

Week 3: Data types, summary statistics, & transformations

ANTH 674: Research Design & Analysis in Anthropology

Professor Andrew Du

Andrew.Du2@colostate.edu

Office Hours: Thursdays, 1:00–4:00pm in GSB 312

1

Statistics vignette

- What is Euler's number & where does it come from?
- $e \approx 2.71828$
- Used as base of natural logarithm
- Fundamental to continuous growth & rate of change



Leonhard Euler

2

Derivation using compound interest

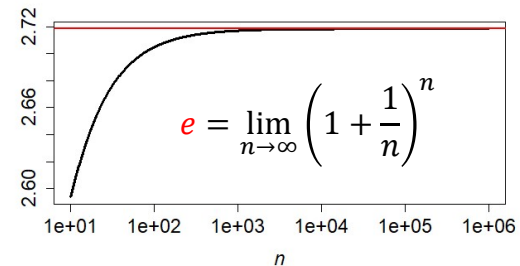
- You have \$1 in your bank, which offers 100% interest every year
 - After one year, \$1 → \$2
- 50% interest twice a year
 - After one year, \$1 → \$1.50 → \$2.25
- 1/12th interest every month
 - After one year, \$1 → \$1.08 → ... → \$2.61
- General formula: $N_0 \left(1 + \frac{1}{n}\right)^n$



3

Derivation using compound interest

- Every day: $1 \left(1 + \frac{1}{365}\right)^{365} = 2.715$
- Every hour: $1 \left(1 + \frac{1}{8760}\right)^{8760} = 2.718$



4

A probability interpretation

- Probability every dropped chocolate is placed in wrong position = $1/e$ (as # chocolates $\rightarrow \infty$)



5

Widely considered the most beautiful equation in math

- Euler's identity

$$e^{i\pi} + 1 = 0$$

6

Lecture outline

1. Different types of data ("qualitative" vs. quantitative)
 1. How are they described & summarized?
 2. How are they plotted (visualizing the distribution)?
2. Different types of data transformations
 1. What are they & what are they used for?
3. Plotting two data types against each other

7

Different data types

How are they described, summarized, and plotted?

8

What is data?

- Wikipedia: **Data** are characteristics or [information](#), usually numerical, that are collected through observation
- Want to learn something from data through analysis and/or visualization (plotting)



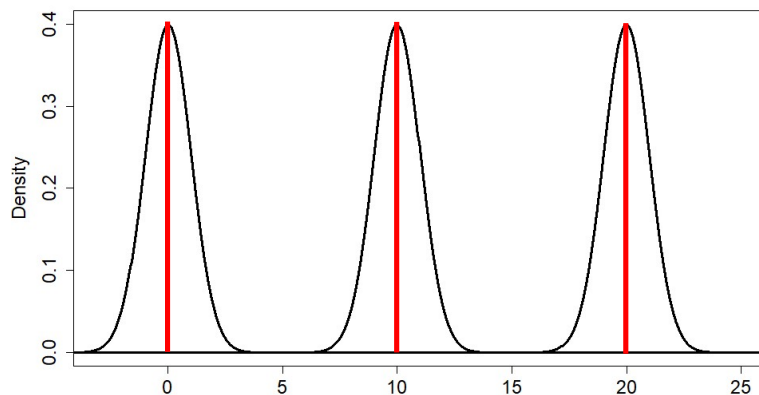
9

Summary statistics

- Used to summarize distribution of data w/ one number (except for multivariate distributions)
 1. Location or central tendency
 2. Spread or variation

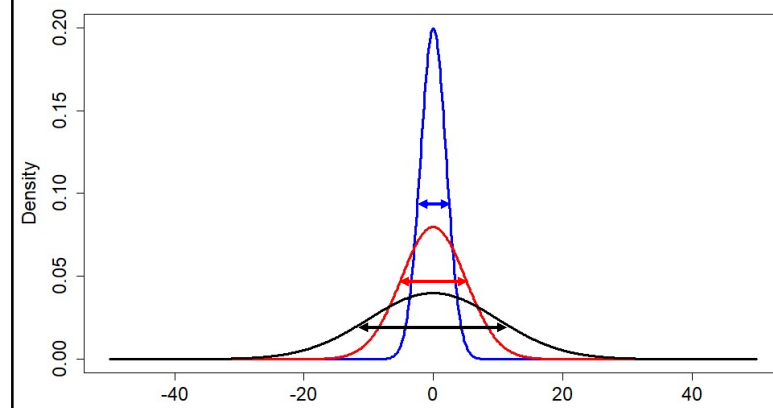
10

Location/central tendency



11

Spread/variation

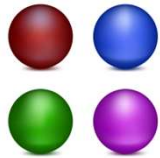


12

Different types of data

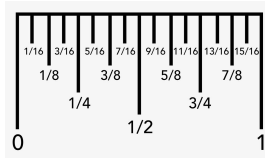
“Qualitative”

- Categorical/nominal
- Ordinal



Quantitative

- Discrete
- Continuous



- Data type tells you which summary statistics, plots, and analyses to use!

13

Questions?



14

“Qualitative” data

15

What is “qualitative” data?

