

Resident Research Course:

Study Design and Analysis – Types of Data

Lecture 1

Dr. Grace Kwong, PhD, P.Stat.

grace.kwong@ucalgary.ca





- Statistics, study design, and critical thinking about bias are important
- Types of sampling methods
- Identify what types of data you are (or want to be) working with and describe it appropriately
 - Graphical presentation
- Statistical analysis plan





- Statistics, study design, and critical thinking about bias are important
- Types of sampling methods
- Identify what types of data you are (or want to be)
 working with and describe it appropriately
 - Graphical presentation
- Statistical analysis plan





- Statistics, study design, and critical thinking about bias are important
- Types of sampling methods
- Identify what types of data you are (or want to be)
 working with and describe it appropriately
 - Graphical presentation
- Statistical analysis plan





- Statistics, study design, and critical thinking about bias are important
- Types of sampling methods
- Identify what types of data you are (or want to be)
 working with and describe it appropriately
 - Graphical presentation
- Statistical analysis plan
 --> have one before you start





- Statistics, study design, and critical thinking about bias are important
- Types of sampling methods
- Identify what types of data you are (or want to be)
 working with and describe it appropriately
 - Graphical presentation
- Statistical analysis plan
 --> have one before you start





Why? Evaluate evidence from research

Example 1:

- This week, you are working in a pediatric ED and you observe that all children who present with COVID-19 are discharged rapidly (mild infection). You talk about it with some other residents working this week and find that they have observed the same thing.
- You compile the data from these children and write a manuscript you conclude that COVID-19 is mild in children.
- The news media contacts you, and you state in your interview that children across Canada should be in school, regardless of transmission trends

What type of study is this? What do you think of the conclusions?





Why? Evaluate evidence from research

Example 2:

- You'd like to examine maternal risk factors for a rare type of congenital neurological anomaly that affects developmental outcomes.
- You enroll all mothers of infants born with this congenital anomaly over the course of one year across all hospitals in Alberta. For your control group, you advertise for volunteer participants (mothers) with healthy births during the same time period at the same hospitals.
- You compare the two groups and find that mothers of children with the congenital anomaly have a lower educational level and were less likely to be 100% compliant to prenatal vitamins.
- You conclude that genetics (IQ-link) and prenatal nutritional deficiencies may interact to cause this congenital anomaly.

What kind of study is this? What do you think of the conclusions?





Why Statistics?

Statistics, following a robust study design, allows one to address questions such as:

Are the effects observed due to chance?

How accurate are these estimates?

Are these conclusions valid?

Considering the # of participants, and how they were enrolled, are these findings trustworthy??





Why Statistics?

Statistics, following a robust study design, allows one to address questions such as:

...p values and results of other significance tests

Are the effects observed due to chance?

... critique statistical analysis plan and study design

Are these conclusions valid?

How accurate are these estimates?

...confidence intervals and other measures



Considering the # of participants, and how they were enrolled, are these findings trustworthy??

... sample selection, size calculation





Why statistics and study design?

- Statistical analysis (correctly done) is a <u>vital</u> part of good research, and is necessary to draw conclusions from any data you've collected
- Your study needs to be well designed and a statistical analysis plan should be worked into it from the beginning

A well-designed study, poorly analyzed, can be rescued by re-analysis, whereas a poorly designed study is beyond redemption, even using sophisticated statistics





- Statistics, study design, and critical thinking about bias are important
- Types of sampling methods
- Identify what types of data you are (or want to be)
 working with and describe it appropriately
 - Graphical presentation
- Statistical analysis plan
 --> have one before you start



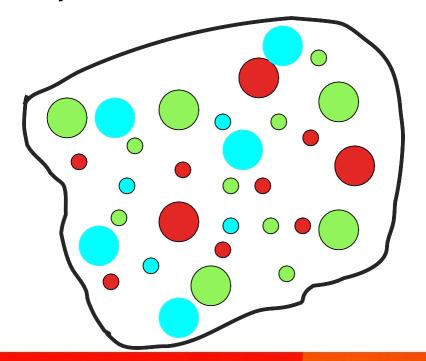


Population vs. Sample

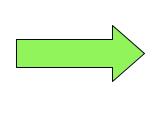
Statistics allow us to draw inferences about a body of data (population) when only a part of the data (sample) is observed

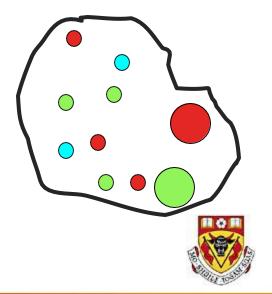
<u>Population</u> – *all* individuals or items under consideration pertaining to a study question <u>Sample</u> – a selected part of the population from which information is obtained.

Everyone in the world with liver cancer



Persons with liver cancer presenting for treatment at one specialized centre in Cairo



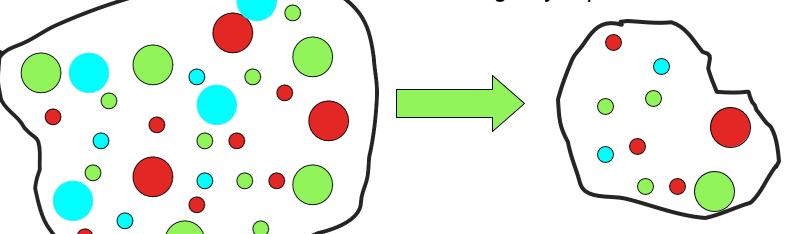




Population vs. Sample

All children with COVID-19

Children with COVID-19 presenting to 50 emergency departments across 14 countries



How might this sample differ from the population?

How representative is this sample?

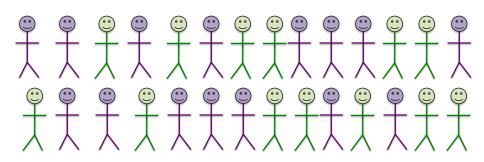




Probabilistic (random) designs

- Simple random sampling
- Systematic sampling
- Cluster sampling
- Stratified sampling

All children with a rare genetic disorder in Alberta





Demographic characteristic 1 (e.g. male or low socioeconomic status)

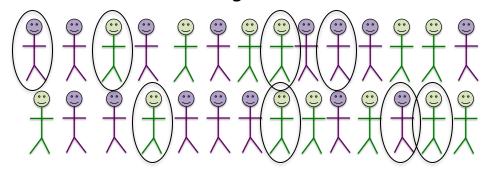






- Probabilistic (random) designs
 - Simple random sampling
 - Systematic sampling
 - Cluster sampling
 - Stratified sampling

All children with a rare genetic disorder in Alberta





Demographic characteristic 1 (e.g. male or low socioeconomic status)

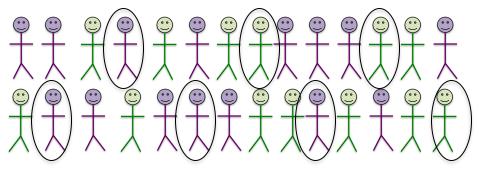






- Probabilistic (random) designs
 - Simple random sampling
 - Systematic sampling
 - Cluster sampling
 - Stratified sampling

