

1) ~~P XV~~ statistics is the science of extracting information from data. 0282-1

2) ^{Ch 1} ~~P 4~~ Individuals are the objects described by a data set. A variable is a characteristic recorded about an individual.

2) * ~~P 4~~ A categorical variable takes on several categories

eg race, hair color, gender
often count the number in each category or find the percentage
adding the categories does not make sense

A quantitative variable takes on numerical values

eg height, # phone calls received
adding the values makes sense

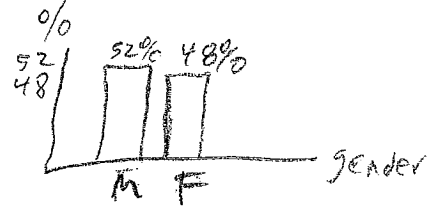
3) ~~P 4~~ The distribution of a variable tells what values it takes and how often.

4) ~~P 6~~ The distribution of a categorical variable lists counts (frequencies) or percents

ex SIU student gender

M	52%
F	48%

5) ~~P 7~~ bar graph vertical axis height = percent or count
horiz axis has categories



A pie chart assigns percents of a circle area to categories and must use all of the categories



Frequency histogram is a graph of the distribution of a quantitative variable.

→ P8-9 a) Divide the range of the data into classes of equal width. Each observation should fall in exactly one class.

ex semester credit hours 40 students

COMMON mistake { 0-4 4-8 8-12 12-16 16-20 }
 which class does 12 go into?
 count

classes	$0 < \text{hours} \leq 4$	1	
	$4 < \text{hours} \leq 8$	2	
	$8 < \text{hours} \leq 12$	14	12 goes in this class
	$12 < \text{hours} \leq 16$	16	12 does not go in this class
	$16 < \text{hours} \leq 20$	7	

b) find the count of observations with each class

c) On the horiz axis mark the scale of the variable. Bar height = count

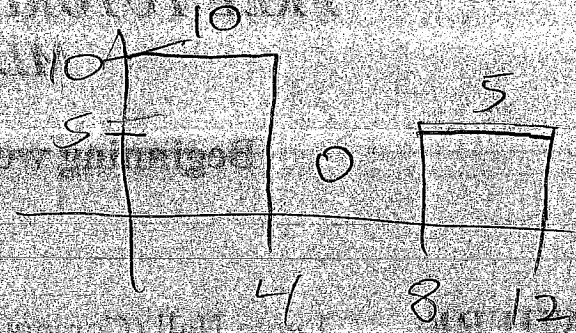


d) label the top of the bar with the count (NOT in box)

between point 8 and 9
 bar graph gaps 0.75

histogram no gaps unless a class
 has a count of 0

$0 <$	≤ 4	10
$4 <$	≤ 8	0
$8 <$	≤ 12	5



8) ^{p282-2} Since the bars have equal width, the area of each bar is proportional to the percentage of individuals in each class.

1st paragraph p10 is not important

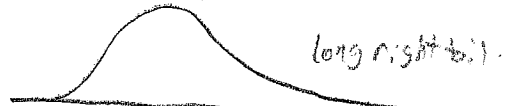
9) ^{p11} To interpret the histograms Look for an overall pattern and deviations from the pattern. shape, center, spread

10) ^{p11} Outliers are observations that fall outside the overall pattern.

11) ^{p12} 3 common shapes:



approx. symmetric
approx. mirror image about midpoint
eg 1.3 p13



right skewed
eg income data
fig 1.4 p13



left skewed
eg scores from easy exam

12) ^{know} Stem plots ^{p15} are for small data sets (p17) Final exam scores

3	0 4	5 6 7 8	six
4	2 6	6	three
5	0 1	2 4 7 8 9 9	two
6	0 1	4 4 4 5 5 5 6 7 8 8 8 9	four
7	0 1	1 2 2 2 3 3 4 5 6 6 8 8 8 8 9	eight
8	1 1 1	2 2 2 2 3 4 4 4 5 6 7 7 7 8 8 8	nine
9	1 2 2		

Stems leaves

stem ten
leaf one

a) leaf = final digit

divide data into groups = stems that contain all but the last digit
eg for exam data 3, 4, ..., 9 stand for scores in 30's, 40's, ..., 90's

b) Place stems in order in a vertical column
(text suggest smallest on top largest on bottom)
Draw a vertical line to separate stems from leaves

c) write each leaf to the right of its stem in increasing order

d) write stem and leaf units on plot
(not in text)

13) p 17 variations on the stem plot

a) If there are too many stems, round the data, then make a stem plot

b) Split the stems eg 7 0-4
7 5-9

one stem for 70, 71, 72, 73, 74

another stem for 75, 76, 77, 78, 79

14) p 17 a time plot plots each

observation, vs the time it was measured

EBay stock price



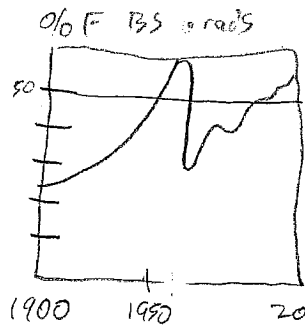
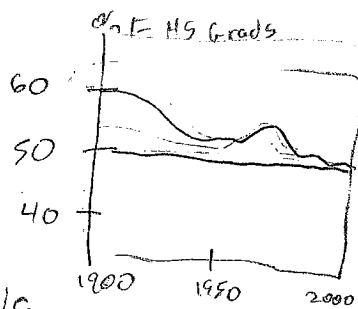
MONTH

19 brands

fat contents of peanut butter

22) 30, 40, 56, 69, 30, 41, 50, 65,

39, 40, 56, 62, 36, 45, 50, 44, 53, 56



ex)

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

Scratch

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

Stem ten
leaf one
 $5(0) + 3(1) = 53$

split stems 0-4

hardly ever use this

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

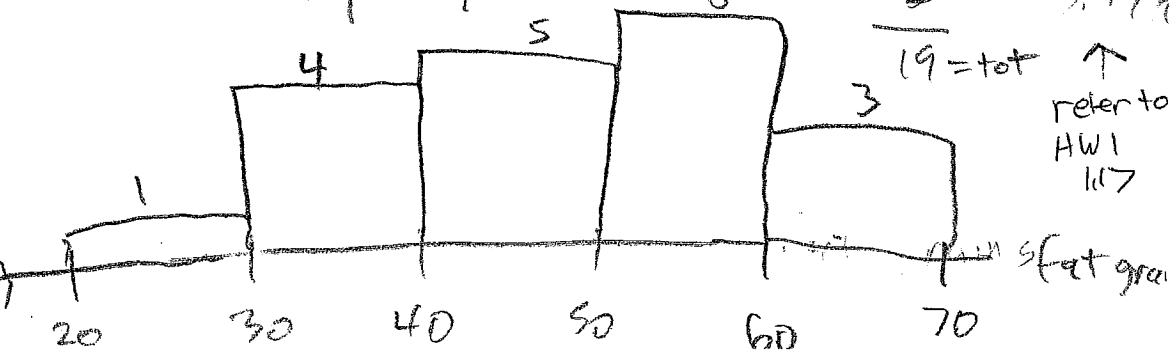
$4(0) + 4(1) = 44$

Stem ten
leaf ones
 $\% = \frac{\text{count}}{\text{tot}} 100\%$

Histogram

	tally	Count	%
20-29		1	
30-39		4	5.26
40-49		5	21.05
50-59		6	26.32
60-69		3	31.58
		19	79%

eg include left hand part not right



15) ^{know} ~~P 27~~ mean $\bar{X} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum x_i$ 3-5 memorize

ex) $\{1, 9, 3, 7\}$ $\bar{X} = \frac{1+9+3+7}{4} = \frac{20}{4} = 5$

Note: could say find $\bar{Y}, \bar{w}, \bar{x}$. The letter should not matter.