- Data assigned to groups, usually based on some qualitative property
1. Categorical/nominal data (unordered)
 2. Ordinal data (ordered)



16

Categorical/nominal data

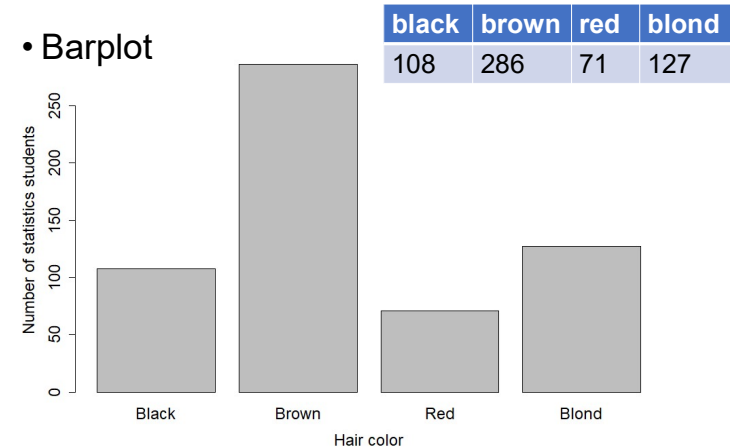
- Unordered qualitative data
- E.g., head/tail (binomial);
basalt/chert/quartz (multinomial)
- Quantified as counts or proportions
- = factors in R



17

Plotting categorical data

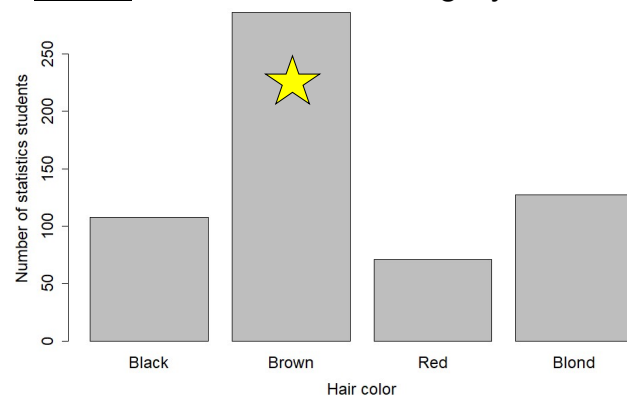
- Barplot



18

Central tendency

- Mode: most common category



19

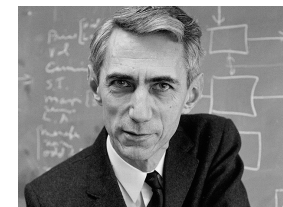
Spread

	black	brown	red	blond
#	108	286	71	127
prop.	0.18	0.48	0.12	0.21

- Information theory measures (most common is Shannon's index)
- Not commonly used

$$-\sum p_i \log(p_i) \quad p_i = \text{proportion of } i^{\text{th}} \text{ category}$$

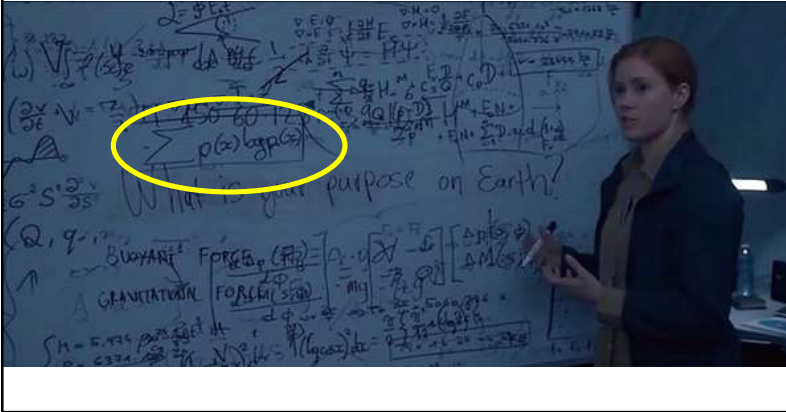
$$-[0.18 \times \log(0.18) + 0.48 \times \log(0.48) + 0.12 \times \log(0.12) + 0.21 \times \log(0.21)] = \mathbf{1.25}$$



Claude Shannon

20

From the movie *Arrival*



21

Ordinal data

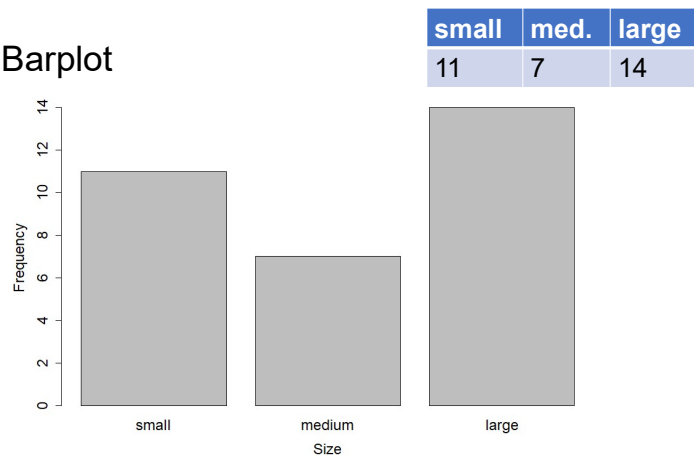


- Ordered qualitative data
- Distance between categories not known
- E.g., small/medium/large; juvenile/adult
- Quantified as counts or proportions
- = ordered factor levels in R

22

Plotting ordinal data

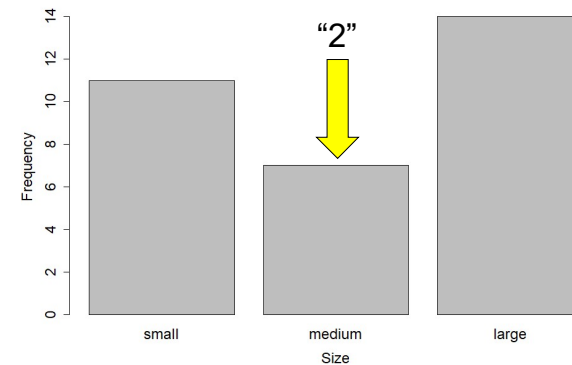
- Barplot




23

Central tendency

- Median: middle value in ordered data
(first convert to ranks: 1, 2, 3)



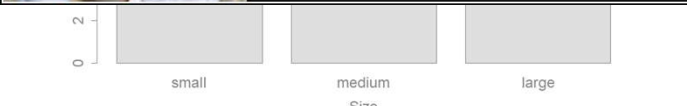
24



The median isn't the message.

