

Introduction to Data Analysis

Kate Davis

January 30, 2015

Introduction to Data Analysis

Statistical Data Analysis is quantitative evaluation of **Numeric Data** and multiple data points with the same unit can be combined using basic arithmetic operations to form a new data point. Height (in mm), weight (in kg), temperature (degrees F), proportions (percent), and monetary values are examples of **continuous** numeric data points. **Discrete** numeric data are whole number data or count data, such as the number of sunspots per month, the number of apples in a bushel, or dice roll values. House numbers, credit scores, and jersey numbers are examples of numbers that are not numeric data points, as none have units nor can these numbers be combined arithmetically to form another numeric data point.

Numeric Data points are numbers that represents value. Generally, each numeric data point has a unit of measure

Continuous Data has an infinite number of possible values within a given range, usually represented by real numbers, percentages or fractions

Discrete Data are data with a finite list of possible values within any given range, and are often integer or count data

Height in Whole Inches

Consider the numeric **Data Set** of **Height in Whole Inches** of our **Population**: students in MA3200 Section 2. Heights would be continuous data, but we have “discretized” this data by rounding to the nearest whole inch. The data, in the original order presented, is:

```
64 70 72 73 69 67 68 66 62 71 66 72 67 74
71 72 67 71 65 65 69 71 69 72 71 68 63 54
```

This set of data has 28 data points. To better evaluate this data, lets sort it. We can begin to see patterns of multiple values, and can quickly see that the lowest or minimum value is 54 inches and the highest or maximum value is 74 inches. The **Range** is 20 inches.

```
54 62 63 64 65 65 66 66 67 67 67 68 68 69
69 69 70 71 71 71 71 71 72 72 72 72 73 74
```

This data set has 14 discrete values for height, fewer than the range of 20 inches. There is a gap in observations between 54 inches and 62 inches, but all other height values in the range are represented.

To gain more knowledge about this dataset, we must describe the **distribution** of values across the measurement range, with a goal of using that information for predictions, estimations and other inferences about the population when a complete **census**

The **statistical distribution** can be estimated or inferred from a data

A **Data Set** is a collection of numeric data points. Each data point within a data set is called an observation, denoted x_i . n denotes the number of observations.

A **Population** is any complete group or set of measure with at least one characteristic in common

The **Range** is difference between the maximum and minimum values of a data set

The Oxford English Dictionary defines **Distribution** as the *way in which something is shared out among a group or spread over an area*

A **Census** is a complete enumeration of every unit, everyone or everything in a population.

A **Statistical Distribution** assigns probabilities to the possible values of a data set

set, and is used to estimate the accuracy of these predictions, estimates, inferences.

Frequency Tables

We can create a **Frequency table** and **histogram** of the data set values. The height data is in whole inches, so we will start with using the integer height value as the class in integer order. A cumulative frequency column is added for additional calculation.

Height	Frequency	CumulativeFrequency
54	1	1
62	1	2
63	1	3
64	1	4
65	2	6
66	2	8
67	3	11
68	2	13
69	3	16
70	1	17
71	5	22
72	4	26
73	1	27
74	1	28

a **Frequency Table** is a summary of data point Frequency by class or interval

a **Histogram** is a chart that displays the distribution of a data set

Table 1: Frequency Table

Measures of Center

To understand more about the distribution of the height in inches of our students, we first examine “centers” of the data: the mode, the median, and the mean.

The **mode** of this dataset is 71 inches with frequency 5 students. The mode is easily found from the frequency table. If there is one clear mode in a distribution, the dataset is said to be *unimodal*. A data set can have more than one mode, or be *multi-modal*.

The **median** of this dataset is 69 inches with frequency of 1 students. If the number n of data points is odd, this is a simple observation of $(n + 1)/2$. If the number of data points is even, the arithmetic average of the nearest two data point values is the median.

The **mean**. The mean is the center that we will use to further examine the “spread” of the values.

$$\bar{x} = \sum_{i=0}^n x_i$$

In our data set, the mean height is 68.17857 inches, which we round to 68.2

The **Mode** of a data set is value that has the highest frequency

The **Median** of a data set is the midpoint of the distribution, or the middle value of the data when sorted in ascending order. The median is the 50th percentile

The **Mean** of a data set refers to the arithmetic mean of the values, denoted \bar{x}

