

Summarizing Data: Measures of Center & Rank

Engineering Statistics
Section 1.3

Josh Engwer

TTU

27 January 2016

PART 0: PRELIMINARIES

NOTATION FOR SAMPLES & POPULATIONS

SAMPLE STATISTICS & POPULATION PARAMETERS

SORTED SAMPLES & ROUNDING NOTATION

Descriptive Statistics

Recall that Statistics consists of two broad branches:

- Descriptive Statistics
- Statistical Inference

The remainder of this chapter focuses squarely on Descriptive Statistics:

Definition

Descriptive Statistics is the organization, summary, visualization and presentation of data that conveys useful information about the data.

Descriptive Statistics involves:

- Data Visualization (Section 1.2)
- Numerical Summaries (this section and the next)

Notation for Samples

For methods & procedures, it's helpful to have consistent notation for samples:

Definition

(Notation for a Single Univariate Sample)

Sample as a whole is denoted by x .

The **sample size** (i.e. # data points) is denoted by n .

Each data point is denoted by a corresponding subscript: $x_1, x_2, \dots, x_{n-1}, x_n$

Definition

(Notation for Two Univariate Samples)

Samples as a whole are denoted by x & y .

The sample sizes are denoted by n & m or n_1 & n_2

The data points are denoted by subscripts:

x_1, x_2, \dots, x_n & y_1, y_2, \dots, y_m OR x_1, x_2, \dots, x_{n_1} & y_1, y_2, \dots, y_{n_2}

For 3+ samples (rare), run thru upper-end of lowercase alphabet as needed:

x, y, z, w, v, u

Notation for Samples (Examples)

- Student Heights (in ft) $x : 6.1, 3.9, 5.6, 4.0, 5.9, 5.9$
 - Sample Size $n_1 = (\# \text{ data points in sample } x) = 6$
 - Data points $x_1 = 6.1, x_2 = 3.9, x_3 = 5.6, x_4 = 4.0, x_5 = 5.9, x_6 = 5.9$

 - Student Weights (in lb) $y : 205, 135, 183$
 - Sample Size $n_2 = (\# \text{ data points in sample } y) = 3$
 - Data points $y_1 = 205, y_2 = 135, y_3 = 183$

 - Student Eye Colors Hazel, Blue, Brown, Hazel
 - Sample Size $n_3 = (\# \text{ data points in sample of categorical data}) = 4$
 - Sample & Data points of categorical data are not labeled.
-

There are two ways to write out a sample:

- As a list of comma-separated values
 - $x : 6.1, 3.9, 5.6, 4.0, 5.9, 5.9$
 - Hazel, Blue, Brown, Hazel
- As a list of space-separated values
 - $x : 6.1 \ 3.9 \ 5.6 \ 4.0 \ 5.9 \ 5.9$
 - Hazel Blue Brown Hazel

Sample Statistics & Population Parameters

Definition

(Sample Statistic)

A **statistic** of a sample is a meaningful characteristic of a the sample. Statistics are denoted by certain "decorations" of the letter for the sample.

Definition

(Population Parameter)

A **parameter** of a population is a meaningful characteristic of the population. Parameters are often (but not always) denoted by lower-case Greek letters.

Definition

(Notation for Populations)

A population itself is never denoted by a letter. However, the size of a finite population is denoted by N . For two finite populations, their sizes are denoted N & M OR N_1 & N_2

Sorted Samples

With discrete numerical data, it's important for some sample statistics that the sample is sorted in ascending order.

As it happens, there's corresponding notation for a sorted sample:

Definition

(Sorted Samples)

Given a sample with n data points $x : x_1, x_2, \dots, x_{n-1}, x_n$

Then the corresponding **sorted sample** is $x : x_{(1)}, x_{(2)}, \dots, x_{(n-1)}, x_{(n)}$

where the data points are sorted in ascending order:

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n-1)} \leq x_{(n)}$$

$x_{(1)}$ is the **smallest data point** in the sample.

