size- α LRT of $H_0 : \mu = \mu_0$ versus $H_A : \mu \neq \mu_0$.

Example 3.21. $\lambda(\vec{x}_n) = \exp\left(-\frac{n}{2\sigma^2}(\bar{x}_n - \mu_0)^2\right)$

$$\Rightarrow \log(\lambda(\vec{x}_n)) = -\frac{n}{2\sigma^2}(\bar{x}_n - \mu_0)^2$$

$$\Rightarrow -2 \cdot \log(\lambda(\vec{x}_n)) = \frac{n}{\sigma^2}(\bar{x}_n - \mu_0)^2$$

By Wilks' theorem, $R = \{\vec{x} \in \mathcal{X}^n : \frac{n}{\sigma^2}(\bar{x} - \mu_0)^2 \geq \chi^2_{(n, 1-\alpha)}\}$

Wald Tests

- **Definition 5.7:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\theta$. For testing $H_0 : \theta = \theta_0$ versus $H_A : \theta \neq \theta_0$, a **Wald test** is a test based on the **Wald statistic**

$$W_n(\mathbf{X}_n) = (\hat{\theta} - \theta_0)^2 I_n(\hat{\theta}),$$

where $\hat{\theta} = \hat{\theta}(\mathbf{X}_n)$ is the (unrestricted) MLE.

"plug-in Fisher information"

- **Theorem 5.12:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\theta$, where the model satisfies the same regularity conditions as in Theorem 5.9. If we test $H_0 : \theta = \theta_0$ versus $H_A : \theta \neq \theta_0$ using $W_n(\mathbf{X}_n)$, then

$$W_n(\mathbf{X}_n) \xrightarrow{d} \chi^2_{(1)}. \quad \text{under } H_0. \quad \underline{\text{Proof: Exercise. }} \square$$

Because $(\hat{\theta} - \theta_0) \sqrt{I_n(\hat{\theta})} \xrightarrow{d} N(0, 1)$ under H_0 (this is usually easier to work with)

Wald Tests: Examples

- **Example 5.25:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$, where $p \in (0, 1)$. Construct an approximate size- α Wald test of $H_0 : p = p_0$ versus $H_A : p \neq p_0$.

$$W_n(\bar{X}_n) = (\hat{p} - p_0)^2 I_n(\hat{p}), \text{ where } \hat{p} = \bar{X}_n.$$

Fisher information: $I_n(p) = \frac{n}{p(1-p)}$.

$$\Rightarrow I_n(\hat{p}) = \frac{n}{\bar{X}_n(1-\bar{X}_n)}$$

So the Wald statistic is $W_n(\bar{X}_n) = \frac{(\bar{X}_n - p_0)^2 \cdot n}{\bar{X}_n(1-\bar{X}_n)}$

So by Theorem 5.12, $R = \left\{ \bar{X} \in \mathcal{X}^n : \frac{(\bar{X} - p_0)^2 \cdot n}{\bar{X}(1-\bar{X})} > \chi^2_{(0.5), 1-\alpha} \right\}$ is the rejection region of an approximate size- α test of $H_0 : p = p_0$ vs $H_A : p \neq p_0$.

Q: By Theorem 5.12, $R' = \left\{ \bar{X} \in \mathcal{X}^n : \left| \frac{(\bar{X} - p_0)}{\sqrt{\bar{X}(1-\bar{X})/n}} \right| > z_{1-\alpha/2} \right\}$ is the rejection region....

Wald Tests: Examples

- *approximate*
Example 5.26: Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$, where $\mu \in \mathbb{R}$. Construct a size- α Wald test of $H_0 : \mu = \mu_0$ versus $H_A : \mu \neq \mu_0$.

$$I(\nu) = \frac{n}{\sigma^2} \quad (\text{Example 5.20, etc}), \quad \hat{\mu}_{\text{MLE}} = \bar{X}_n.$$

$$\text{So } W_n(\bar{X}_n) = \frac{(\bar{X}_n - \mu_0)^2 \cdot n}{\sigma^2} = \left(\frac{\bar{X}_n - \mu_0}{\sqrt{\sigma^2/n}} \right)^2.$$

By Theorem 5.12, $R = \{ \bar{x} \in \mathcal{X}^n : \frac{(\bar{x} - \mu_0)^2 \cdot n}{\sigma^2} \geq \chi^2_{0.5, 1-\alpha} \}$ is the rejection region

& an approximate size- α test & $H_0: \mu = \mu_0$ vs $H_A: \mu \neq \mu_0$.



