

S Stats B Bootcamp

①

sema barlus

02/28/19

SESSION 1

- Assignment (10 pts) (03/04) → 9/10
ch 1-4 - odd # questions

Outline

- distribution of data
- graphs
- measures of central tendency and dispersion
- correlation

① INTRODUCTION

A) what are statistics? → "fact & figures"

- field of study concerned w/ summarizing data, interpreting data, & making decisions based on data

* qty calculated in sample to estimate population

- importance of interpretation!

? statistics includes:

- numerical calculations

- graphs

- interpretations & decisions based on facts and graphs

? ice cream consumption & drowning

? correlation + causation

- third variable involved - summer

B) Importance of statistics:

- credibility

- question the statistics you encounter

? certain preferences?

- need to know more

? evaluating statistical claims:

- statistics presented

- sources of statistical findings

- procedures used to generate the claims

C) Descriptive Statistics

- used to summarize & describe data

- data: info. collected from an expt., survey, historical record, etc.

- don't involve generalizing

↳ inferential statistics

avg salary for occupations in 1999

of unmarried men / 100 unmarried women in US metropolitan areas 1990

(2)

- winning Olympic times for men & women (since 1984)

- can compute desc stats from the data

- ? 2 basic divisions of statistics

- inferential

- descriptive

- ? Desc. stat.

- summarize & describe data

- ? Desc. stats:

- mean of age

- # of ppl

- median

- height

D) INFERENTIAL STATISTICS

- Statistics: rely on sample to draw inferences from population

- sample selection challenges:

- over-representation

- non-random sampling

- non-proportional

- SIMPLE RANDOM SAMPLING

- every member of population has to have an equal chance of being selected into the sample

- selection of 1 member must be independent of every other member "chooses a sample by pure chance"

- * Bias

- SAMPLE SIZE MATTERS

- sampling procedure ^{rather} than results define SRS.

- small sample size may not be rep. of the pop.

- MORE COMPLEX SAMPLING

- when SRS not feasible

(1) Random Assignment

- populations often hypothetical

- random division of sample into 2 groups

- critical for expt validity

Ex - failure to randomize → invalidates experimental findings

- non-random sample → restrict generalizability

(2) Stratified Sampling

- can be used if population has # of distinct strata / groups

1. identify number of sample in each group

2. randomly sample from subgroups w proportion analyses

- ? our data come from populations, but we really care most about population.
- ? A random sample:
 - stratified sampling more info. of pop than random sampling
- ? participant aware for resp. study put into treatment groups on basis of chance
 - Random Assignment
- € ? uncertainty re. conditions about population can be eliminated if
 - obtain data from all members of the pop.
- ? using a Random sample:
 - is no accept uncertainty abt conclusions
 - enables you to calculate statistics
 - is no risk drawing wrong conclusions about population
- ? Random sample
 - equal prob of selection
- ? SRS of student body:
 - ID #'s picked from a table of random #'s
- ? A Biased sample:
 - will likely have group from the population over-represented or under-represented due to systematic sampling factors.

3) sampling framework (con'tnue book)

- SRS is stratified SRS

E) variables

- Independent & dependent variables:
 - var = property or char that takes on value
 - independent var = manipulated by an experimenter (x)
 - determine effect of independent var on dependent var (y)
 - levels of independent variables
 - experimental & control → 2 levels
 - 5 diets → 5 levels
 - # of levels = # of experimental conditions
- Qualitative & Quantitative vars:
 - do not imply numerical ordering → QUALITATIVE
 - Qualitative attribute (e.g. gender)
 - measured numerically (e.g. height) → QUANTITATIVE
- Discrete & continuous vars:
 - DISCRETE - possible values are discrete points on a scale (e.g. # of children)
 - CONTINUOUS - continuous scale (e.g. response time)

④

03.03.11a

F) PERCENTILES

- percent of scores lower than yours (no universally accepted def)

$$\text{Rank} = \frac{P}{100} \times (N+1)$$

P = desired percentile

N = # of entries

R = integer then R, otherwise integer part, decimal part · (R), (R + 1)

- DEFINITIONS

↳ interpolate by multiplying difference

1) "the lowest score that's greater than 65% of the scores"

2) "smallest score greater than or equal to 65% of scores"

3) Rank:

