

- Using the Standard Normal Table
- The Normal Distribution
- Histograms
- Summary Statistics
- Discuss Quizzes/Answer Questions

## Lesson Plan

We shall cover this quickly, so be sure to read the book.

- Visualizations (histograms, other graphics)
- Measures of Dispersion (standard deviation, range, IQR)
- Measures of Central Tendency (mean, median, mode)

are:

Given a collection of data, one needs to find representations of the data that facilitate understanding and insight. Three standard tools for this

### 3. Summary Statistics

The **median** is the middle largest value among the  $n$  observations, if  $n$  is odd. If there are an even number of observations, then we average the two middle-largest values.

The **mode** is just the value that occurs most frequently in the sample. There can be many modes.

$$\cdot \cdot ({}^u X + \cdots + {}^1 X) \frac{u}{1} = \\ {}^1 X \sum_{u=1}^u \frac{u}{1} = \underline{X}$$

The **mean** is just the average of the data. Suppose one has a sample of  $n$  observations with values  $X_1, \dots, X_n$ . Then the mean is just

## 3.1 Measures of Central Tendency

- Example:** Suppose one observes the following data: 1, 0, 2, -2, 1, -2, 5, -1. The mean is  $\underline{X} = \frac{1}{8}(1 + 0 + 2 - 2 + 1 - 2 + 5 - 1) = .5$ . The modes are 1 and -2. The median is the average of 0 and 1, or .5. The single large or small datum will have a large influence on the mean, but not on the median.
- The mean can be pulled in misleading directions if there are outliers. A single large or small datum will have a large influence on the mean, but not on the median.
- An **outlier** is an incorrect or unrepresentative observation that is very different from the others in the sample.

The first formula is better for understanding, and the second is better for

each observation and the mean.  
This is the square root of the average of the squared deviations between

$$\sqrt{\frac{[(\underline{X} - \bar{X})^2 + (\underline{X} - \bar{X})^2 + \dots + (\underline{X} - \bar{X})^2]}{n}} = s$$

To measure how spread out a sample is, we mostly use the **standard deviation**. This is:

## 3.2 Measures of Dispersion

- The majority of observations are usually within 1 sd of the mean. But in some extreme cases it can happen that none of the data are within 1 sd of the mean.
- However, one can prove that:
- At least 75% of the observations must always be within 2 standard deviations of the mean.
  - At least 89% of the observations must always be within 3 standard deviations of the mean.
- This is a result of Chebyshev's Theorem.
- Note Bene:** In some books, one divides by  $n - 1$  rather than  $n$ , but we shall not do that.

The **range** is the largest observation minus the smallest. As a measure of dispersion, it is strongly influenced by outliers.

The **interquartile range** is  $75^{\text{th}} \text{ percentile} - 25^{\text{th}} \text{ percentile}$  of the data minus the  $25^{\text{th}}$  percentile (the median is the  $50^{\text{th}}$  percentile). For our purposes, we shall calculate the  $75^{\text{th}}$  percentile as the median of all observations strictly less than the median, and the  $25^{\text{th}}$  percentile as the median of all observations strictly greater than the median, and the  $25^{\text{th}}$  percentile is as the median of all observations strictly less than the median. (This is not perfectly correct, but is good enough for our purposes.)

The **interquartile range** is not strongly influenced by outliers.

$$\text{standard deviation} = \sqrt{\frac{1}{n}[(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2]}$$

but it is faster to calculate

$$\text{standard deviation} = \sqrt{\frac{1}{n}[(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2]}$$

The standard deviation is

median{0, -2, -2, -1}, or  $1.5 - (-1.5) = 3$ .  
 minus the median of all values less than  $\bar{x}$ , or median{1, 2, 1, 5} minus  
 The interquartile range is the median of all values greater than  $\bar{x}$ ,

The range is  $5 - (-2) = 7$ .

