

STA 522, Spring 2021  
Introduction to Theoretical Statistics I

Lecture 1

01 February, 2021



# Agenda

- ▶ Review random samples
- ▶ Order Statistics
- ▶ Convergence Concepts

## Review: Random Samples

**Definition:** The random variables  $X_1, X_2, \dots, X_n$  are called a **random sample** of size  $n$  from the population  $f(x)$  if  $X_1, X_2, \dots, X_n$  are mutually independent random variables and the marginal pdf or pmf of each  $X_i$  is the same function  $f(x)$ .

**Notation:**  $X_1, X_2, \dots, X_n \sim \text{iid } f$ . Joint pdf/pmf:  
 $f_{X_1, X_2, \dots, X_n}(x_1, \dots, x_n) = f(x_1, \dots, x_n) := \prod_{i=1}^n f(x_i)$

If  $f$  is a member of a parametric family with parameter(s)  $\theta$ , then we may write  $f(x_1, x_2, \dots, x_n \mid \theta) = \prod_{i=1}^n f(x_i \mid \theta)$

Example:  $X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$  with  
 $f(x \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$

## Review: Statistics and Sampling Distributions

**Definition:** Let  $X_1, X_2, \dots, X_n$  be a random sample of size  $n$  from a population and let  $T(x_1, x_2, \dots, x_n)$  be a function (real-valued or vector-valued) whose domain includes the sample space of  $(X_1, X_2, \dots, X_n)$ . The random variable (or vector)  $Y = T(X_1, X_2, \dots, X_n)$  is called a **statistic**. The probability distribution of a statistic is called its **sampling distribution**.

**Note:** A statistic cannot contain a parameter.

Examples:

- (i) sample mean  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ ,
- (ii) sample variance
$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \sum_{i=1}^n X_i^2 - \frac{n}{n-1} \bar{X}^2$$
- (iii) sample standard deviation  $S = \sqrt{S^2}$ .
- (iv) sample minimum, sample maximum, sample quantiles.

**Result (Lemma 5.2.5):** Let  $X_1, X_2, \dots, X_n$  be a random sample from a population, and let  $g(x)$  be a function such that  $E(g(X_1))$  and  $\text{Var}(g(X_1))$  exist. Then

$$E\left(\sum_{i=1}^n g(X_i)\right) = n E(g(X_1))$$

and

$$\text{Var}\left(\sum_{i=1}^n g(X_i)\right) = n \text{Var}(g(X_1)).$$

**Result (Theorem 5.2.6):** Let  $X_1, X_2, \dots, X_n$  be a random sample from a population with mean  $\mu$  and variance  $\sigma^2 < \infty$ . Then

(a)  $E(\bar{X}) = \mu$

(b)  $\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$ ; and

(c)  $E(S^2) = \sigma^2$ .

How to determine the sampling distribution of  $\bar{X}$ ?

- (i) **Using transformations.** Let  $Y = \sum_{i=1}^n X_i$ , so that  $\bar{X} = \frac{1}{n} Y$ . If  $f(x)$  is the pdf of  $Y$ , then the pdf of  $\bar{X}$  is  $f_{\bar{X}}(x) = nf(nx)$ .
- (ii) **Using mfg (Theorem 5.2.7).**  $M_{\bar{X}}(t) = M_Y\left(\frac{t}{n}\right) = [M_X\left(\frac{t}{n}\right)]^n$  where  $M_X(t)$  is the mgf of the underlying population. Then identify the distribution of  $\bar{X}$ .

**Theorem 5.3.1.** Let  $X_1, X_2, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$ . Then

- a.  $\bar{X}$  and  $S^2$  are independent random variables.
- b.  $\bar{X} \sim N(\mu, \sigma^2)$ .
- c.  $(n-1)S^2/\sigma^2 \sim \chi_{n-1}^2$ .

# Order Statistics

**Definition:** The order statistics of a random sample  $X_1, X_2, \dots, X_n$  are the sample values placed in ascending order. They are denoted by  $X_{(1)}, X_{(2)}, \dots, X_{(n)}$  and satisfy  $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ .

## Examples:

- (a) **sample minimum:**  $X_{(1)}$  and **sample maximum:**  $X_{(n)}$  are called the extreme order statistics.
- (b) **sample range:**  $R = X_{(n)} - X_{(1)}$ .
- (c) **sample median:**  $M$  where

$$M = \begin{cases} X_{((n+1)/2)} & \text{if } n \text{ is odd;} \\ \frac{X_{(n/2)} + X_{(n/2+1)}}{2} & \text{if } n \text{ is even.} \end{cases}$$



# Sampling Distributions of Extreme Order Statistics from a Continuous Population

Suppose  $X_1, X_2, \dots, X_n$  is a random sample from a population with continuous cdf  $F$  and pdf  $f$ . Then

1.  $\{X_{(n)} \leq x\} = \{\text{all } X_i \leq x\} = \{X_1 \leq x, \dots, X_n \leq x\}$ . So

$$\begin{aligned} F_{X_{(n)}}(x) &= P(X_{(n)} \leq x) \\ &= P(X_1 \leq x, \dots, X_n \leq x) \\ &= P(X_1 \leq x) \dots P(X_n \leq x) \\ &= F(x) \dots F(x) = [F(x)]^n \end{aligned}$$

Differentiating,  $f_{X_{(n)}}(x) = n f(x) [F(x)]^{n-1}$ .