$x_{(n)}$ is the **largest data point** in the sample.

EXAMPLE: Given sample $x : 5, 4, 8 \implies x_1 = 5, x_2 = 4, x_3 = 8$

Then, the sorted sample is $x : 4, 5, 8 \implies x_{(1)} = 4, x_{(2)} = 5, x_{(3)} = 8$

Rounding Numbers (Compact Notation)

It is convenient to have mathematical notation for **rounding numbers**.

Always Round Down: $\lfloor 3 \rfloor = 3$ $\lfloor 3.1 \rfloor = 3$ $\lfloor 3.5 \rfloor = 3$ $\lfloor 3.9 \rfloor = 3$

Always Round Up: $\lceil 3 \rceil = 3$ $\lceil 3.1 \rceil = 4$ $\lceil 3.5 \rceil = 4$ $\lceil 3.9 \rceil = 4$

Round to Nearest Integer: $\llbracket 3 \rrbracket = 3$ $\llbracket 3.1 \rrbracket = 3$ $\llbracket 3.5 \rrbracket = 4$ $\llbracket 3.9 \rrbracket = 4$

$\lfloor x \rfloor$ is called the **floor function**.

$\lceil x \rceil$ is called the **ceiling function**.

PART I: MEASURES OF CENTER FOR DISCRETE NUMERICAL DATA

MEAN, MEDIAN, TRIMMED MEAN

Mean of a Sample

Definition

(Mean of a Discrete Numerical Sample)

Given a sample with n data points $x : x_1, x_2, \dots, x_n$
Then its **mean**, denoted \bar{x} , is the average of the sample.

$$\bar{x} := \frac{1}{n} \sum_{k=1}^n x_k = \frac{x_1 + x_2 + \dots + x_n}{n}$$

NOTE: \bar{x} is pronounced "x bar"

REMARK: The symbol $:=$ translates to "is defined to be".

Median of a Sample

Definition

(Median of a Discrete Numerical Sample)

Given a sample with n data points

$$x : x_1, x_2, \dots, x_n$$

Then its **median**, denoted \tilde{x} , is the middle value of the sorted sample.

$$\tilde{x} := \begin{cases} x_{([n+1]/2)} & , n \text{ odd} \\ \frac{x_{(n/2)} + x_{(1+[n/2])}}{2} & , n \text{ even} \end{cases} = \begin{cases} \text{Middle data point} & , n \text{ odd} \\ \text{Average of the two} & , n \text{ even} \\ \text{middle data points} & \\ \text{in sorted sample} & \end{cases}$$

NOTE: \tilde{x} is pronounced "x tilde" OR "x twiddle"

For instance, given sample

$$x : 6, 4, 5, 7, 1, 7, 2$$

First, sort the sample:

$$x : 1, 2, 4, \mathbf{5}, 6, 7, 7$$

Then, since sample size ($n = 7$) is odd, the median is:

$$\tilde{x} = x_{([n+1]/2)} = x_{(4)} = \mathbf{5}$$

Median of a Sample

Definition

(Median of a Discrete Numerical Sample)

Given a sample with n data points

$$x : x_1, x_2, \dots, x_n$$

Then its **median**, denoted \tilde{x} , is the middle value of the sorted sample.

$$\tilde{x} := \begin{cases} x_{([n+1]/2)} & , n \text{ odd} \\ \frac{x_{(n/2)} + x_{(1+[n/2])}}{2} & , n \text{ even} \end{cases} = \begin{cases} \text{Middle data point} & , n \text{ odd} \\ \text{Average of the two} & , n \text{ even} \\ \text{middle data points} & \\ \text{in sorted sample} & \end{cases}$$

NOTE: \tilde{x} is pronounced "x tilde" OR "x twiddle"