— *Stephen Jay Gould* —

AZ QUOTES



2

0

small medium large

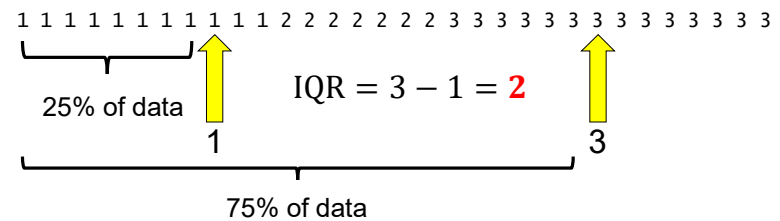
Size

25

small	med.	large
11	7	14

- Interquartile range (IQR): difference between 75th and 25th percentiles (or 3rd and 1st quartiles; median is the 2nd quartile)
- Middle 50% of data

After converting sizes to ranks



26

Questions?



27

Quantitative data

28

What is quantitative data?

- Each data point is a number, and distances have meaning (e.g., 1 vs. 3)

Discrete

- Finite or countable values (e.g., integers)
- In practice, all continuous numbers are discrete due to limited precision of measurements (e.g., 1.21, 1.22, 1.23)

Continuous

- Any value within a *continuous* interval (e.g., 2.4575)

29

Discrete data

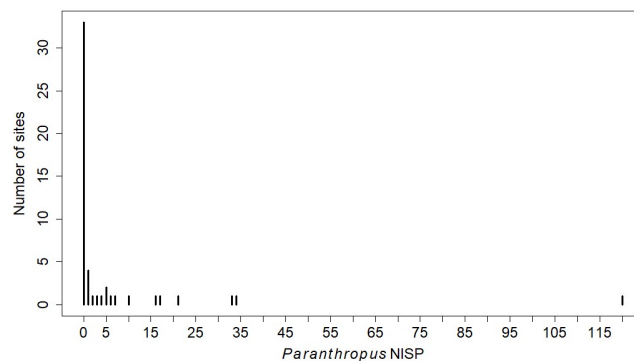
- Finite or countable values
- E.g., count data, anything measured in integers
- Treated as numeric class in R



30

Plotting discrete data

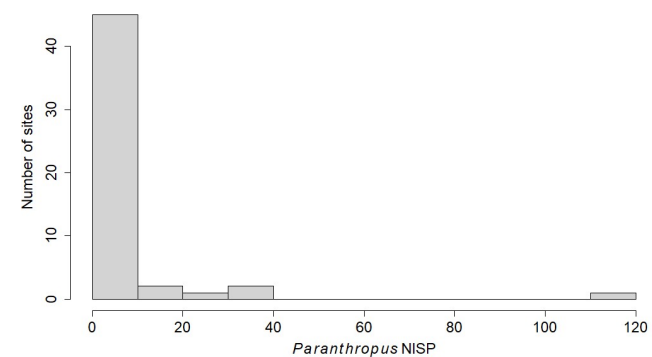
- Line plot (like a barplot w/ more categories)



31

Plotting discrete data

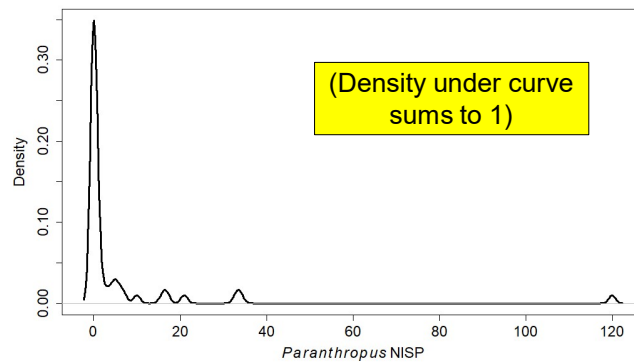
- Histogram



32

Plotting discrete data

- Density plot ("smoothed histogram")



33

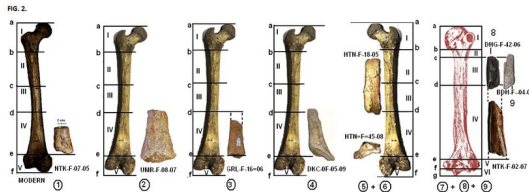
Central tendency & spread

- Same as with continuous data

34

Continuous data

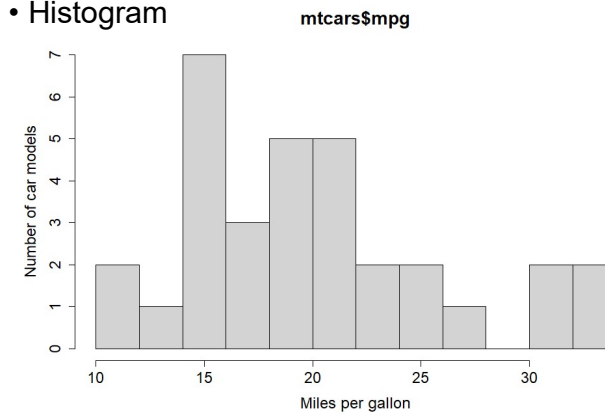
- Any value within a *continuous* interval (e.g., 2.4575)
- E.g., stone tool mass, hominin femur length
- Treated as numeric class in R



35

Plotting continuous data

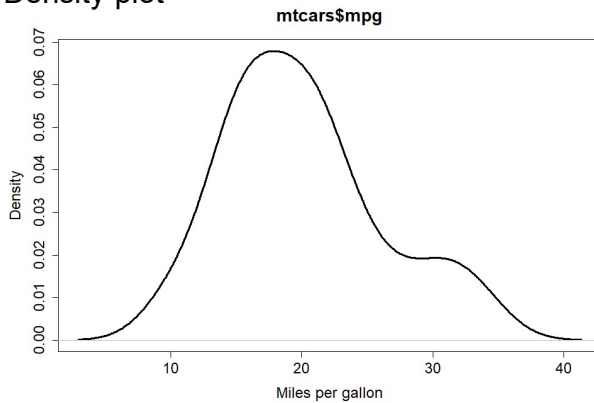
- Histogram



36

Plotting continuous data

- Density plot



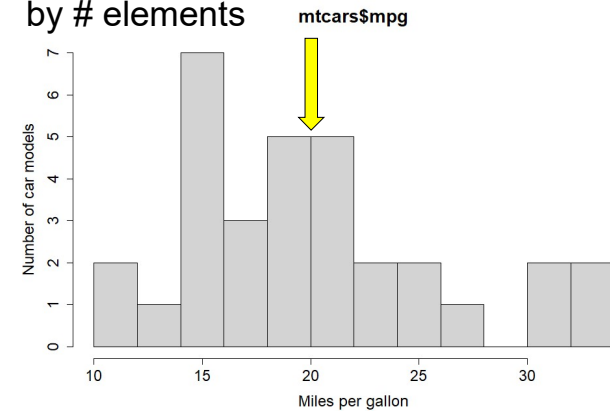
37

Central tendency

$$\frac{\sum x_i}{n}$$

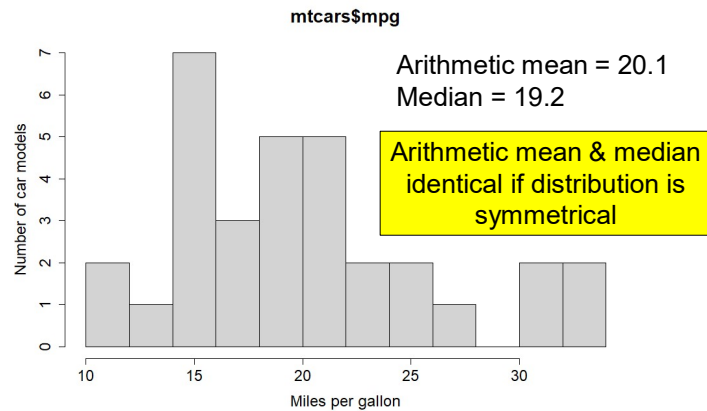
elements $\rightarrow n$ \leftarrow x_i \leftarrow i th element

- Arithmetic mean: sum elements & divide by # elements



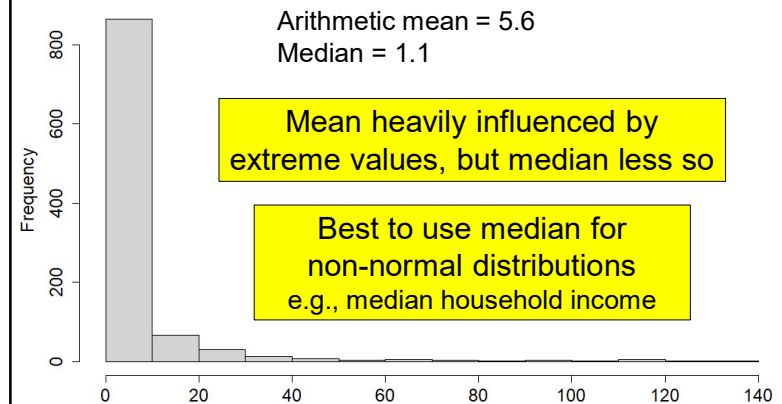
38

Central tendency



39

What if data are non-normal?

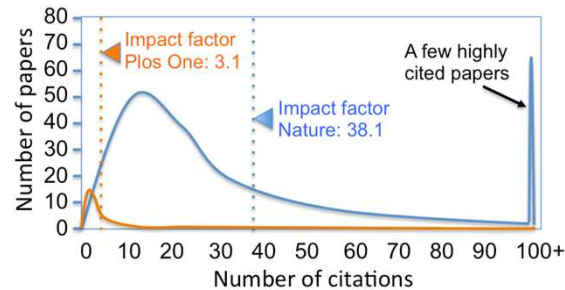


40

An example

90% of *Nature's* 2004 impact factor was based on only 25% of its publications

- Journal impact factor: average # of times articles from journal published in the past two years have been cited in year of consideration



41

Central tendency

$$\left(\prod_{i=1}^n x_i \right)^{\frac{1}{n}}$$

i-th element # elem.

- Geometric mean: n th root of elements multiplied together
E.g., $x \leftarrow c(1, 2, 3)$
 $GM = \sqrt[3]{(1 \times 2 \times 3)}$
- Same as taking arithmetic mean of log-transformed values & calculating antilog

$$GM = \exp\left(\frac{\sum \log(x_i)}{n}\right) \quad \text{E.g., } GM = \exp\left(\frac{\log(1) + \log(2) + \log(3)}{3}\right)$$

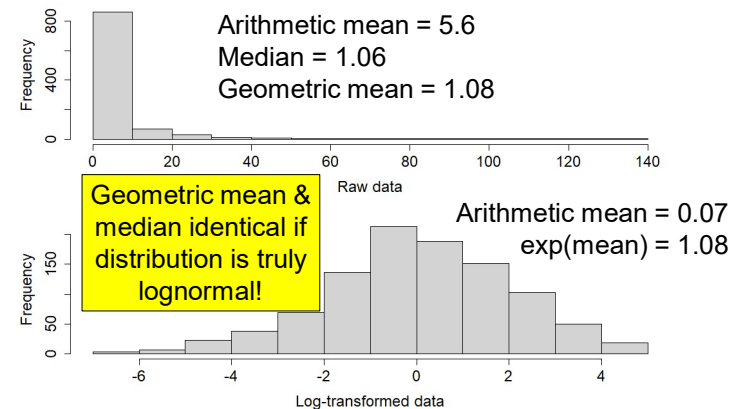
42

Central tendency (geometric mean)

