

# STA261 - Module 5

## Asymptotic Extensions

Rob Zimmerman

University of Toronto

July 30 - August 1, 2024

# 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)
- 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,  $\bar{X}_n \xrightarrow{P} E[X_i]$  (assuming the  $X_i$ 's are iid,  $E[X_i] < \infty$ , etc.)*

# Driving Up the Sample Size

- On the other hand, we have plenty of *limiting* distributions as  $n \rightarrow \infty$
- **Example 5.1:** If  $X_1, X_2, \dots \stackrel{i.i.d.}{\sim} N(\mu, \sigma^2)$ , then  $\bar{X}_n \xrightarrow{p} \mu$  by WLLN and  $\frac{\bar{X}_n - \mu}{\sqrt{\frac{\sigma^2}{n}}} \xrightarrow{d} N(0,1)$  by CLT
- **Example 5.2:** If  $X_n \sim \text{Bin}(n, p_n)$   $\forall n$  and  $n \cdot p_n \xrightarrow{n \rightarrow \infty} \lambda > 0$ , then  $X_n \xrightarrow{d} \text{Poisson}(\lambda)$
- Of course, we never have  $n = \infty$  in real life (STA257 or EXERCISE !)
- 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

Note:  $g$  need not be defined only on  $\mathbb{R}$ ! For example,  $g(\vec{x}) = \sum_{i=1}^n x_i$  totally works (and hence we can use the CMT to prove Slutsky)

- 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

the converse is not true in general; only when  $X=x$  is constant

- Theorem 5.1:** If  $X_n \xrightarrow{p} X$ , then  $X_n \xrightarrow{d} X$ . If  $X_n \xrightarrow{d} x$ , then  $X_n \xrightarrow{p} x$ .

- 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  $X_n + Y_n \xrightarrow{d} X + y$ .

- 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  $g(X_n) \xrightarrow{d} g(X)$ . ("CMT") FYI: also true for a.s. convergence

\*  $X_n \xrightarrow{d} X$  means that  $F_{X_n}(x) \xrightarrow{n \rightarrow \infty} F_X(x)$  whenever  $x$  is a continuity point of  $F_X(\cdot)$  Proofs: STATA (maybe)

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

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

↑ "a.s." = "almost surely"

# Notation Update

- For the rest of this module, we will accentuate statistics of finite samples with the subscript  $n$  (so  $\mathbf{X}$  is now  $\mathbf{X}_n$ , etc.)
- 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 writing "the sequence of sample means  $\{\bar{X}_n\}_{n \geq 1}$  converges in probability to  $\mu$ ", we'll just write " $\bar{X}_n$  converges in probability to  $\mu$ " or simply " $\bar{X}_n \xrightarrow{p} \mu$ "

# 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

by Slutsky's theorem (EXERCISE!)

# Poll Time!

On Quercus: Module 5 - Poll 1

# Asymptotic Unbiasedness

- As in Module 2, we're interested in estimators of  $\tau(\theta)$
- But now we're concerned with their limiting behaviour 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  $\mu$ .

$$\text{Why? } \mathbb{E}_\mu\left[\frac{1}{n+1} \sum_{i=1}^n X_i\right] = \frac{n}{n+1} \mu. \text{ So } \text{Bias}_\mu\left(\frac{1}{n+1} \sum_{i=1}^n X_i\right) = \mu\left(\frac{n}{n+1} - 1\right) \xrightarrow{n \rightarrow \infty} 0.$$



# Consistency

← by the LLN

- $\bar{X}_n \xrightarrow{p} \mu$  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** for  $\tau(\theta)$  if  $W_n \xrightarrow{p} \tau(\theta)$  for every  $\theta \in \Theta$ .
- **Example 5.5:**  $X_1, X_2, \dots \stackrel{iid}{\sim} \text{Exp}(\lambda)$ . Then  $\frac{1}{\bar{X}_n}$  is consistent for  $\lambda^2$ .

Why?  $\bar{X}_n \xrightarrow{p} \frac{1}{\lambda}$  by LLN.