For instance, given sample

$$y : 6, 4, 5, 7, 9, 1, 7, 2$$

First, sort the sample:

$$y : 1, 2, 4, 5, 6, 7, 7, 9$$

Then, since sample size ($n = 8$) is even, the median is:

$$\tilde{y} = \frac{y_{(n/2)} + y_{(1+[n/2])}}{2} = \frac{y_{(4)} + y_{(5)}}{2} = \frac{5+6}{2} = \boxed{5.5}$$

Trimmed Mean of a Sample

Definition

(Trimmed Mean of a Discrete Numerical Sample)

Given a sample with n data points $x : x_1, x_2, \dots, x_n$

Then its $p\%$ **trimmed mean**, $\bar{x}_{tr}(p\%)$, is the mean of the dataset resulting from eliminating the smallest $p\%$ and largest $p\%$ of the sorted sample.

$\bar{x}_{tr}(10\%)$:= Mean of sorted sample x with largest 10% & smallest 10% removed

$\bar{x}_{tr}(25\%)$:= Mean of sorted sample x with largest 25% & smallest 25% removed

Relevant trimming percentages tend to be moderate: between 5% & 25%

$\bar{x}_{tr}(10\%)$ is spoken as "the 10% trimmed sample mean."

REMARK:

For simplicity, the trimming percentage will always evenly divide sample size n .
In other words, the expression $np/100$ will always be an integer.

Otherwise, **interpolation** would be needed which complicates matters!!

Interpolation is formally encountered in **Numerical Analysis**. (MATH 4310)

Mean, Median, Trimmed Means of a Population

The mean, median, and trimmed means can be computed for populations.

Definition

(Notation for Mean, Median, Trimmed Means of a Population)

The **population mean** is denoted by μ . (“mew”)

The **population median** is denoted by $\tilde{\mu}$. (“mew tilde” or “mew twiddle”)

The $p\%$ **trimmed population mean** is denoted by $\mu_{tr(p\%)}$.

The 10% **trimmed population mean** is denoted by $\mu_{tr(10\%)}$.

The 25% **trimmed population mean** is denoted by $\mu_{tr(25\%)}$.

REMARKS:

μ is the lower-case Greek letter mu.

Computing μ , $\tilde{\mu}$, etc for finite populations is not practical due to their enormity.

Computing μ , $\tilde{\mu}$, etc for infinite populations will be encountered in Chapter 4.

Mean, Median and Skewness of a Sample

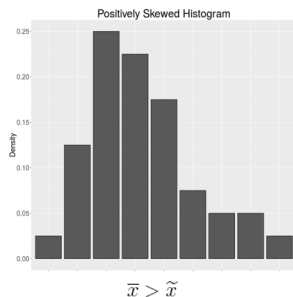
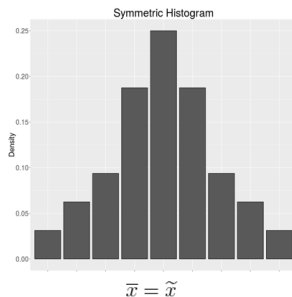
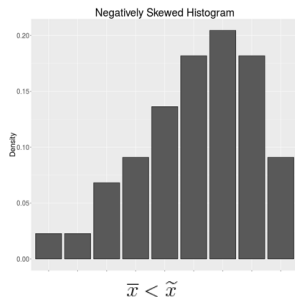
Proposition

Given a sample with n data points

$x : x_1, x_2, \dots, x_n$

Then:

- If $\bar{x} < \tilde{x}$, then the sample is negatively skewed.
- If $\bar{x} = \tilde{x}$, then the sample is symmetric.
- If $\bar{x} > \tilde{x}$, then the sample is positively skewed.