All children with a rare genetic disorder in Alberta



(Every fourth child registered)



Demographic characteristic 1 (e.g. male or low socioeconomic status)

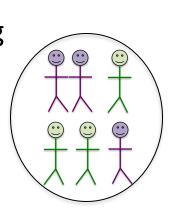


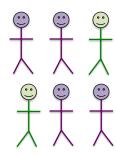


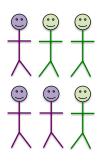


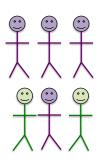
All children with a rare genetic disorder in Alberta

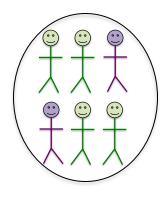
- Probabilistic (random) designs
 - Simple random sampling
 - Systematic sampling
 - Cluster sampling
 - Stratified sampling











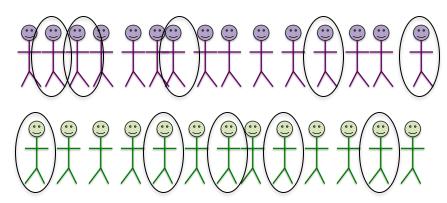




Probabilistic (random) designs

- Simple random sampling
- Systematic sampling
- Cluster sampling
- Stratified sampling

All children with a rare genetic disorder in Alberta





Demographic characteristic 1 (e.g. male or low socioeconomic status)









Convenience sample

Volunteer

Purposive

Snowball sampling

attended clinic A in 2020

Kids who

Kids who's parent saw the online study notice and wanted their child to take part

Kids of lower socio-economic status families

Participants and their friends / acquaintances





Sampling - takeaway

Why should you know sampling strategies?

- For you (medical researcher)
 - Choose appropriate one
 - How representative is your study population?
 - What kind of statistics / analysis will you need?
 - How strong can you state your conclusions
- For knowledge users and stakeholders
 - How seriously should we take these findings?



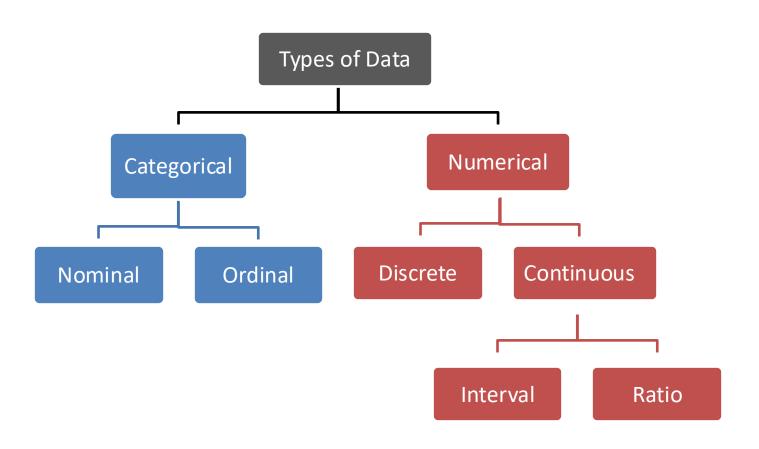


- Statistics, study design, and critical thinking about bias are important
- Types of sampling methods
- Identify what types of data you are (or want to be) working with and describe it appropriately
 - Graphical presentation
- Statistical analysis plan
 --> have one before you start





Types of Data Overview

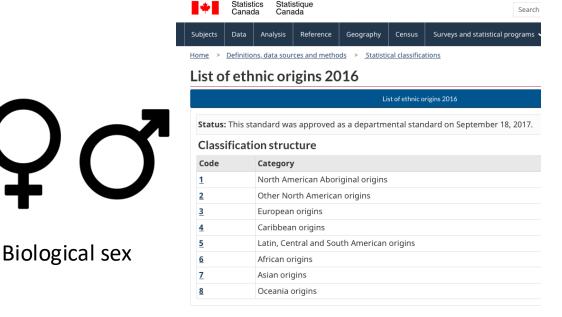






Categorical Data

- Nominal (names) variables:
 - Labels based on an attribute of each element
 - Two or more categories, no ordering



Disease classification (when not related to severity):

E.g. Types of respiratory illness

- Pneumonia
- Asthma
- Emphysema

E.g. location (segment) of liver tumour

E.g. first symptom of COVID-19 experienced

- fever, chills, headache, cough





Categorical Data

No

Ordinal (ordered) variables:

- Have properties of nominal data and can be used to rank or order the observations
- Two or more categories, ordered in magnitude
- Difference between successive not necessarily the same

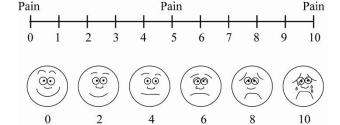
Visual analog scale

Disease classification (related to severity):

E.g. Acute respiratory distress syndrome (ARDS)

- Mild
- Moderate
- Severe





Moderate

Federal tax rates for 2020

- 15% on the first \$48,535 of taxable income, plus
- 20.5% **on the next** \$48,534 of taxable income (on the portion of taxable income over 48,535 up to \$97,069), **plus**
- 26% **on the next** \$53,404 of taxable income (on the portion of taxable income over \$97,069 up to \$150,473), **plus**
- 29% **on the next** \$63,895 of taxable income (on the portion of taxable income over 150,473 up to \$214,368), **plus**
- 33% of taxable income **over** \$214,368

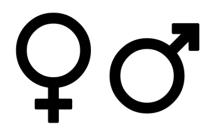


Worst



Categorical Data

- Special case dichotomous (binary) variable
 - Only two possible options
 - May indicate an order but not necessarily



Biological sex (?)

old - young



Presence of fever? Yes



No

Note: Yes / No often turned into a 'dummy' variable (0/1) in analysis



Displaying categorical data

Tables / Summaries

- Frequencies (number of observations that fit in that category) with percentages (proportions)
 - DO NOT CALL IT A RATE





Example: Cardiac arrest

In a study, there are 747 males and 434 females

• 276 males and 195 females had a cardiac arrest

Sex	Frequency	%
Male	747	63.25
Female	434	36.75
Total	Total 1181	

Cardiac arrest	Frequency	%
Yes	471	39.88
No	710	60.12
Total	1181	100.00

Tabulation





Example: Cardiac arrest

N or n Row % Col %	Cardiac arrest	No cardiac arrest	Total
Male	276	471	747
	36.95	63.05	100.00
	58.60	66.34	63.25
Female	195	239	434
	44.93	55.07	100.00
	41.40	33.66	36.75
Total	471	710	1181
	39.88	60.12	100.00
	100.00	100.00	100.00

Cross Tabulation





Displaying categorical data

95% Confidence Intervals (95% CI)

$$\hat{\rho} \pm z \star \sqrt{\frac{\hat{\rho}(1-\hat{\rho})}{n}}$$

 $\hat{\mathbf{p}}$ = proportion \mathbf{z} = value from a standard normal distribution to represent the confidence level you desire (95% = 1.96) \mathbf{n} = sample size

E.g. In the cardiac arrest example, 276 (36.95%) of the 747 male had a cardiac arrest.

