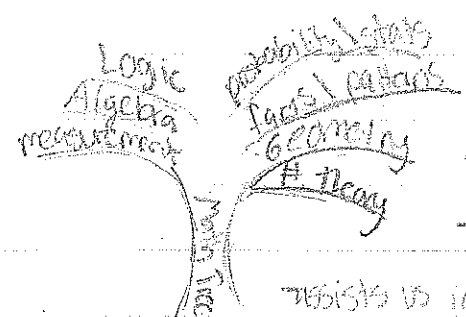


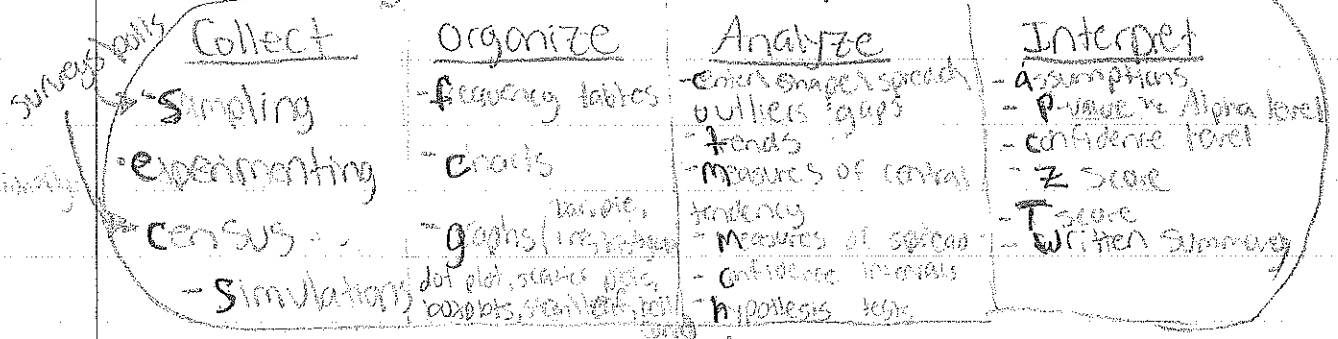
1.1 Notes: Def, Pop, Samples



• Statistical Data ...
 - is found almost everywhere
 - exists in a variety of forms
 - assists us in making informed decisions when faced w/ uncertainty
 \$ w/ bias

I. Definition

A. The collection, organization, Analysis, Interpretation of numerical info.



• Correlation between events \neq causation

II. Caution

- A. Can lie w/ or interpretation of stats
- B. Work diligently to use stats to show entire picture truth by good intentions & skill

III. Population

- A. Is all measurements of interest, entire group of those being measured
- B. Data from all



IV. Sample

- A. Part of population
- B. Smaller section of those being measured
- C. Data from...
- D. Random representative sample (to increase confidence level, ↑ size of random representative sample)
 - list w/ random #s are reliable for getting a RRS⁺
 - categorical (qualitative) or quantitative - data -

1.2 Notes: Random Sample

- Time of test (how long) - ratio
- Time of 1st class - interval
- Score on exam - ratio
- major - Nominal
- course evaluation - Ordinal

• Random # Table: 1) Start w/ a list 2) # all of sample

3) Prop something to land on # 4) Take first #s of sample (take first 3 digits if card 001-126)

5) Discard #s too big & continue on 6) Continue until you have amount you want

*decide if going down or right, etc

• Calc: 1) Go to "math" 2) Go to "PRB" 3) Go to "rand Int(1)" 4) Enter

5) Enter "randInt(1, 126, 3)" 1 = lowest # of sample 3 = how many you want to select
126 = highest # of sample

- stratified is good, systematic could be bad, convenience is terrible, cluster and multistage are good

I. Sampling Approaches

- Convenience - bad, doesn't represent pop, it's easy, statistician feel good, lots o' bias
- Stratified - pop into groups from a list & randomly sample
- Systematic - # everybody, choose a # and then every multiple of #, or #s that ends in #
- Cluster - pop into different geographic areas & randomly select from them
- Multistage - variety of methods

• A representative sample looks like the population

• Pop. consists of all subjects of interest (parameter)

• Sample is part of pop. & produces stats.

• Bias: • selection bias: sample doesn't represent pop.

- Undercoverage: low representation in sample

- convenience sample: easy to get sample that doesn't represent pop

- nonresponse bias: person chosen for sample don't participate

• Response Bias - occurs w/ research

- Social desirability: lies about info.

- leading questions: guides students to a certain answer

Blocking \ Binding Study Guide

2) A placebo treatment is a dummy treatment that can have no physical effect. It is a sham \ fake pill $\frac{1}{2}$ is NO treatment.

3) The placebo effect is that many patients respond favorably to any treatment, even a placebo, because of trust in the doctor and expectations of a cure. The mind! nothing becomes something.

4) Two lurking variables in the original gastric freezing experiment was that the experimental data was misleading because of the placebo effect and the data reflects any special features of that particular study, such as a physician with a soothing manner. No control group.

5) A control group ~~is~~ receives the placebo treatment in order to control the effects of lurking variables on the outcome, the experimental group receives the actual treatment and is the data that was being questioned by @slorder.

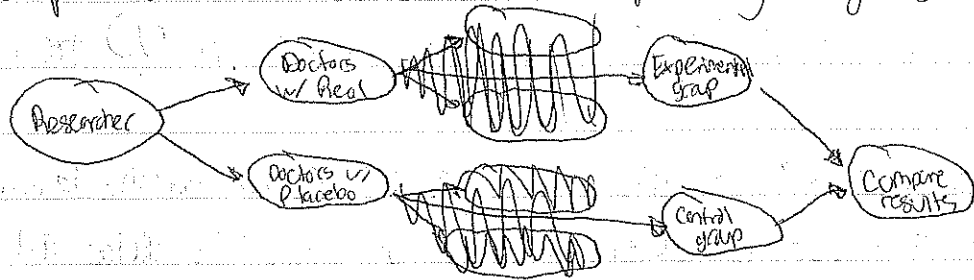


= placebo or other variable
> actual treatment

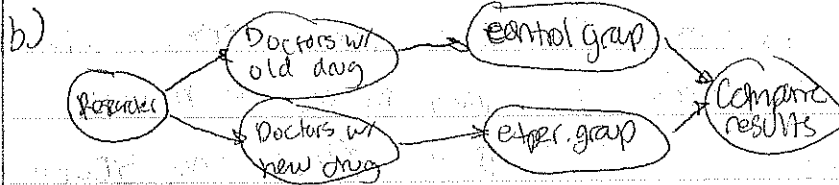
6) Uncontrolled studies give new therapies inflated success because without comparison of treatments, experimental results can be dominated by the details of the experimental arrangement, the selection of subjects, and the placebo effect. The result is often bias.

7) A double-blind experiment is that a researcher wants to see how effective a pill is, a researcher gives doctors a pill, not getting to know if the pill is the real pill or a sugar, placebo pill.

The patients also don't know what pill they are getting.



