

# Lecture 12

10/19/2020

In lec 10  $\frac{\hat{\theta} - \theta}{SE[\hat{\theta}]} \xrightarrow{d} N(0,1) \Rightarrow \frac{\hat{\theta} - \theta}{SE[\hat{\theta}]} \xrightarrow{d} N(0,1)$

We can use this now in our situation:

$$\frac{\hat{\theta}_1 - \hat{\theta}_2}{SE[\hat{\theta}_1 - \hat{\theta}_2]} \xrightarrow{d} N(0,1)$$

$$\Rightarrow \frac{\hat{\theta}_1 - \hat{\theta}_2}{SE[\hat{\theta}_1 - \hat{\theta}_2]} \xrightarrow{d} N(0,1)$$

$$SE[\hat{\theta}_1 - \hat{\theta}_2] = \sqrt{\theta_{\text{shared}}(1 - \theta_{\text{shared}}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

$$\hat{SE}[\hat{\theta}_1 - \hat{\theta}_2] = \sqrt{\hat{\theta}_{\text{shared}}(1 - \hat{\theta}_{\text{shared}}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)} \text{ if } \hat{\theta}_{\text{shared}} \text{ is consistent.}$$

$$\hat{\theta}_{\text{shared}} = \text{avg. over both samples} = \frac{\sum X_{1i} + \sum X_{2i}}{n_1 + n_2}$$

$$\Rightarrow \frac{\hat{\theta}_1 - \hat{\theta}_2}{\sqrt{\frac{\sum X_{1i} + \sum X_{2i}}{n_1 + n_2} \left( 1 - \frac{\sum X_{1i} + \sum X_{2i}}{n_1 + n_2} \right) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \sim N(0,1)$$

e.g.  $H_a: \theta_1 - \theta_2 \neq 0$ ,  $H_0: \theta_1 - \theta_2 = 0$ ,  $\alpha = 5\%$

control:  $n_1 = 81$ ,  $\sum X_{1i} = 27 \Rightarrow \hat{\theta}_1 = 27/81 = 0.333$

experimental:  $n_2 = 79$ ,  $\sum X_{2i} = 12 \Rightarrow \hat{\theta}_2 = 12/79 = 0.152$

$$\hat{\theta}_{\text{shared}} = \frac{27 + 12}{81 + 79} = 0.244$$

$$(\hat{\theta}_1 - \hat{\theta}_2)_{std} = \frac{0.333 - 0.152}{\sqrt{0.244(1-0.244)(\frac{1}{81} + \frac{1}{79})}} = 2.66 \notin [-1.96, +1.96] \Rightarrow \text{Reject } H_0$$

Another (obvious) Wald Test: If  $x_1, \dots, x_n$  iid DGP with mean  $\theta$  and variance  $\sigma^2$  and the estimator  $\hat{\theta}$  is  $\bar{x}$ , then the CLT implies that:

$$\frac{\hat{\theta} - \theta}{\sigma/\sqrt{n}} \xrightarrow{d} N(0,1) \text{ if } \sigma \text{ is known.}$$

If  $\sigma$  is unknown ... I can replace  $\sigma$  with any consistent estimator e.g.  $s, \hat{\sigma}$  &  $\frac{1}{n} \sum (x_i - \theta)^2$ .

$$\Rightarrow \frac{\hat{\theta} - \theta}{s/\sqrt{n}} \xrightarrow{d} N(0,1)$$

Are you allowed to just use the T-test here? Many people just use the T-test here. Technically it's wrong because you need the DGP to be normal iid. But if you use the T-test ... it's "not so bad". (in midterm 2 Q11)

$$H_0: \theta \geq 2, n=30, \bar{x} = 2.57, s = 1.00$$

$$\hat{\theta}_{std} = \frac{2.57 - 2}{\frac{1.00}{\sqrt{30}}} = 3.12 \notin (-\infty, 1.645] \Rightarrow \text{Rej } H_0$$



Another Wald test for two independent samples with unknown variances and you wish to test a difference in means.

$$\frac{\hat{\theta}_1 - \hat{\theta}_2}{\sqrt{\frac{\hat{\sigma}_1^2}{n_1} + \frac{\hat{\sigma}_2^2}{n_2}}} \xrightarrow[\text{from last class}]{d} N(0,1)$$

$$\Rightarrow \frac{\hat{\theta}_1 - \hat{\theta}_2}{\sqrt{s_1^2/n_1 + s_2^2/n_2}} \xrightarrow{d} N(0,1)$$

If you use the Satterthwaite t-test, it "wouldn't be so bad" because unless your population distributions were so very skewed, it should be fine.

Let's use the asymptotic normality of the MLE from (last class) to do a Wald Test. HW1, m has  $\text{dGP: } x_1, \dots, x_n \text{ i.i.d. Gumbel}(\theta, 1)$ . The Gumbel is a rv model for "extreme events" think maximum rainfall per month:

$$\ell'(\theta; x_1, \dots, x_n) = n - e^{-\theta} \sum e^{-x_i} \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \hat{\theta}^{\text{MLE}} = \ln\left(n / \sum e^{-x_i}\right)$$

$$\ell'(\theta; x) = 1 - e^{-\theta} e^{-x} \Rightarrow \ell''(\theta; x) = -e^{-\theta} e^{-x} \quad ?$$

$$\begin{aligned} \ell(\theta) &= E[\ell''(\theta; x)] = E[-e^{-\theta} e^{-x}] = -e^{-\theta} E[e^{-x}] \\ &= -e^{-2\theta}. \end{aligned}$$