What are the lower and upper 95% CIs?

Altogether what does this confidence interval 'range' mean?





Displaying categorical data

95% Confidence Intervals (95% CI)

E.g. In the cardiac arrest example, 276 (36.95%) of the 747 male had a cardiac arrest.

What are the lower and upper 95% CIs?

Altogether what does this confidence interval 'range' mean?





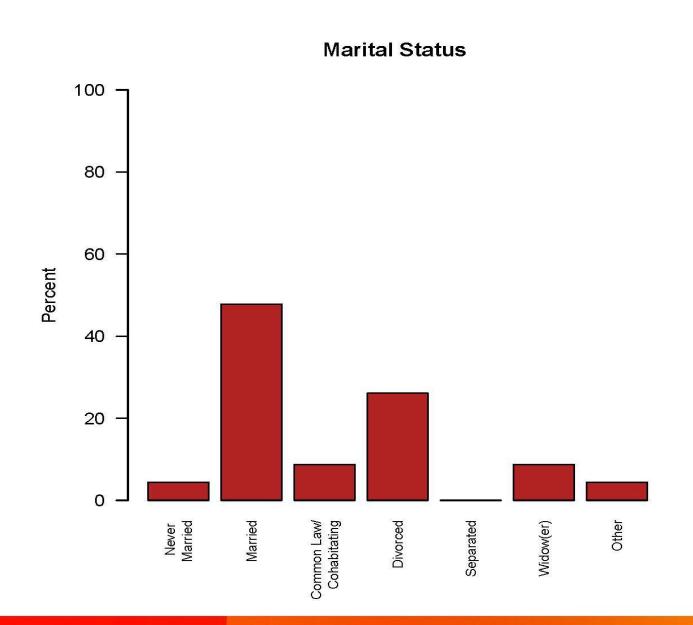
Example: Cardiac arrest

N or n Row % 95 % CI	Cardiac arrest	No cardiac arrest	Total
Male	276	471	747
	36.95	63.05	
	33.5 – 40.4	59.4 – 66.5	
Female	195	239	434
	44.93	55.07	
	40.2 – 49.7	50.3 – 59.8	
Male or Female	471	710	1181
(prev Total)	39.88	60.12	
	37.1 – 42.7	57.3 – 62.9	





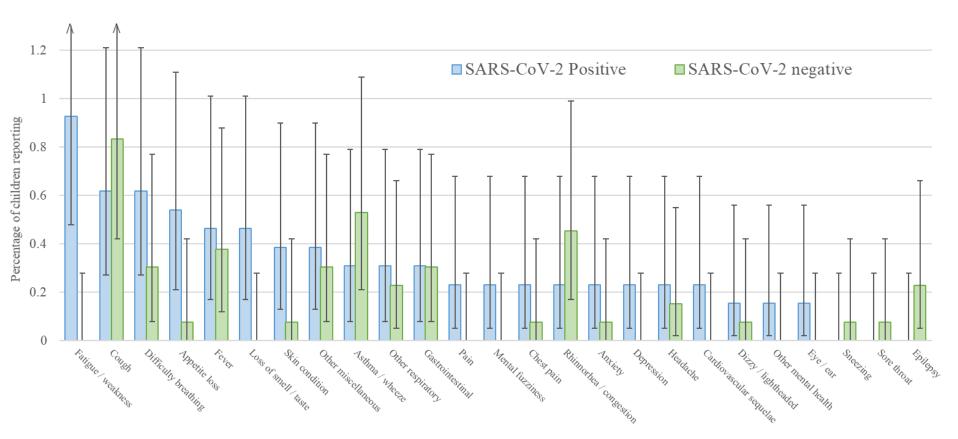
Displaying categorical data: Bar charts







Displaying categorical data: Bar charts



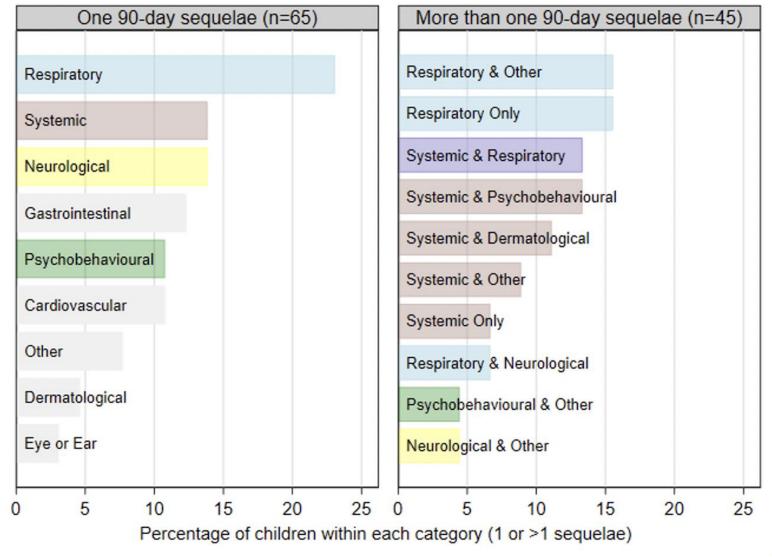
Persistent symptoms in SARS-CoV-2 positive (N=1295) and SARS-CoV-2 negative (N=1321) children **not** hospitalized during the acute phase of SARS-CoV-2

Funk AL, Kuppermann N, Florin TA, Tancredi DJ, Xie J, Kim K, Finkelstein Y, Neuman MI, Salvadori MI, Yock-Corrales A, Breslin KA. Post–COVID-19 conditions among children 90 days after SARS-CoV-2 infection. JAMA network open. 2022 Jul 1;5(7):e2223253-.





