# Statistics 301: Probability and Statistics

Data displays and summary statistics

*Module 2*

*2018*

## Have data will, er. . . ?
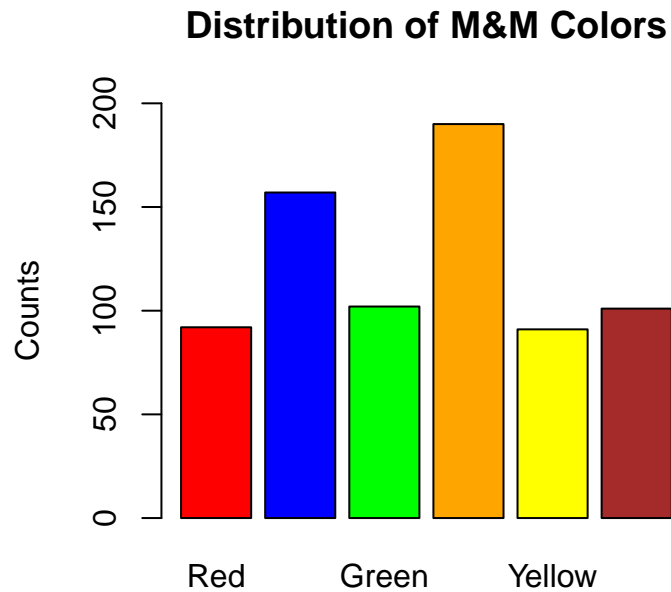
Now what to do with the data?

- Graphs visualize what the data is telling use
- see trends and other features
- knowing what the data looks like in graph form tells us what analysis to use

## Qualitative Graphs

**Bar graph** consist of bars that are separated from each other and are usually rectangular; bars represent the categories
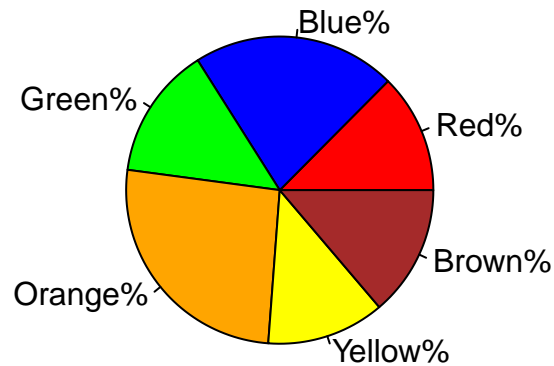
**Pie graph** shows parts of a whole, in pie form. Not very useful and easy to manipulate; not used in this class beyond an example

## Barplot



**Distribution of M&M Colors**

**Pie chart**

## Pie Chart of M&M colors

Blue%

Green%

Red%

Orange%

Brown%

Yellow%

## Quantitative graphs I

**histogram** consists of adjoining rectangles. The horizontal axis ($x$) is the ranges of data values and the vertical axis ($y$) is the frequency (counts) or relative frequency (percents or probability, denoted as $rf$).

**stemplot** (stem-and-leaf plot) comes from the field of exploratory data analysis and is good when the data sets are small. To create the plot, divide each observation of data into a stem and a leaf. The leaf consists of a final significant digit.
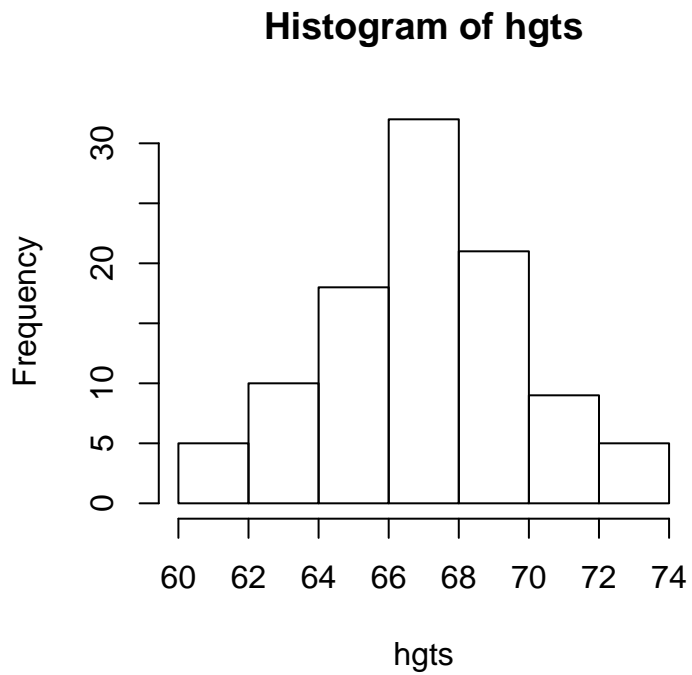
## Quantitative graphs II

**boxplot** also called box-and-whisker plots or box-whisker plots. Gives a good graphical image of the concentration of the data and also show how far the extreme values are from most of the data. A boxplot is constructed from five values: the minimum value ($min$), the first quartile ($Q1$), the median ($Median$), the third quartile ($Q3$), and the maximum value ($max$); the collection of these five summaries are called the "5 Number Summary".