- Used when dealing with data produced by multiplicative processes (e.g., % increases) → lognormal distributions
- E.g., population size, body size, household income, citation numbers
- Can't use when you have zeros or negative numbers

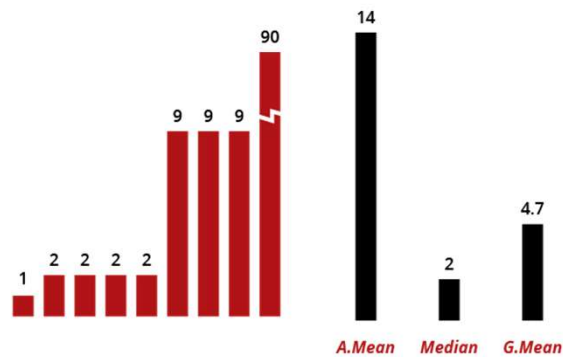
43

Geometric mean (on lognormal distribution)



44

Comparing all three



Order always the same!

45

Questions?



46

Spread

- Variance: measures how far values deviate from the arithmetic mean
- Subtract the mean from each element, square the results, add them up, and divide by number of elements minus 1

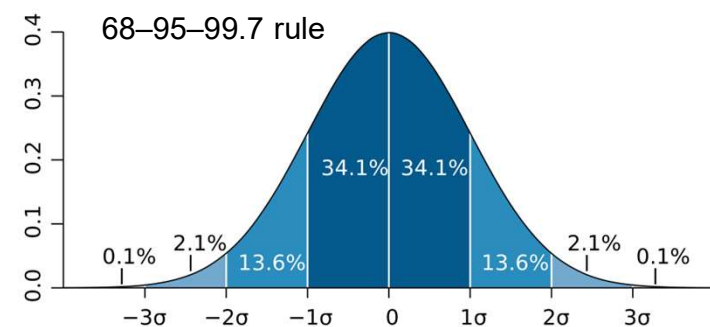
$$\frac{\sum (x_i - \bar{x})^2}{n - 1}$$

*i*th element → x_i Arithmetic mean → \bar{x}
elements → n

- Square-root of variance = standard deviation

47

Spread (standard deviation)

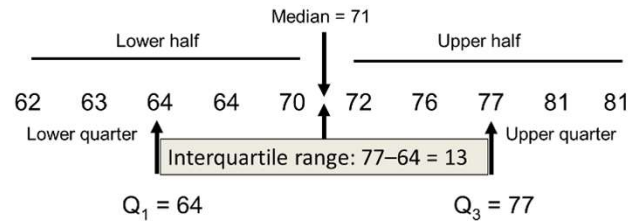


- As with the arithmetic mean, variance and SD are most interpretable for normal distributions

48

What if data are non-normal?

- Interquartile range
- Doesn't depend on arithmetic mean or type of distribution



50

Geometric standard deviation

- Used when you would use the geometric mean (e.g., lognormal data)
- Calculate std. dev. of log-transformed values and take the antilog

$$\exp \left(\sqrt{\frac{\sum [\log \left(\frac{x_i}{GM} \right)]^2}{n}} \right)$$

*i*th element

Geometric mean

elements

51

Summary: which plot to make?

Data type	Plot
<i>Categorical</i>	Barplot
<i>Ordinal</i>	Barplot
<i>Discrete</i>	Line plot Histogram Density plot
<i>Continuous</i>	Histogram Density plot

52

Summary: which statistic to use?

Data type	Location	Spread
<i>Categorical</i>	Mode	Information measures
<i>Ordinal</i>	Median	Interquartile range
<i>Discrete/Continuous</i>		
Normal	Arithmetic mean	- Variance - Standard deviation
Non-normal	- Median - Geometric mean	- Interquartile range - Geometric SD

53

Questions?



54

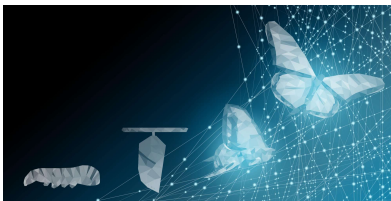
Data transformations

What are they & what are they used for?

55

What is data transformation?

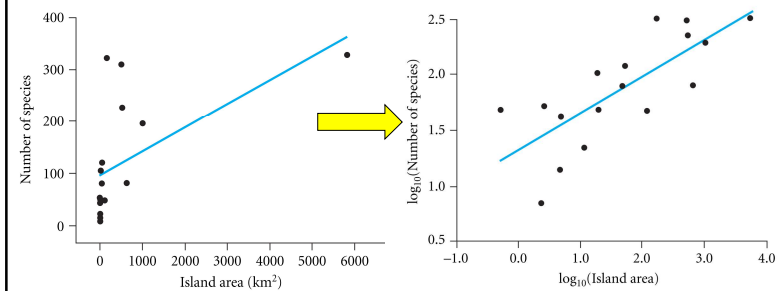
- Applying a mathematical function to data to change its distribution
- Rank order of data maintained (monotonic transformation)



56

Why transform data?

- To make data & results easier to understand and visualize



57

Why transform data?

