

* Summary Statistics (numeric data)

Statistics is the name of an academic discipline, but also a numerical value that summarizes a dataset in some way. The first statistics we will discuss are summarizing the center and spread.

* Measures of Center

1) The sample mean, average, or expected value.

Let our sample be made of data points y_1, y_2, \dots, y_n (n total data). The sample mean is denoted \bar{y} , and is:

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} (y_1 + y_2 + \dots + y_n) = \frac{y_1}{n} + \frac{y_2}{n} + \dots + \frac{y_n}{n} \\ = (\text{sum of all data})/n$$

Notice the mean uses $1/n^{\text{th}}$ of all data to create a "typical value", which means every piece of data is weighted equally.

2) Median / Percentiles / Quartiles

The k^{th} percentile of a dataset is the value for which $k\%$ of the data lies below it, and $(100-k)\%$ of the data lies above it (for k between 0 to 100).

To calculate the k^{th} percentile (denoted $y^{(k)}$):

1) Order the data from smallest to largest

2) Calculate $(\frac{k}{100})(n+1)$ [n = sample size]

2) Use "the rounding rule": If 2) resulted in a whole number, the k^{th} percentile is the ordered number in the $\lceil (\frac{k}{100})(n+1) \rceil^{\text{th}}$ location.

If 2) resulted in a decimal, round $(\frac{k}{100})(n+1)$ up and down, and average the ordered numbers in those two locations.

The median is denoted \tilde{y} , and is the 50^{th} percentile. It is the literal center of the data.

Notice the median uses locations, and is calculated

based off of a maximum of two values (unlike \bar{y})

Quartiles split the data into quarters, and

$Q_1 = 25^{\text{th}}$ percentile = first quartile

$Q_2 = 50^{\text{th}}$ percentile = second quartile = median

$Q_3 = 75^{\text{th}}$ percentile = third quartile

Ex: The survival time of patients with severe chronic heart disease (in months) was:

5, 7, 10, 15, 16, 17, 19, 20, 29, 32, 35 ($n = 11$)

a) Calculate the mean

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{11} (5+7+10+15+\dots+35) = \frac{1}{11} (205) = 18.636$$

b) Calculate the median

\tilde{y} is the 50^{th} percentile, so

1) Data is already ordered (yay)

$$2) \frac{50}{100}(n+1) = 0.50(12) = 6^{\text{th}}$$

3) 6 is a whole #, so $\tilde{y} = 6^{\text{th}}$ value = 17

c) Calculate Q_1 , Q_2 , Q_3 , and the 40^{th} percentile

$Q_1 = 25^{\text{th}}$ percentile, $0.25(n+1) = 0.25(12) = 3^{\text{rd}}$ value

so $Q_1 = 10$

$Q_2 = \tilde{y} = 17$.

$Q_3 = 75^{\text{th}}$ percentile, $0.75(n+1) = 0.75(12) = 9^{\text{th}}$ value,

so $Q_3 = 29$

for $y^{(90)}$, $0.90(12) = 10.8^{\text{th}}$ value, so average 10^{th} and 11^{th} :

$$y^{(90)} = (32 + 35)/2 = 33.5$$



Outliers

Outliers are unusually small or large observations.

There are multiple definitions, but a common one uses a "boxplots" definition.

Boxplot outlier: Any observation that is

- i) Larger than $Q_3 + 1.5(Q_3 - Q_1)$ (upper cutoff)
 ii) Smaller than $Q_1 - 1.5(Q_3 - Q_1)$ (lower cutoff)

Ex: From our previous example,

$$\text{lower cutoff} = Q_1 - 1.5(Q_3 - Q_1) = 10 - 1.5(29 - 10) = -18.5$$

$$\text{upper cutoff} = Q_3 + 1.5(Q_3 - Q_1) = 29 + 1.5(29 - 10) = 57.5$$

So there are no outliers.

* Measures of Spread (numeric data)

We may also be interested in the spread of the data - either overall or from its center.

1) The range of a data set: This is simply the maximum difference in values of the dataset.

$$\text{I.e.: range} = \max\{y_i\} - \min\{y_i\} \\ = (\text{maximum of dataset}) - (\text{minimum of dataset})$$

2) Sample Variance

The variance is the "typical squared deviation from the mean", and denoted s^2 .

A deviation from the mean is: $(y_i - \bar{y})$, so a squared deviation is $(y_i - \bar{y})^2$. The total squared deviations are $\sum_{i=1}^n (y_i - \bar{y})^2$, and the variance is

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n-1} \left[\sum_{i=1}^n y_i^2 - n\bar{y}^2 \right]$$

Fact: We don't use simply $\frac{1}{n-1} \sum (y_i - \bar{y})$ because $\sum (y_i - \bar{y}) = 0$ (this is also why \bar{y} is considered the "center" of the data).

