

STA261 - Module 5

Asymptotic Extensions

Rob Zimmerman

University of Toronto

August 2-4, 2022

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)

$$T(\vec{x}) = \bar{X}_n$$

- 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, $\bar{X}_n \xrightarrow{P} E[X]$ (assuming the X_i 's are 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: $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2) \Rightarrow \bar{X}_n \xrightarrow{d} N(\mu, \frac{\sigma^2}{n})$ by LN, $\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \xrightarrow{d} N(0, 1)$ by CLT
- Example 5.2: $X_n \sim \text{Bin}(n, p_n)$, $p_n \in (0, 1)$. If $n \cdot p_n \xrightarrow{n \rightarrow \infty} \lambda > 0$, then $X_n \xrightarrow{d} \text{Poisson}(\lambda)$
(STA257 or EXERCISE?)
- 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

- 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 $X_n \xrightarrow{p} x$.
the 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 $X_n + Y_n \xrightarrow{d} X + y$.
("UMT")
- 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 true for a.s. convergence!

$X_n \xrightarrow{d} X$ means that $F_{X_n}(x) \xrightarrow{d} F_X(x)$ for all continuity points x of $F_X(\cdot)$

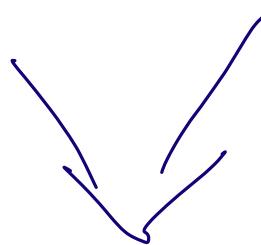
Proofs: STA347
(maybe?)

$X_n \xrightarrow{p} X$ means that $\forall \epsilon > 0, \mathbb{P}(|X_n - X| > \epsilon) \xrightarrow{n \rightarrow \infty} 0$

$X_n \xrightarrow{a.s.} X$ means that $\forall \epsilon > 0, \mathbb{P}(\lim_{n \rightarrow \infty} |X_n - X| > \epsilon) = 0$ *(not used in our course)*

Poll Time!

$$\left. \begin{array}{l} \bar{X}_n \xrightarrow{P} E[X_i] = 3 \text{ by LLN} \\ \Rightarrow \bar{X}_n^2 \xrightarrow{P} q \quad \text{by CMT} \end{array} \right\} \begin{array}{l} \bar{Y}_n \xrightarrow{P} E[Y_i] = 5 \text{ by UW} \\ \Rightarrow \bar{Y}_n^2 \xrightarrow{P} 25 \text{ by CMT} \end{array}$$



$$\bar{X}_n^2 - \bar{Y}_n^2 \xrightarrow{d} q - 25 = -16 \quad \text{by Slutsky}$$

$$\Rightarrow \bar{X}_n^2 - \bar{Y}_n^2 \xrightarrow{P} -16 \quad \text{by Theorem 5.1}$$

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 $\{\bar{X}_n\}_n$ converges in probability to μ ", we'll just write " \bar{X}_n converges in probability to μ " or " $\bar{X}_n \xrightarrow{P} \mu$ ", etc.

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

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

- 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

\backslash
by Slutsky's theorem!

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 $N(\mu, \sigma^2)$ setup, $\frac{1}{n+1} \sum_{i=1}^n X_i$ is asymptotically unbiased for μ .
Why? $E_\theta \left[\frac{1}{n+1} \sum_{i=1}^n X_i \right] = \frac{n}{n+1} \cdot \mu$ so $\text{Bias}_\theta \left(\frac{1}{n+1} \sum_{i=1}^n X_i \right) = \mu \cdot \left(\frac{n}{n+1} - 1 \right) \xrightarrow{n \rightarrow \infty} 0$