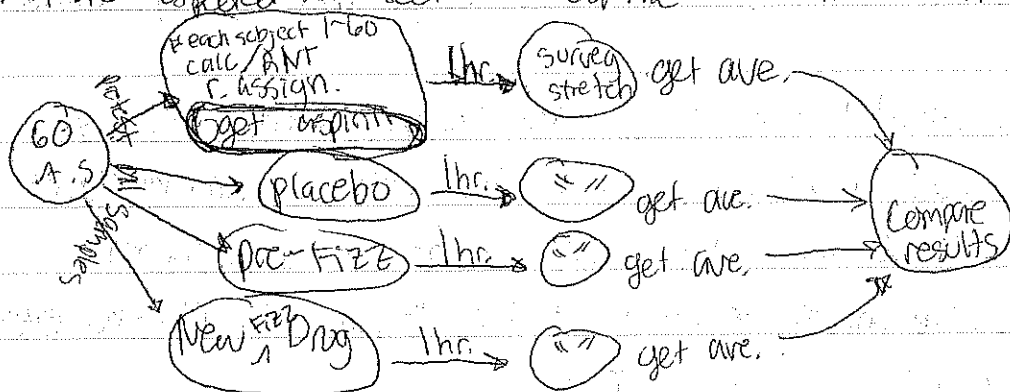
8) a) There is no control group and the data would be inflated.



c) No, this would create bias & affect results.

d) Yes, less chance of bias. Should be double-blind.

9) Block design is the random assignment of units to treatments is carried out separately w/in each block. A block is a group of experimental units or subjects that are similar in ways that are expected to affect the response to the treatments.



maybe have a combo of drugs.

Simulation Notes

I. Steps to making a simulation

- A. Describe Experiment - ex. 64 NFL how many field goals are made; what is 62%
- B. State Assumptions - ex. distance, wind could affect it, but assume the field goal is independent & field goal makes 90%.
- C. Assign Digits on a RNT to simulate - ex. simulate real life kick
- D. Repeat A LOT
- E. State Conclusion that corresponds to real life
 - Can have HS repeated for a simulation
 - written questions are better than verbal questions
 - have questions w/ even #s (-ex. law satisfaction 1 2 3 4) can't put middle
 - space #s evenly; words centered over #s
 - try to have ratio level questions
 - have clear, concise questions; make respondent feel comfortable; assist them in giving truthful, accurate answers (eliminate peers, fears, response bias, telescoping)
 - assure them it's confidential and anonymous

2.2

BASIC Graphs - Visual Representations of Data

I.

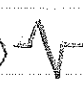
Overview.

- A. Info. can be represented quickly & efficiently
- B. Info. can be presented more attractively
- C. Keep it simple & understandable
- D. Viewers will look for approx 7 seconds
- E. It starts w/ a TITLE
- F. Ends w/ Analysis; look 4 highs/lows/trends & then ask why?

II.

Bar graphs

- A. Allow viewer to see differences in data easily
- B. Good tool for comparing several subjects

- C. Categorical Data on horizontal axis
- D. Label axes
- E. ~~Width~~ increments the same (can have a squiggle \ jump on axis) 
- F. Keep the same width ↳ Does it zoom in for strong comparisons or exaggerate
- G. Can have double bar graphs... USE a key!
- H. Put # inside the bar, NOT on top (atls height if put on top)
- I. Pareto charts - vertical bar graphs arranged tall to small

III. Circle/Pie Graphs

- A. Good for comparing parts of a whole; %S OF 100
- B. %d into slices
- C. Put # of % inside sector or a key
- P. Pieces can't overlap; total can't > 100%; can use "other" section

IV. Line Graphs

- A. Shows change / change in time
- B. When subject is tracked repeatedly
- C. Put a dot as high as data should go
- D. Connects the dots from L to R w/ a straight edge
- E. Can have multiple lines... Make a key (...; etc)
- F. Don't use squiggle in middle of axis
- G. Time series = special line graphs (time = x axis); notes trends

V. Picto Graphs

- A. Rare, use repeated symbols; Use a key; full symbol = amount in key
- B. Keep it simple; symbols same size & same spacing

VI. Dot Plots - newer graphs

- A. Similar to vertical bar graph, but uses dots to build upward
- B. Good for Analysis (center) shape / spread / outliers

VII. Histograms





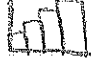
- A. Bar graph in which bars touch
- B. Width of bar has meaning
- C. Data is grouped into classes

- D. Useful w/ large amounts of data (big data)
- E. Overall shape assists us in recognizing patterns
- F. Miss out in weird data cuz everything clumped together
- G. Exact specific data is hidden
- H. How choose # of classes (bars)? Ask "Would info. be hidden? Most? Stand out?" < will nothing stand out?
- I. Usually 5-10; 7 is perf
- J. $\text{Width} = \frac{\text{Big} - \text{small}}{\# \text{ of classes}}$ (ALWAYS round up)
- K. Begin w/ smallest data value, then step by width

VIII. Frequency Tables

- A. Columns to help organize raw data before graphically displaying them
- B. Has class groups; tally marks; frequency #; class midpoints / class marks; relative frequency
- C. First class begins w/ the smallest raw
- D. Put weirds out into last column

IX. Types of Histogram

- A. Normal (Symmetric) 
- B. Rectangular 
- C. Bimodal  (peaks & valleys)
- D. Skewed Right 
- E. Skewed Left 

X. Analysis

- A. Center? B. Shape? C. Spread? (standard deviation) D. outliers?

II. Stem and Leaf Plots

- A. Good for quantitative B. Easy to make C. Groups data
- D. Shape, distribution, still see raw data E. Not good w/ large amount of data F. Back to back G. Can split stems

- H. Max/min & gaps, outliers

Period 1 vs 0000

-	0	8
0	26	89
1	025	
2	15	
3	45	

314 = 34 secs after 60 secs
 -018 = 8 secs before 60 secs

Analysis: Center = 11 median Shape: Right
 Spread = 6-15 secs = median
 Outliers: -8, 31

Skewed Right / Skewed

MPG for Mickey's Carlot

Ex. Toyota Camry 18 mpg	Nissan 22 mpg
General Pkx 24 mpg	Honda 31 mpg
Ford Explorer 23 mpg	Lexus 27 mpg

Domestic	Foreign
0	
6	1
43	27
	4

$3/2/7 = 27$ mpg for foreign & 23 mpg for domestic cars

Measures of Central Tendency

I.

The Averages

A. Mean B. Median C. Mode D. Trimmed Mean

- ID the Beast → Where? When? Predictions → Data Values → 10
- Mean • Median • Mode • 10% trimmed Mean • trend? • Graphs

II.

The Arithmetic Mean

A. $\frac{\text{sum}}{\# \text{ of terms}}$ B. \bar{X} for sample data (" μ " (μ) for population

D. Based on numeric HS E. Easily affected by outliers

III.

The Median

A. Middle value when data ordered small to large

B. Emphasizing position, not numerical value C. Resistant to outliers

D. Half data above median, half below median

IV.

The Mode

A. Most often, most occurring value B. Useful for categorical data

C. Can be bimodal, can have no mode D. Outliers don't affect it

V.