If  $g(x) = \frac{1}{x^2}, x \neq 0$ , then  $g(\bar{X}_n) \xrightarrow{p} g(\frac{1}{\lambda})$  by CMT  
 $= \lambda^2$

$$\Rightarrow \frac{1}{\bar{X}_n^2} \xrightarrow{p} \lambda^2.$$

$X_1, X_2, \dots \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , then

$\frac{\bar{X}_n^2}{S_n^2}$  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  $\sigma^2 > 0$ . Show that the sample mean  $\bar{X}_n$  is consistent for  $\mu$ . Let  $\theta := (\mu, \sigma^2)$ .

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

$$\Rightarrow \forall \varepsilon > 0, \mathbb{P}_\theta(|\bar{X}_n - \mu| < \varepsilon) \xrightarrow{n \rightarrow \infty} 1 \quad \Rightarrow \bar{X}_n \xrightarrow{p} \mu.$$

# 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  $\sigma^2 > 0$ . Show that the sample variance  $S_n^2$  is consistent for  $\sigma^2$ .

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

$$① \xrightarrow{p} 1$$

$$② = \overline{(X_i - \mu)^2} \xrightarrow{p} \mathbb{E}[(X_i - \mu)^2] = \sigma^2 \text{ by WLLN}$$

$$③ = \overline{X_i - \mu}^2 \xrightarrow{p} \mathbb{E}[X_i - \mu]^2 = 0 \text{ by WLLN + CMT}$$

$$\begin{aligned}
 &\xrightarrow{d} 1 \cdot (\sigma^2 + 0) \text{ by Slutsky (x2)} \\
 &= \sigma^2
 \end{aligned}$$

$$\Rightarrow S_n^2 \xrightarrow{p} \sigma^2 \text{ by Theorem 5.1}$$

Actually, we didn't use the  $\mathcal{N}(\mu, \sigma^2)$  distribution anywhere here! We showed that  $S_n^2$  is always consistent for  $\text{Var}(X_i)$  (assuming  $\text{Var}(X_i) < \infty$ )

# 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:**  $X_1, X_2, \dots \stackrel{i.i.d.}{\sim} \text{Bin}(K, p)$ . Then  $\bar{X}_n \xrightarrow{\text{MSE}} \underbrace{K \cdot p}_{= E[X_i]}$ .

Why?  $\text{MSE}_p(\bar{X}_n) = \underbrace{\text{Bias}_p(\bar{X}_n)^2}_{=0 \text{ since } \bar{X}_n \text{ is always unbiased for } E[X_i]} + \text{Var}_p(\bar{X}_n)$   
 $= \text{Var}_p(\bar{X}_n)$   
 $= \frac{1}{n} Kp(1-p) \xrightarrow{n \rightarrow \infty} 0$ . So  $\bar{X}_n \xrightarrow{\text{MSE}} Kp$ .

# Poll Time!

On Quercus: Module 5 - Poll 2

# 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! Hint: use Chebyshev!

(Always remember Chebyshev...)

# 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  $\theta \in \Theta$ , 
$$\text{MSE}_\theta(W_n) = \text{Bias}_\theta(W_n)^2 + \text{Var}_\theta(W_n).$$
$$\begin{array}{ccc} & \downarrow_{n \rightarrow \infty} & \downarrow_{n \rightarrow \infty} \\ & 0 & 0 \end{array}$$

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

By Theorem 5.6,  $W_n$  is consistent for  $\tau(\theta)$ .  $\square$

# 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  $\mu$ .

$$\text{Bias}_0(\bar{X}_n) = \mathbb{E}_0[\bar{X}_n - \mu] = 0.$$

$$\text{Var}_0(\bar{X}_n) = \frac{1}{n} \cdot \text{Var}_0(x_i) \xrightarrow{n \rightarrow \infty} 0.$$

By Theorem 5.7,  $\bar{X}_n$  is consistent for  $\mu$ .

(Also  $\bar{X}_n \xrightarrow{p} \mu$  is exactly what the WLLN says)



# 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

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

(Don't need to memorize!)

- **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  $\sigma^2$ .

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

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

By Theorem 5.7,  $S_n^2$  is consistent for  $\sigma^2$ .

# 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: TONS of examples to play 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  $\lambda$ , by previous stuff. For fixed  $n$ , we know (Module 2) that  $\bar{X}_n$  is the UMVUE of  $\lambda$ , but does that mean we should just completely ignore  $S_n^2$ ?

