

# Lecture 20

## Descriptive Statistics

---

Manju M. Johny

STAT 330 - Iowa State University

# Statistics

---

## Definition: Statistics

A *statistic*,  $T(X_1, \dots, X_n)$  is a function of random variables.

- Start with taking a *simple random sample (SRS)* of size  $n$  from some population/distribution.

$$X_1, \dots, X_n \stackrel{iid}{\sim} f_X(x)$$

- We can then obtain *statistics* based on  $X_1, \dots, X_n$
- Since a statistic is a function  $T(\cdot)$  of random variables, the statistic is also a random variable.
- Thus, the statistic will have its own distribution called the *sampling distribution of the statistic* (more on this later!)

### Definition: Observed Statistics

The *observed statistics*,  $T(x_1, \dots, x_n)$  is the statistic function with observed values plugged in.

- *Descriptive statistics*: Describing what our sample data looks like (graphically or numerically)
- *Inferential statistics*: Use the statistic to infer/learn about the "true" distribution,  $f_X(x)$ , that generated the data.

### Note:

- Use small letters ( $x$ ,  $\bar{x}$ ,  $s^2$ , etc) to represent observations and observed statistics.
- Use capital letters ( $X$ ,  $\bar{X}$ ,  $S^2$ , etc) to represent random variables.

# Mean and Variance

---

## Sample Mean and Variance

Let  $X_1, \dots, X_n \stackrel{iid}{\sim} f_X(x)$  where  $E(X_i) = \mu$  and  $Var(X_i) = \sigma^2$

- **Sample mean** is defined as  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$   
→ estimates the population mean  $\mu$ .
- **Sample variance** is defined as  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$   
→ estimates the population variance  $\sigma^2$   
→ an estimate of the  $Var(X) = E[(X - E(X))^2]$  can be found as  $\frac{1}{n} \sum_{i=1}^n (X_i - (\bar{X}))^2$   
→ typically,  $n$  in the above denominator is replaced with  $n - 1$  to get  $S^2$  (more on this later)
- **Sample standard deviation** is  $S = \sqrt{S^2}$

**Note:** The quantities above are R.V's since they are functions of R.V's  $X_1, \dots, X_n$ .

## Observed Sample Mean and Variance

- To obtain the *observed sample mean* and *observed sample variance*, plug in observed data values  $(x_1, \dots, x_n)$  into sample mean and variance formulas

$$\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2$$

$$s = \sqrt{s^2}$$

**Note:** The quantities above are not random variables since you have plugged in data values. They are values such as 2.4, 100, etc.

# Quantiles

---



# Quantiles

## Definition: Quantiles

The  $q^{\text{th}}$  *quantile* of a distribution,  $f_X(x)$ , is a value  $x$  such that  $P(X < x) \leq q$  and  $P(X > x) \leq 1 - q$ .

This is also called the  $100 \cdot q^{\text{th}}$  *percentile*.

$Q_1 = 0.25^{\text{th}}$  quantile,  $Q_2 = 0.5^{\text{th}}$  quantile (median), and  $Q_3 = 0.75^{\text{th}}$  quantile

## Definition: Quantile Function

The *quantile function* is defined as:

$$F_X^{-1}(q) = \min\{x : F_X(x) \geq q\}$$

The *median* is the  $0.5^{\text{th}}$  quantile (or  $50^{\text{th}}$  percentile)

→ can be written as  $F_X^{-1}(0.5)$

The *sample median* is calculated by:

1. Order sampled values in increasing order:  $X_{(1)}, \dots, X_{(n)}$

- If  $n$  is odd, take the middle value

→ median =  $X_{[\frac{n}{2}]}$

- If  $n$  is even, average the two middle values

→ median =  $\frac{X_{\frac{n}{2}} + X_{\frac{n}{2}+1}}{2}$

**Note:** Since the above values are functions of R.V's, they are R.Vs.

Obtain the *observed sample median* by plugging in the observed values  $(x_1, \dots, x_n)$  from data.

Other sample quantiles we are typically interested in are

- $Q_1 = 0.25^{th}$  quantile
- $Q_3 = 0.75^{th}$  quantile

Many ways to calculate quantiles. Our method for a general  $q^{th}$  sample quantile is ...

1. Compute  $(n + 1) \cdot q$ 
  - If this value is an integer, use  $(n + 1)q^{th}$  ordered value
  - Else, use the average of the 2 surrounding values

## Example

Example 1: A sample  $X_1, \dots, X_n \stackrel{iid}{\sim} f_X(x)$  was taken where  $X_i =$  CPU time for a randomly chosen task. The ordered observed values are 15, 34, 35, 36, 43, 48, 49, 62, 70, 82 (secs)

## Example Cont.

Right now, we're only using these statistics to describe the sample of CPU speeds.

- sample mean and median ( $Q_2$ ) tell us “typical” values
- sample variance tells us how “spread out” / how variable the data are
- $Q_1$  and  $Q_3$  “rank” where values fall in our sample

## Mode, Range, IQR

---

# Mode, Range, and IQR

Other common descriptive statistics to describe the data:

- *Mode*: The most frequent value in our sample. Can have multiple modes in data set
- *Range*:  $\text{Max} - \text{Min} = X_{(n)} - X_{(1)}$   
→ describes the “total” variability of the data
- *Interquartile Range (IQR)*:  $Q_3 - Q_1$   
→ describes the variability of the middle 50% of data

# Robust Statistics

- With all the different options for statistics, how do we choose which ones to use?
  - It depends on your data set
- Statistics that are not affected by extreme values are called *robust statistics*

Example 2: