

# **MATH 503: Mathematical Statistics**

## **Lecture 8: Confidence Intervals, and Bayesian Statistics**

**Reading: CB Sections 7.2.3, 9.1-9.2.1**  
(or HMC Sections 5.4, 7.1, 11.1-11.3)

**Kimberly F. Sellers**

**Department of Mathematics and Statistics**

# ***Today's Topics***

- Confidence intervals
  - “What is it?” and further introduction
  - Interpretation
  - Choosing a large enough sample size
- Bayesian Statistics
  - Subjective Probability
  - Prior and Posterior Distributions
  - Bayesian Point Estimation
  - Additional Terminology

# ***What is a confidence interval?***

An **interval estimate** of a real-valued parameter  $\theta$  is any pair of functions  $L(x_1, \dots, x_n)$  and  $U(x_1, \dots, x_n)$  of a sample that satisfy  $L(\mathbf{x}) \leq U(\mathbf{x}) \quad \forall \mathbf{x} \in \chi$ . If  $X = \mathbf{x}$  is observed, the inference  $L(\mathbf{x}) \leq \theta \leq U(\mathbf{x})$  is made. The random interval  $[L(\mathbf{x}), U(\mathbf{x})]$  is called an **interval estimator**.

# ***What is a confidence interval? (cont.)***

- An alternative to reporting a single estimate (i.e. the point estimate) for a parameter.
- Because of sampling variability, the point estimate is almost never exactly the correct value for the parameter
- Point estimates don't tell us how close we are to the actual parameter
- Instead, confidence intervals are calculated intervals of plausible values for the parameter

# ***Coverage Probabilities & Confidence Coefficients***

- For an interval estimator  $[L(X), U(X)]$  of a parameter  $\theta$ , the **coverage probability** of  $[L(X), U(X)]$  is the probability that the random interval  $[L(X), U(X)]$  covers the true parameter,  $\theta$ , i.e.  $P_{\theta}(\theta \in [L(X), U(X)])$ .
- For an interval estimator  $[L(X), U(X)]$  of a parameter  $\theta$ , the **confidence coefficient** of  $[L(X), U(X)]$  is the infimum of the coverage probabilities, i.e.  $\inf_{\theta} P_{\theta}(\theta \in [L(X), U(X)])$ .

# ***What do you need to compute a confidence interval?***

- a sample of observations (for purposes of discussion, from a normal distribution)
- the sample estimate,  $\hat{\theta}$ , and its standard error
- confidence level, traditionally 95% (other popular choices are 90% or 99%; generally, its  $100(1 - \alpha)\%$ ); the higher the confidence level, the more strongly we believe the value of  $\theta$  falls within the interval

# ***Creating a confidence interval for $\theta$***

- Let  $\theta_0$  denote the true (unknown) value of  $\theta$
- Suppose  $\bar{X}$  estimator of  $\theta_0$  such that

$$\sqrt{n}(\bar{X} - \theta_0) \xrightarrow{d} N(0, \sigma^2),$$

- Assume  $\sigma^2$  known

# ***Confidence interval for $\theta$ (cont.)***

- A  $100(1 - \alpha)\%$  confidence interval (interval estimator) for  $\theta$  of a normal population when the value of  $\sigma$  is known is

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

- The interval estimate is determined by replacing  $\bar{X}$  above with  $\bar{x}$



# ***Examples***

- Confidence interval for  $\mu$  when  $\sigma$  known:

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

- Large-sample interval for  $\mu$  when  $\sigma$  unknown:

$$\bar{x} \pm z_{\alpha/2} \frac{s}{\sqrt{n}}$$

# ***Examples (cont.)***

- Small-sample interval for  $\mu$  when  $\sigma$  unknown:

$$\bar{x} \pm t_{\alpha/2, n-1} \frac{s}{\sqrt{n}}$$

- Large-sample interval for  $p$ :

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}}$$

# ***Example***

A random sample of 100 observations from a normally distributed population possesses a mean equal to 83.2. Assume  $\sigma = 6.4$ . Find a 95% confidence interval for  $\mu$ .

# ***Questions***

- How does the confidence interval change when
  - sample size increases?
  - $\sigma$  increases?
  - $\alpha$  increases (say, to  $\alpha' > \alpha$ )?

## ***Example (cont.)***

Would a 90% confidence interval calculated from the same sample be narrower or wider than the above interval? Why?