- $X_1, X_2, \dots \stackrel{iid}{\sim} N(\mu, \sigma^2).$   $S_n^2$  and  $\hat{\sigma}_{mle}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  are both consistent for  $\sigma^2$ . Which one should we use?

# Asymptotic Normality

- There's a much more useful criterion, but first we need an important CLT-inspired definition

$$T_n = T_n(\bar{X}_n)$$

- **Definition 5.4:** Let  $T_n$  be a sequence of estimators for  $\tau(\theta)$ . If there exists some  $\sigma^2 > 0$  such that

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

FYI: the definition extends to where  $\sqrt{n}$  and  $\tau(\theta)$  are replaced by sequences of constants  $\{b_n\}_{n \geq 1}$  and  $\{c_n\}_{n \geq 1}$ :  $b_n(T_n - c_n) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$

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

Note:  $\tau(\theta)$  is not necessarily the mean of  $T_n$

- 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 some degenerate distribution (i.e., a constant).  $\bar{X}_n \xrightarrow{p} \mu$ , for example.

The distribution of  $\sqrt{n}(\bar{X}_n - \mu)$  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.

$$\sqrt{n}(\bar{X}_n - E[\bar{X}_n]) \xrightarrow{d} N(0, \text{Var}_0(X_i)) \text{ by the CLT}$$

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

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.

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

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

EXERCISE: prove

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

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 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  $T_n$ .
- In other words, if  $\sqrt{n}(T_n - \tau(\theta)) \xrightarrow{d} Y$  for some r.v.  $Y$ , then the asymptotic distribution of  $T_n$  is exactly the distribution of  $Y$
- So if  $T_n$  is an asymptotically normal sequence of estimators for  $\tau(\theta)$  with asymptotic variance  $\sigma^2$ , then its asymptotic distribution is  $\mathcal{N}(0, \sigma^2)$

- **Example 5.14:**  $X_1, X_2, \dots \stackrel{i.i.d.}{\sim} \text{Bin}(k, \theta) \Rightarrow \bar{X}_n$  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  $\mathcal{N}(k\theta, \frac{k\theta(1-\theta)}{n})$ ,"  $\sqrt{n}(\bar{X}_n - k\theta) \sim \mathcal{N}(0, k\theta(1-\theta))$   
 but we **CANNOT** say "for large  $n$ , the distribution of  $\bar{X}_n$  is  $\mathcal{N}(k\theta, \frac{k\theta(1-\theta)}{n})$ "  $\sim$  "approximately distributed as"  
 ... because it's not!

# Poll Time!

On Quercus: Module 5 - Poll 3

# The Delta Method

- If some sequence  $T_n$  is asymptotically normal for  $\theta$  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  $\theta$  to get  $g(T_n) = g(\theta) + g'(\tilde{\theta}_n) \cdot (T_n - \theta)$  for some  $\tilde{\theta}_n$  between  $T_n$  and  $\theta$ .

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

①: Since  $T_n \xrightarrow{P} \theta$  by Slutsky (check!),  $\tilde{\theta}_n \xrightarrow{P} \theta$ . By CMT,  $g'(\tilde{\theta}_n) \xrightarrow{P} g'(\theta)$ .

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

By Slutsky,

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

□



# 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  $g(x) = 1/x$ ,  $x \neq 0$ . Then  $g'(x) = -1/x^2$ ,  $x \neq 0$ .

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

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

So  $1/\bar{X}_n$  has asymptotic distribution  $N(1/\mu, \sigma^2/n\mu^4)$ .

For large  $n$ , the distribution of  $1/\bar{X}_n$  is approximately  $N(1/\mu, \sigma^2/n\mu^4)$ .

# The Delta Method: Examples

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

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

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

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

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

For large  $n$ , the distribution of  $\log(1 - \bar{X}_n)$  is approximately  $N\left(\log(1-\theta), \frac{\theta}{n(1-\theta)}\right)$ .

# The Delta Method: Examples

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

So the asymptotic distribution of  $\tau(\bar{X}_n)$  is  $N(\tau(\mu), \frac{[\tau'(\mu)]^2 \sigma^2}{n})$ .