① Find R of percentile desired

R = 2.25

② integer: then Rth

percentile is # at Rank R

Percentile = 2; Rank = 2

③ R + integer, find pth percentile by interpolation

④ IR = 2

⑤ FR = 0.25

⑥ Find scores w/ Rank IR & Rank IR + 1

= 2, 3 (Rank) (n/scores)

⑦ (0.25)(7.5) + 5 = 5.5

Why? multiply diff b/w scores by FR then add to lower IR

G) LEVELS OF MEASUREMENT

- Types of scales:

◦ measure dependent var depending on type

① NOMINAL SCALES

> names / categorical responses

> no ordering

> lowest level of measurement

② ORDINAL SCALES

> ordered

> allow comparisons of degree w/ c & subjects across diff var

> BUT: diff b/w 2 levels can't be assumed to be same

↳ (diff b/w adjacent scale may not rep egl intervals on underlying scale)
even if numbers are used

③ INTERVAL SCALES

> numerical scales w/ same interpretation → no

> ex. Fahrenheit

> BUT: no true zero point → no ratio comparison.

Measurement
- interval ordinal

④ RATIO SCALES:

- > most informative scale
 - > interval with zero point (call B values into 1)
 - > labels, ordered, diff scale = same meaning, ratio
 - > eq. Kelvin , amt of #

④ psychosocial variant: pain

- > 5-7 pt scale → ordinal } generally "inappropriate"
 - > memory = ratio, but early u hard }

- CONSEQUENCES OF LEVEL OF MEASUREMENT

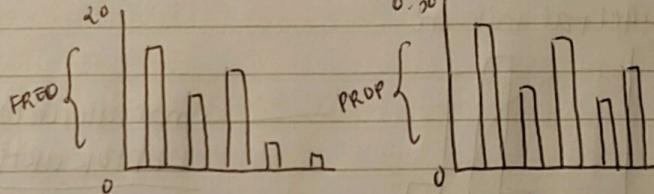
- helps in statistical analysis

H) DISTRIBUTIONS

Distributions of discrete variables

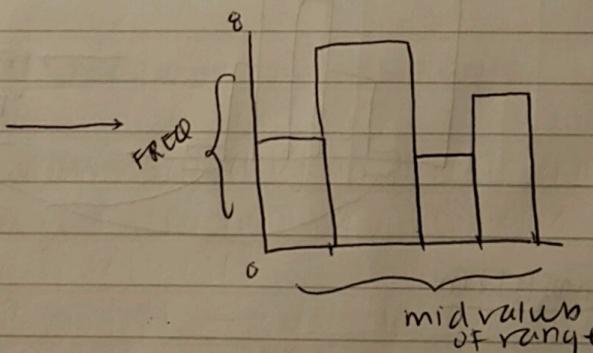
- : Frequency table, freq. dist., A/r prop. dist

COLOR	PREP



- continuous variables:

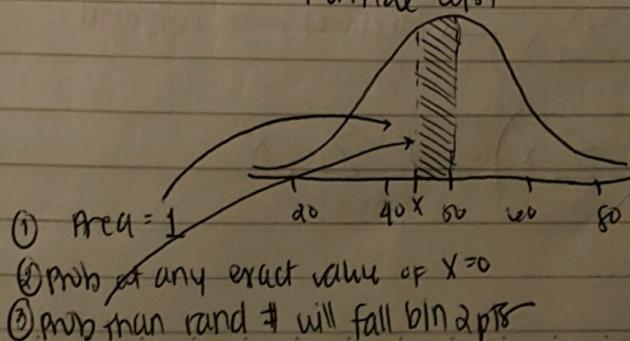
- eq. time to respond → grouped freq dist since entries are unique (ranges) = tabulated



- probability densities

- represent all potential lines at once \rightarrow point dip of const var.