- To make data & results easier to understand and visualize
- To make sure assumptions of statistical methods are not violated

58

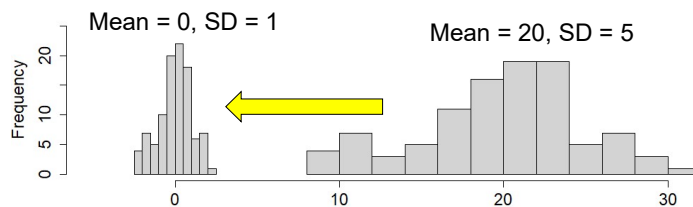
Types of data transformations

1. Centering and scaling
2. Log transformations
3. Square-root transformations
4. Arcsine & logit transformations

59

Centering and scaling

- Transforms data to have mean = 0 & standard deviation = 1 (i.e., Z-scores)
- Subtract the mean & divide by SD



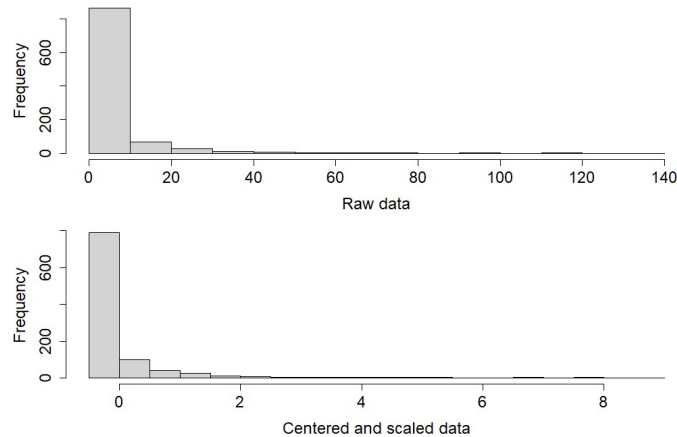
60

Centering and scaling

- Transforms data to have mean = 0 & standard deviation = 1 (i.e., Z-scores)
- Subtract the mean & divide by SD
- Converts data to units of standard deviation, so variables of different units can be compared (e.g., mass & inches)
- Can be used on non-normal data!

61

Centering and scaling



62

Log transformations

- Replaces data with their logarithms
- $b^a = x$; $\log_b x = a$
- Base e (analyses) and 10 (plotting) most common
- Cannot transform zeros and negative numbers



63

Log transformations

- Many statistical & plotting methods deal with **additive/absolute/linear** change
 - E.g., linear regression: $y = a + bx$
 - $1 \rightarrow 2$: **+1**
 - $100 \rightarrow 200$: **+100**
 - $1 \rightarrow 2$: **x2**
 - $100 \rightarrow 200$: **x2**
- } Treated differently w/ **linear** methods
 } But **multiplicative/relative** change the same!

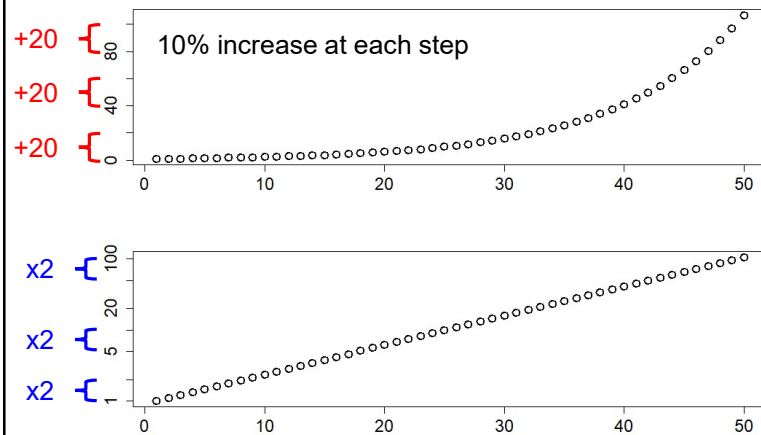
64

Log transformations

- What if you are interested in **multiplicative/relative/proportional/percent** change?
- Log-transformations: **multiplicative** \rightarrow **additive**
- $\log(2) - \log(1) = \log\left(\frac{2}{1}\right) = 0.69$
- $\log(200) - \log(100) = \log\left(\frac{200}{100}\right) = 0.69$
- **Doubling** from 1 to 2 now treated the same as **doubling** from 100 to 200!
- Can now use **linear** methods to investigate **multiplicative** change

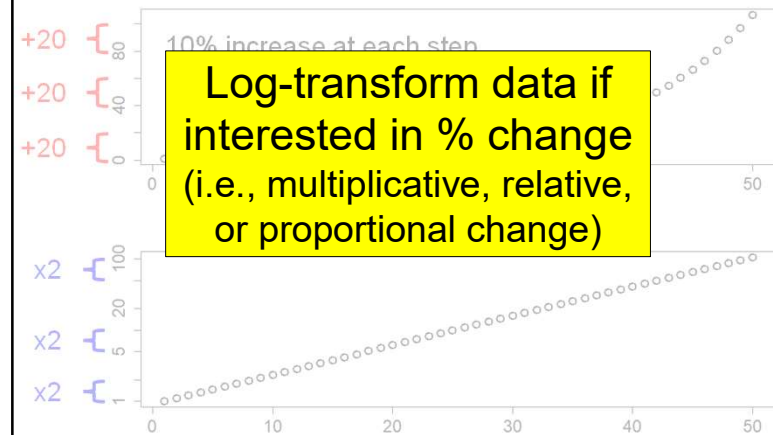
65

Log transformations



66

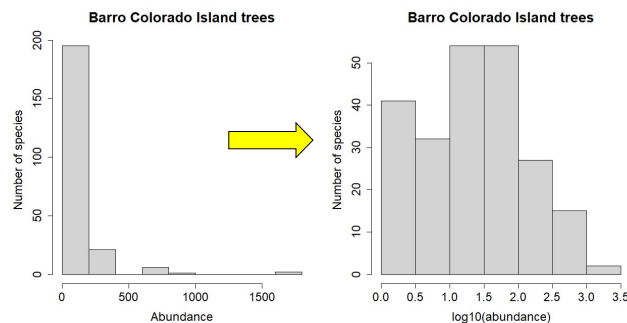
Log transformations



67

Log transformations

- Can aid visualization if data vary over orders of magnitude (spreads out small numbers and squeezes large ones)



68

Other transformations

- Removes dependence between mean & variance of a variable

1. Square-root: \sqrt{x}
 - Used for count data
2. Arcsine: $\arcsin(\sqrt{x})$
 - Used for proportions
3. Logit: $\log\left(\frac{x}{1-x}\right)$
 - Used for proportions

69

Questions?