Consistency

- $\bar{X}_n \xrightarrow{p} \mu$ is the prototypical example of an estimator converging in probability to the “right thing”
by the WLLN
 - We have a special name for this
 - **Definition 5.2:** A sequence of estimators W_n of $\tau(\theta)$ is said to be **consistent** for $\tau(\theta)$ if $W_n \xrightarrow{p} \tau(\theta)$ for every $\theta \in \Theta$.
 - **Example 5.5:** $X_1, X_2, \dots \stackrel{\text{iid}}{\sim} \text{Exp}(\lambda)$, then $\frac{Y}{\bar{X}_n^2}$ is consistent for λ^2 .
Why? $\bar{X}_n \xrightarrow{p} Y_\lambda$ by LN
 $\Rightarrow g(\bar{X}_n) \xrightarrow{p} g(Y_\lambda)$ by CMF, where $g(x) = \frac{Y}{x^2}$
 $\Rightarrow \frac{Y}{\bar{X}_n^2} \xrightarrow{p} \lambda^2$
- $X_1, X_2, \dots \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$. Then
 $\frac{\bar{X}_n^2}{\lambda^2}$ is consistent for $\frac{\mu^2}{\mu^2 + \sigma^2}$
- (EXERCISE)

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

$$\Rightarrow \forall \varepsilon > 0, P(|\bar{X}_n - \mu| > \varepsilon) \xrightarrow{n \rightarrow \infty} 0$$
$$\Rightarrow \bar{X}_n \xrightarrow{P} \mu$$

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 ?

$$\begin{aligned} S_n^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \\ &= \frac{n}{n-1} \left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \right] \\ &= \frac{n}{n-1} \left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 + (\bar{X}_n - \mu)^2 \right] \end{aligned}$$

$\underbrace{\phantom{\sum_{i=1}^n}}_{\textcircled{1}}$ $\underbrace{\phantom{\sum_{i=1}^n}}_{\textcircled{2}}$ $\underbrace{(\bar{X}_n - \mu)^2}_{\textcircled{3}}$

$$\begin{aligned} \textcircled{1} &\xrightarrow{P} 1 \\ \textcircled{2} &= \overline{(X_i - \mu)^2} \xrightarrow{P} \mathbb{E}[(X - \mu)^2] = \sigma^2 \text{ by LLN} \\ \textcircled{3} &= \overline{\bar{X}_n - \mu}^2 \xrightarrow{P} \mathbb{E}[\bar{X} - \mu]^2 = 0 \text{ by LLN+CLT} \end{aligned}$$

$$\xrightarrow{P} 1 \cdot (\sigma^2 + 0) = \sigma^2 \text{ by Slutsky (x2)}$$

$$\Rightarrow S_n^2 \xrightarrow{P} \sigma^2 \text{ by Theorem 5.1 } (\sigma^2 \text{ is a constant!})$$

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:** $X_1, X_2, \dots \stackrel{\text{iid}}{\sim} \text{Bin}(k, p)$, k known, $p \in (0, 1)$. Then $\bar{X}_n \xrightarrow{\text{MSE}} k_p$ because $\text{MSE}_p(\bar{X}_n) = \underbrace{\text{Bias}_p(\bar{X}_n)}_0^2 + \text{Var}_p(\bar{X}_n)$
 $= \text{Var}_p(\bar{X}_n)$
 $= \frac{1}{n} k_p(1-p) \xrightarrow{n \rightarrow \infty} 0 \quad \forall p \in (0, 1).$

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

Poll Time!

$$\text{MSE}_\theta(W_n) = \underbrace{\text{Bias}_\theta(W_n)^2}_{0} + \text{Var}_\theta(W_n)$$
$$= \text{Var}_\theta(W_n)$$

$$\text{LHS} \rightarrow 0 \iff \text{RHS} \rightarrow 0$$
$$\iff \text{Var}_\theta(W_n) \rightarrow 0$$

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.

EXERCISE! Use Chebychev's inequality...

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

$$\downarrow \quad \downarrow \text{by assumption}$$

$0 \quad 0$

$$\Rightarrow \text{MSE}_\theta (W_n) \rightarrow 0$$

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

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

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

You don't
need to know this!