# ***Interpreting a confidence interval***

- Question: how do I interpret, say, a 95% confidence interval?
  - Bayesian perspective: a belief with 95% confidence that the parameter,  $\theta$ , is in the stated interval. We can generalize this idea for any  $100(1 - \alpha)\%$  confidence interval, i.e. we are  $100(1 - \alpha)\%$  sure that parameter,  $\theta$ , is contained in the interval in question.

# ***Interpreting a confidence interval (cont.)***

– Frequentist perspective:

95% of all samples would give an interval that includes the parameter,  $\theta$ , and 5% of all samples would give an erroneous interval. More generally,  $100(1 - \alpha)\%$  of all samples would give an interval that includes the parameter,  $\theta$ , and  $100\alpha\%$  of all samples would give an erroneous interval.

# ***Steps to constructing confidence intervals***

1. Answer the following questions to determine which kind of confidence interval to create:
  - What parameter do I want to estimate?
  - What is my sample size? Large or small?
  - Is the population variance,  $\sigma$ , known or unknown?
2. Determine the appropriate value of  $\alpha$ .
3. Find the corresponding critical value(s) (e.g.  $z_{\alpha/2}$  or  $t_{\alpha/2}$ ) depending on the responses to Step 1.
4. Use the appropriate formula to compute the confidence interval.



# ***Example***

The Gallup Organization conducts annual national surveys on home gardening. Results are published by the National Association for Gardening in National Gardening Survey. A random sample is taken of 25 households with vegetable gardens. The mean size of their gardens is 643 sq. ft. with standard deviation equal to 247 sq. ft.

- a) Determine a 90% confidence interval for the mean size,  $\mu$ , of all household vegetable gardens in the U.S.
  
  
  
  
  
  
  
  
  
  
- b) Interpret your answer from part (a).

# ***Example***

The accounting firm Price Waterhouse periodically monitors the U.S. Postal Service's performance. One parameter of interest is the percentage of mail delivered on time. In a sample of 332,000 mailed items, Price Waterhouse determined that 282,200 items were delivered on time (Tampa Tribune, Mar. 26, 1995). Use this information to estimate with 99% confidence the true percentage of items delivered on time by the U.S. Postal Service. Interpret the results.

# ***Confidence Intervals for Differences of Means***

- Assume samples  $X_1, \dots, X_{n_1}$  and  $Y_1, \dots, Y_{n_2}$  gathered independently s.t. sample sizes have relative order
- Large-sample interval for  $\Delta = \mu_1 - \mu_2$  when  $\sigma_1, \sigma_2$  known:

$$\bar{x} - \bar{y} \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

For  $\sigma_1, \sigma_2$  unknown, replace with  $s_1, s_2$ , respectively

**Note:** this CI is approximate based on asymptotic distribution theory!

# ***Confidence Intervals for Differences of Means***

- $X_1, \dots, X_{n_1} \sim N(\mu_1, \sigma^2)$  and  $Y_1, \dots, Y_{n_2} \sim N(\mu_2, \sigma^2)$ , independently drawn
- Exact confidence interval with common  $\sigma^2$ :

$$\bar{x} - \bar{y} \pm t_{\alpha/2, n-2} \cdot s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

where

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

# ***Confidence Interval for Difference of Proportions***

- $X_1, \dots, X_{n_1} \sim \text{Bern}(p_1); Y_1, \dots, Y_{n_2} \sim \text{Bern}(p_2)$ ,  
independently drawn
- Large-sample interval for  $p_1 - p_2$ :

$$\hat{p}_1 - \hat{p}_2 \pm z_{\alpha/2} \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}}$$

# ***Example***

In order to compare the means of two populations, independent random samples of 400 observations are selected from each population, with the following results:

Sample 1	Sample 2
$\bar{x} = 5275$	$\bar{y} = 5240$
$s_1 = 150$	$s_2 = 200$

Use a 95% CI to estimate the difference between the population means,  $\mu_1 - \mu_2$ .

# ***Sample Size Determination***

- **Idea:** we know the level of confidence we want to have regarding an interval surrounding our unknown parameter of interest,  $\theta$ , and we also have a specified error bound or width for the confidence interval in question.
- **Goal:** to know how large the sample size,  $n$ , needs to be in order to fulfill these specifications.
- **How do we do that?**

# ***Example***

The EPA wants to test a randomly selected sample of  $n$  water specimens and estimate the mean daily rate of pollution produced by a mining operation. If the EPA wants a 95% confidence interval estimate with a bound on the error of 1 milligram per liter (mg/L), how many water specimens are required in the sample? Assume prior knowledge indicates that pollution readings in water samples taken during a day are approximately normally distributed with a standard deviation equal to 5 mg/L.