1, 0, 2, -2, 1, -2, 5, -1

**Example:** Suppose one observes the following data:

range of the  $\underline{Y}$  values?  
 Can you guess formulae for the median, mode, range, and interquartile

$$\cdot \underline{x}ps|a| = \underline{y}ps$$

$$q + \underline{X}^a = \underline{Y}$$

Then

the Centigrade scale.)  
 converting units of measurement, such as changing Fahrenheit data into  
 a new sample  $\underline{Y}_1, \dots, \underline{Y}_n$  where  $\underline{Y}_i = a\underline{X}_i + b$ . (This often arises when  
 Suppose we have  $n$  observations  $\underline{X}_1, \dots, \underline{X}_n$  and we use these to create

### 3.3 Properties of $\underline{X}$ and the $sd$

The histogram shows where sample values are located and where they interval, but not the right.

By convention, the left endpoint of a histogram bar is included in the

In a histogram, the areas under a block represent percentages.

$x$ -value. (This is different from a bar chart.)

The  $x$ -axis gives the sample value, and the  $y$ -axis is the percent per

concentrate.

The histogram shows where sample values are located and where they

## 3.4 The Histogram

This incomplete histogram represents monthly wages for part-time employees.

- What is the height of the missing bar?
- Estimate the percentage of part-time employees who make between \$100 and \$110.
- Does the second bar include people who make exactly \$100/month?
- Does the second bar include people who make exactly \$200/month?

As the sample size goes to infinity ( $n \rightarrow \infty$ ) and as the bin-width of the histogram goes to zero ( $h \rightarrow 0$ ) at the appropriate relative rates, then the histogram becomes smooth.

This limiting smooth curve is called a **distribution**.

The term “Gaussian” refers to Carl Friedrich Gauss. Who was he?

Some limiting histograms are famous and have names. The most famous distribution is the **normal distribution** (a/k/a the Gaussian distribution or the bell-shaped curve).

### 3.5 The Normal Distribution

People believe the normal distribution describes IQ, height, rainfall, measurement error, and many other features. This is only approximately true. But it is a good approximation for features that are the sum of many separate increments.

To make a normal probability plot, find the sample mean  $\bar{X}$  and the sample standard deviation  $s_d$ . For each observation  $X_i$ , find the area under the standard normal curve (from the  $z$ -table in our book, page A-105) to left of  $(X_i - \bar{X})/s_d$ . Plot this area against  $X_i$  for all  $i$ . Perfect normality corresponds to a straight line.

- Make a **normal probability plot**.
- Inspect the histogram.

How can you decide if data are a random sample from a normal distribution?

The standard normal has  $\mu = 0$  and  $\sigma = 1$ .

The  $\mu$  is the mean of the entire population, whereas  $X$  is used to denote the mean of a sample from the population. Similarly,  $\sigma$  is the standard deviation of the entire population, whereas  $s$  is used to denote the standard deviation of a sample.

for  $-\infty < \mu < \infty$  and  $\sigma \geq 0$ .

$$\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) = f(x)$$

equation:

A normal distribution with mean  $\mu$  and standard deviation  $\sigma$  has the

The mean of a normal distribution shows where it is centred. The standard deviation of a normal distribution shows how spread out the normal is.

Any question about a normal distribution can be converted into an equivalent question about the standard normal, and vice-versa.

First we practice reading the book's table (page A-105). Assume you have a population that is standard normal.

What proportion of the population has values between -1.5 and 1.5?

### 3.6 The Standard Normal Distribution

From the table, the proportion is .8664, or 86.64% of the population lies between -1.5 and 1.5. Now go the other way. 80% of the population lies between what two values that are centred at 0?

From the table, the answer is about -1.3 and 1.3. So approximately 80% of the standard normal values are within 1.3 of the mean (recall, the mean is 0).

From the table, the answer is about -1.3 and 1.3. So approximately 80% of the standard normal values are within 1.3 of the mean (recall, the mean is 0).

- What proportion of the population has a value less than -3? (Ans: about .3821)

- What proportion of the population has a value larger than -1? (Ans: about .8413)

- What is the value for which about 90% of the population is smaller? (Ans: about 1.3)

- What is the value of  $z$  such that 25% of a standard normal population is larger than that value? (Ans: about .7)

Some problems to think about in preparation for the quiz on Monday.