- Example 5.10:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f_\theta$, where $\mathbb{E}[X_i] = \mu$ and $\text{Var}(X_i) = \sigma^2$ and $\mathbb{E}[X_i^4] < \infty$. Show that S_n^2 is consistent for σ^2 .

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

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

$\downarrow n \rightarrow \infty \qquad \downarrow n \rightarrow \infty$

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 $\tau(\theta)$
Assignment 5: TOWS & examples to practice with
- 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: $X_1, X_2, \dots \stackrel{iid}{\sim} \text{Poisson}(\lambda), \lambda > 0$.
 \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 UMVUE (from Module 2), but should we just ignore S_n^2 altogether?
 $X_1, X_2, \dots \stackrel{iid}{\sim} N(\mu, \sigma^2)$. S_n^2 and $\hat{\sigma}_{\text{MLE}}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ are both consistent for σ^2 .
Which should we 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! $\bar{X}_n \xrightarrow{d} \mu$, for example.

The distribution of $\ln(\bar{X}_n - \mu)$ as $n \rightarrow \infty$ is more "interesting"

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, $\sqrt{n}(\bar{X}_n - E_p[\bar{X}_n]) \xrightarrow{d} N(0, \text{Var}_p(\bar{X}_n))$

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

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 \bar{X}_n^2 is asymptotically normal.

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

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

Exercise: $\mathbb{E}_{\lambda}[X_i^n] = \frac{n!}{\lambda^n}$.

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

So, \bar{X}_n^2 is asymptotically normal with 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:** $X_1, X_2, \dots \stackrel{\text{iid}}{\sim} \text{Bin}(k\theta), \theta \in (0, 1)$, \bar{X}_n has asymptotic distribution $\mathcal{N}(0, k\theta(1-\theta))$ by Example 5.12.
- We might prefer to speak of the distribution of T_n itself when n is large
We can say "for large n , the distribution of \bar{X}_n approaches $\mathcal{N}(k\theta, \frac{k\theta(1-\theta)}{n})$ "
But CANNOT say "for large n , the distribution \bar{X}_n is $\mathcal{N}(k\theta, \frac{k\theta(1-\theta)}{n})$ " ... because it's not!
 $\sqrt{n}(\bar{X}_n - k\theta) \xrightarrow{\text{approximately distributed as}} \mathcal{N}(0, k\theta(1-\theta))$
 $\Rightarrow \bar{X}_n \xrightarrow{\text{approximately distributed as}} \mathcal{N}(k\theta, \frac{k\theta(1-\theta)}{n})$

Poll Time!

$$\bar{X}_n \xrightarrow{P} \mathbb{E}[X_i] = \theta.$$
$$\Rightarrow \bar{X}_n \xrightarrow{d} \theta.$$

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:

a way to handle the case that $g'(\theta) = 0$.

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 .

$$\Rightarrow \sqrt{n}(g(T_n) - g(\theta)) = \underbrace{g'(\tilde{\theta}_n)}_{\textcircled{1}} \cdot \underbrace{\sqrt{n}(T_n - \theta)}_{\textcircled{2}}$$

①: Since $T_n \xrightarrow{\text{P}} \theta$ by Slutsky, $\tilde{\theta}_n \xrightarrow{\text{P}} \theta$ too

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

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

By Slutsky,

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

□

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 $g(x) = \frac{1}{x}$, so $g'(x) = -\frac{1}{x^2}$ for $x \neq 0$.

By the 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 $N\left(\frac{1}{\mu}, \frac{\sigma^2}{n\mu^4}\right)$.

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 limiting distribution of $\log(1 - \bar{X}_n)$.