2.  $\{X_{(1)} \geq x\} = \{\text{all } X_i \geq x\} = \{X_1 \geq x, \dots, X_n \geq x\}$ . Implies  $F_{X_{(1)}}(x) = 1 - [1 - F(x)]^n$  &  $f_{X_{(1)}}(x) = n f(x) [1 - F(x)]^{n-1}$ .

**Example:**  $X_1, X_2, \dots, X_n \sim \text{iid Uniform}(0, \theta)$ . Find the pdf and the expected value of  $X_{(n)}$ .

$$\text{Here } f(x | \theta) = \frac{1}{\theta} I(0 \leq x \leq \theta) \text{ and } F(x | \theta) = \begin{cases} 0, & x < 0 \\ \frac{x}{\theta}, & 0 \leq x \leq \theta \\ 1, & x > \theta \end{cases}$$

so that

$$\begin{aligned} f_{X_{(n)}}(x | \theta) &= n f(x | \theta) [F(x | \theta)]^{n-1} \\ &= n \left(\frac{1}{\theta}\right) \left(\frac{x}{\theta}\right)^{n-1} I(0 \leq x \leq \theta) \\ &= \frac{n x^{n-1}}{\theta^n} I(0 \leq x \leq \theta) \end{aligned}$$

Find expected value  $E[X_{(n)}] = E[X_{(n)} | \theta]$  using integration:

$$E[X_{(n)}] = \int_{-\infty}^{\infty} x f_{X_{(n)}}(x | \theta) dx = \frac{n}{\theta^n} \int_0^{\theta} x^n dx = \frac{n}{n+1} \theta$$

## Distribution of a general order statistic from a continuous population

**Theorem 5.4.4.** Let  $X_{(1)}, X_{(2)}, \dots, X_{(n)}$  denote the order statistics of a random sample  $X_1, X_2, \dots, X_n$  from a continuous population with cdf  $F(x)$  and pdf  $f(x)$ . The pdf of  $X_{(j)}$  is

$$f_{X_{(j)}}(x) = \frac{n!}{(j-1)!(n-j)!} f(x) [F(x)]^{j-1} [1-F(x)]^{n-j}.$$

**Partial Proof.** Call  $\{X_i \leq x\}$  a “success,”  $\{X_i > x\}$  a “failure.” Define  $Z_i = I(X_i \leq x)$  and  $Y = \sum_{i=1}^n Z_i$ . Note that  $Z_i \sim \text{iid Bernoulli}(F(x)) \implies Y \sim \text{Binomial}(n, F(x))$ . Note that,

$$F_{X_{(j)}}(x) = P(X_{(j)} \leq x) = P(Y \geq j) = \sum_{k=j}^n \binom{n}{k} [F(x)]^k [1-F(x)]^{n-k}$$

The pdf is obtained using differentiation.

## Distribution of a general order statistic from a discrete population

**Theorem 5.4.3.** Let  $X_1, X_2, \dots, X_n$  be a random sample from a discrete distribution with pmf  $f(x_i) = p_i$ , where  $x_1 < x_2 < \dots$  are the possible values of  $X$  in ascending order. Let  $X_{(1)}, X_{(2)}, \dots, X_{(n)}$  denote the order statistics from the sample. Define

$$P_0 = 0$$

$$P_i = p_1 + p_2 + \dots + p_i \quad \text{for } i \geq 1$$

Then

$$P(X_{(j)} \leq x_i) = \sum_{k=j}^n \binom{n}{k} P_i^k (1 - P_i)^{n-k}$$

$$P(X_{(j)} = x_i) = \sum_{k=j}^n \binom{n}{k} \left[ P_i^k (1 - P_i)^{n-k} - P_{i-1}^k (1 - P_{i-1})^{n-k} \right].$$

**Proof:** cdf is similar to the continuous case. The pmf is obtained from the cdf through

$$P(X_{(j)} = x_i) = P(X_{(j)} \leq x_i) - P(X_{(j)} \leq x_{i-1}).$$

**Example:**  $X_1, X_2, \dots, X_n \sim \text{iid Uniform}(0, 1)$ . Find the distribution of the  $j^{\text{th}}$  order statistic, along with its mean and variance.