In this case, it's exact!

OR: By Theorem 5.12, $R' = \left\{ \bar{x} \in \mathcal{X}^n : \left| \frac{\bar{X}_n - \mu_0}{\sqrt{\sigma^2/n}} \right| \geq z_{1-\alpha/2} \right\}$ is the rejection region. —

Hey, it's our two-sided Z-test!

Score Tests

- **Definition 5.8:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\theta$. For testing $H_0 : \theta \in \Theta_0$ versus $H_A : \theta \in \Theta_0^c$, a **score test** (also called a **Rao test** or a **Lagrange multiplier test**) is a test based on the **score statistic**

$$R_n(\mathbf{X}_n) = \frac{[S_n(\hat{\theta}_0 | \mathbf{X}_n)]^2}{I_n(\hat{\theta}_0)},$$

where $\hat{\theta}_0 = \hat{\theta}_0(\mathbf{X}_n) = \underset{\theta \in \Theta_0}{\operatorname{argmax}} L(\theta | \mathbf{X}_n)$ is the restricted MLE under H_0 .

- **Theorem 5.13:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\theta$, where the model satisfies the same regularity conditions as in Theorem 5.9. If we test $H_0 : \theta \in \Theta_0$ versus $H_A : \theta \in \Theta_0^c$ using $R_n(\mathbf{X}_n)$, then

$$R_n(\mathbf{X}_n) \xrightarrow{d} \chi_{(1)}^2. \text{ Under } H_0.$$

Because $\frac{S_n(\hat{\theta}_0 | \vec{X}_n)}{\sqrt{I_n(\hat{\theta}_0)}} \xrightarrow{d} N(0,1) \text{ under } H_0$ (this is usually easier to work with)

Score Tests: Examples

$$H_0 = \{p_0\}$$

- **Example 5.27.** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$, where $p \in (0, 1)$.

Construct an approximate size- α score test of $H_0 : p = p_0$ versus $H_A : p \neq p_0$.

$$\begin{aligned} R_n(\vec{x}_n) &= \frac{[S_n(\hat{p}_0 | \vec{x}_n)]^2}{I_n(\hat{p}_0)} = \frac{[S_n(p_0 | \vec{x}_n)]^2}{I_n(p_0)} \\ &= n^2 \left[\frac{\bar{x}_n}{p_0} - \frac{1 - \bar{x}_n}{1 - p_0} \right]^2 \frac{p_0(1 - p_0)}{n} \\ &= \frac{(\bar{x}_n - p_0)^2}{p_0(1 - p_0)/n} \quad \text{check!} \end{aligned}$$

$$L(p | \vec{x}) = p^{\sum x_i} (1-p)^{n - \sum x_i}$$

$$l(p | \vec{x}) = \sum x_i \log(p) + (n - \sum x_i) \log(1-p)$$

$$S(p | \vec{x}) = \frac{\sum x_i}{p} - \frac{n - \sum x_i}{1-p} = n \left[\frac{\bar{x}_n}{p} - \frac{1 - \bar{x}_n}{1-p} \right]$$

$$S'(p | \vec{x}) = -\frac{\sum x_i}{p^2} + \frac{n - \sum x_i}{(1-p)^2}$$

$$I_n(p) = E_p \left[-\frac{\sum x_i}{p} + \frac{n - \sum x_i}{(1-p)^2} \right] = \frac{n}{p(1-p)}$$

By Theorem 5.13, $R = \{ \vec{x} \in \mathcal{X}^n : \frac{(\bar{x} - p_0)^2}{p_0(1 - p_0)/n} \geq \chi^2_{c(1-\alpha)} \}$ is the rejection region of

an approximate size- α test of $H_0 : p = p_0$ vs $H_A : p \neq p_0$.