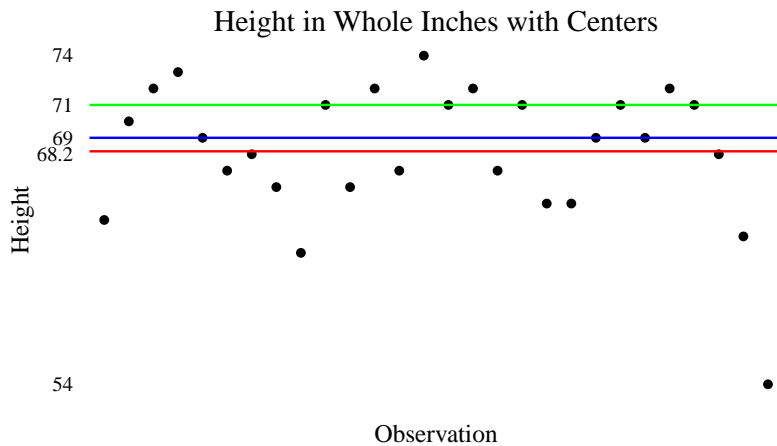


Figure 1: Heights in Observation order with Mode (Green), Median (Blue) and Mean (Red) lines

Distribution Shapes: Histograms, Frequency Polygons, and Ogives

Once we have calculated the frequencies and centers of our datasets we can start to explore the shape and spread of the distribution of values with charts. All charts can be drawn from the frequency table data.

The **histogram** is a view of the overall pattern of the distribution. Histogram bars are evenly sized and each bar represents the same class levels of values, and is centered on the mean of the class. The height of the bar represents the number of observations in that class.

The mode can easily be seen on a histogram, and the median is the vertical line at which there is equal area to the left and to the right in the chart.

A histogram's shape can be symmetric, skewed right with more of the observations on the right or higher values, or skewed to the left with more of the observations on the left or lower values.

A frequency polygon simply displays the frequency for a class, and an ogive, or cumulative frequency polygon displays the cumulative frequency for a class.

A **Histogram** is a visualization of a frequency table.

Measures of Spread

The "spread" of the distribution of a dataset can be quantified by range, first and third quartiles, variance and standard deviation.

The first and third quartiles can be found by examining either the sorted values or the frequency table, and taking the value of the observation at the first quarter and last quarter. For n observations, the

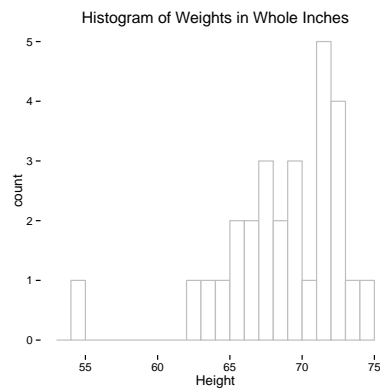


Figure 2: Histograms with Frequency Polygon and Ogive (Cumulative Frequency Polygon). The Height data set is unimodal, skewed right, with out outlier on the left.

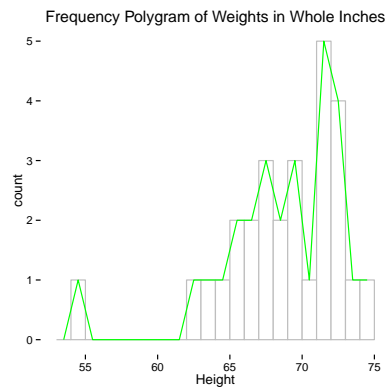


Figure 3: Histograms with Frequency Polygon and Ogive (Cumulative Frequency Polygon). The Height data set is unimodal, skewed right, with out outlier on the left.

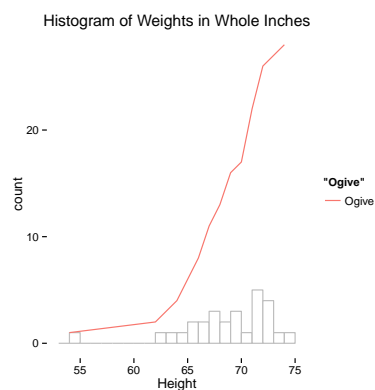


Figure 4: Histograms with Ogive (Cumulative Frequency Polygon).