**scatterplot** is a type of plot or mathematical diagram using Cartesian coordinates to display values for typically two variables for a set of data, an $x$ and $y$. Is useful for identifying trends/associations between two variables.

**time series** also called line plot. Shows the distribution of a varaible over a specified time period.

## Histogram

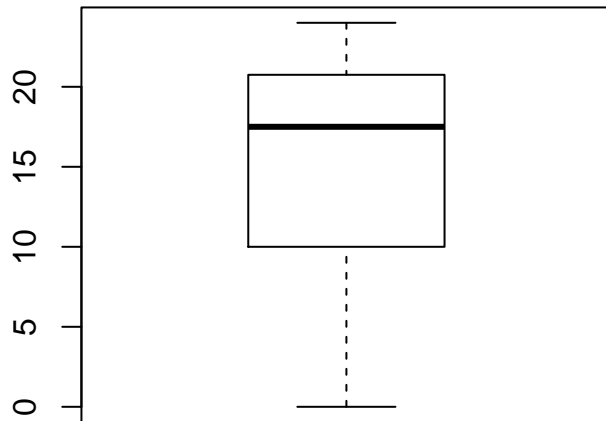**Histogram of hgts**



hgts

## Stemplot

```
The decimal point is 1 digit(s) to the right of the |

0 | 02
0 | 68
1 | 00013
1 | 557889
2 | 0011233344
```
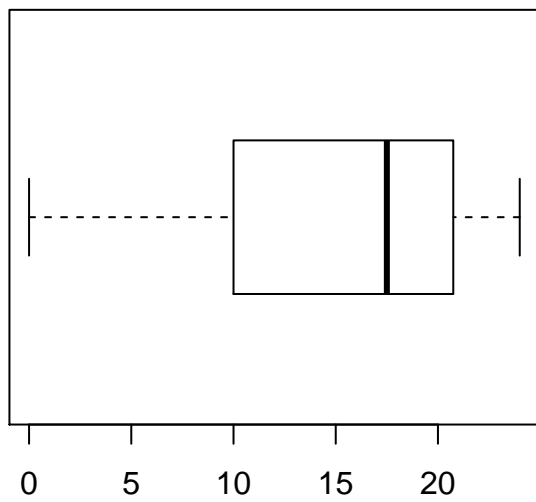
## Boxplot

Vertical boxplot
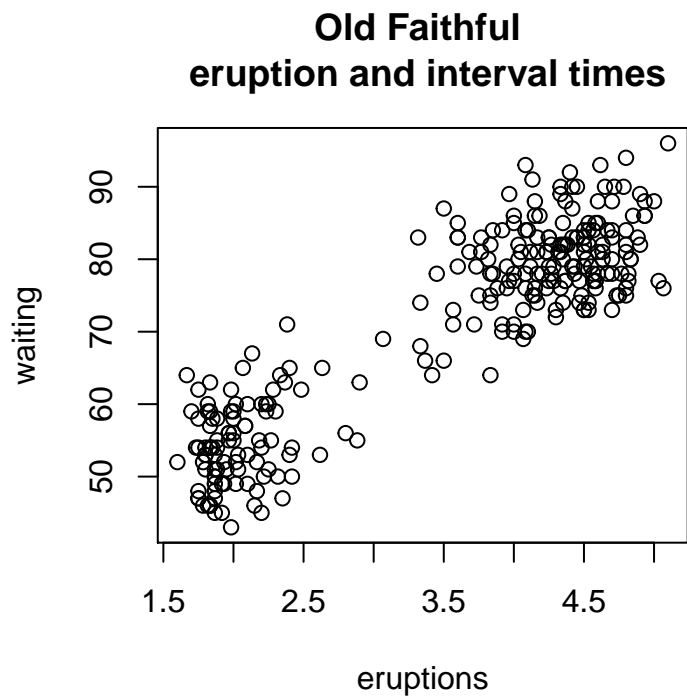
## Hours playing video games per week