16) ^{know} p 30 median M

a) sort the data from smallest to largest

b) count $\frac{n+1}{2}$ obs's from the bottom of the list. (n is the number of obs's)

• If n is odd $M = \frac{n+1}{2}$ th observation

even $M =$ average of the 2 middle obs's
 Note: at least half of the obs's $\geq M$ and at least half are $\leq M$.

ex) $\{1, 9, 3, 7\} \xrightarrow{\text{sort}} \{1, 3, 7, 9\}$

note $\frac{n+1}{2} = \frac{5}{2} = 2.5$

n = size of list = 4

$M = \frac{3+7}{2} = 5$

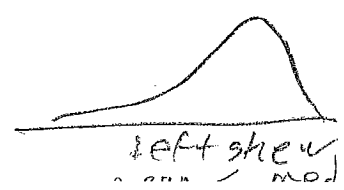
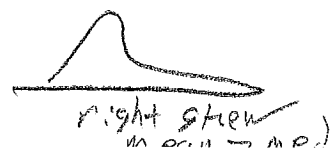
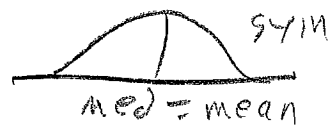
ex) $\{1, 9, 3\} \xrightarrow{\text{sort}} \{1, 3, 9\}$
 \uparrow 2nd obs

$\frac{n+1}{2} = \frac{3+1}{2} = 2$

17) \bar{X} and M are measures of center $n=3!$

18) ^{P 49} If a distribution is (exactly) symmetric

the median = mean



19) the range = largest obs - smallest obs

20) p33 Order the list and find the median. Then the 1st quartile Q_1 is the median of the obs's to the left of the median M while the 3rd quartile Q_3 is the median of the obs's to the right of the median M . The 2nd quartile $Q_2 = M$

ex) n even
50, 59, 72, 81 (ordered)
 $Q_1 = 54.5$ $Q_3 = 76.5$

ex) n odd
50, 59, 72, 81, 97 (ordered)
 $Q_1 = 54.5$ $Q_3 = \frac{81+97}{2} = 89$
 $Q_2 = m = 72$

ex) p43 4, 7, 7, 7, 8, 9, 9
 ↑ ↑ ↑
 Q_1 Q_2 Q_3

21) ^{know} p35 Box plot

- a) order the data; find the five number summary min Q_1 M Q_3 max
- b) Make a box from Q_1 to Q_3 . Draw a line for M .
- c) Draw a line from the edge of the box (Q_1) to min.

Draw a line from the other 4-5 edge of the box (Q3) to the max, label axes

ex exam scores min = 30
(stem plot)

Q1 = med of 1st 35 obs = 18th obs = 60

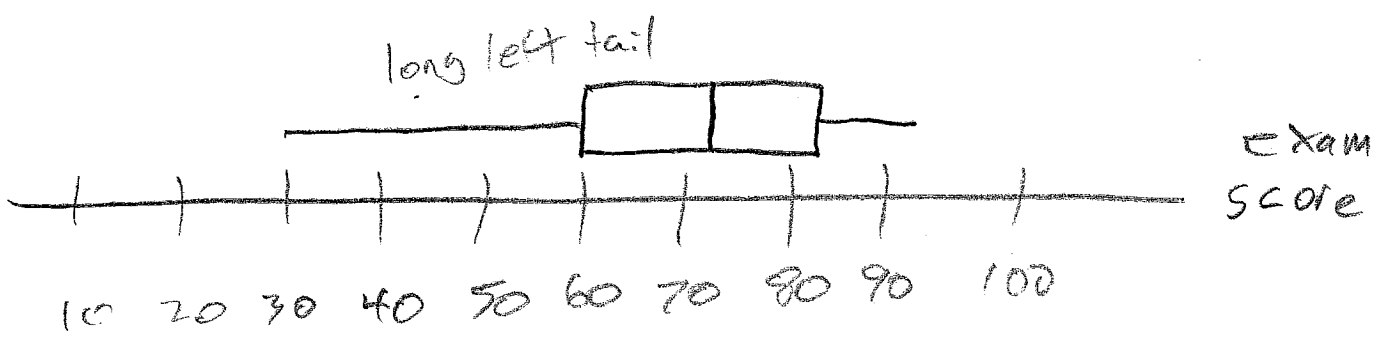
M = 72 (36th obs) Q3 = med of last

35 obs = 54th obs = 82 max = 92

min	Q1	M	Q3	max
30	60	72	82	92

1 2 ... 18 ... 35 36 37 ... 54 ... 11

$\frac{35+1}{2} = 18$ $36+18 = 54$ $71 = 36+35$



22) p38 The standard deviation

MEMORIZE

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

The SD measures spread.

Note: The Variance = $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$

ex) use table for HW, exam, quiz

X	$X - \bar{X}$	$(X - \bar{X})^2$	S
1	-2	4	5
2	-1	1	
3	0	0	
4	1	1	
5	2	4	

$\sum X = 15$

$\bar{X} = \frac{15}{5} = 3$

$0 = \sum (X - \bar{X})$

good check, p39

$\sum (X - \bar{X})^2 = 10$

$= (X_1 - \bar{X})^2 + \dots + (X_5 - \bar{X})^2$

$S = \sqrt{\frac{10}{5-1}} = \sqrt{\frac{10}{4}} = \sqrt{2.5} = 1.581$

common mistakes 1) Forget square root (2.5)
 2) forget to divide by n-1 ($\sqrt{10} = 3.16$)

23) COMMON Q E PROBLEM

Find \bar{X} , median and S for

5, 4, 1, 2, 3 ($\bar{X} = 3, M = 3, S = 1.58$)

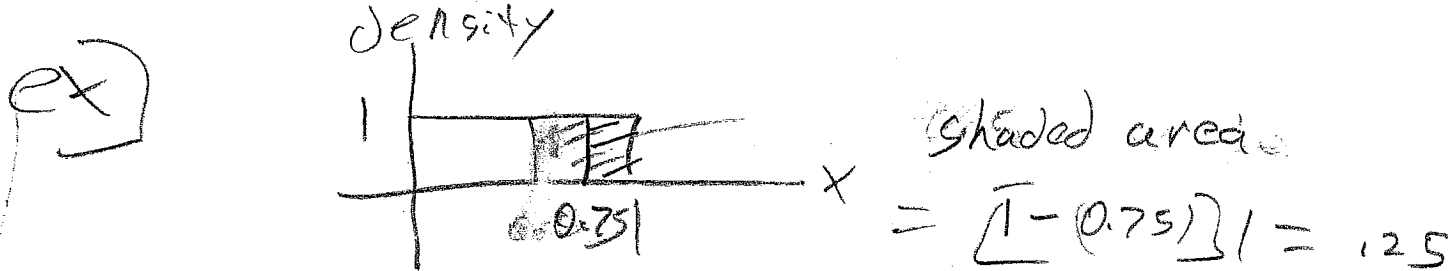
p39 $X_i - \bar{X}$ is called a deviation $\sum (X_i - \bar{X}) = 0$

24) Properties p39 S measures spread about \bar{X}
 $S \geq 0$ and $S = 0$ only when all obs's have the same value. S is larger when there is more spread

S has the same units as the original obs's
 25) n, \bar{X} , S and \bar{X} is the best description of a distribution: always

26) p48 A density curve is never negative and the area under the curve = 1.0

Areas under curve represent proportions of observations



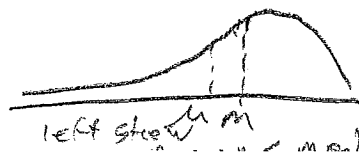
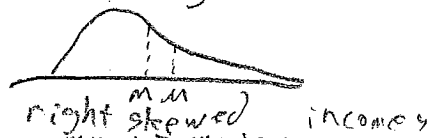
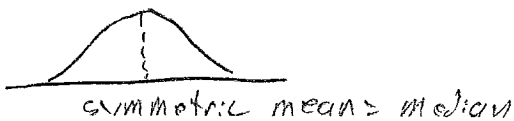
25% of x values are between 0.75 and 1

Density curves are often useful approximations for histograms

27) p49-50 The median of a density curve is the point such that half of the area is to the left and half the area to the right.

