

(6)

03/04/19

③ SUMMARIZING DISTRIBUTIONS

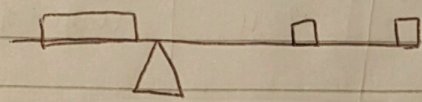
A) CENTRAL TENDENCY:

1) what is Central Tendency?

- location of the center of a distribution → comparing individual scores to dist.
- 3 definitions of central tendency:

① Balance scale:

- point at w/c dist is in balance



← where to place the fulcrum
(geometric middle)
asymmetric distribution

② Smallest Absolute Deviation

- number for w/c the sum of absolute deviation is smallest

③ Smallest Squared Deviation

- target that minimizes the sum of squared deviation ~~deviation~~

* Balance scale summary

- Positive skew: $\text{mean} > \text{median}$
 ↖ balance at mean
- Balanced / unskewed: $\text{mean} = \text{median}$
- Negative skew: $\text{mean} < \text{median}$
 ↖ balance on the mean
- Bimodal dist: bal pt = mean
- if the center of a distribution is defined as the balance pt, then the mean is at the center of the distribution

2) measures of central tendency

① Arithmetic mean

$$M = \frac{\sum x}{N}$$

⊗ geometric mean

$$M = \frac{\sum x}{N}$$

② median

- midpoint of the distribution → same # of scores above & below it
- 50th percentile
- mean of 2 middle #'s

* Abs. diff sum
- is there a formula for this?

* Sqd diff sum
- not really are about the use of dist what is it

③ mode

- most frequently occurring value
- continuous \rightarrow computed from grouped freq dist \rightarrow middle of interval

3) median and mean:

- mean = point on x-axis where a dist would balance
- median = value that minimizes the sum of abs dev
- mean = value that minimizes sum of squared dev.
- symmetric dist: mean = median
- bell shaped dist: mean = median = mode

* mean & median demo:

- depends on the skew of distribution

4) Additional measures of Central Tendency

① TRIMEAN

- ° weighted ave of 25th, 50th, and 75th percentile

$$\text{Trimean} = \frac{(P_{25} + 2P_{50} + P_{75})}{4}$$

° why is it over 4 and not 3? \rightarrow 2 P50

② GEOMETRIC MEAN

- ° multiply all #'s then take nth root of product

the analog of arithmetic mean, geometric mean.

$$\left(\prod x \right)^{\frac{1}{n}}$$

eg. 1, 10, 100

$$1 \times 10 \times 100 = 1000$$

$$\sqrt[3]{1000} \rightarrow 10$$

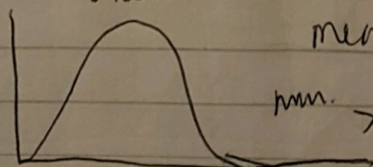
- ° close relationship to logs
- ° geometric mean only makes sense if all the #'s are positive
- ° good for averaging rates

③ TRIMMED MEAN

- ° hybrid of mean & median
- ° remove some of the higher & lower values & compute mean
- ° 10% \rightarrow Remove top 5% and bottom 5%

5) Comparing measures of Central Tendency

- ° symmetric \rightarrow mean = median = trimean = trimmed mean = mode
- ° differences exist w/ skew \rightarrow except bimodal
- ° Positive skew



mean > median

mean > trimean & trimmed mean

mean > geometric mean

03/05/19

B) VARIABILITY

- how much the #s in a dist. differ from each other

1) What is variability?

- "spread" of scores

↳ "variability", "spread", "dispersion"

- measures of variability:

1) RANGE

$$\text{highest score} - \text{lowest score} = \text{range}$$

2) INTERQUARTILE RANGE:

- range of middle 50% of scores in a dist

$$\text{IQR} = P75 - P25$$

- "H-spread" → upper hinge - lower hinge

- SEMI-INTERQUARTILE RANGE:

$$\rightarrow \frac{\text{IQR}}{2}$$

- symmetric dist: ^{mediant}SIQR contains 50% scores in dist.3) VARIANCE

- "how close the scores are to the middle of dist"

$$\sigma^2 = \frac{\sum (X - M)^2}{N}$$

- sample

$$s^2 = \frac{\sum (X - M)^2}{N - 1}$$

sample var used to estimate
- populatn var.M: mean of sample taken from
popul mean μ

- other formula:

$$\sigma^2 = \frac{\sum X^2 - \frac{(\sum X)^2}{N}}{N}$$

$$s^2 = \frac{\sum X^2 - \frac{(\sum X)^2}{N}}{N - 1}$$

④ STANDARD DEVIATION

- square root of variance
- useful when dist is normal or approx normal bec. prop. of dist with given # of SDs from μ can be calculated
- σ = population SD
- s = sample SD