Trimmed Mean

A. Mean that resists the extremes B. Cut off % values from both ends

Measures of Variation Ch. 3.2

average

- It may not represent an entire set of ES well, so a cross reference is the range, the difference ^{measure of spread} between the largest & smallest values of data ^{measure of variation & dispersion}
- 3 measures of variation are range, variance, & standard deviation (σ) ^{pop} or s ^{sample}
- Pros of range: tells us difference between largest & smallest values in a distribution; total span; easy to get
- Cons of range: doesn't tell us how much other values vary from 1 another or from the mean; not important
- The measurement that helps us see how data is different from the mean is the standard deviation
- Divide by $n-1$ sometimes n by N other times because $n-1$ is for sample and N is for population and
- $E(8-5.5)^2$ for ex. 5, this new formula is bad cuz he is using a median and not the mean.

The pop. formulas for mean & standard deviation when compared to the sample formulas are: the difference \pm by N or $n-1$ for standard deviation. Pop. $\frac{\sum x}{N}$ Sample: $\bar{x} = \frac{\sum x}{n}$ Means are same.

$$\sigma = \frac{\sum (x - \mu)^2}{N} \quad s = \frac{\sum (x - \bar{x})^2}{n-1}$$

Ex. 2, 4, 6, 8, 10 of sample.

$\bar{x} = 6$ median = 6 mode = none Range = 10 - 2 = 8 (called total span)

x	$x - \bar{x}$	$(x - \bar{x})^2$
2	2 - 6 = -4	16
4	4 - 6 = -2	4
6	6 - 6 = 0	0
8	8 - 6 = 2	4
10	10 - 6 = 4	16
		$\Sigma = 40$

$$\frac{40}{n-1} = \frac{40}{5-1} = 10 = s^2 \quad \boxed{s = \sqrt{10}}$$

Ex. 19.8, 43.8, 36.1, 52.4, 63.1, 20.7, 46.3

Range = 63.1 - 19.8 = 43.3

$M = 40.31$
 $CV = \frac{\sigma}{M} = \frac{14.62}{40.31} = 36.27\%$

x	$x - \bar{x}$	$(x - \bar{x})^2$	x	$x - \bar{x}$	$(x - \bar{x})^2$
19.8	19.8 - 40.31 = -20.51	420.66	63.1	-22.79	519.39
43.8	43.8 - 40.31 = 3.49	12.18	20.7	-19.61	384.55
36.1	36.1 - 40.31 = -4.21	17.72	46.3	5.99	35.88
52.4	52.4 - 40.31 = 12.09	146.16			
$\Sigma = 1536.64$					

$s = 14.62$

CV (Coefficient of variation) = $\frac{s}{\bar{x}}$ or $\frac{\sigma}{\mu}$ - can compare between populations.

for spread:

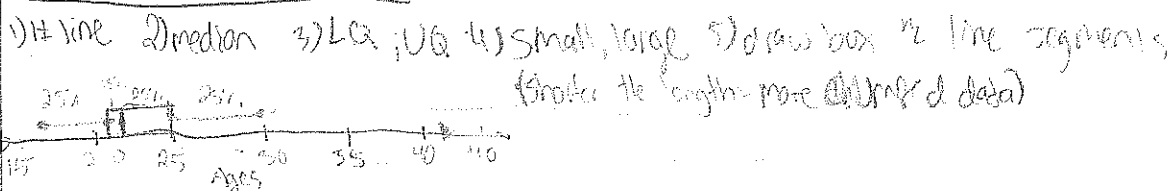
- CV: 0-10% = tight, 13-30% = low, 33-66 = medium, 67-95 = high, 100% or more = ^{data} very high
- Chebyshev's Theorem = $1 - \frac{1}{k^2}$ if $\frac{1}{2} = 75\%$ $k=2, 3, 4$ usually

- need a one sentence conclusion: "At least (75%) of the (students) would fit into the group that (volunteered) from (25.7) to (32.5) (hours each semester)."

3.3:

- Median = 50th Percentile. Percentile = roadmap, placement in data set, compare ranking
- Quartiles = special percentiles, splits data in 4's
- 1st quartile = $Q_1 = 25^{th}$ percentile. 2nd quartile = $Q_2 = 50^{th}$ percentile = median. 3rd quartile = $Q_3 = 75^{th}$ percentile
- Average #s small to large = median. look @ data to left of median & find median of subgroup
- look at data to right of median & find median of this subgroup
- IQR = inner quartile range = $Q_3 - Q_1 = 75^{th} - 25^{th}$ percentile = spread of inner data
- Box and Whisker Plots

to find Q_3 :

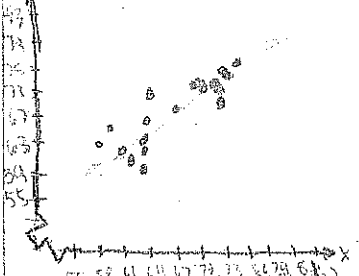


min = 17 $Q_1 = 21$ $Q_2 = 22$ $Q_3 = 25$ max = 31 = 5 # summary
 IQR = 25 - 21 = 4 Upper outlier = $Q_3 + 1.5(IQR) = 25 + 6 = 31$ = no outliers
 Lower outlier = $Q_1 - 1.5(IQR) = 21 - 6 = 15$ = no outliers

Ch. 1 Scatter Plots

Height v Wingspan of Students. (circles, proportional)

Response Variable (Dependent) Height (m)



Analysis: + association
 Trend: as $x \uparrow$, $y \uparrow$
 Longer arms tend to be attached to longer / taller bodies

= go w/ 2 quantities = can color code, have different circle/dot marks, etc

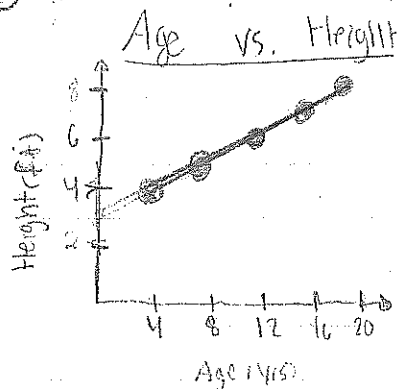
Explaining Variable (Independent) Weight

- can have outliers that are far from line
- dot that follows line but is apart from rest of data (extreme) is called an influential observation



age	x	4	8	12	16	20	6
ht.	y	4	5	6	7	8	3

$y = vx + 3$



This is most likely fake ^{some parents} as the child would be 3 ft. long when it is born. The slope is v . The y-int is 3. At some point, this person will stop growing. As age ↑, height ↑.