"probabilistic density"
normal dist



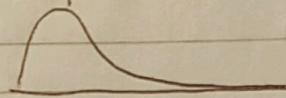
(6)

- shapes of distributions:

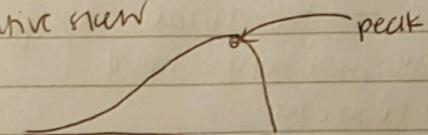
- symmetric

- tails → longer + tail = positive skew (skewed to the right)

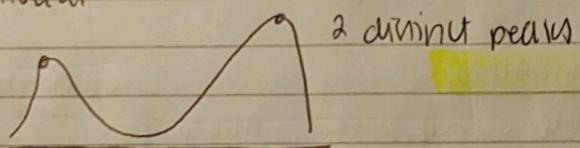
- large / extreme



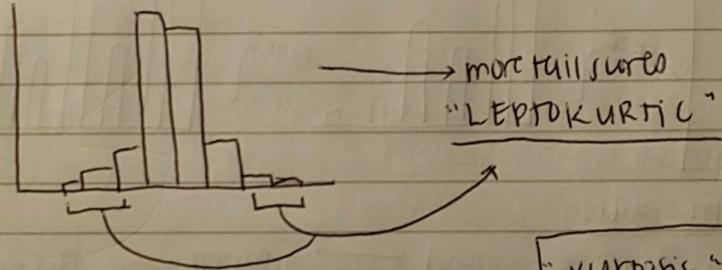
- negative skew



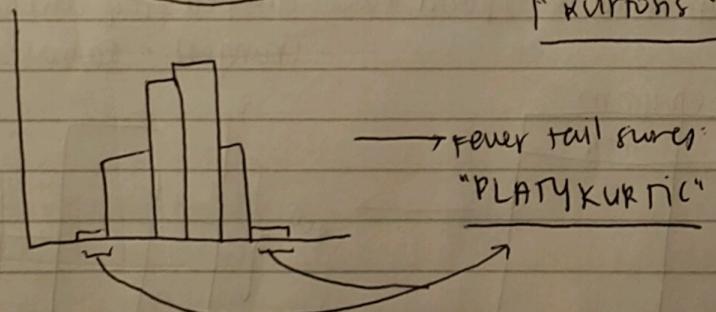
- bimodal dist



- How large fat tails are



"KURTOSIS"



I) SUMMATION NOTATION:

$$\sum_{i=1}^4 x_i = x_1 + x_2 + x_3 + x_4$$

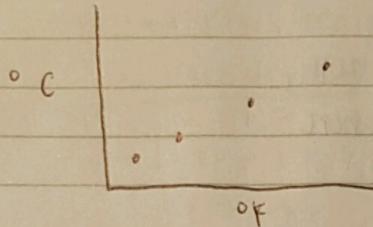
↑ index var

$\sum x$ → all scores to be summed

$$(\sum x)^2 \neq \sum x^2$$

J) LINEAR TRANSFORMATIONS

- Transform data from one measurement scale to another
- Form a straight line = "linear"



K) LOGARITHMS

- log transformation reduces positive skew
- ratio of logs:

◦ opposite of exponents

$$10^2 = 100$$

$$\log_{10}(100) = 2$$

$$\log_{10}(1000) = 3$$

◦ $e = 2.718 \rightarrow$ natural log $\rightarrow \ln(x)$ $\log_e(x)$

◦ changing the base of a log changes the result by a multiplicative constant:

$$\log_{10} \rightarrow \ln \rightarrow \times 2.303$$

$$\ln \rightarrow \log_{10} \rightarrow \div 2.303$$

◦ "anti log" undoes logging \rightarrow change base to N

- Logs and proportional change:

◦ proportional raw changes are equal in log units

◦ arithmetic ops:

$$\log(AB) = \log(A) + \log(B)$$

$$\log(A/B) = \log(A) - \log(B)$$

L) STATISTICAL LITERACY

- case study

(8)

03/03/19

② GRAPHING DISTRIBUTIONS

a) QUALITATIVE VARIABLES:

- graphing: first step in data analysis
- William Playfair 18th c., John Tukey 20th c.
- iMac computer expansion of Apple market share
 - 500 iMac users interviewed
 - no ordering pre-established

① Frequency Tables:

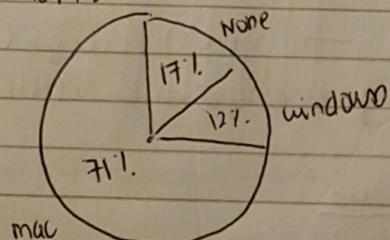
Result	Freq	Rel. Freq
None	x	z
windows		
mac		
Total	y	

proportion of responses in each category

$z = \frac{x}{y}$

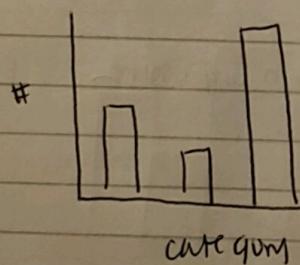
Response
Categ.