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

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

By the delta method, $\sqrt{n}(\log(1-\bar{X}_n) - \log(1-\theta)) \xrightarrow{d} N\left(0, \left(\frac{-1}{1-\theta}\right)^2 \theta(1-\theta)\right)$
 $= N\left(0, \frac{\theta}{1-\theta}\right).$

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

The Delta Method: Examples

- **Example 5.17:** Let $X_1, X_2, \dots, \cancel{X_n} \stackrel{iid}{\sim} f_\theta$ 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, $\sqrt{n}(\tau(\bar{X}_n) - \tau(\mu)) \xrightarrow{d} N(0, [\tau'(\mu)]^2 \sigma^2)$,

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

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

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 $I_n(\theta) = \text{Var}_\theta (S(\theta | \mathbf{X}_n)) = \mathbb{E}_\theta [S(\theta | \mathbf{x})^2]$
- 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

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

$E_\theta[W_n]$ cause unbiased
 $\sqrt{n}\sigma_\theta(W_n)$ cause 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 $\sqrt{n}(W_n - \theta) \xrightarrow{d} N(0, 1/I(\theta))$

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, $\sqrt{n}(\bar{X}_n - \lambda) \xrightarrow{d} N(0, \frac{1}{\lambda^2})$

Let $g(x) = \lambda/x, x \neq 0 \Rightarrow g'(x) = -\lambda/x^2 \Rightarrow g'(\lambda) = -\lambda^2$

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

Now, what's $I(\lambda)$? $I(\lambda|x) = \log(\lambda) - \lambda x$

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

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

$$\Rightarrow I_2(\lambda) = \mathbb{E}_x[-\frac{\partial}{\partial \lambda} S(\lambda|x)] = \frac{1}{\lambda^2}$$

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

So the CRLB is $\frac{[I'(\lambda)]^2}{I_2(\lambda)} = \frac{1}{\frac{1}{\lambda^2}} = \lambda^2$

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

Here $\psi(\lambda) = \lambda$ so $[\psi'(\lambda)]^2 = 1$.

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

$$\Rightarrow l(\lambda|x) = -\lambda + x \cdot \log(\lambda) + c, \text{ where } c \in \mathbb{R} \text{ is free of } \lambda$$

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

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

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

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

$$\text{Var}_\theta(T(\tilde{x})) < \infty \quad \text{and} \quad \frac{1}{\partial \theta} E_\theta[T(\tilde{x})] = \int_{\tilde{x}} \frac{\partial}{\partial \theta} [T(\tilde{x}) \cdot f_\theta(\tilde{x})] d\tilde{x}$$

(ie, we can push $\frac{\partial}{\partial \theta}$ inside the integral)

- 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!

$\text{Unif}(0, \theta)$ does not satisfy $\frac{d}{d\theta} \int_{\chi} = \int_{\chi} \frac{\partial}{\partial \theta}$

because the support $\chi = (0, \theta)$ depends on θ

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 $\text{Supp}(f_\theta) = \{x : f_\theta(x) \neq 0\}$
- ▶ $I_1(\theta) < \infty$ 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 .

Write $\hat{\theta}_n = \hat{\theta}_n(\vec{x}_n)$ for simplicity.

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

$$\ell'(\hat{\theta}_n) \approx \ell'(\theta_0) + (\hat{\theta}_n - \theta_0) \cdot \ell''(\theta_0) \quad \text{with equality as } n \rightarrow \infty$$

Trust me on these! $\Rightarrow 0 \approx \ell'(\theta_0) + (\hat{\theta}_n - \theta_0) \cdot \ell''(\theta_0)$
The regularity conditions make them happen $\Rightarrow \hat{\theta}_n - \theta_0 \approx -\frac{\ell'(\theta_0)}{\ell''(\theta_0)}$

$$\Rightarrow \sqrt{n}(\hat{\theta}_n - \theta_0) \stackrel{?}{=} \frac{\frac{1}{\sqrt{n}} \cdot l'(\theta_0)}{\frac{1}{n} \cdot l''(\theta_0)} \stackrel{①}{\leftarrow} \stackrel{②}{\leftarrow}$$