The mean μ of a density curve is the balance point. See figure 1.6.

The mean is drawn towards the longer tail in skewed histograms.

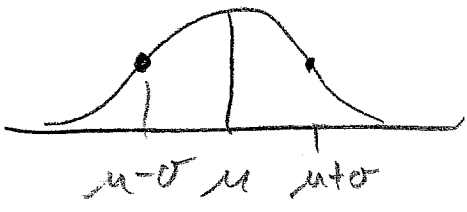


28) p51 $\sigma = \text{SD of density curve}$ 0282 6

\bar{x} estimates μ

s estimates σ

29) p51 the normal curve is a density curve that can approximate many data sets. Given μ and σ , a normal curve can be drawn. A normal curve is symmetric about μ , so μ is the mean and the median.

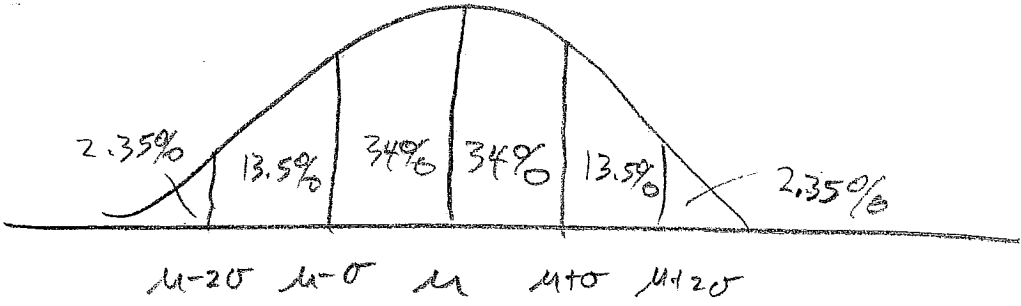


see fig 1.19 p52

30)* p5368 - 95 - 99.7 rule

For normal data

68%	of obs's are within σ of μ
95%	2σ
99.7%	3σ



ex length of rattlesnake is approx normal with $\mu = 8$ in $\sigma = 1$ in. Then $\approx 95\%$ have lengths between 6 in and 10 in. $\approx 2.5\%$ are longer than 10 in.

31) Notation x from $N(\mu, \sigma)$
means x is an obs from a normal dist with mean μ and SD σ .

32) ~~PS2~~
Last paragraph on p 52 is important

3 reasons for normal
many datasets \approx normal
many are not

33) ^{know} PS5 IF x is from a dist with mean μ and SD σ , then

the standardized value of x is the z score of x = $z = \frac{x - \mu}{\sigma}$.

and tells how many SD's x is from μ .
(memorize)

34) PS6 IF X is $N(\mu, \sigma)$, then

$z = \frac{X - \mu}{\sigma}$ is $N(0, 1)$, the standard

normal. Hence one table can be used for every normal curve.

35) ^{PS7-61} The z table in the front of the book can be used to find 3 types of

Probabilities (Areas),

0282 7

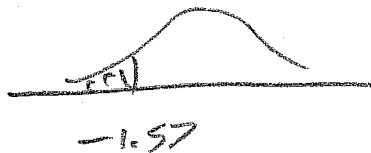
Let Z be $N(0,1)$. The Z table
 tables $P(Z \leq z^*) = P(Z < z^*)$.

I) a) $P(Z < 2.8)$

z^*	.00
2.8	.9974

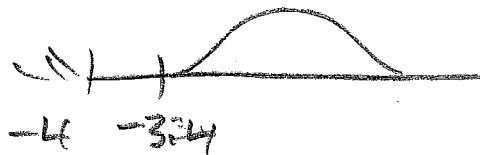


b) $P(Z \leq -1.57)$



z^*	.07
-1.5	.0592

c) $P(Z \leq -4)$



$\approx P(Z \leq -3.49)$

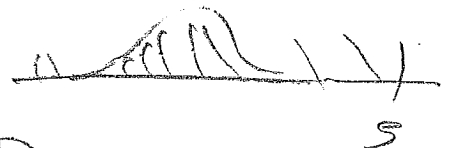
so $P(Z \leq -4) \approx 0$.

d) $P(Z \leq -40) \approx 0$



e) $P(Z \leq 5) \approx P(Z < 3.49)$

so $P(Z \leq 5) \approx 1.0$



f) $P(Z \leq 0) = 0.5$



Key words: less than

$$\text{II } P(Z > z^*) = P(Z \geq z^*) \quad 7.5$$

$$P(Z > z^*) = 1 - \underbrace{P(Z \leq z^*)}_{\text{area from table}}$$

$$\text{a) } P(Z \geq -1) = 1 - P(Z < -1) = \frac{\text{area}}{-1}$$

$$= 1 - \frac{\text{area}}{-1} = 1 - 0.1587 = 0.8413$$

$$\text{b) } P(Z \geq 1.763) = P(Z \geq 1.76)$$

$$= 1 - P(Z < 1.76) = 1 - \overset{\uparrow}{\text{closest}} 0.9608 = 0.0392$$



$$\text{c) } P(Z \geq 0) = 0.5$$

$$\text{d) } P(Z > 1.645)$$

$$= 1 - P(Z < -1.645) = 1 - \frac{0.0505 + 0.0495}{2} = 1 - 0.05 = 0.95$$



\swarrow 1.64 and 1.65 are equally close to -1.645

key words! greater than

III $P(a \leq z \leq b) = P(a < z \leq b) = P(a = z < b)$
 $= P(a < z < b)$
 $= P(z \leq b) - P(z \leq a)$



a) $P(-1 \leq z \leq 1) = .8413 - .1587 = .6826$

b) $P(2.297 \leq z \leq 4.093) \approx P(z \leq 4.09) - P(z < 2.26)$
 $\approx 1 - .9881 = 0.0119$

Key words "between."