Displaying categorical data: Bar charts





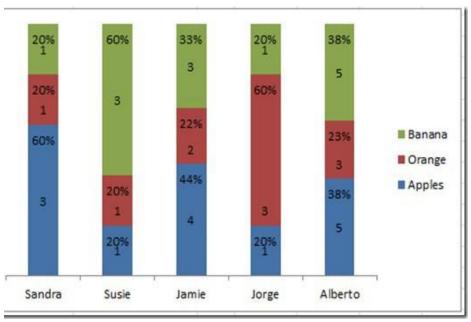
Continued... persistent symptoms in SARS-CoV-2 positive children

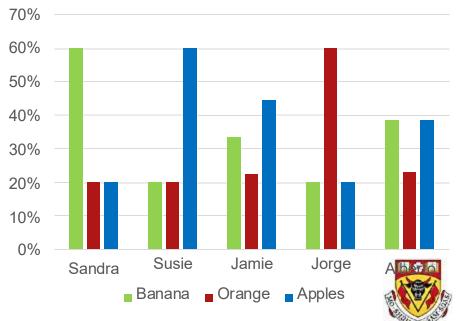




Don't overuse stacked bar graphs

- Use a graph to make comparison easy → not repeat all numbers (that's what tables are for!)
- In stacked bars you can not easily tell different heights of different bars

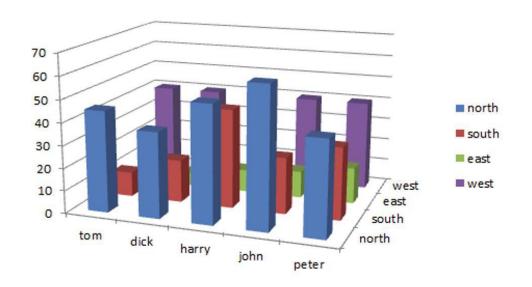


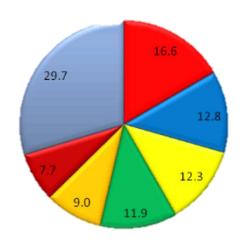




What not to do!

- Don't add "dimensions" that are not real dimensions and which distort the graph → keep it simple & clear
- Use pie charts: slices are not easily comparable and numbers are often added... Hence it is just a colored disarrayed table!









Numerical Data

Discrete

- Can only take on certain numerical values
- E.g. Counts of events
 - Number of children hospitalized with COVID-19
 - Number of concussions in an adolescent football player
 - Number of children per household in Southern Alberta
- Discrete variables that would not be a 'count'?





Numerical Data

Continuous

Interval data

- Ordered, difference between any two values is meaningful
- Can theoretically be any value within a range
- No true 'zero' measurement that indicates the lack of characteristic

• Examples:

- Date (with/without time) = 19/10/2020 09:04:56.864
- Temperature in Celsius/Fahrenheit = -14.8 Celsius
- pH level = 7.4





Numerical Data

Continuous

Ratio data

- Interval data <u>BUT</u> a 'true zero' does exist
 - Because of this, a ratio of two data points is meaningful
 - 10 kg is twice as heavy as 5 kg (ratio)
 - 10 °C is **not** twice as hot as 5 °C (interval)

Examples:

- Viral load in copies/ml
- Height in cm
- Temperature in Kelvin (with its absolute zero that represents no temperature)
- Time to an event (e.g. survival after transplant)



Numerical → Categorical

Interval and ratio data can be converted into categorical data

- Examples:
 - Precise age → < 1 year, 1-<3 years, 3-<5 years
 - Precise height and weight → specific BMI value → BMI category
- Lose some precision / detail, however, is a 1 unit increase in your numerical data meaningful?





Describing numerical data - overview

- Data distributions and their features
 - Probability distributions
 - Measures of central tendency mean, median, mode
 - Measures of dispersion
 - Skewness

Graphical representation of numerical data





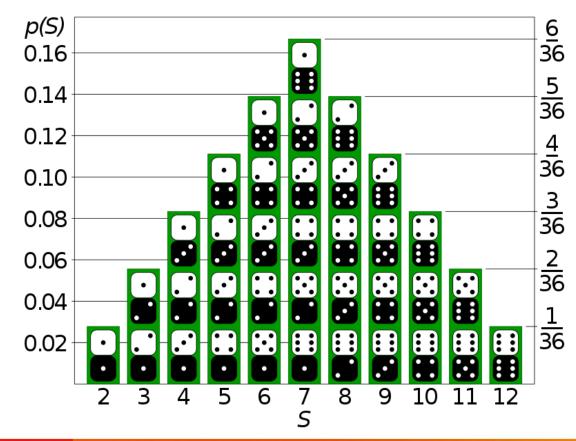
Probability Distributions

 In 'probability distributions', each possible event gets/has a certain amount of probability (chance) of happening.

Probability mass function →

Retrieved at:

https://commons.wik imedia.org/wiki/File: Dice_Distribution_(b ar).svg





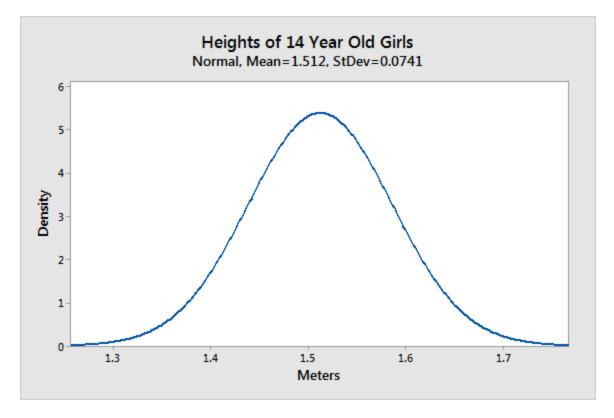


Probability Distributions

 In 'probability distributions', each possible event gets/has a certain amount of probability (chance) of happening.

Normal (Gaussian) distribution →

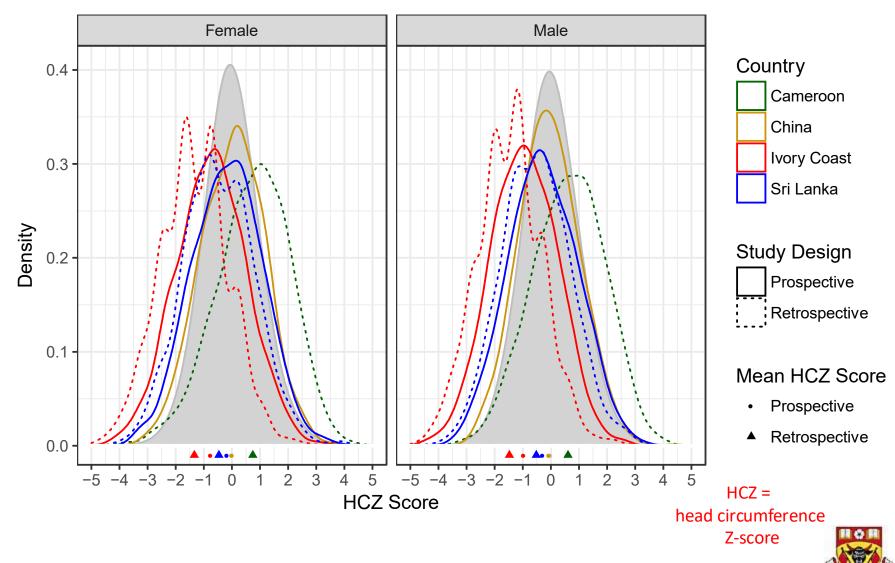
Retrieved at: https://statisticsbyji m.com/basics/norma l-distribution/