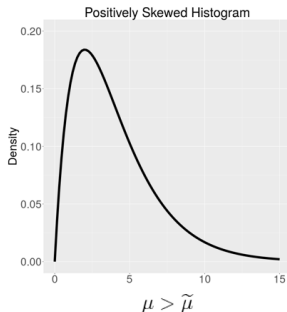
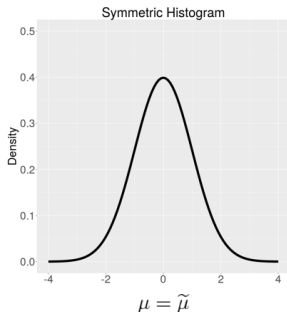
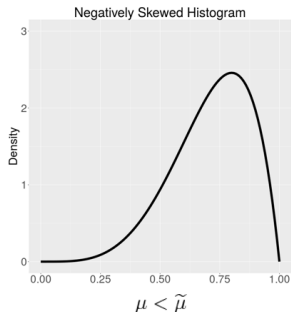


Mean, Median and Skewness of a Population

Proposition

Given a population with mean μ and median $\tilde{\mu}$. Then:

- If $\mu < \tilde{\mu}$, then the population is negatively skewed.
- If $\mu = \tilde{\mu}$, then the population is symmetric.
- If $\mu > \tilde{\mu}$, then the population is positively skewed.



Measures of Center & Their Sensitivity to Outliers

The sample mean, median, and trimmed mean are examples of **statistics**. The sample mean, median, and trimmed mean are all **measures of center**. This means they indicate "central locations" of the sample in unique ways.

A question that's important for many situations is:

HOW SENSITIVE ARE THESE MEASURES OF CENTER TO **OUTLIERS**??

- The mean, \bar{x} , is extremely sensitive to outliers.
- Lightly-trimmed means (e.g. $\bar{x}_{tr(5\%)}$) are largely sensitive to outliers.
- Heavily-trimmed means (e.g. $\bar{x}_{tr(25\%)}$) are largely insensitive to outliers.
- The median, \tilde{x} , is almost completely insensitive to outliers.

So, if a subdivision has all its houses priced in the \$100,000's and, later, a ten million-dollar house is built there, then:

- the mean house price will increase substantially...
- ...but the median house price will only increase slightly.

PART II: MEASURES OF RANK FOR DISCRETE NUMERICAL DATA

PERCENTILES, QUARTILES, HINGES

Percentiles of a Sample

Definition

(Percentile of a Discrete Numerical Sample)

Given a sample with n data points $x : x_1, x_2, \dots, x_n$

Then its p -th **percentile**, denoted $x_{p/100}$, is the smallest data point such that $p\%$ of the sample is less than or equal to that data point:

$$x_{p/100} := x_{(\lceil np/100 \rceil)} = \left(\left\lceil \frac{np}{100} \right\rceil \right)\text{-th data point in sorted sample}$$

e.g. (37% of sample x) $\leq x_{0.37} \equiv$ (37th percentile of sample x)

e.g. (98% of sample y) $\leq y_{0.98} \equiv$ (98th percentile of sample y)

REMARK: The symbol \equiv translates to "represents" OR "is represented by".

Software packages (e.g. MATLAB, R, SPSS, SAS, Minitab) may define percentiles slightly differently.

Quartiles of a Sample

Definition

(Quartiles of a Discrete Numerical Sample)

Given a sample with n data points $x : x_1, x_2, \dots, x_n$ Then:

(1) $x_{Q1} := x_{0.25} \equiv 1^{st}$ **quartile** of sample x

- i.e. (25% of sample x) \leq (1^{st} quartile of sample x)

(2) $x_{Q2} := x_{0.50} \equiv 2^{nd}$ **quartile** of sample x

- i.e. (50% of sample x) \leq (2^{nd} quartile of sample x)
- 2^{nd} quartile, x_{Q2} , is never used since it's exactly or very close to median, \tilde{x} .

(3) $x_{Q3} := x_{0.75} \equiv 3^{rd}$ **quartile** (75^{th} percentile) of sample x

- i.e. (75% of sample x) \leq (3^{rd} quartile of sample x)

REMARK: Software packages (e.g. MATLAB, R, SPSS, SAS, Minitab) may define quartiles slightly differently.