OR: By Theorem 5.13, $R' = \{ \vec{x} \in \mathcal{X}^n : \left| \frac{\bar{x} - p_0}{\sqrt{p_0(1 - p_0)/n}} \right| \geq z_{1-\alpha/2} \}$ is the rejection region of...

Score Tests: Examples

- **Example 5.28:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \frac{\sigma^2}{n})$, where $\mu \in \mathbb{R}$. Construct an approximate size- α score test of $H_0 : \mu = \mu_0$ versus $H_A : \mu \neq \mu_0$.

$$S(\mu | \bar{x}) = \frac{n(\bar{x} - \mu)}{\sigma^2}$$

$$\text{so } R_n(\bar{x}_n) = \frac{[S(\mu_0 | \bar{x}_n)]^2}{I_{\mu_0}(\mu_0)}$$

$$\Rightarrow S'(\mu | \bar{x}) = -\frac{n}{\sigma^2}$$

$$\Rightarrow I_{\mu_0}(\mu) = n/\sigma^2$$

$$= \frac{n^2(\bar{x}_n - \mu_0)^2}{\sigma^4} \cdot \frac{\sigma^2}{n}$$
$$= \left(\frac{\bar{x}_n - \mu_0}{\sqrt{\sigma^2/n}} \right)^2$$

By Theorem 5.13, $R = \{ \bar{x} \in \mathcal{X}^n : \frac{(\bar{x} - \mu_0)^2 \cdot n}{\sigma^2} \geq \chi^2_{\alpha/2, 1-\alpha} \}$ is the rejection region

& an approximate size- α test & $H_0: \mu = \mu_0$ vs $H_A: \mu \neq \mu_0$.



In this case, it's exact!

OR: By Theorem 5.13, $R = \{ \bar{x} \in \mathcal{X}^n : \left| \frac{\bar{x} - \mu_0}{\sqrt{\sigma^2/n}} \right| \geq z_{1-\alpha/2} \}$ is the rejection region & ...

Hey, it's our two-sided Z -test again!

The Trinity of Tests

- The LRT, the Wald test, and the score test form the backbone of classical hypothesis testing
- Observe that under H_0 , all three tests are asymptotically equivalent (i.e., all three test statistics all converge in distribution to a $\chi^2_{(1)}$)
- For this reason, the three tests are sometimes collectively referred to as the **trinity of tests**
- Although asymptotically equivalent, the speed of convergence to $\chi^2_{(1)}$ can be quite different for each one – for small n , they can be quite different in terms of power and other “small-sample” properties
 - Fact: $\ell(\theta|x) = a\theta^2 + b\theta + c$ for some $a, b, c \in \mathbb{R}$, then all three give the same result (1982)
- One might tell you to reject H_0 while another might not!

Approximate Confidence Intervals

- Using any of the asymptotic tests to test $H_0 : \theta = \theta_0$ versus $H_A : \theta \neq \theta_0$, it's sometimes possible to invert any of the test statistics to obtain an approximate $(1 - \alpha)$ -confidence interval for θ
- Out of the three, the LRT is usually the hardest to invert into an actual interval, and the Wald statistic is usually the easiest
- In practice, you can always try to use numerical solvers when the algebra doesn't work
- For Wald and score intervals, the standard recipe is to take the square root of the test statistic and compare it to $\mathcal{N}(0, 1)$

Approximate Confidence Intervals: Examples

- **Example 5.20:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$, where $p \in (0, 1)$. Construct an approximate 95% confidence interval for p based on the Wald statistic.

Example 5.20: $P_p \left(\frac{|\bar{X}_n - p|}{\sqrt{\bar{X}_n(1-\bar{X}_n)/n}} < z_{1-\alpha/2} \right) \approx 1-\alpha$

$$= P_p \left(-z_{1-\alpha/2} < \frac{p - \bar{X}_n}{\sqrt{\bar{X}_n(1-\bar{X}_n)/n}} < z_{1-\alpha/2} \right)$$
$$= P_p \left(\bar{X}_n - z_{1-\alpha/2} \cdot \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}} < p < \bar{X}_n + z_{1-\alpha/2} \cdot \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}} \right)$$

$\Rightarrow \left(\bar{X}_n - z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}, \bar{X}_n + z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}} \right)$ is an approximate $(1-\alpha)\text{-CI}$ for p

- This confidence interval shows up everywhere in polling (and is a staple of introductory Statistics classes); its half-length is called the **margin of error**

Almost always with $\alpha = 0.05 \Rightarrow z_{1-\alpha/2} \approx 1.96$

Approximate Confidence Intervals: Examples

- **Example 5.30:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$, where $p \in (0, 1)$. Construct an approximate 95% confidence interval for $\log\left(\frac{p}{1-p}\right)$ based on the Wald statistic.