② pie charts:



- slice proportional to rel. frequency
- relative freq. $\times 100$
- good for small # catag.

③ Bar charts:



- Comparing distributions:

- Bar charts \rightarrow diff bln dist (oriented horiz.)
- Graphical mistakes to avoid:

- 3D bar charts
- lie factor = size of effect shown in a graph : size of min data (Tukey) $\rightarrow > 1.05$ or < 0.95 = Bad
- setting baseline $\neq 0$
- using line graph for qual vars

B) QUANTITATIVE VARIABLES

- measured on numeric scale

D) Stem & Leaf Display

- Best for small / moderate Amts of data

- # of rough down passy

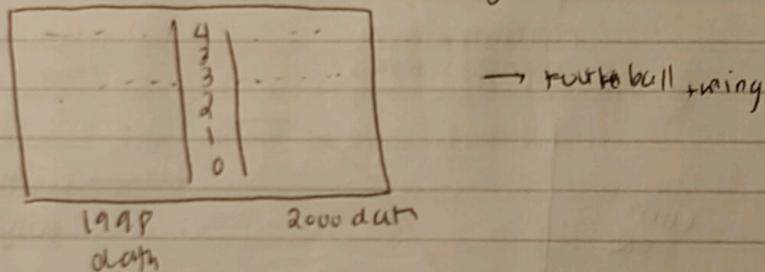
3	2 3 3 7
2	0 0 0 1 4 2 3
1	8 7 6 5
0	6 9

↑ ↓
tens ones

- to clarify shape of dist.

- can also split

- Back-to-Back Stem & Leaf Display



→ run & ball training

E) CUMULATIVITY

- effects of pruning

- number of per cluster \approx b

- useful Mr n $/ N = 200$

→ plan accuracy

F) Histograms

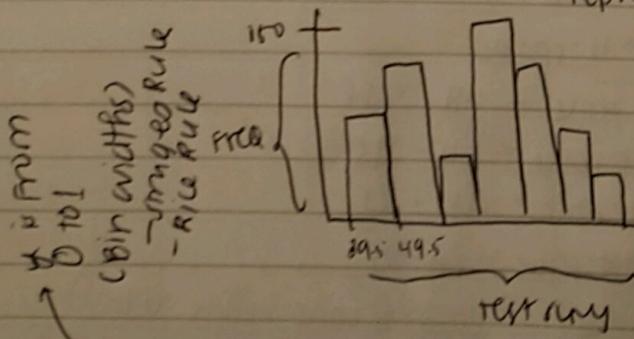
- shape of distribution \rightarrow large # of observations

- steps:

1. Create Freq. Table (intervals) — class intervals

↳ place intervals midway or b/n 2 #s

↳ class freqs represented by bars



→ discrete

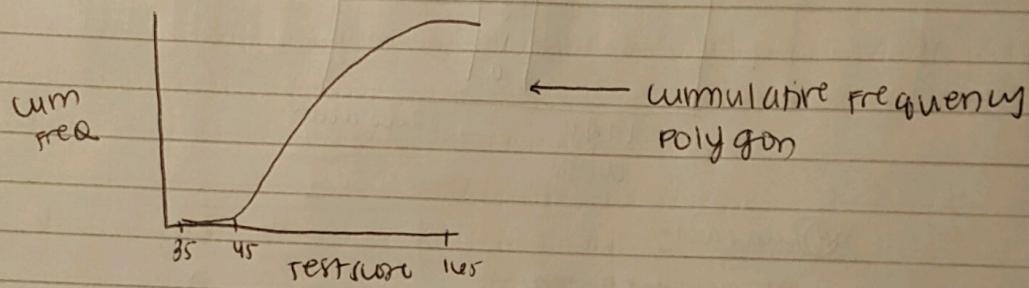
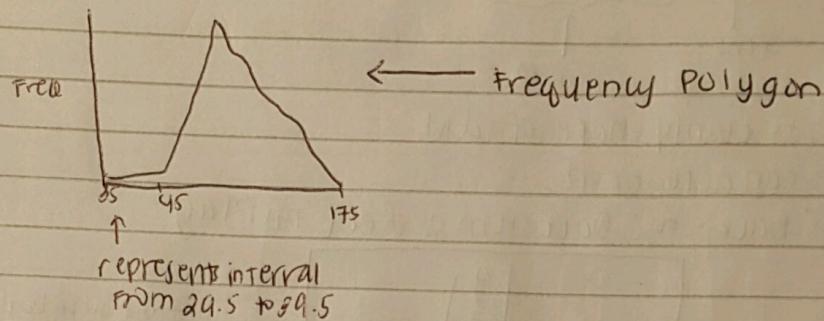
- for continuous, can bucket into whole #s