Here  $f(x) = I(0 < x < 1)$  and  $F(x) = x$  for  $0 < x < 1$ . Therefore

$$\begin{aligned} f_{X_{(j)}}(x) &= \frac{n!}{(j-1)!(n-j)!} f(x) [F(x)]^{j-1} [1 - F(x)]^{n-j} \\ &= \frac{n!}{(j-1)!(n-j)!} x^{j-1} (1-x)^{n-j} I(0 < x < 1) \\ &= \frac{\Gamma(n)}{\Gamma(j)\Gamma(n-j+1)} x^{j-1} (1-x)^{n-j} I(0 < x < 1) \end{aligned}$$

This shows that  $X_{(j)} \sim \text{Beta}(j, n-j+1)$ . From this, we can deduce that

$$E[X_{(j)}] = \frac{j}{n+1}$$

and

$$\text{Var}[X_{(j)}] = \frac{j(n-j+1)}{(n+1)^2(n+2)}.$$

## Joint Distribution of Order Statistics

**Theorem 5.4.6.** Let  $X_{(1)}, X_{(2)}, \dots, X_{(n)}$  denote the order statistics of a random sample  $X_1, X_2, \dots, X_n$  from a continuous population with cdf  $F(x)$  and pdf  $f(x)$ . The joint pdf of  $X_{(i)}$  and  $X_{(j)}$ ,  $1 \leq i < j \leq n$ , is

$$f_{X_{(i)}, X_{(j)}}(u, v) = c f(u) f(v) [F(u)]^{i-1} [F(v) - F(u)]^{j-1-i} [1 - F(v)]^{n-j}$$

for  $-\infty < u < v < \infty$ , where  $c = \frac{n!}{(i-1)!(j-1-i)!(n-j)!}$ .

Joint distribution pdf of all the order statistics from a continuous population:

$$f_{X_{(1)}, \dots, X_{(n)}}(x_1, \dots, x_n) = \begin{cases} n! f(x_1) \dots f(x_n), & -\infty < x_1 < \dots < x_n < \infty \\ 0, & \text{otherwise} \end{cases}$$

**Example:** Let  $X_1, X_2, \dots, X_n \sim \text{iid uniform}(0, a)$ ,  
 $X_{(1)}, X_{(2)}, \dots, X_{(n)}$  denote the order statistics. Find the joint pdf  
of the sample range  $R = X_{(n)} - X_{(1)}$  and the mid-range  
 $V = \frac{X_{(1)} + X_{(n)}}{2}$ . Hence find the marginal pdf of  $R$ .

First obtain the joint pdf of  $X_{(1)}$  and  $X_{(n)}$ :

$$\begin{aligned} f_{X_{(1)}, X_{(n)}}(x_1, x_n) &= \frac{n(n-1)}{a^2} \left( \frac{x_n}{a} - \frac{x_1}{a} \right)^{n-2} I(0 < x_1 < x_n < a) \\ &= \frac{n(n-1)(x_n - x_1)^{n-2}}{a^n} I(0 < x_1 < x_n < a) \end{aligned}$$

Solve for  $X_{(1)}, X_{(n)}$  to obtain  $X_{(1)} = V - R/2$  and  
 $X_{(n)} = V + R/2$ . Jacobian of this transformation is -1.

Support of  $(R, V)$ :

$$\begin{aligned}0 < x_1 < x_n < a \\ \implies 0 < v - r/2 < v + r/2 < a \\ \implies 0 < r < a, \quad r/2 < v < a - r/2\end{aligned}$$

The joint pdf of  $(R, V)$  is

$$f_{R,V}(r, v) = \frac{n(n-1) r^{n-2}}{a^n}; \quad 0 < r < a, \quad r/2 < v < a - r/2$$

The marginal pdf of  $R$  is

$$f_R(r) = \int_{r/2}^{a-r/2} \frac{n(n-1) r^{n-2}}{a^n} dv = \frac{n(n-1) r^{n-2} (a-r)}{a^n}; \quad 0 < r < a$$

It is easy to see that  $\frac{R}{a} \sim \text{Beta}(n-1, 2)$  distribution.

**HW:** find the marginal pdf of  $V$ .



# Convergence Concepts

What happens to sample statistics, particularly  $\bar{X} = \bar{X}_n$ , when the sample size  $n \rightarrow \infty$ ?

For a real sequence  $(a_n)_{n=1}^{\infty}$  defining convergence is somewhat straightforward:  $(a_n)_{n=1}^{\infty}$  is said to converge to a point  $a$  if  $\lim_{n \rightarrow \infty} a_n = a$ .

How to define convergence of random variables?

- ▶ convergence in probability
- ▶ convergence in almost sure sense
- ▶ convergence in distribution (or law)
- ▶ convergence in mean [later]