# ***Example***

According to estimates made by the General Accounting Office, the Internal Revenue Service (IRS) answered 18.3 million telephone inquiries during a recent tax season, and 17% of the IRS offices provided answers that were wrong. These estimates were based on data collected from sample calls to numerous IRS offices. How many IRS offices should be randomly selected and contacted in order to estimate the proportion of IRS offices that fail to correctly answer questions about gift taxes with a 90% confidence interval of width .06?

# ***What is subjective probability?***

- A bet
- All probability axioms and associated results hold under this philosophical construct!

# ***Useful Tools***

- The **conditional probability** of  $A$  given that the event  $B$  has occurred is

$$P(A|B) = \frac{P(A \cap B)}{P(B)},$$

where  $P(B) > 0$ .

- Note: by cross-multiplying both sides, we get “the **multiplication rule**”:

$$P(A \cap B) = P(A | B) \cdot P(B)$$

# ***Useful Tools (cont.)***

## ***The Law of Total Probability***

$$P(B) = P(B | A)P(A) + P(B | A')P(A')$$

## ***Bayes' Theorem/ Bayes' Law***

$$P(A | B) = \frac{P(B | A)P(A)}{P(B | A)P(A) + P(B | A')P(A')}$$

# ***Prior & Posterior Distributions***

- $X|\theta \sim f(x|\theta), \theta \in \Omega$
- $\Theta \sim \pi(\theta)$ ; called prior distribution
- Suppose  $X_1, X_2, \dots, X_n$  random sample drawn from population indexed by  $\theta$  with pdf  $f(x|\theta)$ . Then, sampling distribution is

$$f(\mathbf{x}|\theta) = \prod_{i=1}^n f(x_i|\theta)$$

and joint pdf of  $\mathbf{X}' = (X_1, X_2, \dots, X_n)$  and  $\Theta$  is

$$f(\mathbf{x}, \theta) = f(\mathbf{x}|\theta)\pi(\theta)$$

# ***Prior & Posterior Dists. (cont.)***

- Conditional pdf of  $\Theta$ , given sample  $X$ :

$$\pi(\theta | \mathbf{x}) = \frac{f(\mathbf{x}, \theta)}{m(\mathbf{x})} = \frac{f(\mathbf{x}|\theta)\pi(\theta)}{m(\mathbf{x})},$$

where  $m(\mathbf{x}) = \int_{-\infty}^{\infty} f(\mathbf{x}|\theta)\pi(\theta)d\theta$

- This is posterior pdf of  $\theta$  given  $x$
- **What does this mean???**

# ***Example***

Let  $X_i \mid \theta \sim \text{Poisson}(\theta)$  iid, and  $\Theta \sim \text{Gamma}(\alpha, \beta)$  where  $\alpha, \beta$  known. Find the posterior distribution of  $\theta \mid \mathbf{x}$ .

# Notes

- $\pi(\theta \mid \mathbf{x}) \propto f(\mathbf{x} \mid \theta) \cdot \pi(\theta)$
- If a sufficient statistic  $Y = u(\mathbf{X})$  exists, then  $\pi(\theta \mid \mathbf{x}) \propto f(u(\mathbf{x}) \mid \theta) \cdot \pi(\theta)$ , thus

$$\pi(\theta \mid y) \propto f(y \mid \theta) \cdot \pi(\theta)$$



# ***Example***

Let  $X_i \mid \theta \sim \text{Bernoulli}(\theta)$  iid  $i = 1, \dots, n$ , and  $\Theta \sim \text{Beta}(\alpha, \beta)$  where  $\alpha, \beta$  known. Find the posterior distribution of  $\theta \mid \mathbf{x}$ .