For large  $n$ , the distribution of  $\tau(\bar{X}_n)$  is approximately  $N(\tau(\mu), \frac{[\tau'(\mu)]^2 \sigma^2}{n})$   
i.e., the distribution of  $\sqrt{n}(\tau(\bar{X}_n) - \tau(\mu))$  as  $n \rightarrow \infty$

# 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  $I_n(\theta) = \text{Var}_\theta (S(\theta | \mathbf{X}_n))$
- 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 same conditions as the CRLB itself

# 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

$$\frac{W_n - \tau(\theta)}{\sqrt{[\tau'(\theta)]^2/nI_1(\theta)}} \xrightarrow{d} \mathcal{N}(0, 1)$$

- 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

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

# 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  $\lambda$ .

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

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

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

Now, what's  $I_1(\lambda)$ ?  $l(\lambda|x) = \log(\lambda) - \lambda x$

$$\Rightarrow S(\lambda|x) = 1/\lambda - x$$

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

$$\Rightarrow I_1(\lambda) = \mathbb{E}_x[-\frac{\partial}{\partial \lambda} S(\lambda|x)] = 1/\lambda^2.$$

Same thing! So we conclude that  $1/\bar{X}_n$  is asymptotically efficient for  $\lambda$ .

So the CRLB is  $\frac{(\psi'(\lambda))^2}{I_1(\lambda)} = \frac{1}{1/\lambda^2} = \lambda^2$

# Asymptotic Efficiency: Examples

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

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

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

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

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

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

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

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

$\psi(\lambda) = \lambda$

Since the asymptotic variance of  $\bar{X}_n$  is equal to the CRLB for unbiased estimators of  $\lambda$ , we conclude that  $\bar{X}_n$  is asymptotically efficient for  $\lambda$ .

# 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  $\theta_0$
- We'll show that the MLE is asymptotically efficient, under certain “regularity conditions”



# Regularity Conditions

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

$$\textcircled{1} \text{Var}_\theta(T(\tilde{X}_n)) < \infty \quad \forall \theta \in \Theta \quad \textcircled{2} \quad \frac{d}{d\theta} \mathbb{E}_\theta[T(\tilde{X}_n)] = \int \frac{\partial}{\partial \theta} [T(\tilde{x}) \cdot f_\theta(\tilde{x})] d\tilde{x}$$

- Such conditions are generically referred to as *regularity conditions*, and they're used to rule out various pathological counterexamples and edge cases (i.e., we can push the derivative inside the integral)
- 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!

## On Quercus: Module 5 - Poll 4

Unif(0,  $\theta$ ) does not satisfy  $\frac{d}{d\theta} \int_{\mathcal{X}} \dots = \int_{\mathcal{X}} \frac{\partial}{\partial \theta} \dots$

because the support  $\mathcal{X} = (0, \theta)$  depends on  $\theta$ .

# 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  $\theta_0$  based on a sample of size  $n$ . Suppose the following regularity conditions hold:
  - ▶  $\Theta$  is an open interval (not necessarily finite) in  $\mathbb{R}$
  - ▶ The log-likelihood  $\ell(\theta | \mathbf{x}_n)$  is three times continuously differentiable in  $\theta$
  - ▶ The support of  $f_\theta$  does not depend on  $\theta$
  - ▶  $I_1(\theta) < \infty$  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  $\theta_0$ .

Write  $\hat{\Theta}_n = \hat{\Theta}_n(\bar{\mathbf{x}})$  for simplicity.

*Proof (sketch).* Take a Taylor series of  $\ell'(\hat{\Theta}_n | \bar{\mathbf{x}})$  around  $\theta_0$ . For large  $n$ , we get

$$\begin{aligned} \ell'(\hat{\Theta}_n | \bar{\mathbf{x}}) &\approx \ell'(\theta_0 | \bar{\mathbf{x}}) + (\hat{\Theta}_n - \theta_0) \cdot \ell''(\theta_0 | \bar{\mathbf{x}}) \text{ with equality as } n \rightarrow \infty \\ \Rightarrow 0 &\approx \ell'(\theta_0 | \bar{\mathbf{x}}) + (\hat{\Theta}_n - \theta_0) \cdot \ell''(\theta_0 | \bar{\mathbf{x}}) \\ \Rightarrow \hat{\Theta}_n - \theta_0 &\approx \frac{-\ell'(\theta_0 | \bar{\mathbf{x}})}{\ell''(\theta_0 | \bar{\mathbf{x}})} \end{aligned}$$