Probability Distributions



Zika-virus related microcephaly surveillance across 4 countries



Central tendency or location parameters

- The value that observations are "centered" around
- Mean average value
 - Sum of values divided by the number in the sample
- Median value at which 50% of the observations have accumulated
 - To calculate: rank data, and find the value for which 50% of observations fall above and 50% below, or use rank formula
- Mode most frequently <u>observed</u> value







Weight (in kg) of 10 five year old boys

Data: 14, 14, 15, 17, 17, 18, 18, 18, 19, 20





Weight (in kg) of 10 five year old boys

- Data: 14, 14, 15, 17, 17, 18, 18, 18, 19, 20
- Mean: (14+14+15+17+17+18+18+18+19+20) /10 = 17
- Median (via rank formula):
 - Index i = (pth percentile*# observations)/100
 - ❖ 50th percentile*10 observations/100 = 5
 - If i is not an integer (e.g. 3.5, 46.5) the next value is pth percentile
 - If i is an integer (e.g. 4, 40), average ith and ith +1
 - Our i is 5 (average 5th and 6th values), median = 17+18 / 2 = 17.5
- Mode: most frequent value
 - '18' appears the most (3 times), mode = 18





Example continued...

Weight (in kg) of 10 five year old boys

Data: 14, 14, 15, 17, 17, 18, 18, 18, 32, 40

- Median (via rank formula):
- Mode: most frequent value
 - ❖ '18' appears the most (3 times), mode = 18

Changed!

Stayed the same!





Using Measures of Central Tendency

Mean

- When data aren't skewed
 - If symmetric distribution, mean ~ median
- Lack of outliers / extreme values
 - Highly influenced by extreme values

Median

- When data are skewed or when there are outliers
 - Robust to extreme values

Mode

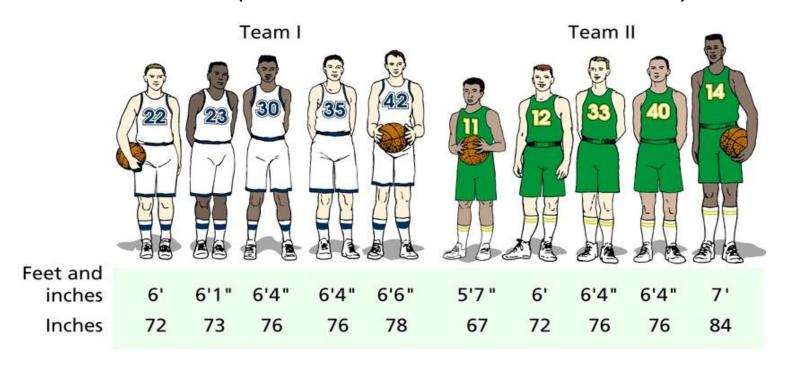
- Use with nominal data (e.g. most common ethnicity)
- Use to check for bimodal distribution





Dispersion

 How spread out the observations are in the distribution (sometimes around the "centre").

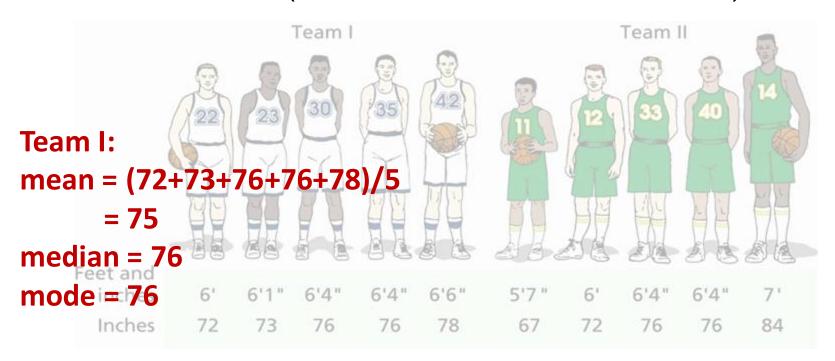






Dispersion

 How spread out the observations are in the distribution (sometimes around the "centre").

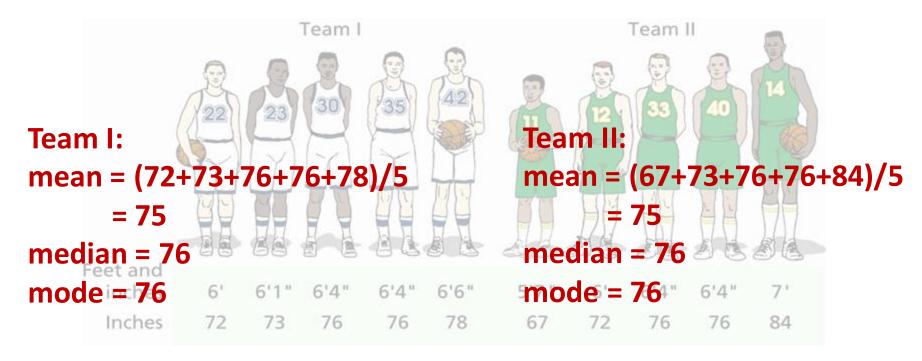






Dispersion

 How spread out the observations are in the distribution (sometimes around the "centre").

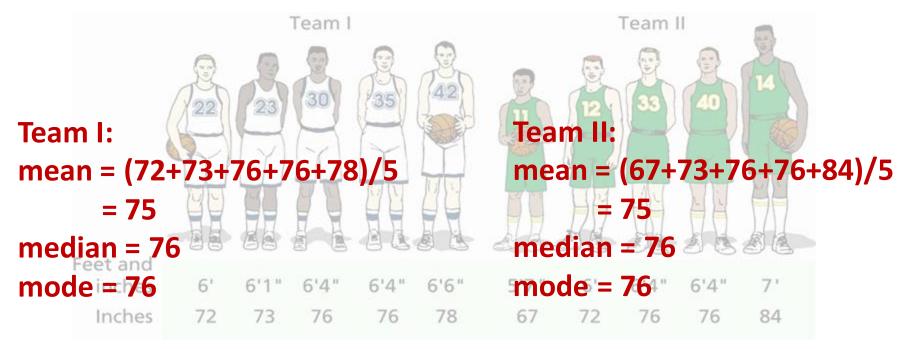






Dispersion

 How spread out the observations are in the distribution (sometimes around the "centre").



The "data sets" have the same Mean, Median, and Mode yet clearly differ!