# ***Bayesian Point Estimation: Definitions and Framework***

- $X_1, X_2, \dots, X_n$  iid with pdf  $f(x; \theta)$ ,  $\theta \in \Omega$ .
- $Y = u(X_1, X_2, \dots, X_n)$ : statistic on which point estimate of  $\theta$  is determined
- $\delta(y)$ : decision rule, determining value of point estimate at  $\theta$ 
  - Numerical point estimate is a decision
- $\mathcal{L}(\theta, \delta(y))$ : loss function, reflecting severity of difference between  $\theta$  and  $\delta(y)$

# ***Loss Functions***

- Examples include:
  - Squared-error loss:  $\mathcal{L}(\theta, \delta(y)) = [\theta - \delta(y)]^2$
  - Absolute-value error loss:  $\mathcal{L}(\theta, \delta(y)) = |\theta - \delta(y)|$
  - “Goal-post” loss:
$$\mathcal{L}(\theta, \delta(y)) = \begin{cases} 0, & \text{if } |\theta - \delta(y)| \leq a \\ b, & \text{otherwise} \end{cases}$$
- Don't have to be symmetric
- Goal of point estimation: find  $\delta(y)$  s.t.  $\mathcal{L}(\theta, \delta(y))$  is minimum

# ***Risk Function***

- **Problem:**  $\theta$  unknown, thus minimizing  $\mathcal{L}(\theta, \delta(y))$  impossible

- Introduce risk function:

$$R(\theta, \delta) = E[\mathcal{L}(\theta, \delta(y))] = \int_{-\infty}^{\infty} \mathcal{L}(\theta, \delta(y)) f_Y(y; \theta) dy$$

- **Goal:** select  $\delta(y)$  s.t.  $R(\theta, \delta)$  minimum for all  $\theta \in \Omega$ 
  - **Problem:**  $\delta(y)$  s.t.  $R(\theta, \delta)$  minimum does not minimize  $R(\theta', \delta)$
  - **Resolution:** restrict decision function to certain class, or order risk functions

# ***Minimax Principle***

- If the decision function given by  $\delta_0(y)$  is such that, for all  $\theta \in \Omega$ ,

$$\max_{\theta} R[\theta, \delta_0(y)] \leq \max_{\theta} R[\theta, \delta(y)]$$

for every other decision function  $\delta(y)$ , then  $\delta_0(y)$  is called a minimax decision function.

# ***Example***

Let  $X_1, X_2, \dots, X_{25} \sim N(\theta, 1)$  iid,  $-\infty < \theta < \infty$ ,  $Y = \bar{X}$ , and  $\mathcal{L}(\theta, \delta(y)) = [\theta - \delta(y)]^2$ . Consider two decision functions  $\delta_1(y) = y$ , and  $\delta_2(y) = 0$  for  $-\infty < y < \infty$ .

1. Find the associated risk function associated with each decision function.
2. Which is better if you:
  - (a) restrict the choice of  $\delta(y)$  to satisfy  $E[\delta(y)] = \theta$ ?
  - (b) want  $\delta(y)$  to satisfy the minimax principle?

# ***Example***

Let  $X_1, X_2, \dots, X_n$  denote a random sample from a distribution that is  $\text{Bin}(1, \theta)$ ,  $0 \leq \theta \leq 1$ . Let  $Y = \sum_{i=1}^n X_i$ , and let  $\mathcal{L}(\theta, \delta(y)) = [\theta - \delta(y)]^2$ . Consider decision functions of the form  $\delta(y) = by$ , where  $b$  does not depend upon  $y$ . Prove that  $R(\theta, \delta) = b^2 n \theta (1 - \theta) + (bn - 1)^2 \theta^2$ .

# ***Notes***

- $Y = u(\mathbf{X})$  statistic  $\Rightarrow \delta(Y)$  statistic
- If  $\delta(Y)$  satisfies  $E[\delta(y)] = \theta$ ,  $\mathcal{L}(\theta, \delta(y)) = [\theta - \delta(y)]^2$ , and minimizes  $R(\theta, \delta)$ , then  $\delta(y)$  is MVUE



# ***How does this change with knowledge from random sample?***

- Goal: determine point estimator for  $\theta$ , i.e. select decision function  $\delta$  s.t.  $\delta(\mathbf{x})$  is predicted value of  $\theta$  when computed  $\mathbf{x}$  and posterior distribution  $\pi(\theta \mid \mathbf{x})$  known
- In other words, choose  $\delta(\mathbf{x})$  s.t. minimize
$$E\{\mathcal{L}(\Theta, \delta(\mathbf{x})) \mid \mathbf{X} = \mathbf{x}\} = \int_{-\infty}^{\infty} \mathcal{L}(\theta, \delta(\mathbf{x})) \pi(\theta \mid \mathbf{x}) d\theta \quad (1)$$
- $\delta(\mathbf{x})$  called a Bayes estimator of  $\theta$
- Can generalize to find Bayes estimator of  $\ell(\theta)$  by replacing  $\theta$  with  $\ell(\theta)$  in  $\mathcal{L}(\theta, \delta(\mathbf{x}))$  of Equation (1)