70

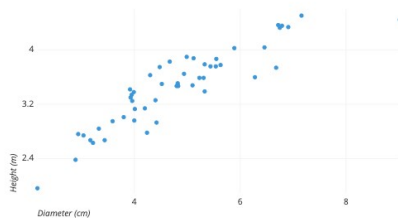
Plotting two data types against each other



71

Why do this?

- Want to see how two variables are related to each other
- Good way to visually describe your two variables



72

First, a note on ordinal & discrete data when plotting

- Ordinal treated as categorical (maintaining order of categories)
- Discrete treated as continuous
- When ordinal data have many categories, can mimic discrete data (e.g., rank abundance of species)
- When discrete data are too few, can mimic ordinal data (e.g., 3, 4, & 5 number of forward gears in mtcars)

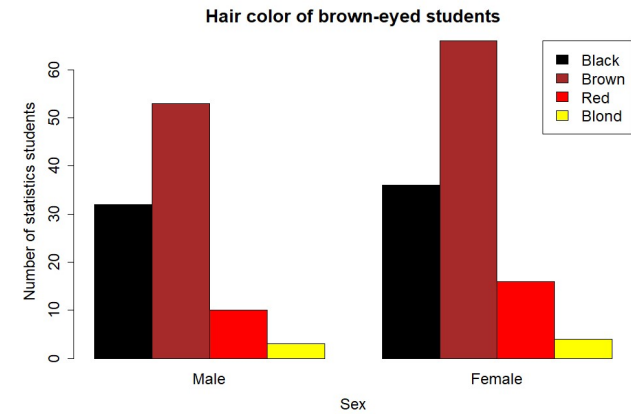
73

Which plot to make

		X-axis	
		Categorical	Continuous
Y-axis	Categorical	Bar plot	Box plot Violin plot
	Continuous	Box plot Violin plot	Scatter plot

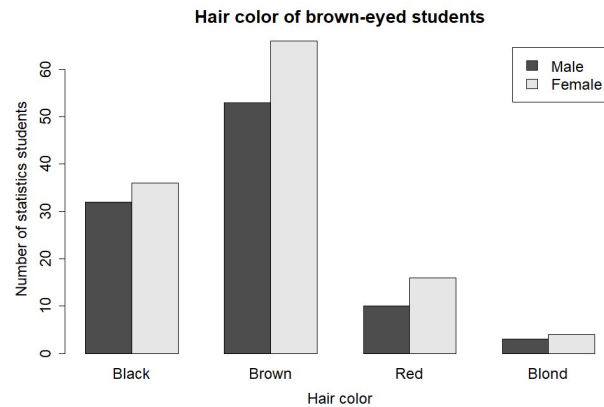
74

Barplots (categorical vs. categorical)



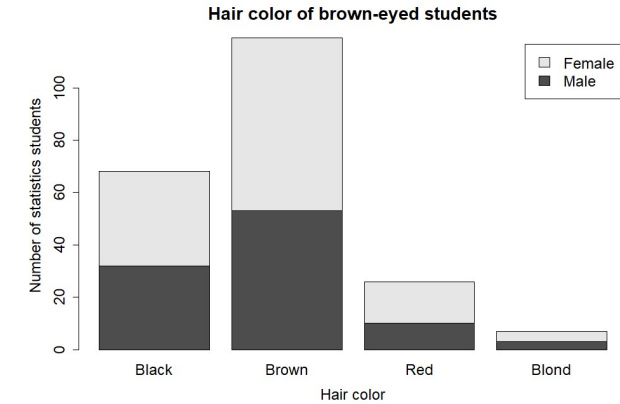
75

Barplots (categorical vs. categorical)



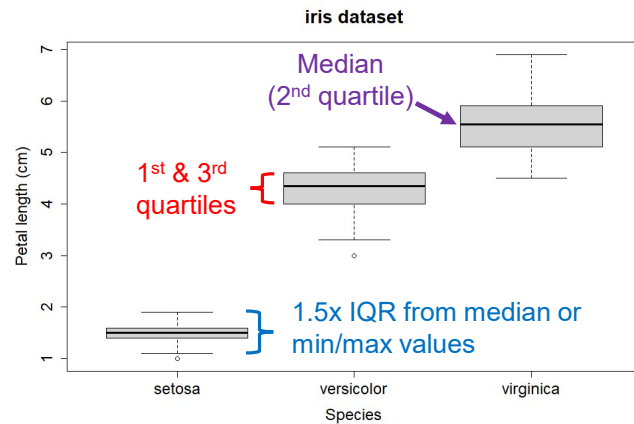
76

Barplots (stacked) (categorical vs. categorical)



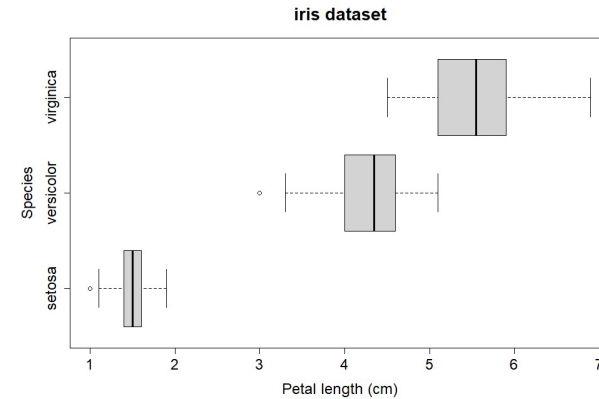
77

Boxplots (categorical vs. continuous)