Measures of dispersion

Range

Difference between largest and smallest values

Interquartile range → IQR

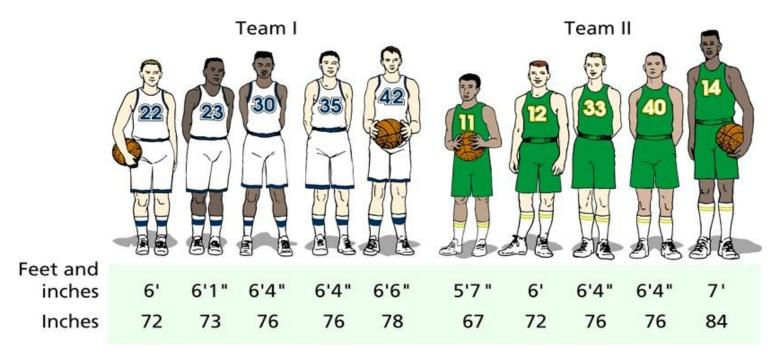
- ❖ Difference between the third quartile (75th percentile) and first quartile (25th percentile)
- Central 50% of the data

■ Standard deviation → S or SD

- ❖ Square root of the sample variance (S²⁾
 - (Sum of squared deviations from mean) / n-1
- Most commonly used, measures spread around the mean ... & highly affected by mean







Range: 72-78

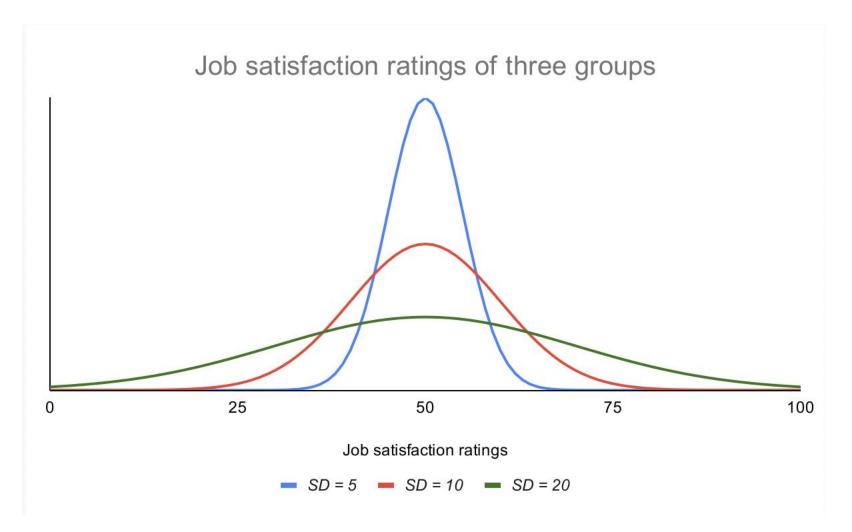
IQR: 73 – 76

Range: 67-84

IQR: 72 – 76







Largers vs smaller standard deviations





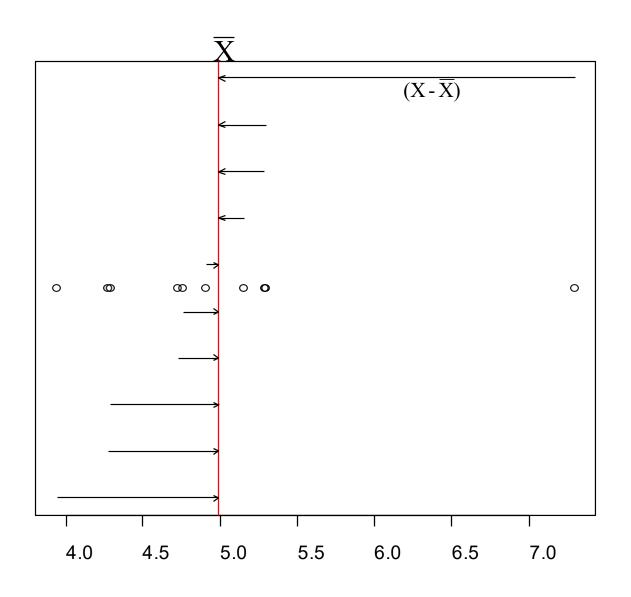
Calculating the variance and standard deviation (X = total cholesterol)

Subject (i)	X
1	3.94
2	4.73
3	5.30
4	5.16
5	4.91
6	5.29
7	4.29
8	7.30
9	4.76
10	4.27





Calculating the variance and standard deviation (X = total cholesterol)







Calculating the variance and standard deviation (X = total cholesterol)

i	X	X	(X -X)	$(X - X)^2$
1	3.94	4.99	-1.05	$1.\overline{11}$
2	4.73	4.99	-0.26	0.07
3	5.30	4.99	0.31	0.09
4	5.16	4.99	0.17	0.03
5	4.91	4.99	-0.08	0.01
6	5.29	4.99	0.30	0.09
7	4.29	4.99	-0.70	0.50
8	7.30	4.99	2.31	5.31
9	4.76	4.99	-0.23	0.06
10	4.27	4.99	-0.72	0.53
Σ	49.95			7.79

$$\overline{X} = \frac{1}{n} \prod_{i=1}^{n} X_i = 4.99$$

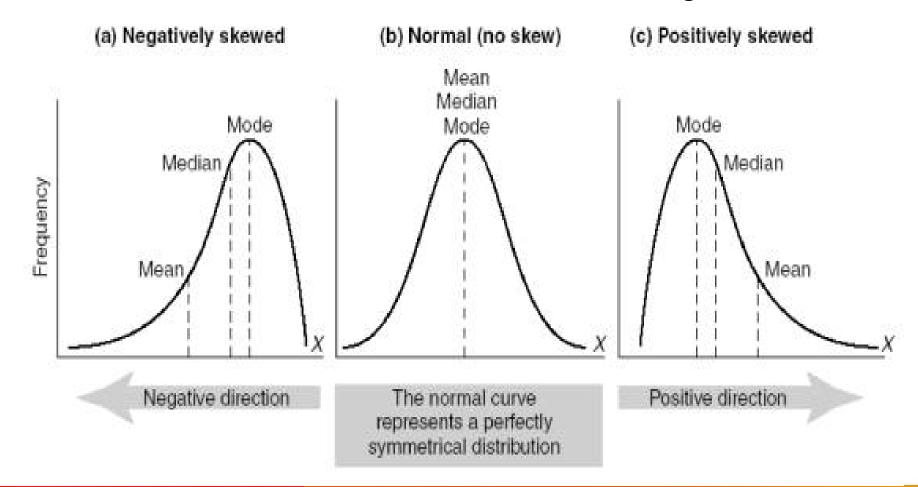
$$S^{2} = \frac{\sum_{i=1}^{n} (X_{i} - \bar{X})^{2}}{n-1} = 0.87 \qquad S = \sqrt{0.87} = 0.93$$

$$S = \sqrt{0.87} = 0.93$$





- **Skewness** (measured by the *skewness* statistic)
 - Are observations symmetrically distributed? Do they tend to extend farther to the left or to the right?