Boxplot

Horizontal boxplot

## Hours playing video games per week

Scatterplot

**Old Faithful
eruption and interval times**



Time series (line plot)

**CO2 1957–1997**

Multiple time series

## CO2 1957–1997

Ukraine

4000000

1000000

2003    2005    2007    2009

year

Contour plot

## The Topography of Maunga    Height (meters)

Meters West

600

500

400

300

200

100

100    400    700

Meters North

190
180
170
160
150
140
130
120
110
100
90

## Measures of Location I

**percentiles**: divide ordered data into hundredths; useful for comparing values, especially with large populations. Ex: unemployment rates, SAT scores

**quartiles**: divide ordered data into quarters
- $Q1$: quartile 1 refers to the $25^{th}$ percentile (25% of the data is less than $Q1$ and 75% of the data is more than $Q1$)
- $M$: the *median* refers to the $50^{th}$ percentile, the center-most value of the ordered dataset (50% of the data is less than $M$ and 50% of the data is more than $M$).
- $Q3$: quartile 3 refers to the $75^{th}$ percentile (75% of the data is less than $Q3$ and 25% of the data is more than $Q3$)

## Measures of Location II

**mean**: the mathematical average. $\mu$ represents the population mean and $\overline{X}$ represents the sample mean. It is calculated by the sum of the observation values divided by the number of observations

**minimum, maximum**: the smallest ($min$) and largest ($max$), respectively value of the dataset

**mode**: the most frequently occurring observation(s); can have more than 1, can also have 0

**population size**: the number of elements in a population; denoted as $N$ (always upper case)

**sample size**: the number of observations in a dataset; denoted as $n$ (always lower case)

## Formulas I

**Median**: values would be calculated the same way for the population as for the sample. We mostly deal with samples. The data must be ordered before calculating the values

$M = \frac{n+1}{2}^{th}$ observation (if $n$ is odd); $M = avg\left(\frac{n}{2}, \frac{n}{2} + 1\right)^{th}$ observations (if $n$ is even)

$Q1 = $ median of lower half of ordered dataset
$Q3 = $ median of upper half of ordered dataset
$min$ the smallest value
$max$ the largest value
$mode$ the observation(s) that occur most frequently $N$ population size (most often not known or difficult to find)
$n$ the sample size; the number of observations in a dataset

## Examples

Dataset: {1,11,6,7,4,8,9,10,6,8,3,2,10,1} $\rightarrow$ order it: {1,1,2,3,4,6,6,7,8,8,9,10,10,11}
$n = 14$

$M = avg\left(\frac{n}{2}, \frac{n}{2} + 1\right)^{th} = avg\left(\frac{14}{2}, \frac{14}{2} + 1\right) = avg(7^{th}, 8^{th})$ observations. $M = avg(6,7) = \frac{6+7}{2} = 6.5$. Note that 6.5 is not an original observation (and it does not have to be)

$Q1$: median of the lower half of the dataset (not including $M$): {1,1,2,3,4,6,6} so $M = 3$ (the $4^{th}$ observation)
$Q3$: median of the upper half of the dataset (not including $M$): {7,8,8,9,10,10,11} so $M = 9$ (the $4^{th}$ observation)
$min = 1$, $max = 11$, $mode = 1, 6, 8, 10$ (there are multiple *modes*), $N$ is unknown/not given, $n = 14$

## The Mean

Population values (parameters) are usually denoted with letters from the Greek alphabet, while sample values (statistics) are usually denoted with English letters ($\pm$ a few extra symbols here and there); and there are, of course, a few exceptions

Sample mean: $\overline{X} = \frac{\sum x_i}{n} =$ where $x_i$ are the values of each observation in the sample

Example of sample mean: $\overline{X} = \frac{\sum x_i}{n} = \frac{1+1+\cdots+11}{14} = 6.143$

## Measures of Variability

**IQR (interquartile range)**: shows the "spread" of the middle 50% of the data; also used to help identify outliers

**outlier**: a data point that is significantly different than the other data points. Some could be due to data entry errors, some are unique and usually require more investigation

**range**: the difference between the *max* and *min* values; shows entire "spread" of the data

**variance**: the average *squared* distance each data point is from its mean

**standard deviation**: the average distance each data point is from its mean; is used most often as the main measure of variability (thus its name)

## Formulas II

$IQR = Q3 - Q1$

A value is considered a potential outlier if it is $< Q1 - IQR(1.5)$ or if $> Q3 + IQR(1.5)$