(this is where those regularity conditions are needed - trust me on this!)

$$\rightarrow \sqrt{n}(\hat{\theta}_n - \theta_0) \approx \frac{-\frac{1}{\sqrt{n}} l'(\theta_0 | \vec{x})}{\frac{1}{n} l''(\theta_0 | \vec{x})} \quad \textcircled{1}$$

$$\begin{aligned} \textcircled{1} \quad -\frac{1}{\sqrt{n}} l'(\theta_0 | \vec{x}) &= -\frac{1}{\sqrt{n}} S(\theta_0 | \vec{x}) \\ &= -\frac{1}{\sqrt{n}} \sum_{i=1}^n S(\theta_0 | x_i) \\ &= \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n S(\theta_0 | x_i) - 0 \right) \\ &= \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n S(\theta_0 | x_i) - \mathbb{E}_{\theta_0} [S(\theta_0 | x_i)] \right) \\ &= -\sqrt{n} \left( S(\theta_0 | X) - \mathbb{E}_{\theta_0} [S(\theta_0 | x_i)] \right) \\ &\xrightarrow{d} -N(0, \text{Var}_{\theta_0}(S(\theta_0 | x))) \text{ by the CLT} \\ &= -N(0, I_1(\theta_0)) \end{aligned}$$

$$\begin{aligned} \textcircled{2} \quad \frac{1}{n} l''(\theta_0 | \vec{x}) &= \frac{1}{n} \frac{\partial^2}{\partial \theta^2} S(\theta | \vec{x}) \Big|_{\theta=\theta_0} \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} S(\theta | x_i) \Big|_{\theta=\theta_0} \\ &= \frac{\frac{\partial^2}{\partial \theta^2} S(\theta | X) \Big|_{\theta=\theta_0}}{\mathbb{E}_{\theta_0} \left[ \frac{\partial^2}{\partial \theta^2} S(\theta | X) \Big|_{\theta=\theta_0} \right]} \\ &= -I_1(\theta_0) \end{aligned}$$

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

So  $\hat{\theta}_n$  is asymptotically efficient! Consistency follows from Slutsky's theorem (again). "□"

# 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  $\sigma^2$  is known. Find the asymptotic distribution of the MLE of  $\mu$ . ( $\hat{\mu}_n = \bar{X}_n$ ).

Check the conditions of Theorem 5.9:

①  $\mathcal{F} = (-\infty, \infty) \subseteq \mathbb{R}$  is open in  $\mathbb{R}$  ✓

②  $l'''(\mu|x) = 0$ , which is continuous in  $\mu$  ✓

③  $f_\mu(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} > 0 \quad \forall x \in \mathbb{R}$

so the support of  $f_\mu$  is  $\mathbb{R}$ , which doesn't depend on  $\mu$  ✓

④  $I(\mu) = \frac{1}{\sigma^2} < \infty \quad \forall \mu \in \mathbb{R}$  ✓

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

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

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

$$l''(\mu|x) = -\frac{1}{\sigma^2} \Rightarrow I_1(\mu) = \mathbb{E}_\mu[-(-\frac{1}{\sigma^2})] = \frac{1}{\sigma^2}$$

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

# Asymptotically Efficient MLEs: Examples

- **Example 5.21:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim}$  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  $1/p$ .

# 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  $\theta$  is not asymptotically normal.

$$\hat{\Theta}_n = X_{(n)}.$$

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

But ...  $\mathbb{P}_\theta(Y_n \leq y)$

$$= \mathbb{P}_\theta(\theta - X_{(n)} \leq y/\sqrt{n})$$

$$= \mathbb{P}_\theta(X_{(n)} \geq \theta - y/\sqrt{n})$$

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

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

$$\xrightarrow{n \rightarrow \infty} \begin{cases} 1, & y/\theta \geq 0 \\ 0, & y/\theta < 0 \end{cases} = \mathbb{1}_{y \geq 0}$$

**EXERCISE:** sometimes different scalings of  $T_n - \theta$  give us interesting results (e.g., if  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , then  $\mathbb{1} \cdot (\bar{X}_n - \mu) \xrightarrow{d} \mathbb{1}_{x \leq 0}$  but  $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2)$ .)  
In the  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Unif}(0, \theta)$  case, what - if anything - does  $n\left(\frac{n+1}{n} X_{(n)} - \theta\right)$  converge in distribution to?