30) ~~P57761~~ Forwards know for Q E Final:
 Finding probabilities when X is $N(\mu, \sigma)$
 with a picture
 step 1) state problem, 2) standardize, 3) draw picture
 4) use z table

i) $P(X \leq b) = P(X \leq b) = P\left(z \leq \frac{b-\mu}{\sigma}\right)$ table
↓

ii) $P(X > a) = P(X \geq a) = P\left(z \geq \frac{a-\mu}{\sigma}\right) = 1 - P\left(z \leq \frac{a-\mu}{\sigma}\right)$

iii) $P(a \leq X \leq b) = P(a < X < b) = P(a = X < b) = P(a < X \leq b)$
 $= P\left(\frac{a-\mu}{\sigma} < z \leq \frac{b-\mu}{\sigma}\right) = P\left(z \leq \frac{b-\mu}{\sigma}\right) - P\left(z \leq \frac{a-\mu}{\sigma}\right)$ table

ex) IQ test scores are $\approx N(\mu=100, \sigma=15)$


a) Find chance $85 < X < 115$

step 1

$$\begin{array}{c} | \quad | \quad | \\ \hline 85 \quad 100 \quad 115 \end{array} X = P(85 < X < 115)$$

step 2 standardize

$$\frac{85-100}{15} = -1 \quad \frac{115-100}{15} = 1$$

step 3  $Z = P(-1 \leq Z \leq 1)$

step 4 $.8413 - .1587 = .6826$

b) Find chance ^{randomly selected} IQ score is greater than 107.


step 1

$$\begin{array}{c} | \quad | \\ \hline 100 \quad 107 \end{array} X = P(X \geq 107)$$

step 2 standardize $X=107$

$$Z = \frac{107-100}{15} \approx 0.47$$

step 3

 $Z = P(Z > 0.47)$

step 4 $= 1 - \begin{array}{c} \leftarrow \text{table} \\ \text{table} \\ \hline .47 \end{array} = 1 - .6808 = .3192$

common mistake standardize, then

write down table value (eg .6808 instead of .3192,

37) P 61 Backwards Normal calculation P 82 09
 Finding an X value know for
 @ Exam Final
 given a proportion

step 1 state problem with a picture step 2 use table

step 3 unstandardize

ex} For the IQ scores find the score needed to be in the highest 90%



Find x^* so that $P(X \leq x^*) = 0.1$



Find the largest prob ≤ 0.1 and the smallest prob ≥ 0.1 .
 Then z^* corresponds to the closer

z^*	0.08	0.09
-1.2	0.1003	0.0985

closer so $z^* = -1.28$

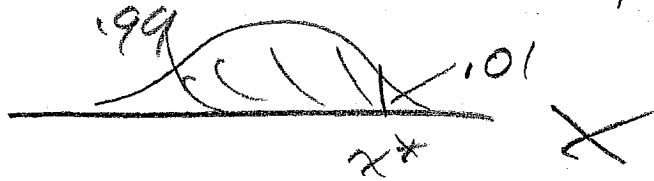
step 3 unstandardize $z^* = \frac{x^* - \mu}{\sigma}$

so $x^* = \mu + z^* \sigma$ P 72

$x^* = 100 + (-1.28)(15) = 100 - 19.8 = 80.2$

Find the IQ score such that 99% of the scores are smaller, (0.5)

Step 1



Step 2

z^*	.02	203
2.3	.9898	.9901 closer

$$z^* = 2.33$$

Step 3 unstandardize

$$x^* = \mu + z^* \sigma$$

$$= 100 + 15(2.33/15) = 134.95$$

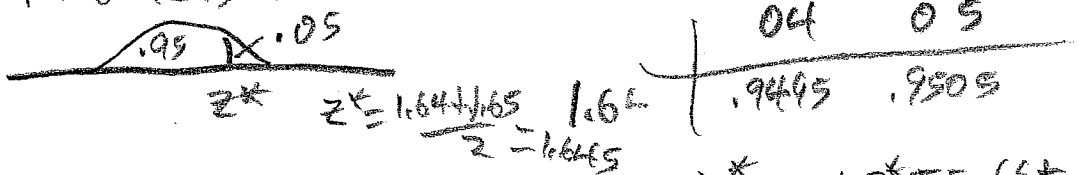
check standardize x^*

$$\frac{134.95 - \mu}{\sigma} = \frac{134.95 - 100}{15} = 2.33$$



ex) Suppose women heights $\approx N(\mu = 66 \text{ in } \sigma = 3 \text{ in})$

Find height so that 5% are taller



$$x^* = \mu + z^* \sigma = 66 + 3(1.645) = 70.93$$

1) p 80-85 a response variable

= dependent variable Y . measures an outcome of a study. An explanatory variable = independent variable X attempts to explain the observed outcomes.

~~p 80~~ Sometimes have 2 variables, but they are not ^{response} explanatory

ex) predict Y for given values of X

eg Y = college Fresh GPA X = SAT score

Y is the response X the explanatory variable

2) p 81, 135 association \neq causation

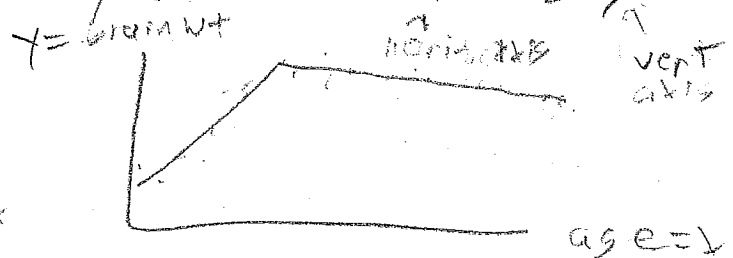
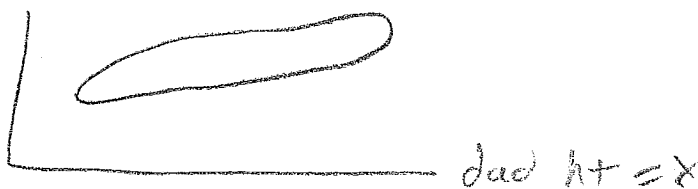
2 variables X and Y can be related, but this does not mean that changes in X result in changes of Y

ex X = shoe size Y = reading score of kids $H-O$

bigger X 's go with bigger Y 's

increase in age results in increase of shoe size and reading score

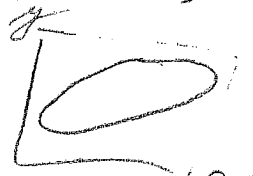
3) p 82 A scatter plot is a plot of X vs Y
 eldest son ht = Y
 dad ht = X



4) p84 Look for an overall pattern especially linear (football shaped) plots and for outliers

5) p84 2 variables are positively associated when above ave values of one variable tend to accompany above ave values of the other and below ave values also tend to occur together

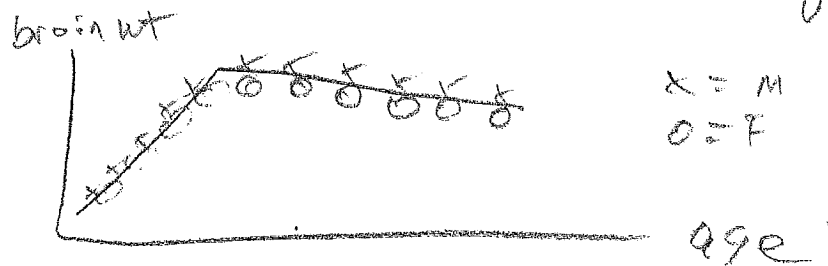
2 variables are negatively associated when above ave values of one variable tend to accompany below ave values of the other and vice versa.



p85 Strong relationships make accurate predictions possible.