$range = max - min$

Sample variance: $s^2 = \frac{\sum (x_i - \overline{X})^2}{n-1} = \frac{(x_1 - \overline{X})^2 + (x_2 - \overline{X})^2 + \cdots + (x_n - \overline{X})^2}{n-1}$ where $n - 1 = df$, the degrees of freedom (more to come on $df$ later)

Sample variance: $s = \sqrt{\frac{\sum (x_i - \overline{X})^2}{n-1}} = \sqrt{s^2}$

*Note*: $\sigma^2$, $\sigma$, $s^2$, and $s$ MUST $\geq 0$

## More Examples

Dataset: $\{1,11,6,7,4,8,9,10,6,8,3,2,10,1\} \rightarrow$ order it: $\{1,1,2,3,4,6,6,7,8,8,9,10,10,11\}$
$n = 14$, the 5# summary is: $\{1,3,6.5,9,11\}$ $(min, Q1, M, Q3, max)$

$IQR = Q3 - Q1 = 9 - 3 = 6$; $IQR(1.5) = 6(1.5) = 9$; lower outlier boundary: $= Q1 - IQR(1.5) = 3 - 9 = -6$ and upper outlier boundary: $= Q3 + IQR(1.5) = 9 + 9 = 18$. Since no observations are outside those boundaries, there are no potential outliers in the dataset

## Continued Example

$range = max - min = 11 - 1 = 10$

$s^2 = \frac{\sum (x_i - \overline{X})^2}{n-1} = \frac{(1-6.143)^2 + (1-6.143)^2 + (2-6.143)^2 + \cdots + (11-6.143)^2}{13} = 11.824$

$$s = \sqrt{\frac{\sum (x_i - \overline{X})^2}{n-1}} = \sqrt{s^2} = \sqrt{11.824} = 3.439$$

With no idea of population size nor population values, we cannot compute the parameters (population values)

## Describing distributions I

**symmetric**: if a vertical line can be drawn at some point in the histogram such that the shape to the left and the right of the vertical line are mirror images of each other. In symmetric distributions, the mean and median are the same (approximately equal) and the mode(s) are generally in the center as well

**left (or negative) skew**: when it looks like the graph is "pulled" to the left (fewer observations to the left than the right); mode(s) generally on right side

**right (or positive) skew**: when it looks like the graph is "pulled" to the right (fewer observations to the right than the left); mode(s) generally on left side
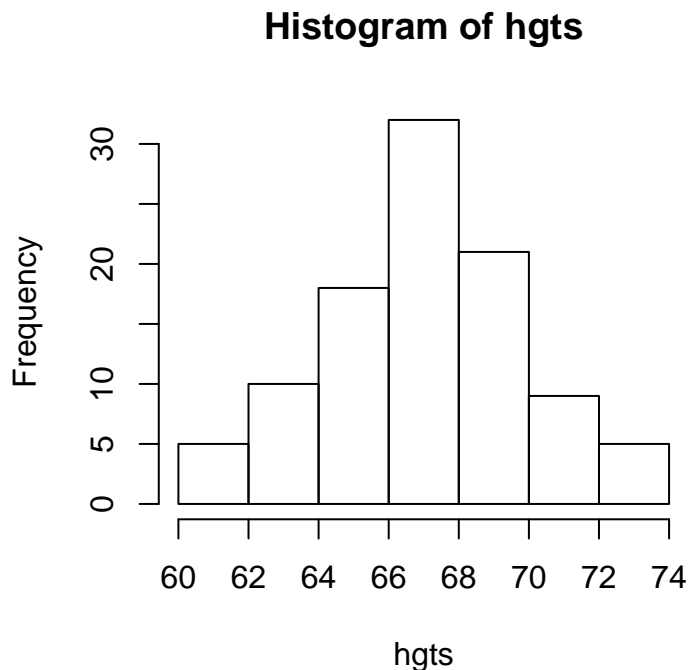
## Describing distributions II