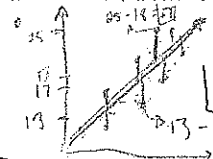
- Measures of Variation: S^2 (variance) - standard deviation - CV (coefficient of variation)
- IQR (interquartile range) - range - Chebychev's

extrapolation

Things in future might not be same in future, so don't extend line beyond parameters of data, safe to assume true if only slightly off from line (at yrs) but be sure to say nothing is certain

Linear Regression & The Coefficient of Determination

$y = a + bx$
 y-int: a , slope: b



residual: vertical distance from line to point
 $y - y_p$
 Line = Least Squares Regression Line / Line of best fit

$b = \frac{n \sum xy - (\sum x)(\sum y)}{n \sum x^2 - (\sum x)^2}$

n = # of lines
 x = x values
 y = y values

$a = \bar{y} - b\bar{x}$

\bar{x} = mean of x values
 \bar{y} = mean of y values

- put (x, y) on graph always! w/ a different mark to be above lowest value
- ins = influential observations, oas = outliers

Analysis for Scatterplot: +/- association, trend, outliers, C (H word), r^2

- r = how close points are to line = correlation coefficient; # between -1 & +1
- perfect correlation = -1 or +1 (sign depends on + or - slope/association)
- no linear correlation = 0
- r^2 = coefficient of determination (always between 0 & 1)

$r = \frac{n \sum xy - (\sum x)(\sum y)}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}$

Since r is [close] to 1, we have an indication of a strong + linear correlation between x variable & y variable

As x increases/decreases, y increases/decreases. Interpretation of r or r^2 = (% of variation in y) is explained by the Least-Squares Regression line. % variation in x remains unexplained (one conjectures)

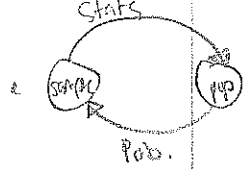
CORRELATION DOES NOT MEAN CAUSATION!

- Probability - a measure of chance/likelihood of an event happening between $0 < P < 1$
- Prob Notation: $P(A)$ = prob of event A "P of A" - Prob. fund: # of desired results / # of total possible results
- Odds = ratio of favorable / not favorable (f')
- $P(A')$ = complement = $1 - P(A)$; $P(A) + P(A') = 1$; sample space = listing of all possible outcomes

4.2

4.1

5.1



-10-

5.2

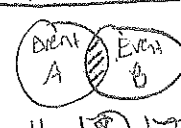
- Stats - sample is known • Draw conclusions about pop. based on results of an sample
- Prob: pop. is known; ask "what are chances... what is likelihood in this try?" we are guessing about sample
- Stat reasons from sample toward pop. • Prob reasons from pop. toward sample
- 1- Pepsi • 2- R^{1st} • 3- Coke $S = C, P, R, C, R, P, R, C, P, R, C$ • expected # = who guesses correctly = probability \times # of people in class = $\frac{1}{6} \cdot 24 = \frac{1}{3} \cdot 32 = 4$ (whole #)

COMPOUND EVENTS - two or more events that ^{happen together} occur one after the other

- dependency = when things change \rightarrow 1 affects the other.
- Independent: $P(A \& B) = P(A) \times P(B)$ • Dependent: $P(A \& B) = P(A) \cdot P(B|A)$ $P(B|A) = P(B)$



Or = 1 of other = sample space



And = Both events = sample space

$P(A \cup B) = Or = +$
 $P(A \cap B) = \& = \times$

- $P(\text{King} \cup \text{Ace}) = P(\text{King}) + P(\text{Ace}) = \frac{4}{52} + \frac{4}{52} = \frac{8}{52} = \frac{2}{13}$ & disjoint / mutually exclusive = no overlap
- $P(\text{Heart} \cup \text{King}) = P(\text{Heart}) + P(\text{King}) - P(\text{King of Heart}) = \frac{13}{52} + \frac{4}{52} - \frac{1}{52} = \frac{16}{52} = \frac{4}{13}$ & non-disjoint / not mutually exclusive = overlap

12 black Sox $P(\text{matching Sox}) = P(B|B) + P(B|A) + P(A|B) + P(A|A) + P(W|S) + P(S|W)$ ($= 0!$)
 6 blue Sox $= P(BL) \cdot P(BL|BL) + \dots$
 8 Brown Sox $= (\frac{12}{36} \cdot \frac{11}{35}) + (\frac{6}{36} \cdot \frac{5}{35}) + (\frac{6}{36} \cdot \frac{7}{35}) + (\frac{6}{36} \cdot \frac{4}{35}) + (\frac{4}{36} \cdot \frac{3}{35}) + (\frac{1}{36} \cdot \frac{0}{35})$

- 5 Grey Sox
- 4 Red Sox
- 1 Striped Sox

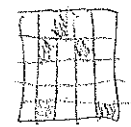
$P(\text{W win}) = P(BNSNSNSNS) = P(S \cdot S / S \cdot B / SS \cdot S / SSS \dots)$

conditional probability - probability that an event B will occur given the knowledge that an event A has already occurred. $P(B \text{ given } A)$

tree diagram = more attractive sample space (flows)

$P(B|A) = P(A \cap B) / P(A)$

3 Red Δ 2 Red \circ 4 Blue Δ 4 Blue \circ $Z = P(\Delta \text{ or } \circ) = \frac{7}{13} + \frac{4}{13} - \frac{4}{13} = \frac{11}{13}$

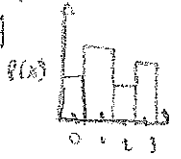


5.3

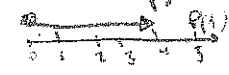
6.1

Discrete Random Variable - quantitative observations that are countable (finite)

- 2, -1, 0, 1, 2 ... fractions & decimals - ex: # of students who earned an A
- histogram & bars area of all bars added together = 1 (are entire sample space present no gaps mutually exclusive no overlaps)



Continuous Random Variable - countless, infinite observations

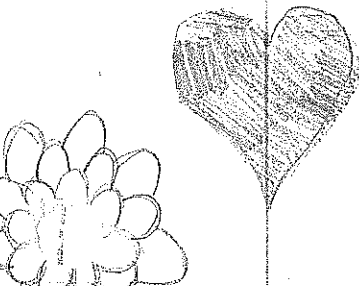


-ex: temp in room (has decimals & fractions) & time

Probability distribution - all probabilities + to 1. Has a mean & standard deviation

$\mu = \text{mean} = \sum x \cdot P(x) = \text{expected value} = \text{balance} / \text{central point} = \sum x p(x)$

$\sigma = \text{St. dev.} = \sqrt{\sum (x - \mu)^2 \cdot P(x)}$



23.2 - Ch. 6.1: 10) $P(60+ \text{ years-old}) = P(60-69) + P(70-79) + P(80+) = .25 + .091 + .018 = \boxed{.359}$

X	24.5	34.5	44.5	54.5	64.5	74.5	84.5
P(x)	.057	.097	.145	.247	.25	.091	.018
xP(x)	1.394	3.347	9.61	15.91	16.13	6.79	15.5

$\mu = \sum xP(x)$

added together = $\mu = \sum = \boxed{53.72}$

X	X - μ	(X - μ) ²	(X - μ)²	(X - μ) ² P(X)
24.5	-29.2	852.64	852.64	48.6
34.5	-19.22	369.408	369.408	35.83
44.5	-4.2	17.64	17.64	16.50
54.5	.8	.64	.64	1.869
64.5	10.78	116.21	116.21	27.05
74.5	20.8	432.64	432.64	39.37
84.5	30.78	947.41	947.41	17.053

$\sum (x - \mu)^2 P(x) = 18.6 + 35.83 + 16.5 + 1.869 + 27.05 + 39.37 + 17.053$

$= 146.272$

$\sigma = \sqrt{146.272} = 12.09$

$\sigma = 13.66 \text{ years old}$

$CV = \frac{\sigma}{\mu} \cdot 100 = \frac{13.66}{53.72} \cdot 100 = \boxed{25.43\%}$

2.2

Binomial Experiments - Bernoulli Experiments

- fixed # of trials
- outcomes (success) is
- trials = independent

Jacob Bernoulli was a Swiss mathematician who studied binomial experiments in the late 1600s extensively.