6) p84 an outlier is an observation that falls outside the overall pattern. (see p18, too)

7) p88 use colors ^{or symbols} to show categorical variables in a scatterplot of 2 quantitative variables



8) p89 The strength of a linear relationship of 2 quantitative variables x and y is measured by the correlation r .

$$r = \frac{1}{n-1} \sum \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right) = \frac{1}{n-1} \sum z_x z_y \quad 0.28211$$

where \bar{x} and s_x are the mean and SD of x
 while \bar{y} and s_y are the mean and SD of y ,

ex	X husband	Y wife ht	$\frac{x-\bar{x}}{s_x} = z_x$	$\frac{y-\bar{y}}{s_y} = z_y$	product $z_x z_y$
	72	66	1.1858	0	0
	68	64	-0.3953	-0.9535	0.3769
	70	66	0.3953	0	0
Common exam problem delete 5 numbers	68	65	-0.3953	-0.4767	0.18843
	71	70	0.7905	1.9069	1.50240
	65	65	-1.5810	-0.4767	0.75366
	$\sum x = 414$	$\sum y = 396$	0.0000	0.0000	2.82639
	$\bar{x} = \frac{414}{6} = 69$	$\bar{y} = \frac{396}{6} = 66$	always true except for rounding		11
			$72 - 69 = 3$		$\sum z_x z_y$

$$s_x = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2} = \sqrt{\frac{1}{5} [3^2 + (-1)^2 + 1^2 + (-1)^2 + 2^2 + (-4)^2]}$$

$$= \sqrt{\frac{32}{5}} = \sqrt{6.4} = 2.530$$

$$s_y = \sqrt{\frac{1}{n-1} \sum (y_i - \bar{y})^2} = \sqrt{\frac{1}{5} [0^2 + (-2)^2 + 0^2 + (-1)^2 + 4^2 + (-1)^2]}$$

$$= \sqrt{\frac{22}{5}} = \sqrt{4.4} = 2.0976$$

3rd column 1st row $\frac{72 - 69}{2.530} \approx 1.1858$

4th column 1st row $\frac{65 - 66}{2.0976} \approx -0.4767$

2nd row 1st column $\frac{66 - 69}{2.530} = -1.1858$
 2nd row 2nd column $\frac{64 - 66}{2.0976} = -0.9535$

Since $\sum z_x z_y = 2.82639,$

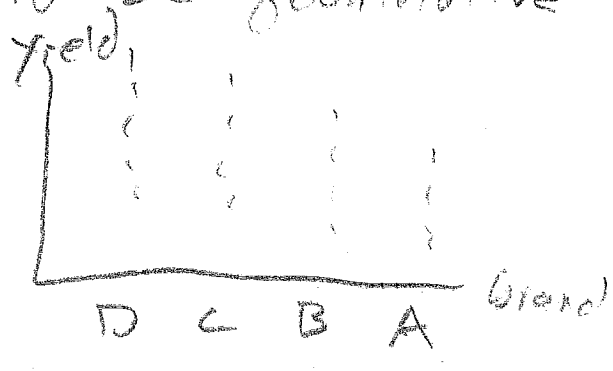
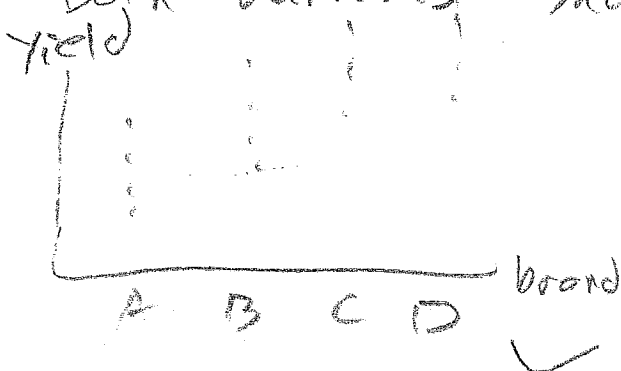
11.5

$$r = \frac{2.82639}{5} = 0.565$$

exam Q: get table with 4 entries missing a z_x , z_y and products

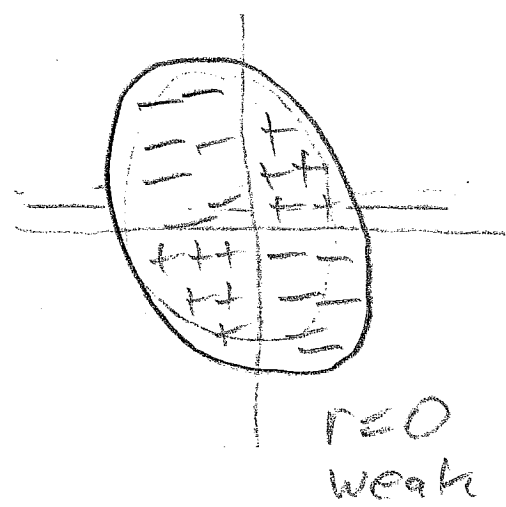
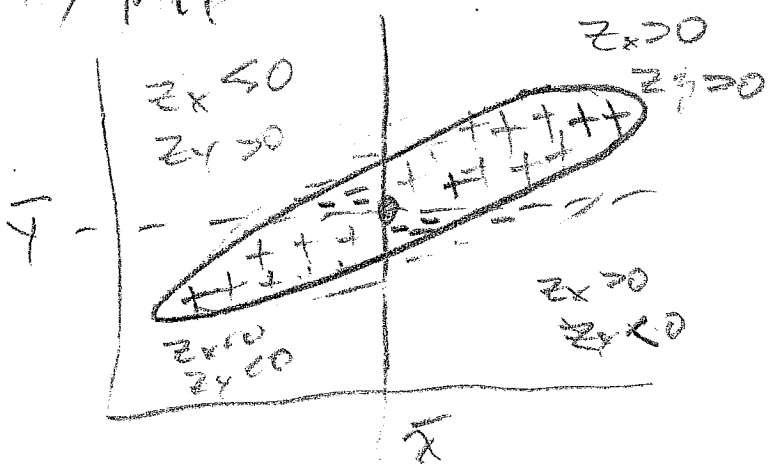
9) p98 Both variables should be quantitative

ex



positive and negative association doesn't make sense

10) p99



$r < 0$
weak

$r > 0$ strong

11) p99 properties of r when scatterplot is (football shaped) linear

$\text{corr } x, y = \text{corr } y, x$, unit free

$$-1 \leq r \leq 1$$

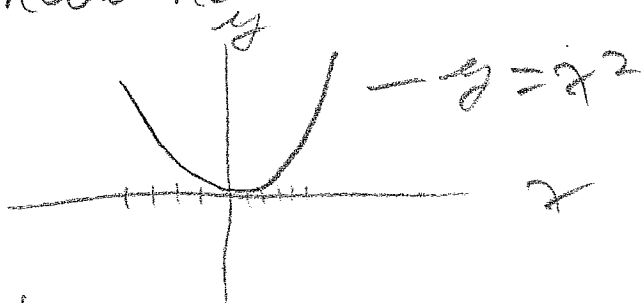
positive assoc if $r > 0$
neg $<$

r far from 0 means strong association
ie knowing x helps a lot in predicting y

$r = \pm 1$ if the points in the scatterplot
fall exactly on a line

adding or subtracting the same # to all values of x or y does not change r

12) If the scatterplot is non linear, r
should not be used



perfect nonlinear
relationship, but
 $r = 0$ is possible

§ 2.3

13) P106 The regression line describes how
the response variable y changes as the
explanatory variable x changes.