$$\frac{\hat{\theta}^{MLE} - \theta}{\sqrt{\frac{I(\theta)^{-1}}{n}}} = \frac{\hat{\theta}^{MLE} - \theta}{\frac{e^{\theta}}{\sqrt{n}}} = \frac{\ln\left(\frac{n}{\sum e^{-x}}\right) - \theta}{e^{\theta}/\sqrt{n}} \sim d, N(0,1)$$

$$X_1 = 2.15, X_2 = 1.91, \overset{y_3}{3.66}, \overset{y_4}{4.85}, \overset{y_5}{3.03}, \overset{y_6}{1.03}, \overset{y_7}{3.58},$$

$$n = 7$$

$$\hat{\theta}^{MLE} = 2.26 \quad \text{Test } H_0: \theta > 2, \alpha = 5\%$$

$$\hat{\theta}_{STD}^{MLE} = \frac{2.26 - 2}{\frac{e^{\theta}}{\sqrt{7}}} = \frac{0.26}{2.79} = 0.09 \in (-\infty, 1.645]$$

$\Rightarrow$  Retain  $H_0$

There are three goals of statistical inference

(1) Point Estimation

Goal here is to provide a best guess,  $\hat{\theta}$  at the value of  $\theta$ . You don't know if your specific guess is good, is close, is bad, is far. How do we ask the question "is it good / bad"? We imagined  $\hat{\theta}$  coming from a distribution  $\hat{\theta}$ , the "sampling distribution". There are properties about the "sampling distribution". There are properties about the sampling distribution e.g. some good properties are unbiasedness, consistency, low MSE, low risk (for general loss functions).

(2) Testing

Goal here is to test a theory about a specific  $\theta$ . We used hypothesis testing. What makes a 'good test'? One property is "power". There are other properties we didn't discuss.



### (3) Confidence Sets

The goal here is to create a set of values for  $\theta$  that you're "confident in". The approach we use here is the "confidence interval".

Def

an "interval estimate" are two statistics

$w_L(x_1, \dots, x_n)$  &  $w_U(x_1, \dots, x_n)$  s.t.  $w_L < w_U$  for all data sets

combined in an interval:  $[w_L(x_1, \dots, x_n), w_U(x_1, \dots, x_n)]$

e.g.  $[1.789, 2.463]$

and of course, the "interval estimator" is

$[w_L(x_1, \dots, x_n), w_U(x_1, \dots, x_n)]$   
which is a "random interval".

Def

An interval estimator has "coverage probability"

$P(\theta \in [w_L(x_1, \dots, x_n), w_U(x_1, \dots, x_n)] | \theta)$  An illustration

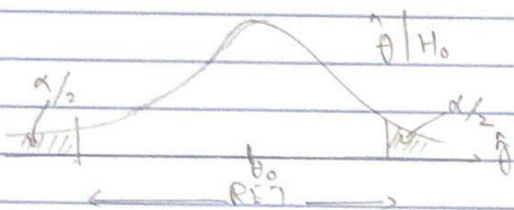
	$\theta$	
Dataset 1 :	$[w_L, w_U]$	The coverage prob. is computed over every dataset. For these 4 datasets, the coverage prob. would be $3/4 = 75\%$ .
Dataset 2 :	$[w_L, w_U]$	
Dataset 3 :	$[w_L, w_U]$	
Dataset 4 :	$[w_L, w_U]$	

We define the "confidence interval" with coverage probability  $1 - \alpha$  for parameter  $\theta$  as this interval estimate and interval estimator (depending on context). Denoted  $C_{\theta, 1-\alpha}$ .

Given  $\alpha$ , how do we find the confidence interval? Let's begin with the DGP iid normal mean  $\theta$ , variance  $\sigma^2$  & variance known & the estimator  $= \bar{X}$ .

Consider testing:

$H_a: \theta \neq \theta_0$  vs.  $H_0: \theta = \theta_0$  @ size  $\alpha$ .



$$z_{1-\alpha/2} = F_z^{-1}(1-\alpha/2)$$

$$P(\hat{\theta} \notin \text{REJ} | H_0) = 1 - \alpha.$$

$$P(\hat{\theta} \in [\theta_0 - z_{1-\alpha/2} \cdot \sigma/\sqrt{n}, \theta_0 + z_{1-\alpha/2} \cdot \sigma/\sqrt{n}] | \theta = \theta_0) \\ = P(\hat{\theta} - \theta_0 \in [-z_{1-\alpha/2} \cdot \sigma/\sqrt{n}, +z_{1-\alpha/2} \cdot \sigma/\sqrt{n}] | \theta = \theta_0)$$

$$= P(\theta_0 - \hat{\theta} \in [-z_{1-\alpha/2} \cdot \sigma/\sqrt{n}, +z_{1-\alpha/2} \cdot \sigma/\sqrt{n}] | \theta = \theta_0)$$

$$= P(\theta_0 \in [\underbrace{\hat{\theta} - z_{1-\alpha/2} \cdot \sigma/\sqrt{n}}_{W_L}, \underbrace{\hat{\theta} + z_{1-\alpha/2} \cdot \sigma/\sqrt{n}}_{W_U}] | \theta = \theta_0)$$

$$= P(\theta_0 \in [W_L(X_1, \dots, X_n), W_U(X_1, \dots, X_n)] | \theta = \theta_0) \text{ Since valid } \forall \theta_0, \dots$$

$$\Rightarrow CI_{\theta, 1-\alpha} \doteq [\hat{\theta} - z_{1-\alpha/2} \cdot \sigma/\sqrt{n}, \hat{\theta} + z_{1-\alpha/2} \cdot \sigma/\sqrt{n}]$$

We constructed our first confidence interval by "inverting the test".