Hinges of a Sample

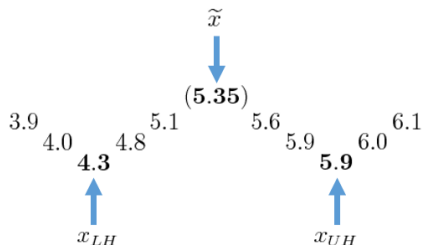
Definition

(Hinges of a Discrete Numerical Sample)

Given a sample with n data points $x : x_1, x_2, \dots, x_n$ Then:

- (1) its **lower hinge**, x_{LH} , is the median of the lower half of sorted sample.
- (2) its **middle hinge**, x_{MH} , is exactly the median of entire sample: $x_{MH} = \tilde{x}$
- (3) its **upper hinge**, x_{UH} , is the median of the upper half of sorted sample.

For instance, given sample $x : 3.9, 4.0, 4.3, 4.8, 5.1, 5.6, 5.9, 5.9, 6.0, 6.1$



Hinges of a Sample

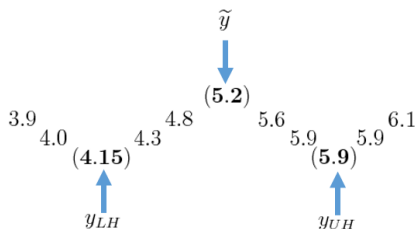
Definition

(Hinges of a Discrete Numerical Sample)

Given a sample with n data points $x : x_1, x_2, \dots, x_n$ Then:

- (1) its **lower hinge**, x_{LH} , is the median of the lower half of sorted sample.
- (2) its **middle hinge**, x_{MH} , is exactly the median of entire sample: $x_{MH} = \tilde{x}$
- (3) its **upper hinge**, x_{UH} , is the median of the upper half of sorted sample.

For instance, given sample $y : 5.9, 3.9, 5.9, 4.8, 5.6, 4.0, 6.1, 4.3$



Hinges of a Sample

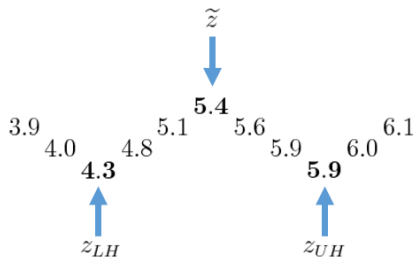
Definition

(Hinges of a Discrete Numerical Sample)

Given a sample with n data points $x : x_1, x_2, \dots, x_n$ Then:

- (1) its **lower hinge**, x_{LH} , is the median of the lower half of sorted sample.
- (2) its **middle hinge**, x_{MH} , is exactly the median of entire sample: $x_{MH} = \tilde{x}$
- (3) its **upper hinge**, x_{UH} , is the median of the upper half of sorted sample.

For instance, given sample $z : 3.9, 4.0, 6.1, 4.3, 5.9, 4.8, 5.1, 6.0, 5.6, 5.4, 5.9$



Textbook Logistics for Section 1.3

- Difference(s) in Terminology:

TEXTBOOK TERMINOLOGY	SLIDES/OUTLINE TERMINOLOGY
Measures of Location (Mean, Median, ...)	Measures of Center
Measures of Location (Percentiles, Hinges, ...)	Measures of Rank

- Difference(s) in Notation:

CONCEPT	TEXTBOOK NOTATION	SLIDES/OUTLINE NOTATION
Arbitrary Trimming Percentage	$100\alpha\%$	$p\%$
Trimmed Mean	$\bar{x}_{tr(12.5)}$	$\bar{x}_{tr(12.5\%)}$

- Ignore "Categorical Data and Sample Proportions" section (pg 34)
 - Sample proportions were encountered with freq. tables & histograms.
 - This numerical look at sample proportions will be encountered in Ch6.

Fin.