$$\begin{aligned} ① -\frac{1}{\sqrt{n}} \cdot l'(\theta_0) &= -\frac{1}{\sqrt{n}} S(\theta_0 | \vec{x}_n) \\ &= -\frac{1}{\sqrt{n}} \sum_{i=1}^n S(\theta_0 | x_i) \\ &= \sqrt{n} \left(-\frac{1}{n} \sum_{i=1}^n S(\theta_0 | x_i) - 0 \right) \\ &= \sqrt{n} \left(-\frac{1}{n} \sum_{i=1}^n S(\theta_0 | x_i) - \overbrace{\mathbb{E}_{\theta_0}[S(\theta_0 | x)]}^{=0} \right) \\ &= \sqrt{n} \left(\overline{-S(\theta_0 | x)}_n - \mathbb{E}_{\theta_0}[S(\theta_0 | x)] \right) \end{aligned}$$

$$\xrightarrow{d} N(0, \text{Var}_{\theta_0}(-S(\theta_0 | x))) \text{ by the CLT} \\ = N(0, I_1(\theta_0))$$

$$\begin{aligned} ② \frac{1}{n} l''(\theta_0) &= \frac{1}{n} \frac{\partial^2}{\partial \theta^2} S(\theta | \vec{x}_n) \Big|_{\theta=\theta_0} \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} S(\theta | x_i) \Big|_{\theta=\theta_0} \\ &= \overline{\frac{\partial^2}{\partial \theta^2} S(\theta | x)} \Big|_{\theta=\theta_0} \\ &\xrightarrow{?} \mathbb{E}_{\theta_0}[\frac{\partial^2}{\partial \theta^2} S(\theta | x) \Big|_{\theta=\theta_0}] \\ &= -I_2(\theta_0) \end{aligned}$$

By Slutsky's theorem, $\sqrt{n}(\hat{\theta}_n - \theta_0) \stackrel{?}{=} \frac{\frac{1}{\sqrt{n}} \cdot l'(\theta_0)}{\frac{1}{n} \cdot l''(\theta_0)} \xrightarrow{+} \frac{1}{I_2(\theta_0)} \cdot N(0, I_1(\theta_0)) = N(0, 1/I_2(\theta_0))$.

$\therefore \hat{\theta}_n$ is asymptotically efficient! Consistency follows from Slutsky. "D"

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 !

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 . ($\hat{\mu}_n = \bar{X}_n$).

Check the conditions & Theorem 5.9:

- ① $\Theta = (-\infty, \infty) \subseteq \mathbb{R}$ is open in \mathbb{R} ✓
- ② $\ell'(\mu | x) = 0$ is continuous in μ ✓
- ③ $f_\mu(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} > 0 \quad \forall x \in \mathbb{R}$ ✓
- ④ $I_{\ell''}(\mu) = \frac{1}{\sigma^2} < \infty \quad \forall \mu \in \mathbb{R}$ ✓

- Θ is an open interval (not necessarily finite) in \mathbb{R}
- The log-likelihood $\ell(\theta | x_n)$ is three times continuously differentiable in θ
- The support of f_θ does not depend on θ $\text{Supp}(f_\theta) = \{x : f_\theta(x) \neq 0\}$
- $I_{\ell''}(\theta) < \infty$ for all $\theta \in \Theta$

$$\ell(\mu | x) = c - \frac{(x-\mu)^2}{2\sigma^2} \quad \text{where } c \text{ is free of } \mu$$

$$\ell'(\mu | x) = \frac{x-\mu}{\sigma^2}$$

$$\ell''(\mu | x) = -\frac{1}{\sigma^2} \leftarrow \ell'(\mu | x) \Rightarrow I_{\ell''}(\mu) = -\frac{1}{\sigma^2} \leftarrow \frac{1}{\sigma^2}$$

$$\ell'''(\mu | x) = 0$$

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

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

$\ell' \downarrow I_{\ell''}(\mu)$. 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 $1/p$.