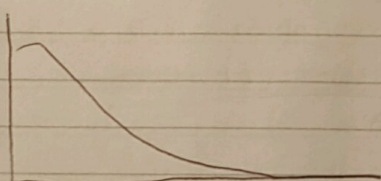
⑤ VARIABILITY SIMULATION:

⑥ VARIANCE ESTIMATION SIMULATION

3) shapes of distributions:

- numerical indexes for measure of shape

① SKEW



- large positive skew
- mean > median

- Pearson numerical index of skew:

$$\frac{3(\text{mean} - \text{median})}{\sigma}$$

σ

- third moment about the mean:

$$\frac{\sum (x - \mu)^3}{\sigma^3}$$

② KURTOSIS:

$$\left[\frac{\sum (x - \mu)^4}{\sigma^4} - 3 \right]$$

⑧ normal distribution \rightarrow to define no kurtosis

⑧ comparing distributions simulation

3) effects of linear transformations:

- if a vari X has mean μ , SD of σ and var σ^2 , then new vari Y created using linear transformation

$$Y = bX + A$$

will have mean $b\mu + A$

SD $b\sigma$

var $b^2\sigma^2$

4) Variance sum Lam I

$$\sigma_{\text{sum}}^2 = \sigma_m^2 + \sigma_f^2$$

\uparrow \uparrow
 var(males) var(females)

$$\sigma_{\text{diff}}^2 = \sigma_m^2 + \sigma_f^2$$

$$\sigma_{x \pm y}^2 = \sigma_x^2 + \sigma_y^2 \longrightarrow \text{VARIANCE SUM LAM}$$

• only applying when var are independent

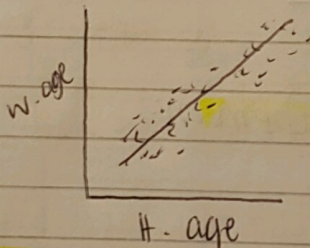
5) Statistical Literacy:

④ DESCRIBING BIVARIATE DATA:

- Data set of 2 vars = "bivariate"

A) INTRODUCTION TO BIVARIATE DATA

- summarizing data similar to / analogous to univariate
- eg. age
 - histogram, μ , σ
 - scatter plot



"positive association" - relationship

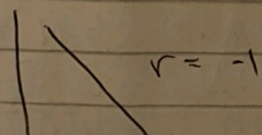
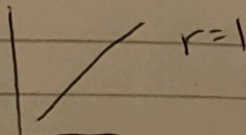
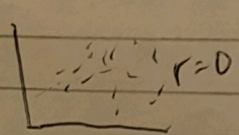
- cluster along line
"linear"

B) VALUES OF THE PEARSON CORRELATION

- Pearson product-moment correlation coefficient
- measure of strength of the linear relship b/w 2 vars.
 - ↳ if not linear, doesn't adequately report relship
- ρ = population } σ
 r = sample
- -1 to 1 range

\uparrow perfect negative relationship
 \uparrow perfect positive relationship

0 = no linear relationship



C) PROPERTIES OF PEARSON'S r

⊗ guessing correlation demo

- symmetric: "the corr of X with Y = corr of Y with X"
- unaffected by linear transformations
- range from -1 to 1

D) COMPUTING PEARSON'S r

$x, y \rightarrow$ deviation scores from the mean (means of $x, y = 0$)

$X, Y \rightarrow$ variables for w/c to get corr of

$xy \rightarrow$ high value for sum of $\sum xy$ means + corr, and low val for $\sum xy$ means negative corr

$$r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

$$\sqrt{\sum x^2 \sum y^2}$$

"r should be the same despite linear transformation"

$$r = \frac{\sum xy - \frac{\sum x \sum y}{N}}{\sqrt{\left(\sum x^2 - \frac{(\sum x)^2}{N}\right) \cdot \left(\sum y^2 - \frac{(\sum y)^2}{N}\right)}}$$

$$\left\{ \begin{array}{l} \text{alternative} \\ \text{formula.} \end{array} \right.$$

X) Range Restriction Demonstration

E) VARIANCE SUM LAW II

- when X and Y are not assumed to be independent - aka. when X and Y are correlated.

POPULATION

$$\sigma_{X \pm Y}^2 = \sigma_X^2 + \sigma_Y^2 \pm 2\rho\sigma_X\sigma_Y$$

corr b/n X and Y in pop.

SAMPLE

$$s_{X \pm Y}^2 = s_X^2 + s_Y^2 \pm 2r s_X s_Y$$

F) STATISTICAL LITERACY