(i.e., the asymptotic distribution of  $Y_n$  is degenerate at 0, so it's 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  $\lambda(\mathbf{X}_n) = \frac{L(\theta_0|\mathbf{X}_n)}{L(\hat{\theta}_n|\mathbf{X}_n)}$ , where  $\hat{\theta}_n = \hat{\theta}(\mathbf{X}_n)$  is the unrestricted MLE of  $\theta$  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  $H_0 : \theta \in \Theta_0$  versus  $H_A : \theta \in \Theta_0^c$  using  $\lambda(\mathbf{X}_n)$ , then

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

under  $H_0$ .

Proof: EXERCISE!

# Poll Time!

On Quercus: Module 5 - Poll 5

# Approximate LRTs: Examples

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

$$\text{Example 3-23} \Rightarrow \lambda(\vec{X}_n) = \left(\frac{p_0}{\bar{X}_n}\right)^{\sum X_i} \left(\frac{1-p_0}{1-\bar{X}_n}\right)^{n-\sum X_i}$$

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

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

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

is the rejection region of an approximate size- $\alpha$  test of  $H_0: \theta = \theta_0$  vs  $H_A: \theta \neq \theta_0$ .

# 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- $\alpha$  LRT of  $H_0 : \mu = \mu_0$  versus  $H_A : \mu \neq \mu_0$ .  $\sigma^2$  known

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

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

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

In fact, it's an exact size- $\alpha$  test! Why? Compare to a Z-test...

# 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**

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

where  $\hat{\theta}_n = \hat{\theta}_{\text{MLE}}(\mathbf{X}_n)$  is the usual 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

$$W_n(\mathbf{X}_n) \xrightarrow{d} \chi_{(1)}^2$$

under  $H_0$ .

Proof: EXERCISE !

# Wald Tests: Examples

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

$$W_n(\vec{X}_n) = (\hat{p}_n - p_0)^2 \cdot I_n(\hat{p}_n), \text{ where } \hat{p}_n = \bar{X}_n. \quad \text{See Slide 50}$$

$$\text{What's the Fisher information? } I_n(p) = \frac{n}{p(1-p)} \Rightarrow I_n(\hat{p}_n) = \frac{n}{\bar{X}_n(1-\bar{X}_n)}$$

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

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

OR:  $R' = \{ \vec{x} \in \mathcal{X}^n : \left| \frac{\bar{x} - p_0}{\sqrt{\bar{x}(1-\bar{x})/n}} \right| > z_{\alpha/2} \}$  is the rejection region of an approximate size- $\alpha$  test of  $H_0: p = p_0$  vs  $H_A: p \neq p_0$ .

EXERCISE: does  $R = R'$ ?

# 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- $\alpha$  Wald test of  $H_0 : \mu = \mu_0$  versus  $H_A : \mu \neq \mu_0$ .  $\hat{\mu}_n = \bar{X}_n$

From Example 5.20,  $I_n(\mu) = n/\sigma^2$ , so  $W_n(\bar{X}_n) = \frac{(\bar{X}_n - \mu_0)^2}{\sigma^2/n} = \left( \frac{\bar{X}_n - \mu_0}{\sqrt{\sigma^2/n}} \right)^2$ .

By Theorem 5.12,  $R = \{ \bar{x} \in \mathcal{X}^n : \frac{(\bar{x} - \mu_0)^2}{\sigma^2/n} > \chi_{(1), \alpha}^2 \}$  is the rejection region of an approximate (exact, in this case!) size- $\alpha$  test of  $H_0: \mu = \mu_0$  vs  $H_A: \mu \neq \mu_0$ .

OR:  $R' = \{ \bar{x} \in \mathcal{X}^n : \left| \frac{\bar{x} - \mu_0}{\sqrt{\sigma^2/n}} \right| > z_{\alpha/2} \}$  is the rejection region of an approximate (exact) size- $\alpha$  test of  $H_0: \mu = \mu_0$  vs  $H_A: \mu \neq \mu_0$ .