EXERCISE !

Use the delta
method for $\hat{1}/\hat{p}$.

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.

$$\hat{\theta}_{MLE} = X_{(n)}.$$

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

$$\text{But... } P_\theta(Y_n \leq y)$$

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

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

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

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

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

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

Not a normal random variable!

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 under H_0 ,

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

No proof!
(it's very similar to the proof & Theorem 5.9)

Poll Time!

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

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

$$\Rightarrow -2 \cdot \log(\lambda(\vec{x})) \in (0, \infty)$$

Always positive!

Approximate LRTs: Examples

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

$$\text{Example 5.23} \Rightarrow \lambda(\bar{X}_n) = \left(\frac{p_0}{\bar{X}_n} \right)^{\bar{X}_n} \left(\frac{1-p_0}{1-\bar{X}_n} \right)^{n-\bar{X}_n}$$

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

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

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

is the rejection region for an approximate size- α test of $H_0: \theta = \theta_0$ vs $H_A: \theta \neq \theta_0$, when n is large.

Approximate LRTs: Examples

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

$$\begin{aligned}\text{Example 3.21} \Rightarrow \lambda(\bar{x}_n) &= \exp\left(-\frac{n}{2\sigma^2}(\bar{x}_n - \mu_0)^2\right) \\ \Rightarrow \log(\lambda(\bar{x}_n)) &= -\frac{n}{2\sigma^2}(\bar{x}_n - \mu_0)^2 \\ \Rightarrow -2 \cdot \log(\lambda(\bar{x}_n)) &= \frac{n}{\sigma^2}(\bar{x}_n - \mu_0)^2\end{aligned}$$

By Wilks' theorem, $R = \{\bar{x} \in \mathcal{X}^n : \frac{n}{\sigma^2}(\bar{x} - \mu_0)^2 \geq \chi^2_{(n, 1-\alpha)}\}$ is the rejection region of an approximate size- α test of $H_0: \mu = \mu_0$ vs $H_A: \mu \neq \mu_0$.

In fact, it's exact! (i.e., not just approximate). Why?

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**

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

"plug-in Fisher information"

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

- **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

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

Proof: EXERCISE!

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 $H_A : p \neq p_0$.

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

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

$$\text{So } W_n(\bar{x}) = \frac{(\bar{X}_n - p_0)^2 \cdot n}{\bar{X}_n(1-\bar{X}_n)} \xrightarrow{\text{under } H_0} \chi^2_{(1)}, \text{ under } H_0 \text{ by Theorem 5.12.}$$

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

Or $R' = \left\{ \bar{x} \in \mathcal{X}^n : \left| \frac{\bar{x} - p_0}{\sqrt{\bar{x}(1-\bar{x})}/n} \right| > z_{1-\alpha/2} \right\}$ is the rejection region...
because $\chi^2_{(1)} \stackrel{d}{=} N(0, 1)^2$

Wald Tests: Examples

- **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 an approximate size- α Wald test of $H_0 : \mu = \mu_0$ versus $H_A : \mu \neq \mu_0$.

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

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

By Theorem 5.12, $R = \{ \bar{x} \in \mathcal{X}^n : \left| \frac{\bar{x} - \mu_0}{\sqrt{\sigma^2/n}} \right|^2 > \chi_{\alpha/2}^2 \}$ is the rejection region & an approximate size- α test $\leftarrow H_0: \mu = \mu_0$ vs $H_A: \mu \neq \mu_0$.

(In this case, it's also exact!)