The sort of problems which have exactly 2 possible outcomes is now called binomial / Bernoulli Experiments.

The central problem of a binomial experiment is finding the probability of r successes out of n trials.

Binomial Experiments don't work with dependent situations. Only independent!

If asked what car is recommended to buy, n (# of trials) = # of teachers asked and # of outcomes possible = as many models of cars there are

This is NOT a binomial experiment, cuz 72 possibilities

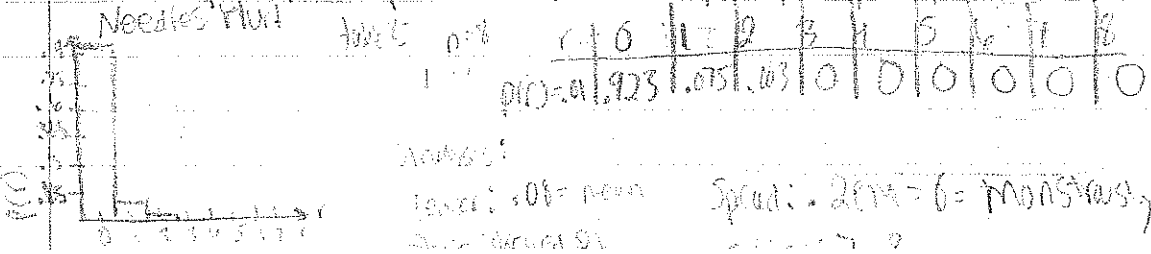
$P(r) = C_{n,r} p^r q^{n-r}$
n = # of trials, p = prob. of success, q = 1 - p = prob. of failure

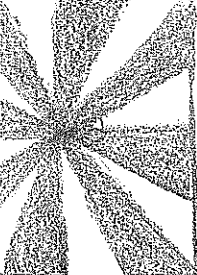
Ex. $P(r=6) = C_{10,6} (.75)^6 (.25)^{10-6} = .1460$

There exists a 14.60% chance that we will get 6 out of 10 people would be concerned about the security of the email.

If $r \geq 6$, then add probabilities of r=6, r=7, r=8, r=9, r=10.

15.4.10



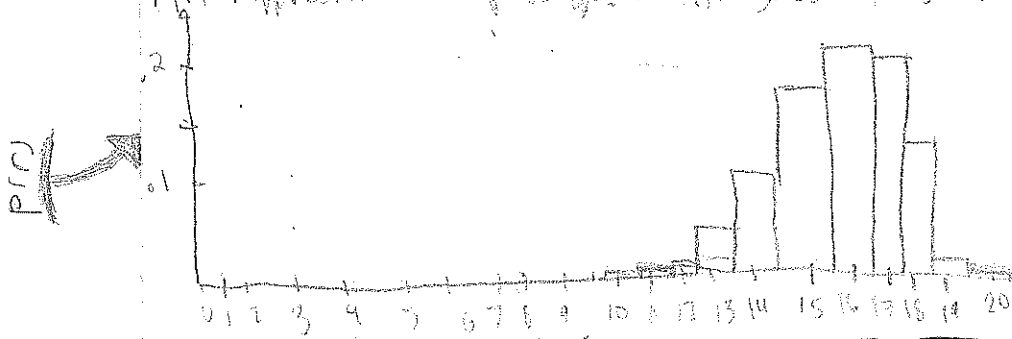


⇒

Assume people learn independently.
 $p = 0.8$ grade 2 700 = student takes class

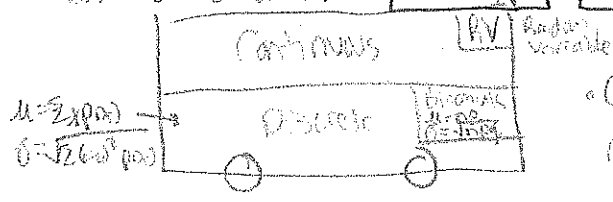
S = Pass F = Fail $n = 20$ $r = 0-20$

Art Appreciation



• $\text{Center} = \mu = (20)(0.8) = 16$
 • $\text{Type} = \text{skewed left}$
 • $\text{Spread} = \sqrt{npq} = 1.79 = \sigma$; low
 • $\text{differs } 0-8$

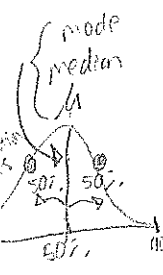
For sum of binomials: $\mu = np$ $\sigma = \sqrt{npq}$



• Calculator: 2nd VARS \Rightarrow binomPDF (just 1 r) or binomCDF (arg. of 5) \Rightarrow binomPDF (n, p, r) or binomCDF (n, p, starting r)

not a bad
 after all,

- Law of Large #s: Mean result of a large # of independent trials comes close to the true mean of the distribution
- Myth of small #s = true mean will occur with few # of trials
- $\sigma^2(x+y) = \sigma^2 x + \sigma^2 y$ Variance of $x+y$ = Variance of x + Variance of y



- High/Low? trend? why? Analysis = bar, pie, picto, pareto, time plot, line
- CSSO Analysis = histogram, stemplot, dot & whisker, dotplot
- For Scatter Analysis = r - association, trends, influential observations, outliers, correlation ($r = \#$ of words) \wedge $r^2 = \text{coefficient of determination} = 49.73\%$ of the variation is explained by the LSRL \wedge the variation in x .

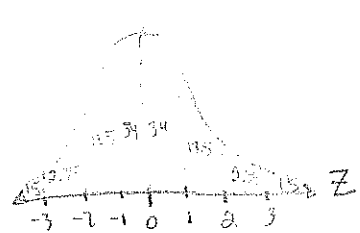
stratified = group people based on characteristics \wedge cluster = group people based on location

h. 7.