Displaying numerical data

Histograms

 Bars representing frequency (or % of data observed) for specific values or groups of values

Boxplots

- Three number summary with outlier identification
 - 25% (lower quartile), 50% (median), 75% (upper quartile)
- Good for visually comparing two or more groups
- Other scatter plots (for correlation, not today)





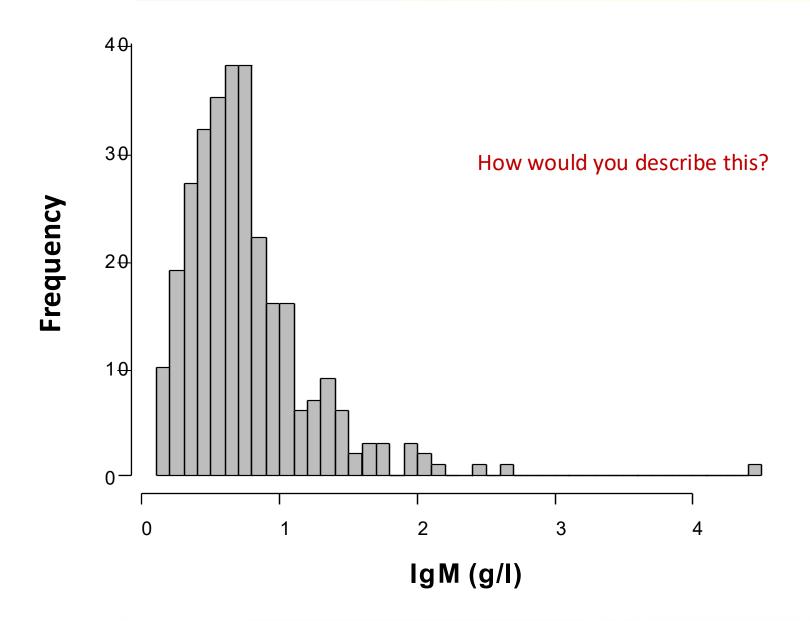
IgM data for 298 healthy children (age 2-6 yrs)

Concentration (g/l)	Frequency	Concentration (g/l)	Frequency
0.1 to < 0.2	3	1.6 to < 1.7	2
0.2 to < 0.3	7	1.7 to < 1.8	3
0.3 to <0.4	19	1.8 to < 1.9	3
0.4 to < 0.5	27	1.9 to < 2	0
0.5 to < 0.6	32	2 to < 2.1	3
0.6 to < 0.7	35	2.1 to < 2.2	2
0.7 to <0.8	38	2.2 to < 2.3	1
0.8 to < 0.9	38	2.3 to < 2.4	0
0.9 to < 1	22	2.4 to < 2.5	0
1 to < 1.1	16	2.5 to < 2.6	1
1.1 to < 1.2	16	2.6 to < 2.7	0
1.2 to < 1.3	6	2.7 to <2.8	1
1.3 to < 1.4	7		
1.4 to < 1.5	9		
1.5 to < 1.6	6	4.5 to < 4.6	1





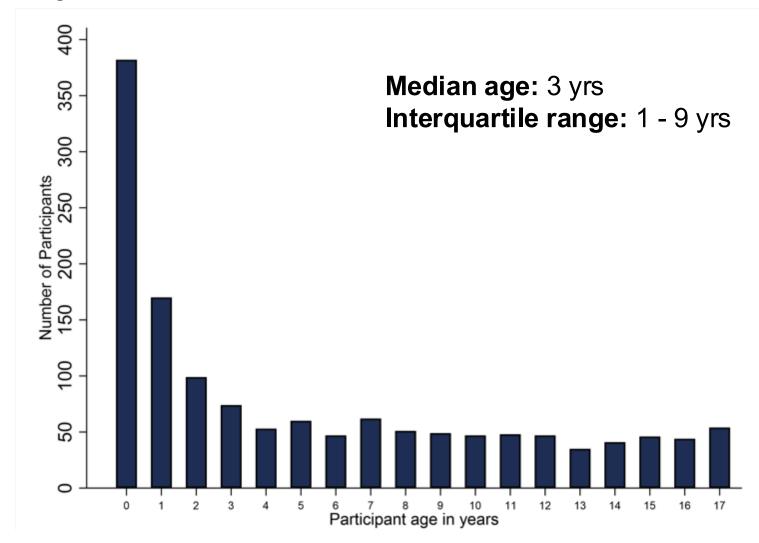
Histogram of IgM (g/l) in 298 healthy children







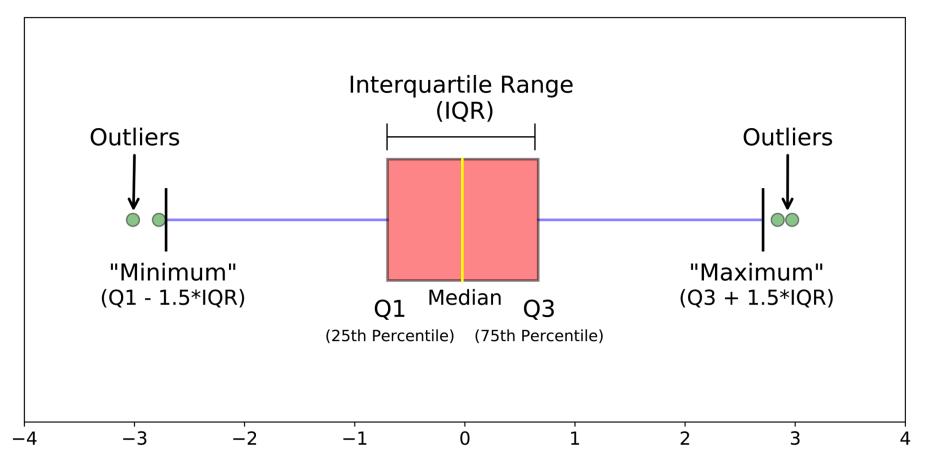
Ages of 1398 children with COVID-19 enrolled in the PERN-COVID-19 cohort study







The Box and Whisker Plot: Anatomy



Note: Points above and below the "whiskers" are considered to be outliers. These outliers are not representative of the general distribution of values





Boxplot Example: CTFC

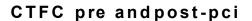
In coronary angiography, the corrected TIMI frame count (CTFC) is the number of frames required for dye to reach a standardized distal landmark.

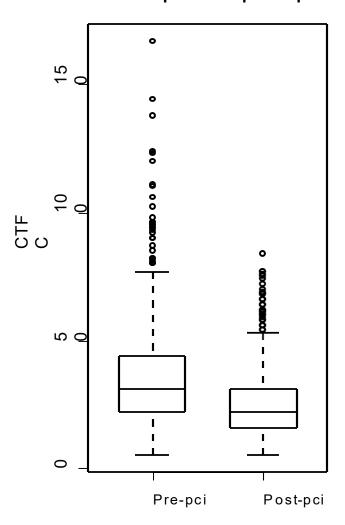
- The following boxplots show the CTFC before and after angioplasty in approximately 1000 patients.
 - Note the skewness in the distributions.