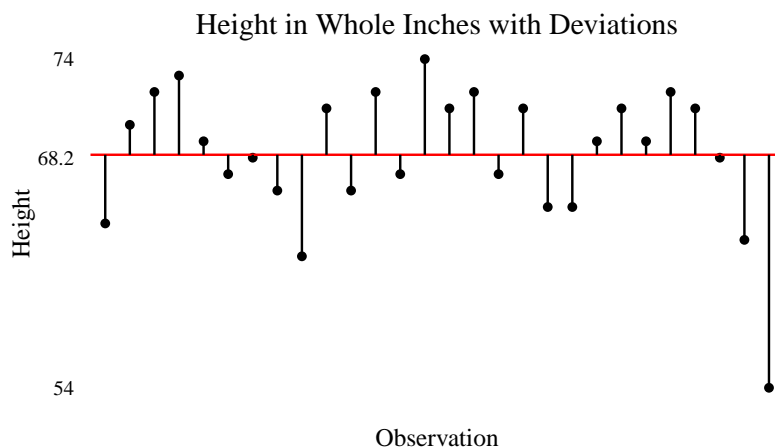
first quartile is the value $(n + 1)/4$ th entry and the third quartile is the value at the $(n + 1) * 3/4$ th entry, and similar to median, the second quartile, if the calculated entry value is not a whole number, the arithmetic mean of the nearest two observation values determines the mean.

In our height dataset of 28, the first quartile is the 7th observation, 66, and the third quartile is the 21th observation, 71.

```
54 62 63 64 65 65 66
66 67 67 67 68 68 69
69 69 70 71 71 71 71
71 72 72 72 72 73 74
```

We would like to measure the **Deviation** from the mean. The deviations from the mean are both positive and negative.

The **Deviation** is the amount by which a single measurement differs from a fixed value, such as the mean.



The deviations are both positive and negative, and the sum of the deviations is zero, so this statistic alone is not suitable for further analysis. If we square the deviations, the sum is no longer zero; in fact, the sum of the squared deviations is the **Variance**. The variance of our dataset is 478.10714 square inches. To get back to our original unit of inches, we take the square root of the variance, 21.86566, or **Standard Deviation**, denoted s . Variance is often denoted s^2 .

$$Var(x) = \sum_{i=1}^N (x_i - \bar{x})^2$$

$$s(x) = \sqrt{Var(x)}$$

The standard deviation, mean and quartiles are used to create a **boxplot**.

The **Variance** is a measure of variability or spread

The **Standard Deviation**, denoted s , is the standard measure of spread used in statistical analysis.

Height	Deviation	DeviationSq
64	-4.18	17.46
70	1.82	3.32
72	3.82	14.60
73	4.82	23.25
69	0.82	0.67
67	-1.18	1.39
68	-0.18	0.03
66	-2.18	4.75
62	-6.18	38.17
71	2.82	7.96
66	-2.18	4.75
72	3.82	14.60
67	-1.18	1.39
74	5.82	33.89
71	2.82	7.96
72	3.82	14.60
67	-1.18	1.39
71	2.82	7.96
65	-3.18	10.10
65	-3.18	10.10
69	0.82	0.67
71	2.82	7.96
69	0.82	0.67
72	3.82	14.60
71	2.82	7.96
68	-0.18	0.03
63	-5.18	26.82
54	-14.18	201.03

Table 2: Deviations

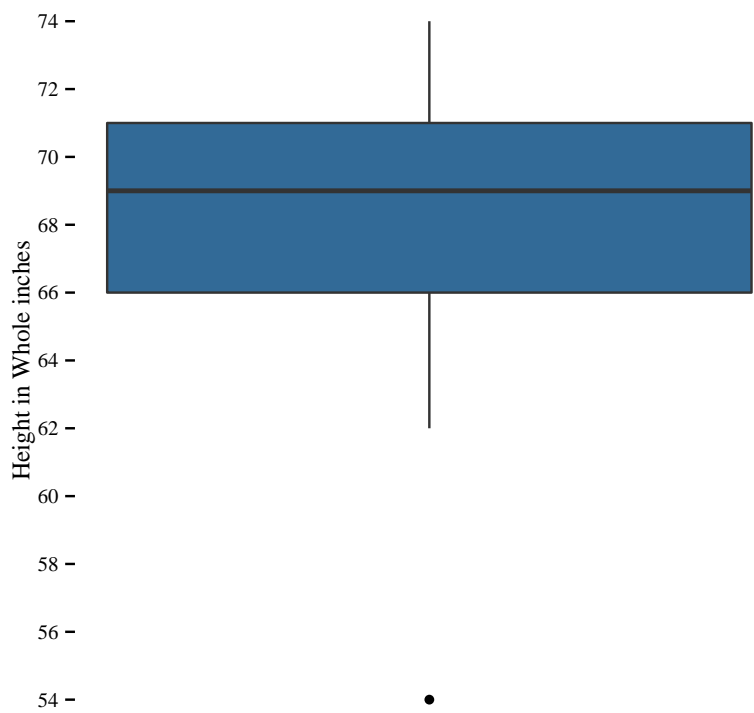


Figure 5: Boxplot