- can also be based on relative frequencies rather than actual frequencies

↳ proportion rather than #

3) Frequency Polygons → comparing distributions (overlaid)

- shape of distributions
- similar to histograms but for sets of data
- cumulative freq dist
* steps:
 - ① choose class interval
 - ② x-axis = score values
 - ③ middle of int w/ tick, label w/ mid value
 - ④ y-axis = freq of classes
 - ⑤ connect the pts



4) BOX PLOTS:

- identifying outliers & comparing distributions
- * steps:

1. 25th, 50th, 75th percentiles
2. Check p. 94 for terminology
Upper hinge, lower hinge
75th 25th

$$H\text{-spread} = UH - LH$$

$$\text{Step} = 1.5 \times H\text{-spread}$$

$$\text{Upper inner fence} = UH + 1.5$$

$$\text{Lower inner fence} = LH - 1.5$$

$$\text{Upper outlier} = UH + 2.5$$

$$\text{Lower outlier} = LH - 2.5$$

Upper Adjacent - largest value below UIF

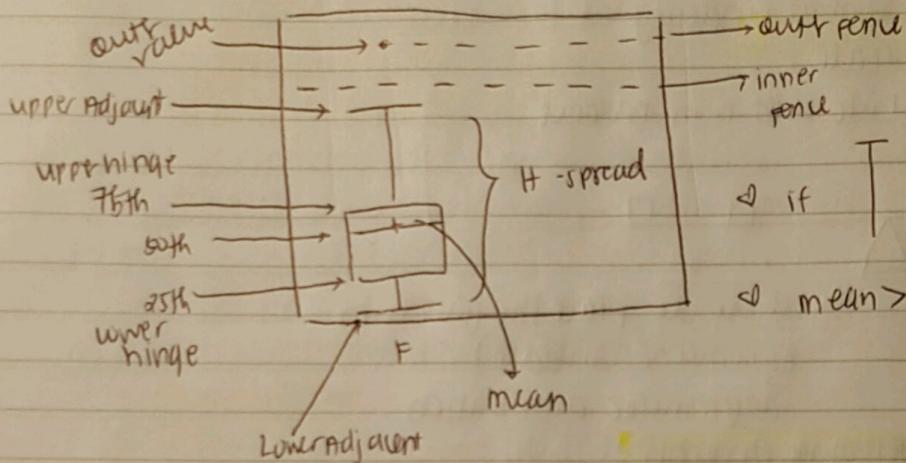
Lower Adjacent - smallest " Above LIF

Outlier value - beyond IF but not OF

Far Out value = value beyond it

terminology

- whisker \rightarrow VH & HH to VA & LA.
- mean = +
- DV = 0
- PDV = *



\curvearrowleft if T longer than I , then positively skewed
 \curvearrowleft mean $>$ median \rightarrow + skewed

- variation on box plots

a) Bar charts:

- frequency counts
- %, increases
- change over time
- compare means of diff experimental conditions = box plot

b) Line graphs:

- Bar graph w/ top of bars represented by pts joined by line (rest of the bar is suppressed)
- appropriate when both x and y axis display ordered (rather than qualitative) variables
- changes over time

c) Dot plots:

- random information

c) STATISTICAL LITERACY:

Eyecolor	Freq	Rel Freq
Blue	10	10/24
Brown	"	"/24
Grey	1	1/24
Green	4	4/24
Total	24	24/24

\rightarrow relative freq table