Fact: The denominator is $(n-1)$ because then s^2 more closely estimates the overall population variance.

Most people do not interpret the variance, since its units are $(\text{units of } y)^2$

3) Sample Standard Deviation, denoted s

$$s = \sqrt{s^2} \quad (\text{deviation})$$

which is interpreted as "The typical distance of a data point to its mean"

Ex: Continuing our example:

d) Calculate s and interpret it.

$$\sum(y_i^2) = 5^2 + 7^2 + 10^2 + \dots + 35^2 = 4795$$

$$s^2 = \frac{1}{n-1} (4795 - (11)(18.636)^2) = 97.469$$

So

$$s = \sqrt{97.469} = 9.873$$

A typical deviation of survival time in months from the mean is 9.873



Graphs for numerical data

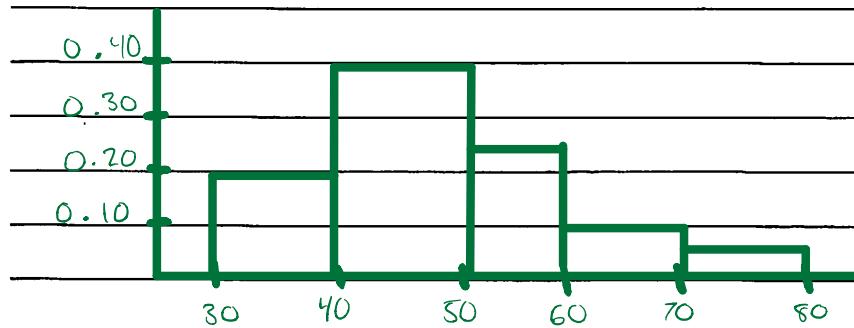
I will only go over the most common graphs, which are the ones you will learn to do in R.

Definition: The frequency of a datapoint is the number of times the data point occurred.

Definition: The relative frequency is the frequency divided by n (always between 0 and 1).

1) Histograms - Histograms take numerical data and selects groups of bins (ranges, typically equi-width). These give the widths of the boxes. For each interval, we also calculate the frequency / relative frequency of data points in the interval. This is the height of the boxes.

Ex: Weights of border collies (a type of dog) in lbs.



From here we can tell the interval that is most/least common, and the shape of the data (the distribution).

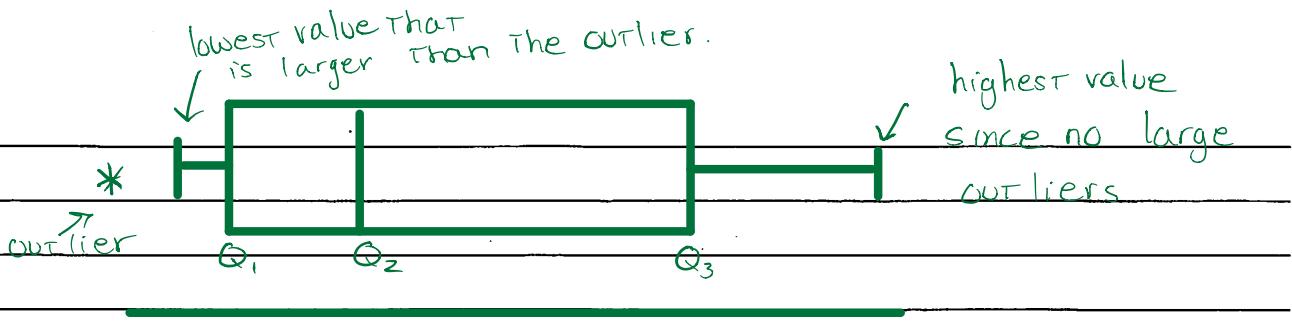
Definition: The Five Number Summary

This is simply five numbers that tell us a lot about a data set: M_{in} , Q_1 , Q_2 , Q_3 , Max

They are also found in our next graph.

2) Boxplot - We draw a box around lines at Q_1 , Q_3 , and put another line at Q_2 . Then, we identify any outliers, and draw "whiskers" from the box to the highest and lowest value that are not outliers. Lastly, we mark outliers with *'s or o's.