14) P109 A line has the form $y = a + bx$ where
 $b = \text{slope}$, $a = \text{intercept}$. The least squares
line picks a and b so that

$\sum (y - \hat{y})^2$ is minimized where

$$a = \bar{y} - b\bar{x}, \quad b = r \frac{s_y}{s_x}$$

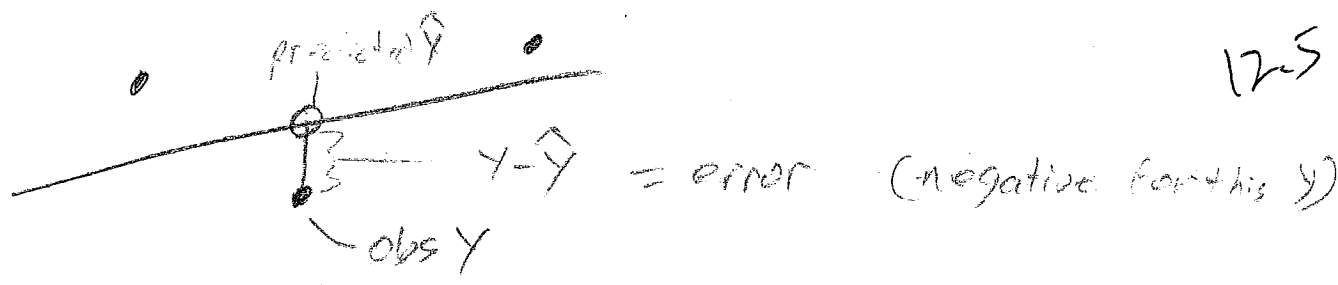
$y - \hat{y} = \text{vertical error}$

$r = \text{corr. coeff.}$

$s_y, s_x = \text{SDs}$

$$\hat{y} = a + bx$$

$\hat{y} = \text{predicted value}$



15) LS line $\hat{Y} = a + bX$, $a = \bar{Y} - b\bar{X}$

$b = r \frac{s_Y}{s_X}$. know for Q Exam Final

Given $\bar{Y}, \bar{X}, r, s_Y, s_X$, you should be able to get the LS slope b and intercept a .

16) p123 the slope b is the amount \hat{Y} changes when X changes by one unit

17) To plot the LS line take one of the smallest values of X and find \hat{Y} largest "

These 2 points determine the line.

18) Extrapolation is using regression to predict y for x outside of the range of the observed values of the explanatory variable. Such predictions are often very bad.

ex) Predict wt from ht
 $\bar{X} = 70$ in $s_X = 3$ in
 $\bar{y} = 162$ lb $s_Y = 30$ lb
 min ht = 57 max ht = 79
 men 18-24
 $x = ht$ $y = wt$
 $r = 0.47$

$b = 0.47 \frac{30}{3} = 4.7 \frac{lb}{in}$

$$a = 7 - 4x = 7 - 4(20) = -167 \quad 0282 \quad 13$$

$$Q = a + bx = -167 + 4.7x \text{ is the LS line}$$

Predict y if i) $x=20$: $\hat{y} = -167 + 4.7(20) = 162$ 165 ↓

ii) $x=23$: $\hat{y} = -167 + 4.7(23) = 176.1$

iii) $x=60$: $\hat{y} = -167 + 4.7(60) = 115$

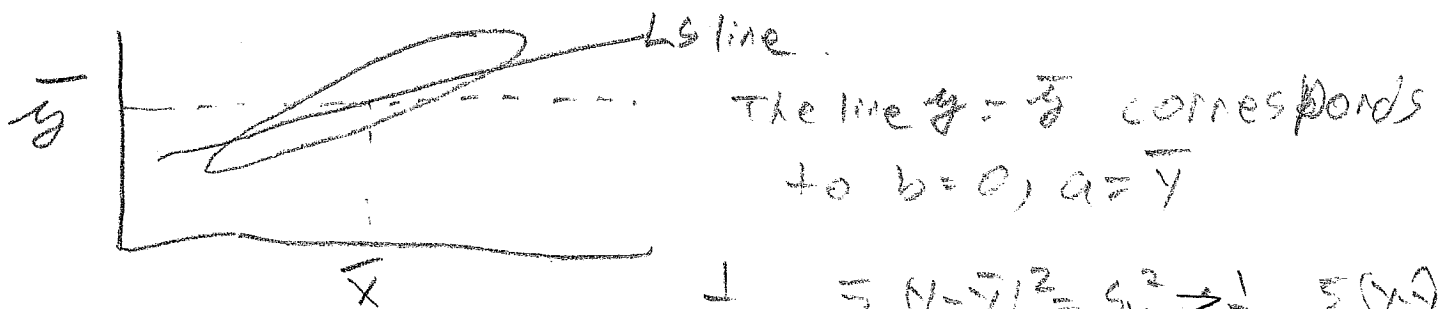
iv) $x=1$: $\hat{y} = -167 + 4.7(1) = -162.3$

19) p113 units of x and y matter

Slope = $b = \frac{r s_y}{s_x}$ means a change in 1SD in x results in a change of r SD in y
 ($-1 \leq r \leq 1$)

The LS line always passes through (\bar{x}, \bar{y}) .

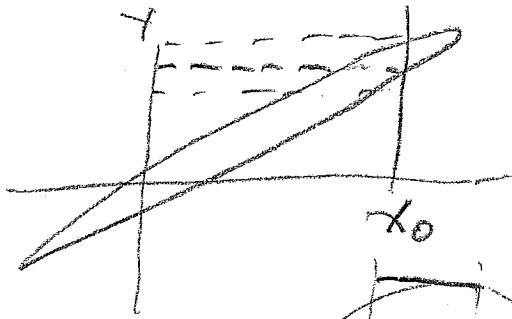
20) p113 The square of the correlation, r^2 , is the fraction of the variation of y (about \bar{y}) that is explained by the LS line



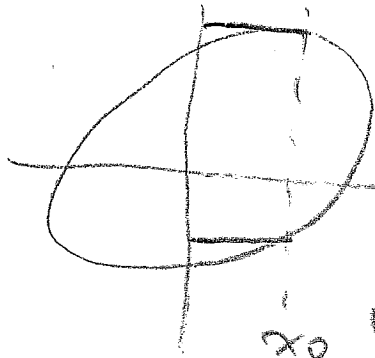
$$\frac{1}{n-1} \sum (y - \bar{y})^2 = s_y^2 \approx \frac{1}{n-1} \sum (y - \hat{y})^2$$

If there is a linear relationship r^2 small means that to predict a new value of y , \bar{x} is about as good as LS, i.e. knowing x LS minimizes this

does not help much in predicting y . If r^2 large means the LS line is much better than \bar{y} for predicting a new y value; knowing x helps a lot in predicting y .



$r^2 \approx .9$ If r^2 is close to 1, for fixed x_0 , spread in y is small (compared to s_y).



r^2 close to 0, for fixed x_0 , spread in y is large (close to s_y)

x_0 $r^2 \approx .05$

EX) $y = \frac{2}{3}x + 32$

$y = F_0$

$x = C_0$

$x = \frac{3}{2}(y - 32)$

suppose

$y = 10$
 $x = \frac{10 - 32}{2} = -11$

50

90

then $\hat{y} = \frac{150}{3} = 50$

$\bar{y} = 10, 50, 90$
at $x = \frac{110}{2} = 55$

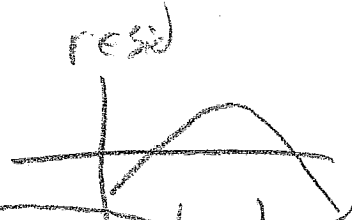
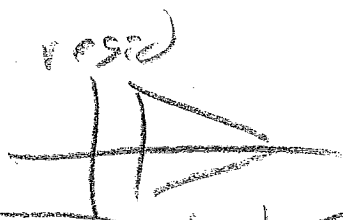
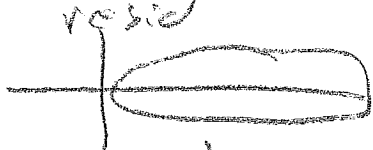
21) p116

residual = $y - \hat{y}$

ANS a residual plot is a scatterplot of the residuals vs x or of the residuals vs \hat{y}

want to see residuals scattered about resid = 0 line with no pattern

Fact $\sum (y - \hat{y}) = 0$

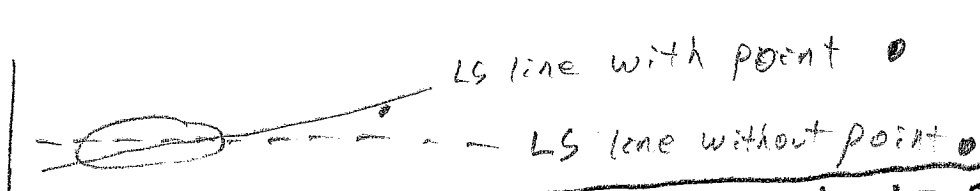


22) Does slope is influential if solution had ...

margin = ...