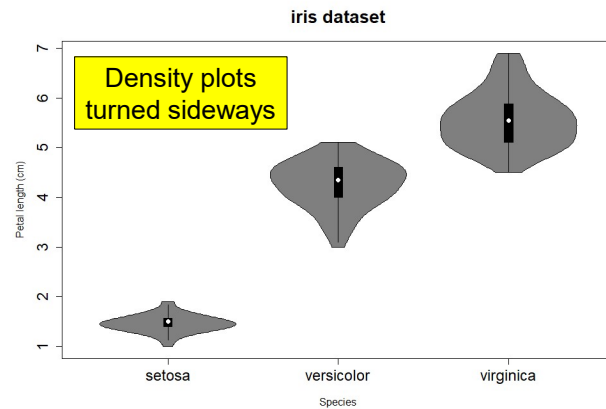
78

Boxplots (continuous vs. categorical)



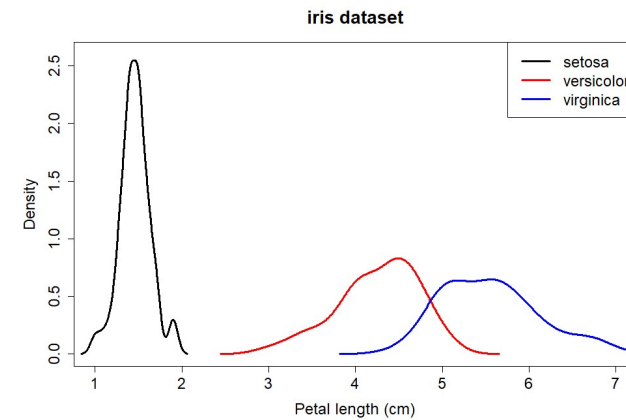
79

Violin plots (categorical vs. continuous)



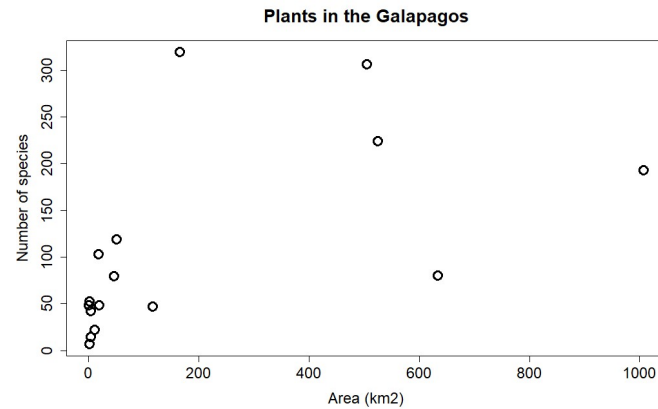
80

Density plot works too! (categorical vs. continuous)



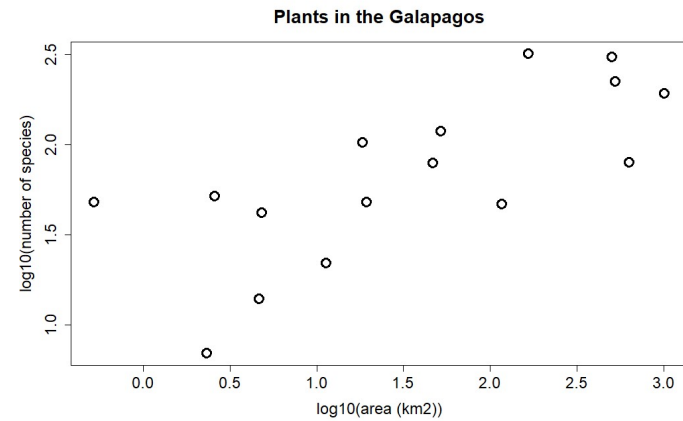
81

Scatter plots (continuous vs. continuous)



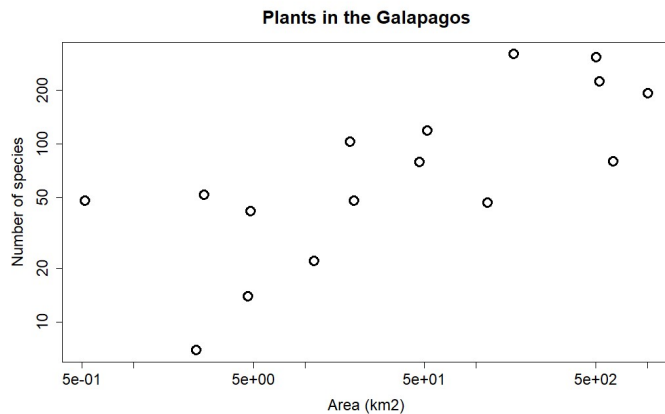
82

Scatter plots (log-transformed variables) (continuous vs. continuous)



83

Scatter plots (log-transformed axes) (continuous vs. continuous)



84

Plotting summary

- These are the general rules of plotting (R will actually automatically make these plots according to your variable type)
- BUT, use best judgment for showing what *YOU want* to show (according to your research question)
- Data visualization is very important: want to convey your data and results as clearly & effectively as possible!

85

Questions?



86

Summary

- There are four main data types: categorical, ordinal, discrete, & continuous
- Data type tells you which summary statistics and plots to use
- Data transformations aid visualization and interpretation & help data satisfy statistical assumptions
- How to plot & compare two variables depends on data type

87