Ex:

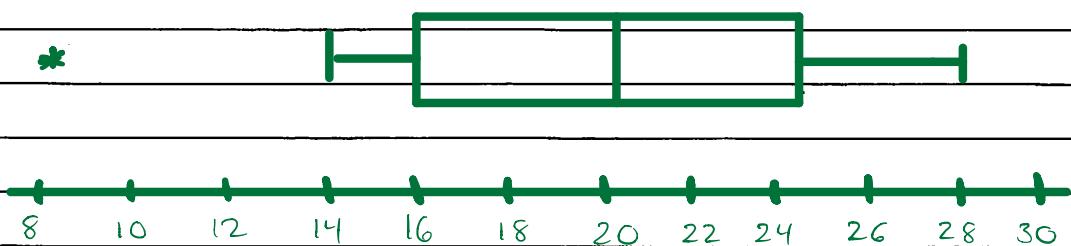


UNITS OF Y

Boxplots can also be vertical instead.

Hourly Wage in \$

Ex:



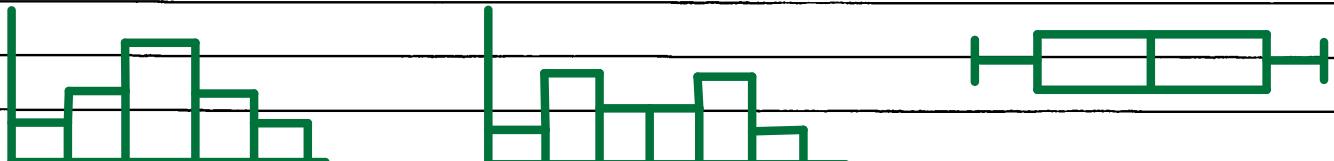
- What value has 25% of the data above it? Q₃, or ~ 24
- What is the range of the data? $(\text{max} - \text{min}) = 28 - 8 = 20$
- What values contain the middle 50% of the data (with 25% on either side)? (16, 24)

*

Distributions of Data (Shapes)

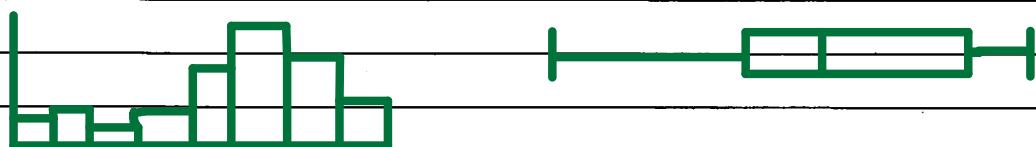
There are three main shapes of data.

1) Symmetric data: One half is a mirror image of the other.



Examples of \approx symmetric distributions: IQ, weight, height

2) Left Skewed: Most of the data is to the right, with a long left "tail".



3) Right Skewed: Most of the data is to the left, with a long right "tail".





Categorical Data

For summarizing categorical data, we really only have freq. and relative freq.

Ex: Grade: A B C D

Freq: 15 20 10 5

Rel. Freq: 0.30 .40 .20 .10
 $= \frac{15}{50}$ $\frac{20}{50}$ $\frac{10}{50}$ $\frac{5}{50}$

If we had a second categorical variable, we could also calculate relative freq. by group

Ex: A B C D | Relative Freq for M only:

(24) Male	8	10	2	4	A	B	C	D
-----------	---	----	---	---	---	---	---	---

(26) Female	7	10	8	1	$\frac{6}{24} \approx .33$	$\frac{10}{24} \approx .42$	$\frac{2}{24} \approx 0.08$	$\frac{1}{24} \approx 0.17$
-------------	---	----	---	---	----------------------------	-----------------------------	-----------------------------	-----------------------------

Relative Freq. Overall

M	$\frac{8}{50} = .16$	$\frac{10}{50} = .20$	$\frac{2}{50} = .04$	$\frac{4}{50} = .08$
---	----------------------	-----------------------	----------------------	----------------------

F	.14	.20	.16	.02	$\frac{7}{26} = 0.27$	$\frac{10}{26} = .38$	$\frac{8}{26} = .31$	$\frac{1}{26} = .04$
---	-----	-----	-----	-----	-----------------------	-----------------------	----------------------	----------------------