22) p124 an observation is influential 14

if deleting it and then refitting the LS line results in a marked change in the fit,



Minitab
 \downarrow Know for X
 \downarrow Exam final \downarrow
 gas used = $1.094 + 0.188 \text{ days}$

23) Minitab output p2.4

Minitab output p110
 pred
 constant $1.0892 = a$
 $X \rightarrow$ D days $0.1888 = b$
 pred y if
 $x = 20$
 $y = 1.0892 + 0.1888(20)$
 $= 4.9652$

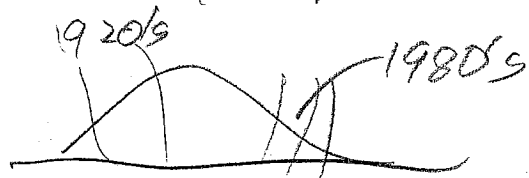
24) p130 r and LS can be very misleading if the relationship between x and y is nonlinear, eg if outliers are present.

25) p132 p xxvi
A lurking variable is a variable that has an important effect on the relationship among the variables in the study but is not included among the variables studied.

(skip ex) shoe size = x y = reading score of kids 2-8
 age of child is a lurking variable

(skip ex) NY Times Feb 24 1998
 some IQ tests have been given over and over again eg every year by the military. Each year they are standardized to have mean 100 and SD 15. A researcher got year by year results and the raw scores

The IQ scores have been going up each year.



Conclusion

generation X is far more intelligent than their grandparents generation

X Year IQ test taken in military

Y = test score

possible lurking variables:

- 1) literacy rate increased
- 2) military more selective, no draft
- 3) Standardized tests are converted to IQ's so schools teach a lot more IQ test problems than they did.

ex)

Raw SAT, ACT scores decreased from 1950 to 1990

Conclusion educators ability decreased

lurking 1) 50% to 10% students took exam

now top 2/3 take it

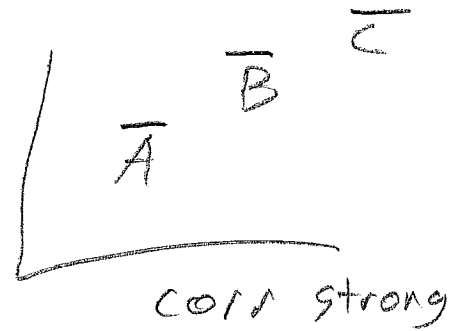
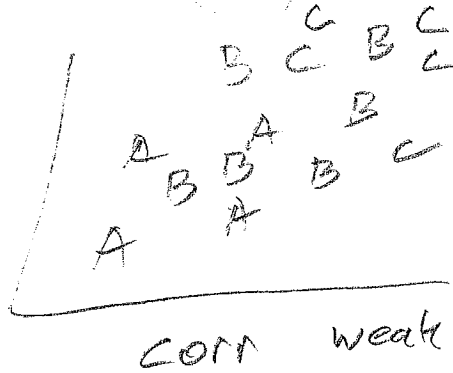
2) Immigration demographics have changed, Perhaps fewer % of immigrants have English as a native language,

2.6) PIB Often averages are computed
eg Y = ave income men 35-54 9 regions of USA
X = ave education by 9 regions

but among individuals

0282 15

$$\text{corr}(X, Y) \approx 0.4$$



Correlations based on averages are usually higher than correlations based on individuals (the correlation that is usually of interest).

skip
27)

P135

Association does not imply causation,
eg shoe size reading score kids 2-8

ch 3

1) P168

A population is the entire group of individuals we want information about. A sample is the subset of the population actually studied.

2) P168 The design of the sample is the method used to choose the sample from the pop.

3) P168 P XXVIII A voluntary response sample consists of people who choose themselves by responding to a general appeal.
The sample chooses the individuals easiest to reach

ex) Magazines often print surveys
eg on sexual behavior, and ask readers
to respond.

Problem! People who respond tend to differ
from the pop as a whole

ex Ann Landers pxxviii 10000 letters 70%
said kids not worth it. This biases
the sample.
A statistically designed poll said 91%
would have kids again.

ex) 1936: The Literary Digest, it sent questionnaires
had correctly picked the winner in
the past 5 presidential elections.

to 10 million people for Landon vs FDR
presidential election. They got the names

from phone books, club lists, etc.

The digest poll said Landon would win overwhelmingly
In 1936 on 25% of homes had phones. 57% to 43%

Gallup took a sample of 3000 people
from the same lists and predicted the
Digest results before the Digest publication
(off by 1%), with a sample of 50000

Gallup predicted	FDR 56%	Landon 44%
Actual result	62%	38%
Digest	43%	57%

After the election, Digest's cover said
"Boy are we red in the face."

~~The Digest poll was not a voluntary response sample since they chose who got the questionnaires. The sample had a much larger proportion of wealthy than the pop as a whole. Before 1936 rich and poor had similar voting behavior.~~

In 1936 the poor esp 9 million unemployed were for FDR the rich (eg 11 million phone owners) were for Landon.

Of those who got questionnaires, only 20% responded. Respondents tend to differ from non respondents.

4) p174 A probability sample gives each member of the population a known nonnegative probability (chance) of being selected.

5) p180 The result from a probability sample, eg the proportion of Republicans that will vote for Bush, will follow some density curve, often the normal distribution. eg probability sample of 1000 Republicans $\mu = .627$ say they will vote for Bush $\approx N(\mu = .627, \sigma = .015)$
So 95% of samples of size 1000 will give a proportion between $\mu \pm 2\sigma$ or $.607$ to $.657$

With a probability sample, you can give a guarantee on how accurate the result is.

P180 This guarantee is called the "margin of error".
b) A probability method uses impersonal chance to select the sample.

b) Polling companies such as Gallop now use probability samples (the design).

P181 Many surveys have terrible designs

a) voluntary response sample

b) samples of convenience: ask the people easiest to reach (the interviewer chooses the sample)

ex) ask SIU student passing through the student center if they are Rep or Dem

ex) At a conference on Leukemia give questionnaires to Doctors regarding improvements in treatment. (Interviewer chose all people attending conference although many went to other)

P169 The design is biased if it systematically favors certain outcomes.

In ex1) Too many women get asked

In ex2) doctors who attend conferences differ from doctors who do not.

Digest poll got too many rich.

ex Send person out to 10 different neighborhoods, have person interview people at 3 homes per neighborhood. Wealthier are over represented, so

are under represented

0282 17

9) p171 The simplest probability sample is a simple random sample (SRS).

For a SRS of size n from a pop of size M , the prob that any member of the pop gets in the sample is $\frac{n}{M}$. That is, each person has an equal chance of being selected.

Idea! put everyones name in a hat, shuffle the names, then draw n out,

ex) pop of size 10000 sample size $n=1000$

Put # 1-10,000 with each name. Use random numbers to select the 1000,

chance any person is selected = $\frac{1000}{10000} = \frac{1}{10}$.