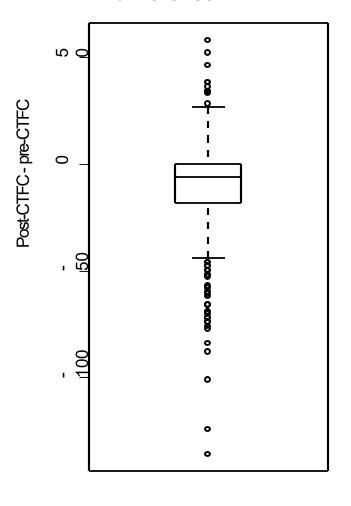


Boxplot Example: CTFC



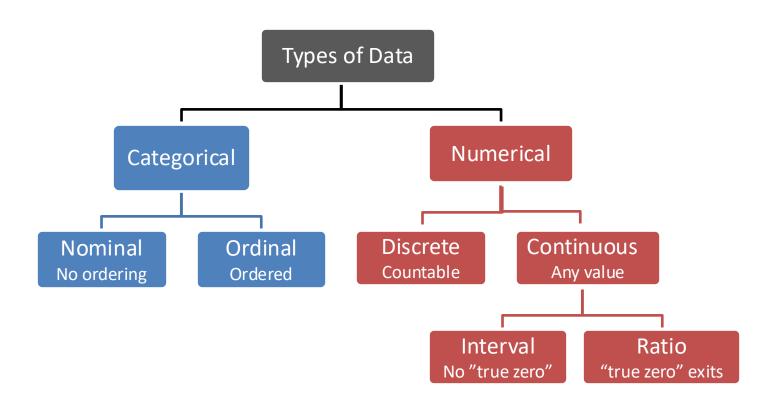


difference in CTFC



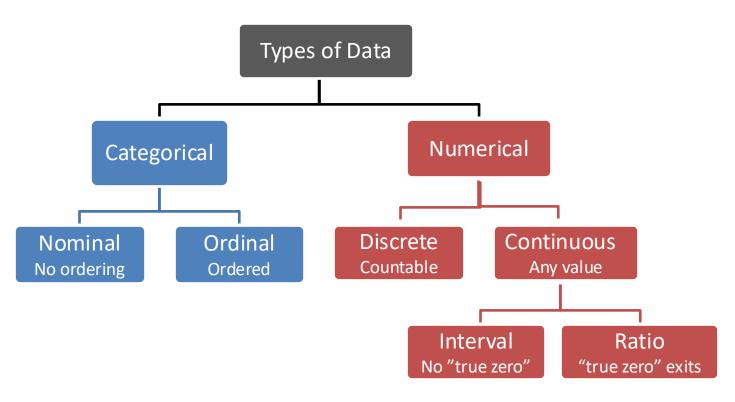








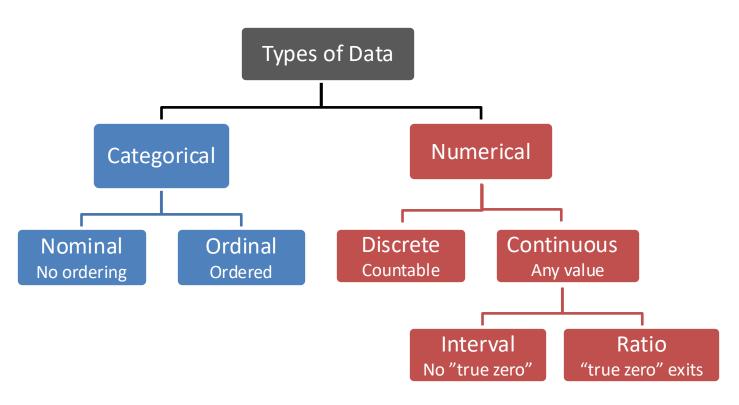




- Difference between 2 categories not equal
- Frequency counts with % (proportions)
- Mode
- Bar charts





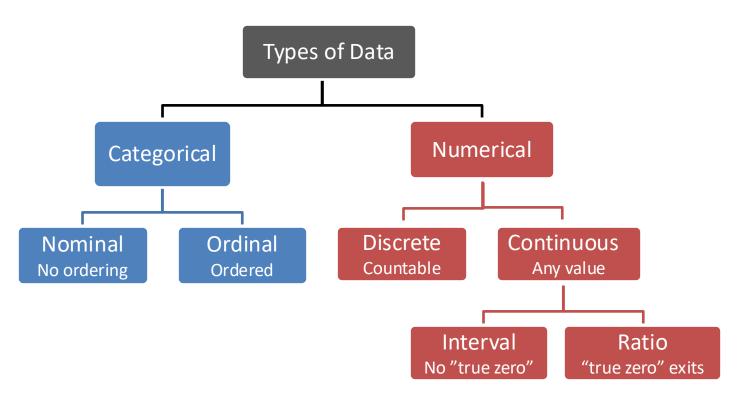


- Difference between 2 categories not equal
- Frequency counts with % (proportions)
- Mode
- Bar charts

- Difference between values is meaningful
- Mean, median
- SD, Range/IQR
- Skewness
- Histograms, box plots







- Difference between 2 categories not equal
- Frequency counts with % (proportions)
- Mode
- Bar charts

- Difference between values is meaningful
- Mean, median
- SD, Range/IQR
- Skewness
- Histograms, box plots

Graphics should <u>help</u> readers understand the data





Lecture objectives

- Statistics, study design, and critical thinking about bias are important
- Types of sampling methods
- Identify what types of data you are (or want to be)
 working with and describe it appropriately
 - Graphical presentation
- Statistical analysis plan
 --> have one before you start





Your statistical analysis methods

- Definitions
 - E.g. What is a 'severe outcome' for a child with COVID-19
- State descriptive statistics used
 - Keep in mind is your sample representative?
- State statistical test or methods used for comparison of groups (if any) in order of when done
- Cite references for complex or uncommon statistical tests used to analyze the data
 - Check similar papers, consult with a designated study (or department) statistician and epidemiologist





Your statistical analysis methods

- Supplementary appendix for more detail, complex methods, 'personalization' of models
- Describe how you checked the assumptions of your statistical methods
 - E.g. is your quantitative data normally distributed?
 - E.g. do participants differ by an unadjusted variable (e.g. site)?
- Provide a sample size calculation
 - Not done at time of report writing, should have been done (or declared 'convenience') at time of proposal





Types of Measures

- The way (scale of measurement) in which your data is collected, will influence statistical analysis
 - E.g. blood pressure high, normal, low vs exact value in mm Hg
 - E.g. viral presence (yes/no) vs viral load in CT (cycle threshold) values
 - Challenges: e.g. quantification of tissue damage? socioeconomic status?
 - Consider which type of measure is appropriate for answering your research question
 - Proportion of treatment failures? Number of diarrheal episodes per patient? Change in quality of life pre and post intervention.





Your statistical analysis methods

In reality, prior to undertaking statistical analysis (as designated in your plan)

- → Careful graphical and exploratory analysis
 - Histograms, scatter plots, univariate analysis, etc
 - Aids in data cleaning and the decision about the appropriateness of planned methods
 - It is not preferred but OK to change your statistical analysis plan! All changes need to be stated and justified in detail.