From Example 5.29,

$$\begin{aligned} 1-\alpha &\approx P_p\left(\bar{X}_n - z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}} < p < \bar{X}_n + z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}\right) \\ &= P_p\left(\log\left(\frac{\bar{X}_n - z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}}{1 - \bar{X}_n + z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}}\right) < \log\left(\frac{p}{1-p}\right) < \log\left(\frac{\bar{X}_n + z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}}{1 - \bar{X}_n - z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}}\right)\right) \\ \text{So } \left(\log\left(\frac{\bar{X}_n - z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}}{1 - \bar{X}_n + z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}}\right), \log\left(\frac{\bar{X}_n + z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}}{1 - \bar{X}_n - z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}}\right)\right) \text{ is an} \\ &\text{approximate } (1-\alpha)\text{-CI for } \log\left(\frac{p}{1-p}\right). \end{aligned}$$

Approximate Confidence Intervals: Examples

- **Example 5.31:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$, where $\lambda > 0$. Construct an approximate 95% confidence interval for λ based on the Wald statistic.

MLE: $\hat{\lambda} = \bar{X}_n$.

$$l(\lambda | \bar{x}) = -n\lambda + \sum x_i \cdot \log(\lambda) + c, \text{ where } c \text{ is free of } \lambda$$

$$\Rightarrow S(\lambda | \bar{x}) = -n + \frac{\sum x_i}{\lambda}$$

$$\Rightarrow S'(\lambda | \bar{x}) = -\frac{\sum x_i}{\lambda^2}$$

$$\Rightarrow I_n(\lambda) = -E_\lambda \left[-\frac{\sum x_i}{\lambda^2} \right] = \frac{n}{\lambda}$$

$$\Rightarrow I_n(\hat{\lambda}) = \frac{n}{\bar{X}_n}$$

$$\begin{aligned} \text{So } W_n(\bar{X}_n) &= \frac{(\bar{X}_n - \lambda)^2 \cdot n}{\bar{X}_n} \\ &= \frac{(\bar{X}_n - \lambda)^2}{\bar{X}_n / n} \end{aligned}$$

$$\text{So } 1 - \alpha \approx P_\lambda \left(-z_{1-\alpha/2} < \frac{\lambda - \bar{X}_n}{\sqrt{\bar{X}_n/n}} < z_{1-\alpha/2} \right)$$

$$= P_\lambda \left(\bar{X}_n - z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n}{n}} < \lambda < \bar{X}_n + z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n}{n}} \right)$$

$$\text{So } \left(\bar{X}_n - z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n}{n}}, \bar{X}_n + z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n}{n}} \right)$$

is an approximate $(1-\alpha)$ -CI for λ .

When the Fisher Information Causes Problems...

- When f_θ is too complicated to allow for exact $(1 - \alpha)$ -confidence intervals, it's standard practice to use Wald intervals and score intervals
- But there might be another problem: calculating the Fisher information $I(\cdot)$.
- In real-life multiparameter models, $I_n(\theta)$ is a matrix and is often impossible to work out directly, which makes calculating $I_n(\hat{\theta}_0)$ or $I_n(\hat{\theta})$ futile
- When this happens, people like to swap $I_n(\cdot)$ with $J_n(\cdot)$ in the Wald and score statistics ... but is this justified ???
- Yes — it can be shown $J_n(\hat{\theta}_n)$ is a consistent estimator of $I_n(\theta_0)$
- Moreover, in a famous 1978 paper, Efron and Hinkley showed empirically that $J_n(\hat{\theta})$ is superior to $I_n(\hat{\theta})$ Optional reading, if you're curious