CONTINUOUS Random Variables - 7.1

- Normal curve - area under curve = 1, curve is bell-shaped w/ highest point over mean (μ), curve is symmetrical about vertical line through μ , curve approaches but doesn't cross horizontal axis, inflection points occur @ $\mu \pm \sigma$
- Empirical rule = 68% of data will fall w/in $\mu \pm \sigma$, 95% w/in $\mu \pm 2\sigma$, 99.7% w/in $\mu \pm 3\sigma$
- Need $P(14 \leq X \leq 35) = ?$ There exists a 47.5% chance that...
- Z score = standard deviation distance from $\mu \rightarrow Z = \frac{(X - \mu)}{\sigma}$
- raw score $X = Z\sigma + \mu$ (derive from Z score equation)

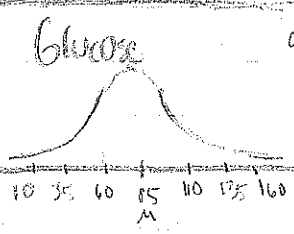
ex. $\mu = 4.8$ $\sigma = 3$ $P(X < 4.2)$ $Z = \frac{4.2 - 4.8}{3} = -0.2$ \rightarrow $Z < -2.06$ \Rightarrow $1.28 = \frac{X - 4.8}{3} \rightarrow 5.184 < X$



Standard Normal Curve →
 # (6) Draw curve left of $z = -0.47$ & find area



7.3 #25



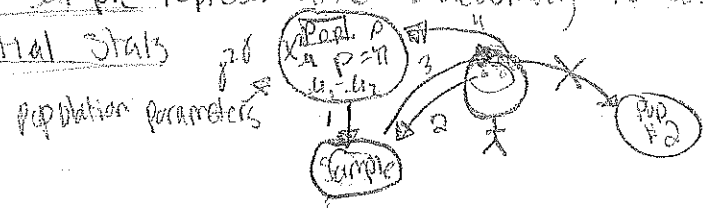
- a) $P(x > 60) \Rightarrow z = \frac{60-85}{25} = -1.00 \rightarrow 1 - .1587 = .8413$
- b) $P(x < 110) \Rightarrow z = \frac{110-85}{25} = 1.00 \rightarrow .8413$
- c) $P(60 < x < 110) \Rightarrow .8413 - .1587 = .6826$
- d) $P(x > 140) \Rightarrow z = \frac{140-85}{25} = 2.20 = .9861 \rightarrow 1 - .9861 = .0139$

There exists a 1.39% chance that an adult under 50 years old will have > 140 milligrams of glucose in their blood.

7.4

Sampling Distributions - review of basic methods & terms

- Population - all measurements of interest
- Sample - subset of measurements from population
- Random Sample - representative & accurately reflects population
- Inferential Stats



- 1) Get a random representative sample
 - 2) Observe & collect data
 - 3) Analyze results
 - 4) Apply back to pop
- Inferential stats can not reasonably be applied to a different population than the one in which sample was taken

Sample Statistics

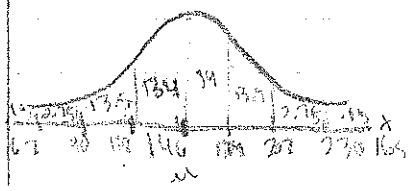
- Statistic - numerical descriptive measure of a sample $\hat{p} = \text{stat} = \frac{x}{n}$ & \bar{x}, s, s^2, \dots
- Use Sample Stats to make inferences (conclusions, decisions) about population parameters
- Ch. 1-4: Descriptive Stats - Organize, Summarize #s
- Ch. 5-7: Probability Theory & Distributions
- Ch. 8-11: Inferential Stats - Methods of using sample to obtain reliable info. about pop. (margin of error)

Binomial Distribution

- $P(x = n) = (1-p)^{n-1} p$ $n = \# \text{ we want}$ $p = \text{prob of success}$ $P(x > n) = (1-p)^n$ $u = 1/e$
- $P(x = 17) = ?$ $p = .25$ $q = .75$ $P(x = 17) = (1-.25)^{16} \cdot .25$ $P(x = 17) = .007506$
- $P(x > 17) = ?$ $P(x > 17) = (1-.25)^{17} = .007517$
- $P(x \leq 17) = 1 - P(x > 17) = 1 - .007517 = .992483$

-Ex.

- A female ≈ 146 lbs $\sigma \approx 28$ lbs ≈ 1100 lbs chair broke \rightarrow caused pain \rightarrow suing for \$1M
- Normal cdf (lower value, upper value, μ, σ)
- Never use normal pdf
- inv Norm (Area, μ, σ) if given % or Area





7.5

- Standard error = $\sigma_{\bar{x}} = \sigma / \sqrt{n}$ $z = \frac{\bar{x} - \mu}{\sigma_{\bar{x}}}$
- Theorem 7.1: For a normal distribution: 1) \bar{x} is a normal distribution
- 2) $\mu_{\bar{x}} = \mu$ 3) $\sigma_{\bar{x}} = \sigma / \sqrt{n}$
- Central Limit Theorem: If $n \geq 30$, then Theorem 7.1 applies for any type of distribution

8.1

- $E = \text{error} = t_c \cdot \sigma / \sqrt{n}$ (z score \times standard error); $\bar{x} - E < \mu < \bar{x} + E$ = point estimate
- sentence: If we took 100 samples of $n=4$, we expect to capture population mean μ (19 times) out of 100.
- $S \approx \sigma$ when n is big (≥ 30)

7.6

- If $np > 5$ and $nq > 5$, then r has a binomial distribution that is approx. normal
- If $np > 5$ and $nq > 5$, then the random variable $\hat{p} = r/n$ can be approximated by a normal random variable x w/ $\mu_{\hat{p}} = p$ and $\sigma_{\hat{p}} = \sqrt{pq/n}$ (proportions and %) not ave.

8.2

- $E \approx t_c \cdot s / \sqrt{n}$ "degrees of freedom = $n-1$ " When σ is unknown
- Ex. $n=4$ d.f. = 3 $c=90\%$ $\bar{x}=33$ $s=5$
- $E = t_c (s/\sqrt{n}) = 2.353 (\frac{5}{2}) = 5.8825$ 27.12 < μ < 38.88
- If we took 100 samples of size $n=4$, we expect to capture the population mean μ of brain cells of dead science teachers 90 times out of 100.
- W.S. Gossett - 1908 chemist @ Guinness brewery - found t_c .
- As n / d.f. increases, t_c decreases and goes toward z_c .
- Ex. $n=14$ d.f. = 13 $\bar{x} = 1.93$ $s = .38$ $c = 90\%$
- $E = t_c (s/\sqrt{n}) = 1.771 (.38/\sqrt{14}) = .1799$ 1.25 < μ < 1.61
- If we took 100 samples of size $n=14$, we expect to capture the population mean μ of mL of depth perception 90 times out of 100.
- $n = (z_c \cdot \sigma / E)^2$ (round up to nearest whole #)
- t curve is less tall and wider than z score if p estimate is known reg. (un) this to get in