Relative Freq. for F only

A	B	C	D
---	---	---	---

A	B	C	D
---	---	---	---

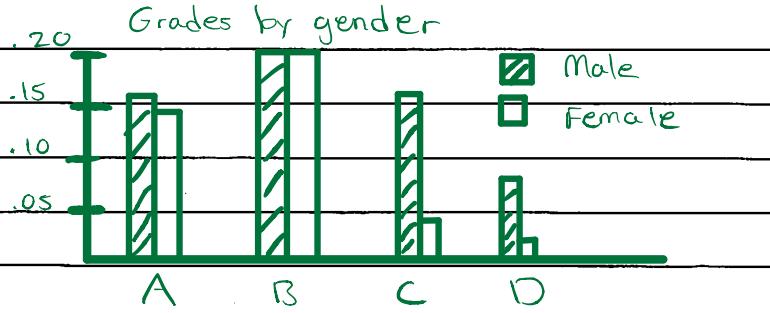
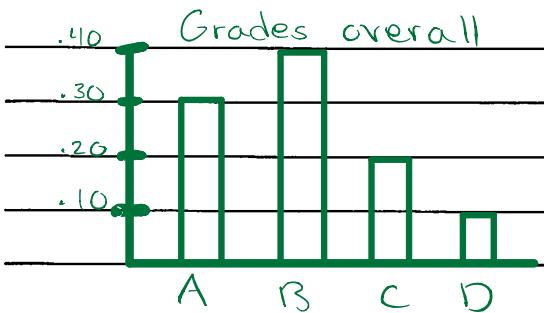


Graphs for categorical data

I will go over two main graphs, That can be used for one or two categorical variables.

1) Bar Plots - Essentially histograms for categorical data.

The main difference is the bars are separated



Notice that: i) The y-axis changed between the two graphs
 ii) The second plot uses the "overall" relative freq, not the "by group" ones.

It is relatively easy to quickly compare grades by gender, or the different categories.

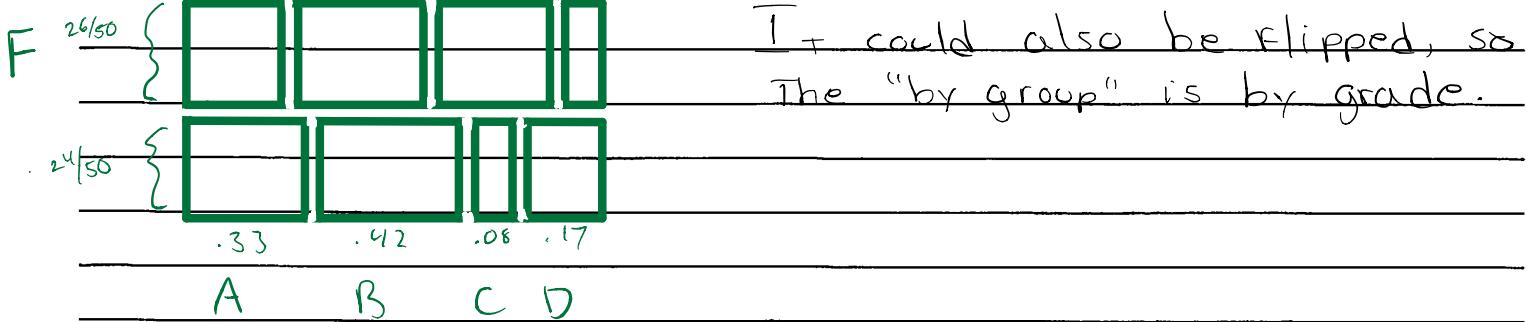
2) Mosaic Plots: A plot of boxes, where the height and width of boxes represent relative freq by group.

One variable:



The C box is $\frac{1}{2}$ of the B, and the D is $\frac{1}{4}$ of B etc.

Two variables: Two rows of the above, but the height of each row represents the proportion of either grades / gender, and boxes are based on rel freq. by group (M/F)



You do not have to draw any plots by hand, but may have to interpret them.

Now we move on to starting probability Theory. First, we should distinguish between a population value, or a sample value.

Definition: A parameter is the value of a statistic based on population data. Almost always not measurable.

Definition: A sample statistic is the value of a statistic based on sample data.

We use sample statistics to estimate population parameters.

We also use different notation for each:

Parameter	Sample Statistic
Mean	μ
Std. Dev	s
Variance	σ^2
Proportion	p
	\bar{y}
	s^2
	\hat{p}

* Probability Theory

The probability of an event is a number between 0 and 1 which represents how likely the event is to occur.

Notation: $P\{\text{A}\}$ = The prob. of an event A.

Three basic prob. facts:

- (i) $0 \leq P\{\text{A}\} \leq 1$ for any event A
- (ii) If $P\{\text{A}\} = 0$, the event A does not occur
- (iii) If $P\{\text{A}\} = 1$, The event A always occurs

There are two main ways to calculate probabilities

1) Logic Based Calculations-

This is where we assume all possible events are either equally likely, or their probabilities are known, based on a logical or natural assumption. We also assume

We know all possible outcomes in advance.