**unimodal**: one main mode in the data set
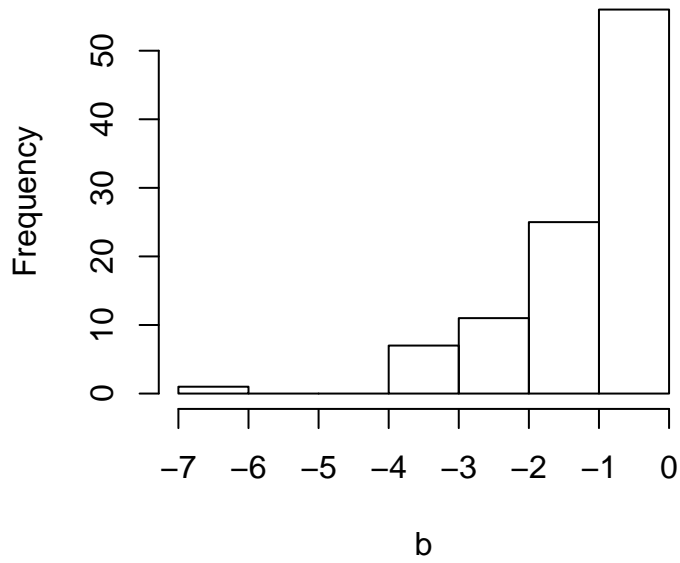
**bimodal**: two modes (more than two is multimodal)
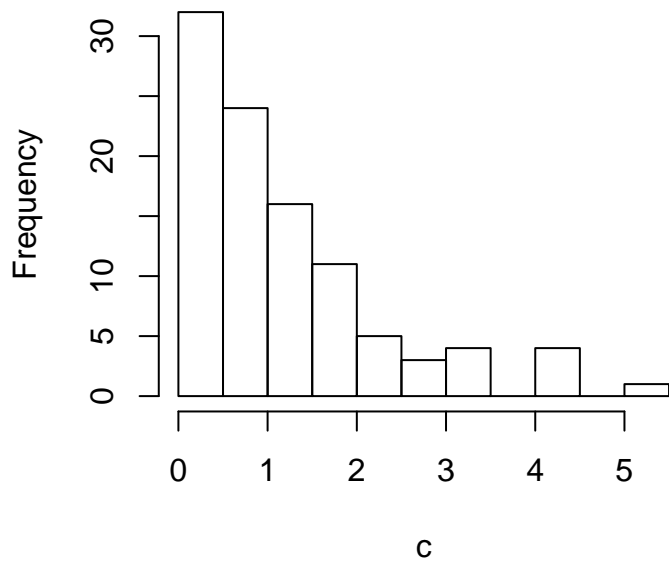
## Symmetric

```
hist(hgts)
```

**Histogram of hgts**



**Left skewed**

```
hist(b)
```

## Histogram of b



Right skewed

```
hist(c)
```

## Histogram of c

## Unimodal

```r
hist(hgts)
```

**Histogram of hgts**



## Bimodal

```r
with(faithful,hist(eruptions))
```

**Histogram of eruptions**

## Appropriate summaries for different distributions

When looking at data that has a non-symmetric distribution (skewed, bimodal, multimodal, . . . ), the best statistics to use would be the $5\#$ summary with the $IQR$, with a boxplot.

When looking at data that has a symmetric distribution (especially the "normal" distribution – bell curve), the best statistics to use would be the mean and standard deviation, with a histogram (usually – boxplots are ok too when $n$ is small).

## Empirical Rule

The Empirical Rule (ER) is appropriate *only* with symmetric distributions. ER states:

68% of observations are within the interval $\overline{X} \pm 1s$
95% of observations are within the interval $\overline{X} \pm 2s$
99.7% of observations are within the interval $\overline{X} \pm 3s$

## ER example

Example: $\overline{X} = 15$, $s = 2$
68% of observations are within the interval $\overline{X} \pm 1s = 15 \pm 2 = (13, 17)$
95% of observations are within the interval $\overline{X} \pm 2s = 15 \pm 2(2) = (11, 19)$
99.7% of observations are within the interval $\overline{X} \pm 3s = 15 \pm 3(2) = (9, 21)$

## Tchebysheff's Theorem

Tchebysheff's Theorem (TT) is appropriate with any distribution (but if symmetric, use ER). It states:

At least $1 - \frac{1}{k^2}\%$ of observations are within $\overline{X} \pm ks$

When $k = 1$, the percent is 0, so skip to $k = 2$ standard deviations
When $k = 2$, the percent is $1 - \frac{1}{2^2}\% = 75\%$ of the observations are within $\overline{X} \pm 2s$
When $k = 3$, the percent is $1 - \frac{1}{3^2}\% = 89\%$ of the observations are within $\overline{X} \pm 3s$

## TT example

Example: (same as ER example)
$k = 2 : 75\%$ of the observations are within $\overline{X} \pm 2s = 15 \pm 2(2) = (11, 19)$
$k = 3 : 89\%$ of the observations are within $\overline{X} \pm 3s = 15 \pm 3(2) = (9, 21)$

Note that TT is more conservative than ER