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

Hey! It's our old friend, the 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

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

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

Score Tests: Examples

$$\hat{H}_0 = \left\{ p_0 \right\}$$

- **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 $H_A : p \neq p_0$.

$$\begin{aligned}
 R_n(\bar{x}_n) &= \frac{S(\hat{p}_0 | \bar{x}_n)^2}{I_n(\hat{p}_0)} = \frac{S(p_0 | \bar{x}_n)^2}{I_n(p_0)} \\
 &= n^2 \left(\frac{\bar{x}_n}{p_0} - \frac{1 - \bar{x}_n}{1 - p_0} \right)^2 \cdot \frac{p_0(1-p_0)}{n} \\
 &= \frac{(\bar{x}_n - p_0)^2}{p_0(1-p_0)/n}
 \end{aligned}$$

$$\begin{aligned}
 L(p | \bar{x}) &= p^{\sum x_i} (1-p)^{n-\sum x_i} \\
 \ell(p | \bar{x}) &= \sum x_i \cdot \log(p) + (n - \sum x_i) \cdot \log(1-p) \\
 S(p | \bar{x}) &= \frac{\sum x_i}{p} - \frac{n - \sum x_i}{1-p} = n \left(\frac{\bar{x}_n}{p} - \frac{1 - \bar{x}_n}{1-p} \right) \\
 S'(p | \bar{x}) &= -\frac{\sum x_i}{p^2} - \frac{n - \sum x_i}{(1-p)^2} \\
 I_n(p) &= -E_p \left[\frac{-\sum x_i}{p^2} - \frac{n - \sum x_i}{(1-p)^2} \right] \\
 &= \frac{np}{p^2} + \frac{n-np}{(1-p)^2} \\
 &= n \left(\frac{1}{p} + \frac{1}{1-p} \right) \\
 &= \frac{n}{p(1-p)}
 \end{aligned}$$

By Theorem 5.13, $R = \left\{ \bar{x} \in \mathcal{X}^n : \frac{(\bar{x} - p_0)^2}{p_0(1-p_0)/n} \geq \chi_{0.05, n-1} \right\}$ is the rejection region for an approximate size- α test of $H_0 : p = p_0$ vs $H_A : p \neq p_0$.

Score Tests: Examples

- **Example 5.28:** 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 score test of $H_0 : \mu = \mu_0$ versus $H_A : \mu \neq \mu_0$.

EXERCISE !

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: if $\ell(\theta|x) = a\theta^2 + b\theta + c$ for some $a, b, c \in \mathbb{R}$, then all three tests are equivalent for finite n (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 5.29:** Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$, where $p \in (0, 1)$. Construct an approximate $(1 - \alpha)$ -confidence interval for p based on the Wald statistic.

Example 5.25:

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

- 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 see $\alpha=0.05 \Rightarrow z_{1-\alpha/2} \approx 1.96$

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 $(1 - \alpha)$ -confidence interval for $\log\left(\frac{p}{1-p}\right)$ based on the Wald statistic.

From Example 5.29,

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

is an approximate $(1-\alpha)$ -CI for $\log\left(\frac{p}{1-p}\right)$.

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 $(1 - \alpha)$ -confidence interval for λ based on the Wald statistic.

$$\hat{\lambda}_{MLE} = \bar{X}_n$$

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

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

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

$$\Rightarrow I_n(\lambda) = -E_\lambda[\sum x_i / \lambda^2] = n/\lambda$$

$$\Rightarrow I_n(\bar{X}_n) = n/\bar{X}_n$$

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

1 - $\alpha \approx P_\lambda \left(-z_{1-\alpha/2} < \frac{\lambda - \bar{X}_n}{\sqrt{\bar{X}_n/n}} < z_{1-\alpha/2} \right)$
 $= P_\lambda \left(\bar{X}_n - z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n}{n}} < \lambda < \bar{X}_n + z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n}{n}} \right)$
So $(\bar{X}_n - z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n}{n}}, \bar{X}_n + z_{1-\alpha/2} \sqrt{\frac{\bar{X}_n}{n}})$
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 actually justified ???
- Yes — it can be shown that $J_n(\hat{\theta}_n)$ 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...~