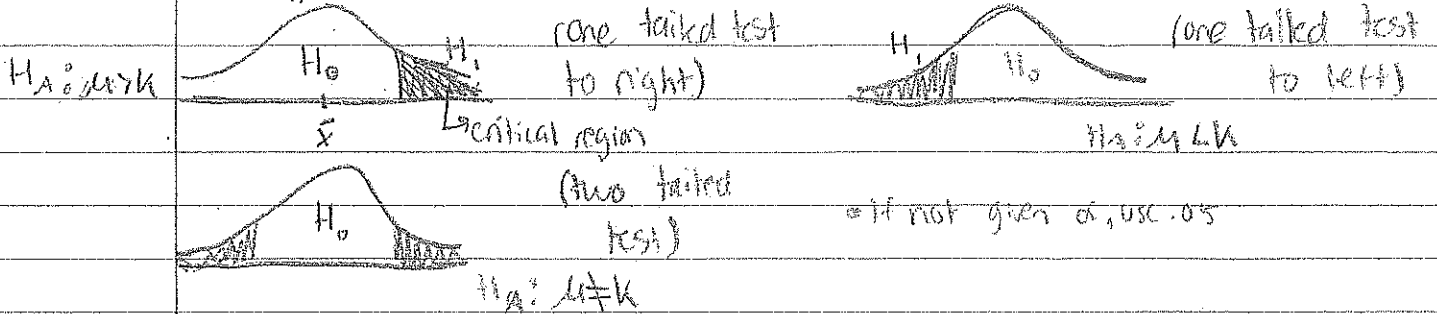
8.3

- To estimate proportion p : $\hat{p} - E < p < \hat{p} + E$ $E = z_c \cdot \sqrt{pq/n}$ $n = \frac{1}{p(1-p)} (\frac{z_c}{E})^2$ if p est. is known
- If $np > 10$ & $nq > 10$ → then distribution is normal
- A normal approximation to the Binomial is justified ("Anatbig") say this!
- Sentence: If we took 1000 samples of size $n=4$, we expect to capture population proportion p of x (950) times out of 1000.

9.1

- Hypothesis: an assumption or belief about a parameter
- Hypothesis testing: procedure based on sample info. by which I accept/reject hypothesis
- Null hypothesis: H_0 ; the hypotheses we are testing. "No change. No difference."
- Alternate hypoth.: H_1 or H_a ; hypothesis to be accepted if null is rejected. "ours is better"

Alt. hypothesis will always have 1 of 3 looks:



- A type I error is worst - when one rejects null when null is still true
- A type II error is also bad - when one accepts null when it no longer should be

Error table:

	"Accept"	"Reject"
H_0 T	no error (circled)	Type I error
H_0 F	Type II error	did right thing (circled)

Level of significance (α) = P(willing to make a type I error)
 $\alpha = .01$ - willing to risk Type I error 1/100 times

Beta = β = P(making a Type II error)

Power of a test ($1 - \beta$) = P(bottom Rt. of table)

- As $\uparrow n$, we \uparrow power

4 Ingredients to a Great Statistic:

- Null hypothesis - Alt. hypothesis - sample data ($\bar{x} \rightarrow z, t$) - critical value

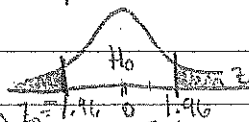
EX: #12 a) $H_0: \mu = 30$ b) $H_A: \mu \neq 30$ c) $H_A: \mu > 30$ d) $H_A: \mu < 30$

Conclusion: Accept / Reject H_0 , Accept / Reject H_A , at what level of sig.
 : There exists sufficient statistical evidence to suggest that ... (H_0 / H_A)

EX: $\mu = 130^\circ F$ $\sigma = 8$ $\bar{x} = 131.08^\circ F$ $\sigma = 1.5 F$ $\neq 130$

$H_0: \mu = 130$ The company claims that the fire sprinkler is activated at a temp. of 130°

$H_A: \mu \neq 130$



$$z = \frac{131.08 - 130}{1.5/\sqrt{8}} = 6.148$$

conc: H_0 (H_A) $\alpha = .05$ $z_c =$

: There exists sufficient statistical evidence to suggest that the sprinkler is indeed activate @ a temp $\neq 130^\circ F$.

P-value: Assuming H_0 is true, probability that test statistic will take on values as extreme as / more extreme than observed test statistic

- smaller the P-value computed from sample data, the stronger the evidence against H_0

Stats

- P-value = probability of chance; low P-value = results not due to random chance alone
 → P-value for 2-tailed test = $2 P(Z > |z_{\alpha/2}|)$

- If P-value $\leq \alpha$, then reject H_0 . If P-value $> \alpha$, then don't reject H_0 .

- There exists insufficient statistical evidence to suggest that...
 if keep on accepting H_0 say this

9.3
 Testing a
 population p

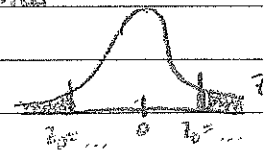
- $Z = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}}$ 'dubio' ✓ that $np > 10$ & $nq > 10$ & $n < \frac{1}{10} N$ then write ANOVA IS and assure there is more than 10x sample size
- looking for true population percentage p (could also be π)
- state "assuming this is a good, representative random sample" for ch. 8 p 9

10.1
 Hypothesis Testing
 w/ Dependent
 Groups

- 2 samples are dependent if each data value in 1 sample can be paired in meaningful way w/ corresponding data value in other sample
- involve paired data samples that help us draw conclusions about the difference of 2 groups
- frequently happens in before/after sits in which some item is measured prior to and after treatment (test before learning @ beginning of year, 1 after test)
- sometimes there is a matching link, a natural match for collecting pairs even w/out a before/after test. ♥
- advantages: reduce danger of extra factors/variables except for characteristics we want? reduces variance
- \bar{d} (d bar) is mean difference between paired data. S_d is sample standard deviation. (after - before) always 0
- test is same as a t distribution. $H_0: \mu_d = 0$ (no difference before and after).
 $H_A: \mu_d \neq 0$. $\bar{d} \rightarrow t$ score = $\frac{\bar{d} - \mu_d}{(S_d/\sqrt{n})}$
- state why dependent

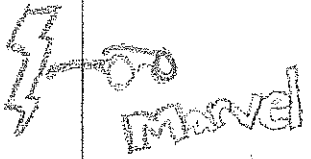
10.2
 Two Diff.
 Populations

- 2 samples are independent if each sample is completely unrelated
- $H_0: \mu_1 - \mu_2 = 0$ or $H_0: \mu_1 = \mu_2$ (Identify both μ) $H_A: \mu_1 \neq \mu_2$
- Large samples - difference of averages
- if both $n \geq 30 \rightarrow CLT \rightarrow$ Normal



$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{(\sigma_1)^2}{n_1} + \frac{(\sigma_2)^2}{n_2}}} = \frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{\frac{(\sigma_1)^2}{n_1} + \frac{(\sigma_2)^2}{n_2}}}$$

- Confidence interval: $(\bar{x}_1 - \bar{x}_2) - E < \mu_1 - \mu_2 < (\bar{x}_1 - \bar{x}_2) + E$ where
 $E = Z_{\alpha/2} \sqrt{\frac{(\sigma_1)^2}{n_1} + \frac{(\sigma_2)^2}{n_2}}$



CLAIMS BUT THE

$$t = (\bar{x}_1 - \bar{x}_2) / s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \quad s = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}} \quad d.f = n_1+n_2-2$$

$$E = t_c \cdot s \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

Pool standard deviation only if standard deviation pool =

10.3

$n\bar{p} > 10$ & $n\bar{q} > 10$ for both populations \rightarrow ANATBEL

$n < .1N$ for both pops. = sample = random, independent, & representative

$H_0: p_1 = p_2$ or $\pi_1 = \pi_2$ $z = \hat{p}_1 - \hat{p}_2 / \sqrt{(\hat{p}_1\hat{q}_1/n_1) + (\hat{p}_2\hat{q}_2/n_2)}$

