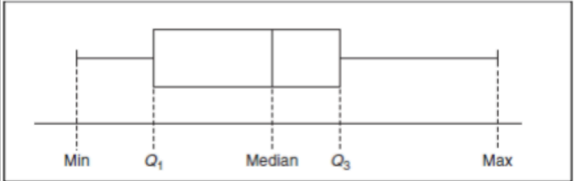


<p>Population - entire collection of objects or individuals about which information is desired</p> <p>Sample - part of the population that is selected for analysis</p> <ul style="list-style-type: none"> - Watch out for sample bias <p>Simple random sampling - every possible sample of a certain size has the same chance of being selected</p>	<p>Mean - the point of balance on a seesaw; the arithmetic average of the data</p> <ul style="list-style-type: none"> - susceptible to extreme values (outliers) by moving towards extreme values <p>Median - in an ordered vector, the median is the middle number</p> <ul style="list-style-type: none"> - not affected by extreme values <p>Quartiles - split the ranked data into 4 equal groups</p>	<p>Variance - the average distance, squared</p> $s_x^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1}$ <ul style="list-style-type: none"> - squaring gets rid of the negative value problem <p>Standard deviation - shows variation about the mean</p> $s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n - 1}}$ <ul style="list-style-type: none"> - highly affected by outliers - has same units as original data 																		
<p>Observational study - there can always be lurking variables affecting results</p> <p>Experimental study - lurking variables can be controlled; can give good evidence for causation (if internal validity is high)</p>	<p>Summary measures</p> <table border="1"> <thead> <tr> <th colspan="3">Describing data numerically</th></tr> <tr> <th>Center & location</th><th>Other Measures of location</th><th>Variation</th></tr> </thead> <tbody> <tr> <td>Mean</td><td>Quartiles</td><td>Range</td></tr> <tr> <td>Median</td><td></td><td>Interquartile Range</td></tr> <tr> <td></td><td></td><td>Variance</td></tr> <tr> <td></td><td></td><td>Standard Deviation</td></tr> </tbody> </table>	Describing data numerically			Center & location	Other Measures of location	Variation	Mean	Quartiles	Range	Median		Interquartile Range			Variance			Standard Deviation	<p>Linear transformations</p> $Y = a + bx$ <p>a: shifts data by a b: changes scale</p> <p>Linear transformations change the centre and spread of data.</p> <ul style="list-style-type: none"> - $Var(a + bX) = b^2 Var(X)$ - $Average(a + bX) = a + b[Average(X)]$
Describing data numerically																				
Center & location	Other Measures of location	Variation																		
Mean	Quartiles	Range																		
Median		Interquartile Range																		
		Variance																		
		Standard Deviation																		
<p>Box and Whisker Plot</p> 	<p>Range = $X_{\max} - X_{\min}$</p> <ul style="list-style-type: none"> - disadvantages: ignores the way in which data are distributed; sensitive to outliers <p>Interquartile range (IQR): 3rd quartile - 1st quartile</p> <ul style="list-style-type: none"> - not used often - not affected by outliers 	<p>Effects of Linear Transformations:</p> <ul style="list-style-type: none"> - $mean_{\text{new}} = a + b * mean$ - $median_{\text{new}} = a + b * median$ - $stdev_{\text{new}} = b * stdev$ - $IQR_{\text{new}} = b * IQE$ <p>z-score: new data set will have mean of 0, and variance of 1</p> $z = \frac{x - \bar{x}}{S}$																		

<p>Detecting outliers</p> <ul style="list-style-type: none"> - Classic outlier detection $z = \left \frac{x - \bar{x}}{s} \right \geq 2$ - The Boxplot Rule $X < Q1 - 1.5(Q3 - Q1)$ <p>or</p> $X > Q3 + 1.5(Q3 - Q1)$ 	<p>Skewness</p> <ul style="list-style-type: none"> - measures the degree of asymmetry exhibited by the data <ul style="list-style-type: none"> - negative values = skewed left - positive values = skewed right - if skewness < 0.8 = no need to transform the data 	<p>Measures of association</p> <ul style="list-style-type: none"> - Covariance Covariance > 0 = larger x, larger y Covariance < 0 = larger x, smaller y $cov_{xy} = \frac{1}{n-1} \sum (x - \bar{x})(y - \bar{y})$ <p>Units = Units of x · Units of y Covariance is only +, -, or 0 (can be any number)</p>
<p>Central Limit Theorem</p> <ul style="list-style-type: none"> - as n increases, \bar{x} (sample mean) should get closer to μ (population mean) - mean of $\bar{x} = \mu$ - variance of $\bar{x} = \sigma^2/n$ $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ <p>if population is normally distributed, n can be any value</p> <p>any population, n needs to be ≥ 30</p> $Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$	<p>Probability</p> <ul style="list-style-type: none"> - measure of uncertainty - all outcomes have to be exhaustive (all options possible) and mutually exhaustive (no 2 outcome can occur at the same time) <p>Rules</p> <ul style="list-style-type: none"> - Probabilities range from $0 \leq Prob(A) \leq 1$ - The probabilities of all outcomes must add up to 1 - The complement rule = A happens or A doesn't happen $P(\bar{A}) = 1 - P(A)$ $P(A) + P(\bar{A}) = 1$ - Addition rule: $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$ 	<p>Correlation - measures strength of a linear relationship between two variables</p> $r_{xy} = \frac{cov_{xy}}{(s_x)(s_y)}$ <p>correlation is between +1 and -1 sign: direction of relationship absolute value: strength of relationship (-0.6 stronger than +0.4)</p> <p>0 - 0.2: very weak 0.2 - 0.4: weak to moderate 0.4 - 0.6: medium to substantial 0.6 - 0.8: very strong 0.8 - 1.0: extremely strong</p> <p>correlation doesn't imply causation correlation of a variable with itself is 1</p>