Ex: A coin. A coin has two sides, and we assume both sides are equally likely when we flip a coin. Thus,

$$\Pr\{\text{Head}\} = \Pr\{\text{Tail}\} = 1/2$$

In general, a probability is calculate by

$$\Pr\{A\} = \frac{\#\text{ of events with } A}{\#\text{ of all possible events}}$$

using the logic-based definition. This requires a lot of assumptions.

2) Relative Frequency Calculations

These are usually based on samples, and

$$(\# \text{ of times } A \text{ occurred in the sample})$$

$$\hat{\Pr}\{A\} = \frac{n}{N}$$

Notice I used $\hat{\Pr}\{A\}$. This is because $\hat{\Pr}\{A\}$ estimates $\Pr\{A\}$ (which would be the "true" probability).

We can show that as you sample more people, or as $n \rightarrow \infty$, $\hat{\Pr}\{A\} \rightarrow \Pr\{A\}$.



Probability Rules

There are several fundamental rules of probability that we will use throughout the course:

1) $0 \leq \Pr\{A\} \leq 1$

2) If we have n possible outcomes of an experiment, A_1, A_2, \dots, A_n , then $\sum_i \Pr\{A_i\} = 1$ (all probabilities add to 1)

3) The probability that A does not occur is denoted $\Pr\{A^c\}$, where $\Pr\{A^c\} = 1 - \Pr\{A\}$

④

Definition: The union of two events A and B is denoted $\{A \text{ or } B\}$, and is the event where either A occurs, B occurs, or both occur. "One or the other or both".

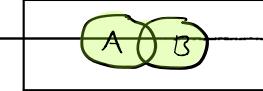
Definition: The intersection of two events A and B is denoted $\{A \text{ and } B\}$, and is the event where both

A and B occur at the same time. Note: $\Pr\{A \text{ and } B\} = \Pr\{B \text{ and } A\}$

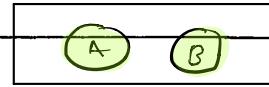
Definition: A and B are disjoint if $\Pr\{A \text{ and } B\} = 0$,
i.e. They cannot occur at the same time.

Rules continued...

4) For any two events A and B,
 $\Pr\{A \text{ or } B\} = \Pr\{A\} + \Pr\{B\} - \Pr\{A \text{ and } B\}$



5) For disjoint events A and B,
 $\Pr\{A \text{ or } B\} = \Pr\{A\} + \Pr\{B\}$



Ex: Plants were given three treatments, and their survival after one year was recorded.

	A	B	C	
Survived	13	18	10	41
Died	17	12	30	59
	30	30	40	100

Let S = plant survival. Find the following:

a) The probability a plant survives.

$$\hat{\Pr}\{S\} = (\# \text{ of plants who } S) / n = 41/100 = 0.41$$

b) The probability a plant is not in treatment A.

$$\hat{\Pr}\{A^c\} = 1 - \Pr\{A\} = 1 - \frac{30}{100} = 70/100 = 0.70$$

c) The probability they did not survive, or they had treatment B, or both.

$$\begin{aligned} \hat{\Pr}\{S^c \text{ or } B\} &= \hat{\Pr}\{S^c\} + \hat{\Pr}\{B\} - \hat{\Pr}\{S^c \text{ and } B\} \\ &= 59/100 + 30/100 - 12/100 = 77/100 = 0.77 \end{aligned}$$

We may also be interested in the prob. of survival by group. Or, "out of the A plants, what was the prob. of survival?" This is a conditional probability, because it subsets the data into only one group.

Definition: The conditional probability of A given B is the prob that A occurred, out of only the events where B occurred. It is denoted $\{A|B\}$.

Rules (continued..)

6) $P_r\{A|B\} = P_r\{A \text{ and } B\} / P_r\{B\}$ or

$$P_r\{B|A\} = P_r\{B \text{ and } A\} / P_r\{A\}$$

(Which gives (using algebra))

7) $P_r\{A \text{ and } B\} = P_r\{A|B\} P_r\{B\}$ or

$$P_r\{A \text{ and } B\} = P_r\{B|A\} P_r\{A\}$$

Ex (cont.)

d) Out of only the A plants, the prob of survival.

$$P_r\{A \text{ and } S\} = 13/100, P_r\{A\} = 30/100, \text{ so that}$$

$$P_r\{S|A\} = P_r\{S \text{ and } A\} / P_r\{A\} = (13/100) / (30/100) = 13/30 = 0.433$$

e) If a plant survived, what is the prob. it was an A?

$$P_r\{A|S\} = P_r\{A \text{ and } S\} / P_r\{S\} = (13/100) / (41/100) = 13/41 = 0.317$$

Notice that $\hat{P}_r\{A|S\} \neq \hat{P}_r\{S|A\} !!$