$\bar{p} = r_1 + r_2 / (n_1 + n_2)$ $\bar{q} = 1 - \bar{p}$

11.2

asks if population fits given pattern

$H_0: \chi^2 = 0$ the pop. fits the pattern/distribution, the $\chi^2 > 0$ pop. doesn't fit a different pattern

$\chi^2 = (O-E)^2 / E$ $E = \% \text{ of } n$ $d.f. = \# \text{ of categories} - 1$ $n = 137$

Ex.	2005	2015	E	(O-E)	(O-E) ²	(O-E) ² /E	p-value
V	4%	3	4% * 137 = 5.48	-2.48	6.15	1.122	calc: 2 dpl VARS
B	65%	77	89.05	-12.05	145.2025	1.6306	χ^2 cdf (0, upper lim = int
Safety	13%	0	17.81	-8.81	77.6161	4.3580	value, degrees of freedom)
Meds	12%	41	16.44	24.56	603.1936	36.7006	
ot	6%	7	8.22	-1.22	1.4884	0.18107	$\chi^2 = 43.982$

$H_0: \chi^2 = 0$ Distribution of workers' wants of 2015 fits the distribution of workers' wants of 2005

$H_a: \chi^2 > 0$ $d.f. = 4$

There exists sufficient statistical evidence to suggest that 2015 feelings are not fitting following 2005 feelings.

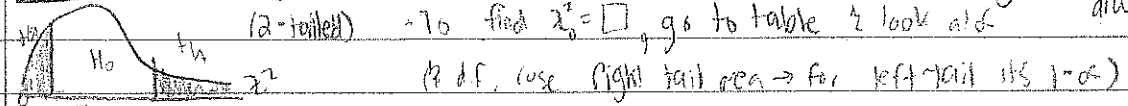
11.1 205

Assumptions: Each cell has ^{expected} $n_{ij} > 5$ sample is random, representative, & independent.

$n < .1N$

11.3

$\chi^2_{\text{variance}} = (n-1)S^2/\sigma^2$ ($d.f. = n-1$) $H_0: \sigma^2 = k$ $H_1: \sigma^2 \neq k$ Only χ^2 that isn't always $>!$



To find $\chi^2_0 = \square$, go to table & look at $d.f.$ (use right tail area \rightarrow for left tail it's $1-\alpha$)

$\chi^2_0 = L$ $\chi^2_0 = U$ $H_0: \sigma^2 = k$ the variance remains the same

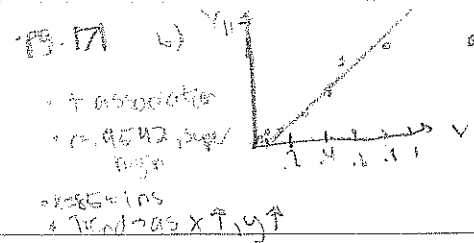
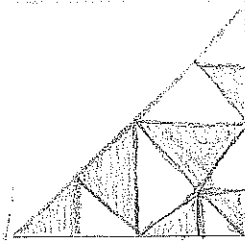
There exists insufficient statistical evidence to suggest that the variation in x is (indeed) (more than, less than, \neq) than claimed by (...)

Assumptions - random, rep, independent sample; $n < .1N$

perfect-linear correlation:		$r = 1$ (corr.)	no linear corr:	
-----------------------------	--	-----------------	-----------------	--

Pg. 164 #15

$r^2 =$ coefficient of determination \rightarrow between 0 & 1 $\rightarrow r^2\%$ of variation in y is explained by variation in x & LSRL. $1 - r^2\%$ remains unexplained.



• hi 2 good fit \Rightarrow yes sig.
 $\hat{\beta} = 17.11K - 2.48$

IDK

residual = $y - \hat{y}_0$

- + association
- r = 0.9542
- r = 0.95 = ins
- trend \rightarrow as $x \uparrow, y \uparrow$

• $r =$ linear coefficient \rightarrow between -1 and 1; how lined are dots = determinant of correlation

Analysis for graph = ins/out, trends, & association, correlation ($r = 1 \neq$ word)

↳ influential obs. \rightarrow as $x(f), y(f)$

representat

11.4

POP	Sample	Assumptions: - set (x, y) of ordered pairs is random & normal distributions
M, μ, σ^2	$S, \bar{x}, s^2, \hat{\mu}, \hat{\sigma}^2$	(majority in middle)
$P = \pi, \delta, \sigma^2$	$\hat{p}, \hat{s}^2, \hat{x}, \hat{y}$	
$\mu_1, \mu_2, \sigma_1, \sigma_2$	$r, \hat{\beta}, \hat{\sigma}^2 = a + bx$	
$y = \mu + \beta x$	$\hat{y} = a + bx$	

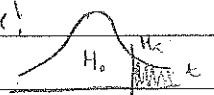
• To test β : $H_0: \beta = 0$ $x \& y$ have no linear relationship correlation. $H_a: \beta \neq 0$ [Curve]

$t = r\sqrt{n-2} / \sqrt{1-r^2}$ d.f. = $n-2$

• To test β : $H_0: \beta = 0$ Pop. slope of $_ = 0$ & LSRL slope is flat ($_$) $H_a: \beta \neq 0$

• $t = b - B / Se$ / $\sqrt{\sum x^2 / n - (\sum x)^2}$ t -value for $P > \beta$ will be the same!

• Calc. Lin Reg T test & input values in stats \rightarrow test



• Standard Error = on ave. residual \rightarrow $Se = \sqrt{\sum (y - \hat{y})^2 / n - 2}$ \rightarrow fit/representative of x, y data

• Pattern = LSRL isn't good fit \rightarrow why aren't linear \rightarrow perhaps $x \& y$ have

$Se = \sqrt{(\sum (y - \hat{y})^2) - a \sum y - b \sum xy} / n - 2$ d.f. = $n - 2$

• Confidence interval: $y_0 - E \leq y \leq y_0 + E$ $E = t_{\alpha} Se \sqrt{1 + \frac{1}{n} + (n(x - \bar{x}))^2 / n \sum (x - \bar{x})^2}$

• if $\alpha = 0$ -value, still reject null (H_0)

Review

• Graph for discrete = histogram; $\mu = \sum (x \cdot p(x))$; $\sigma = \sqrt{\sum (x - \mu)^2 \cdot p(x)}$

• binomial = 2 outcomes w/ prob. being $_$, $n =$, ind. ; $\mu = np$ $\sigma = \sqrt{npq}$

• geometric = $_$ $\mu = 1/p$ \rightarrow until have success

$\mu_{x+y} = E(x) + E(y)$ $\sigma_{x+y}^2 = (\sigma(x) + \sigma(y)) \rightarrow \sigma_{x+y} = \sqrt{(\sigma(x))^2 + (\sigma(y))^2}$