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  $H_0 : \theta = \theta_0$  versus  $H_A : \theta \neq \theta_0$ , a **score test** (also called a **Rao test** or a **Lagrange multiplier test**) is a test based on the **score statistic**

$$R_n(\mathbf{X}_n) = \frac{[S_n(\theta_0 | \mathbf{X}_n)]^2}{I_n(\theta_0)}.$$

- **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  $H_0 : \theta = \theta_0$  versus  $H_A : \theta \neq \theta_0$  using  $R_n(\mathbf{X}_n)$ , then

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

under  $H_0$ .

Equivalently, 
$$\frac{S_n(\theta_0 | \mathbf{X}_n)}{\sqrt{I_n(\theta_0)}} \xrightarrow{d} N(0, 1).$$

# Score Tests: Examples

- **Example 5.27:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim}$  Bernoulli ( $p$ ), where  $p \in (0, 1)$ . Construct an approximate size- $\alpha$  score test of  $H_0 : p = p_0$  versus  $H_A : p \neq p_0$ .

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

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

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

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

$$S'(p | \vec{x}) = n \left( \frac{-\bar{x}}{p^2} - \frac{-(1-\bar{x})}{(1-p)^2} \right)$$

$$I_n(p) = -\mathbb{E}_p \left[ n \left( \frac{\bar{X}_n}{p^2} - \frac{1 - \bar{X}_n}{(1-p)^2} \right) \right] = \frac{n}{p(1-p)}$$

By Theorem 5.13,  $R = \{ \vec{x} \in \mathcal{X}^n : \frac{(\bar{x} - p_0)^2}{p_0(1-p_0)/n} > \chi_{(2), \alpha}^2 \}$  is the rejection region of an approximate size- $\alpha$  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- $\alpha$  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
  - One might tell you to reject  $H_0$  while another might not!
- FYI: 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  $n$  (proved in 1982)

# 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  $\theta$
- 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  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim}$  Bernoulli ( $p$ ), where  $p \in (0, 1)$ . Construct an approximate  $(1 - \alpha)$ -confidence interval for  $p$  based on the Wald statistic.

$$\text{Example 5.25, } 1 - \alpha \approx P_p \left( \frac{|\bar{X}_n - p|}{\sqrt{\bar{X}_n(1 - \bar{X}_n)/n}} < z_{\alpha/2} \right) \text{ when } n \text{ is large}$$

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

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

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

- 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**

In practice you almost always see  $\alpha = 0.05$  (thanks, Fisher...), whence  $z_{\alpha/2} \approx 1.96$

# Approximate Confidence Intervals: Examples

- **Example 5.30:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim}$  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, since  $p \mapsto \log\left(\frac{p}{1-p}\right)$  is a monotone increasing bijection

$$1 - \alpha \approx \mathbb{P}_p \left( \bar{X}_n - z_{\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}} < p < \bar{X}_n + z_{\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}} \right)$$
$$= \mathbb{P}_p \left( \log \left( \frac{\bar{X}_n - z_{\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}}{1 - \bar{X}_n + z_{\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}} \right) < \log\left(\frac{p}{1-p}\right) < \log \left( \frac{\bar{X}_n + z_{\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}}{1 - \bar{X}_n - z_{\alpha/2} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}} \right) \right)$$

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

# Approximate Confidence Intervals: Examples

- **Example 5.31:** Let  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim}$  Poisson ( $\lambda$ ), where  $\lambda > 0$ . Construct an approximate  $(1 - \alpha)$ -confidence interval for  $\lambda$  based on the Wald statistic.

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

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

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

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

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

$$\ell(\lambda | \vec{x}) = -n\lambda + \sum x_i \cdot \log(\lambda) + c, \text{ where...}$$

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

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

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

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



# When the Fisher Information Causes Problems...

- When  $f_\theta$  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!*
- 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  $J_n(\bar{X}_n)$  is a consistent estimator of  $I_n(\theta_0)$*
- Moreover, in a famous 1978 paper, Efron and Hinkley showed empirically that  ~~$J_n(\bar{X}_n)$~~  is superior to  $I_n(\hat{\theta})$   
*Optional reading, if you're curious...*

$$J_n(\bar{X}_n)$$