Confidence intervals - tells us how good our estimate is - Want high confidence, narrow interval - As confidence increases, interval also increases	Margin of Error - The confidence interval is given by $\hat{p} \pm 1.96\sqrt{\frac{\hat{p}\hat{q}}{n}}$ - The standard form of any confidence interval is estimate \pm (margin of error)	Hypothesis testing - Null hypothesis - H ₀ , a statement of no change and is assumed true until evidence indicates otherwise - Alternative hypothesis: H _a , is a statement that we are trying to find evidence to support - Type I error: reject the null hypothesis when the null hypothesis is true - Type II error: do not reject the null hypothesis when the alternative hypothesis is true																					
Effect size - <i>Cohen's d</i> $d = \frac{\mu_{treatment} - \mu_{control}}{s}$ where d is Cohen's d, μ is a mean, and s is the standard deviation	Standard Error of the Mean (S.E.M.) $\sigma_M = \frac{\sigma}{\sqrt{(n)}}$ where σ_M is the S.E.M., σ is the standard deviation and n is the sample size.																						
<table><tr><th></th><th>Hypothesis testing (p-values)</th><th>Effect size</th><th>Power</th></tr><tr><th>Type of conclusion</th><td>Threshold based, yes/no</td><td>Gradient, small to large</td><td>Gradient, small to large</td></tr><tr><th>Sample size</th><td>Sensitive</td><td>Not sensitive</td><td>Sensitive</td></tr><tr><th>Alpha level</th><td>Sensitive</td><td>Not sensitive</td><td>Sensitive</td></tr><tr><th># of tails</th><td>Sensitive</td><td>Not sensitive</td><td>Sensitive</td></tr></table>				Hypothesis testing (p-values)	Effect size	Power	Type of conclusion	Threshold based, yes/no	Gradient, small to large	Gradient, small to large	Sample size	Sensitive	Not sensitive	Sensitive	Alpha level	Sensitive	Not sensitive	Sensitive	# of tails	Sensitive	Not sensitive	Sensitive	Example of Type I and Type II errors - Null hypothesis - H ₀ , a statement of no change and is assumed true until evidence indicates otherwise - Alternative hypothesis: H _a is a statement that we are trying to find evidence to support - Type I error: reject the null hypothesis when the null hypothesis is true - Type II error: do not reject the null hypothesis when the alternative hypothesis is true Methods of Hypothesis testing - Confidence intervals - Test statistics - p-values - C.I. and p-values always safe to do because don't need to worry about size of n (can be bigger or smaller than 30)
	Hypothesis testing (p-values)	Effect size	Power																				
Type of conclusion	Threshold based, yes/no	Gradient, small to large	Gradient, small to large																				
Sample size	Sensitive	Not sensitive	Sensitive																				
Alpha level	Sensitive	Not sensitive	Sensitive																				
# of tails	Sensitive	Not sensitive	Sensitive																				

Cheat Sheet

R Functions (supplement)

For more information on each function, use the question mark and the name of the function in R

Studio, e.g.:

?sum()

Check data frame type

is.integer(), is.character(),

is.logical()

as.integer(), as.character(),

as.logical()

class()

Glimpse into the data

names()

head()

summary()

Formatting the data

sort()

rle()

table()

round()

Basics

length()

sum()

Central Tendency

median()

mean()

Variability

min(), max()

range()

quantile()

IQR()

var()

cor()

sd()

Graphs

plot()

hist()

boxplot()

barplot(

 meansVector,

 ylim = c(0,0),

 names.arg = c("list","labels"),

 main = "title"

)

Arrows on graphs

arrows(

 middle, mean-sem,

 middle, mean+sem,

 length = 0.05,

 angle = 90,

 code=3

)

Multiple graph facets

par(mfrow=c(1,2))

Other

unique()

sqrt()

scale()

Subsetting Vectors

variable[variable == "value"]

variable[variable > 0]

Subsetting Data Frames

df[df\$col == 0,]

df[df\$col > 0,]

Probability

pnorm()

Programming Concepts

for(i in 1:10) {

 your code goes here

}

myfunc <- function(variable) {

 your code goes here

}

if(condition == TRUE) {

 do something

}