# ***How does this change ...? (con't.)***

- $E\{\mathcal{L}(\Theta, \delta(x)) \mid X = x\}$  is r.v. that is a function of  $X$
- Expected risk =

$$\begin{aligned} & \int_{-\infty}^{\infty} \left\{ \underbrace{\int_{-\infty}^{\infty} \mathcal{L}(\Theta, \delta(x)) \pi(\theta \mid x) d\theta}_{E\{\mathcal{L}(\Theta, \delta(x)) \mid X=x\}} \right\} m(x) dx \\ &= \int_{-\infty}^{\infty} \left\{ \underbrace{\int_{-\infty}^{\infty} \mathcal{L}(\Theta, \delta(x)) f(x \mid \theta) dx}_{R(\theta, \delta), \text{ for every } \theta \in \Theta} \right\} \pi(\theta) d\theta \end{aligned}$$

- $\delta(x)$  minimizing  $E\{\mathcal{L}(\Theta, \delta(x)) \mid X = x\}$  for every satisfactory  $x \Rightarrow$  it minimizes expected risk

# ***Estimators Associated with Loss Functions***

- For squared-error loss function,

$$\mathcal{L}(\theta, \delta(y)) = [\theta - \delta(y)]^2,$$

the conditional distribution mean of  $\Theta$ , given  $X = x$ , is Bayes estimator

- For absolute-value loss function,

$$\mathcal{L}(\theta, \delta(y)) = |\theta - \delta(y)|,$$

the conditional distribution median of  $\Theta$ , given  $X = x$ , is Bayes estimator

# ***Example***

Let  $X_i \mid \theta \sim \text{Binomial}(1, \theta) = \text{Bernoulli}(\theta)$  iid,  $i = 1, \dots, n$ , and  $\Theta \sim \text{Beta}(\alpha, \beta)$  where  $\alpha, \beta$  known. Find the Bayes estimator of  $\theta$  using a squared-error loss function.

# ***Example***

Let  $X_i \mid \theta \sim N(\theta, \sigma^2)$  iid,  $i = 1, \dots, n$ , where  $\sigma^2$  known, and  $\Theta \sim N(\theta_0, \sigma_0^2)$  where  $\theta_0, \sigma_0^2$  known. Find the Bayes estimator of  $\theta$  using a squared-error loss function.

# ***Example***

Let  $Y_1$  and  $Y_2$  be statistics that have trinomial distribution with parameters  $n, \theta_1$ , and  $\theta_2$ . Here  $\theta_1$  and  $\theta_2$  are observed values of the random variables  $\Theta_1$  and  $\Theta_2$ , which have a Dirichlet distribution with known parameters  $\alpha_1, \alpha_2$ , and  $\alpha_3$ . Show that the conditional distribution of  $\Theta_1$  and  $\Theta_2$  is Dirichlet and determine the conditional means  $E(\Theta_1 \mid y_1, y_2)$  and  $E(\Theta_2 \mid y_1, y_2)$ .

# ***Bayesian Terminology***

- A class of prior pdfs for the family of distributions with pdfs  $f(x | \theta)$ ,  $\theta \in \Omega$  is said to be a conjugate family of distributions if the posterior pdf of the parameter is in the same family of distributions as the prior.
- Let  $\mathbf{X}' = (X_1, \dots, X_n)$  be a random sample from the distribution with pdf  $f(x | \theta)$ . A prior  $\pi(\theta) \geq 0$  for this family is said to be improper if it is not a pdf but the function  $\pi(\theta | \mathbf{x}) \propto f(\mathbf{x} | \theta) \cdot \pi(\theta)$  can be made proper.

# ***Bayesian Terminology (cont.)***

- A noninformative prior is a prior that treats all values of  $\theta$  uniformly equal
  - Note: Continuous noninformative priors are often improper
  - Example:  $X \mid (\theta_1, \theta_2) \sim N(\theta_1, \theta_2)$  and  $\theta_1$  has prior  $h(\theta_1) = 1, -\infty < \theta_1 < \infty$