Every sample of size 1000 out of 10000

(an enormous #) is equally likely.

ex) pop size = 7 a b c d e f g table B line 130 sample of size 2

ex) pop = 31 want sample of size 5

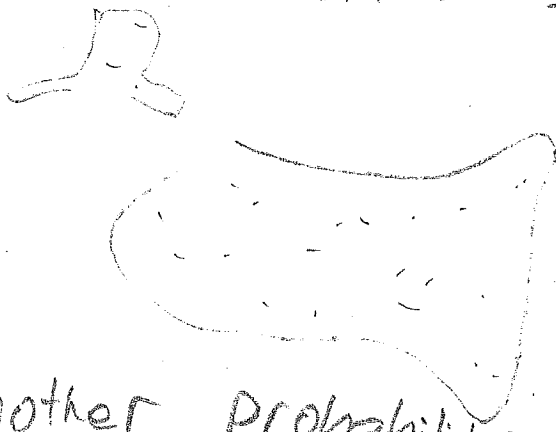
table B: Start at line 130 a_1, a_2, \dots, a_{31}
01 02 31

69 05 16 48 17 87 ~~40~~ 95 ~~84~~ 53
40 64 89 87 20 19

persons $a_5, a_6, a_7, a_9, a_{10}$ selected

10) Problem with SRS. They are too expensive if the pop is big and spread out.

eg SRS of 1000 USA voters for 2000 pres election

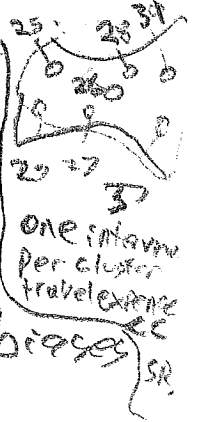


11) P174 Another probability sample is a stratified random sample. First divide pop into groups of similar individuals called strata, eg by age race gender income and education. Take a SRS from each strata to get the full sample.

The design yields clusters of nearby households. One interviewer can visit each cluster, so travel expense is much lower than a SRS of households.

12) PTSA probability method called a multistage sample is used if interviewers are sent out. Idea: Regions → Precincts → Households → Pop centers → wards. eg current population survey → Precincts → Households → Pop centers → wards. ask every voter in the selected households.

13) Random digit dialing. Regions are determined by area code, then use random phone #'s.



14) Probability samples eliminates many biases but there are still problems

15) p178 Under coverage

Some groups
0282
TB

eg homeless, prison inmates,
tend to be underrepresented (eg get
sample proportions and compare to
census proportions)

16) p178 Non response individual chosen for
sample can't be contacted or won't cooperate
(non respondents differ from respondents)

want survey response rate to be high
eg > 80% (Digest only had 20%
response rate)

17) p178 response bias

respondents may say what they think the
interviewer wants to hear

eg attitude of interviewer!

18) wording of question

Do you use tide?

What detergent do you use?

Show me the brand of detergent you use.

19) Census controversy

Census wants to use prob methods to

20) ^{for probability} The accuracy ^{method} of a sample depends ^{is} on the sample size.

(provided that the sample size is a small percentage of the pop size), ^{and if both samples are from the same pop} $\frac{\sigma_1}{\sqrt{n_1}} = \frac{\sigma_2}{\sqrt{n_2}}$ means sample 1 is more accurate ^{if pops differ}

ex) Find proportions of women in
SRS of 2500 Carbondale residents
SRS 2500 Chicago residents

The results will have the same accuracy.

ex) Ask everyone in class whether R or not.

Get 100% accuracy if you ask everyone.

ex) .01% Alaskans are polled
and .01% CA are polled
which sample is more accurate
or are the accuracies the same?
soln CA sample is much larger so
more accurate

4) p217 A phenomenon is random if individual outcomes are uncertain ²⁵ but in the long run follows a distribution like the normal. eg flip coin, toss dice

ex) $y = x^2$ if $x = 2$ $y = 4$ not random

grade on exam
spin of roulette wheel

5) p217 The probability of an outcome of a random phenomenon is the proportion of times the outcome would occur in a long series of ^{independent} experiments.

Figure 1: 4 DD Plots

eg flip fair coin prob Heads = $\frac{1}{2}$

6) p218 independent means the outcome of one trial (expt) must not influence the outcome of any other trial

eg) ^{new} manufacturing employee learns so # of mistakes decreases. If we count her mistakes per day the trials are not independent. trial

7) p218 simulation ← ch10

ex) generate many samples with a computer for each sample generated compute a statistic eg \hat{p} 's make a histogram of $\hat{p}_1, \dots, \hat{p}_{1000}$. Find mean and SD of $\hat{p}_1, \dots, \hat{p}_{1000}$ etc.

ex) simulate fair coin Prob H = 0.5 ← ch10
assign 0-4 to H 5-9 T toss die $P(2) = \frac{1}{6}$

10/10/11: 1 9 2 2 3 9 5 0
H T H H H T T H

$$\hat{p} = \frac{5}{8} = 0.6250$$

12/31/11: 5 4 5 8 0 9 1 5
T H T T H T H T H T $\hat{p} = \frac{3}{8} = 0.375$ etc.

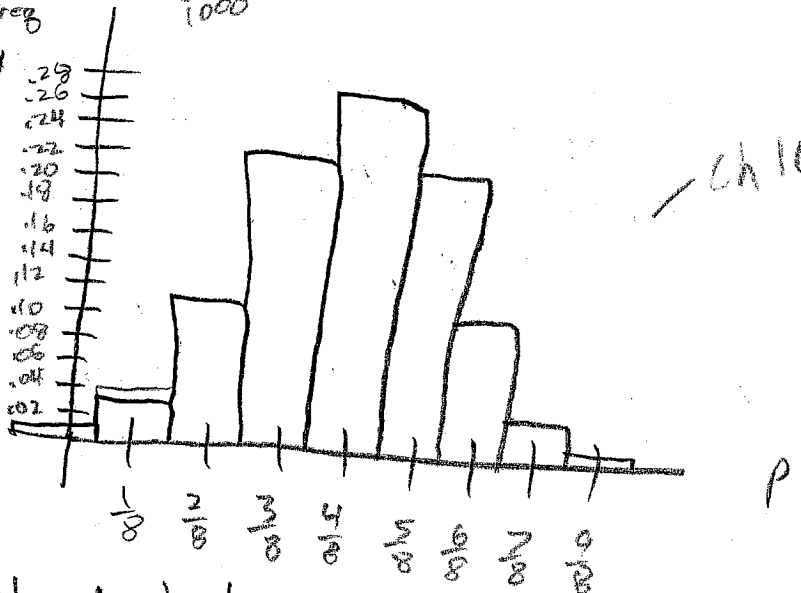
Generate 1000 samples

heads COUNT freq

$\frac{4}{1000} = .004$ etc

0282 20

0	4	.004
1	31	.031
2	109	.109
3	219	.219
4	273	.273
5	219	.219
6	110	.110
7	31	.031
8	4	.004
<hr/>		
	1000	



\$4.2

\$4.2 8) p 220 prob model sample space S

= set of all possible outcomes.

An event A is a subset of the sample space S.

A probability model consists of S and a way

of assigning probabilities to events.

ex) toss coin once $S = \{H, T\}$ $P(H) = P(T) = \frac{1}{2}$

ex) genders: 2 child families $S = \{bb, bg, gb, gg\}$

↑ 1st ↑ 2nd

event 2 boys $\{bb\}$

prob .5 | b .39 so each of the 4 outcomes

has prob = $\frac{1}{4}$,

Prob [at least one boy] = $P\{bb, bg, gb\} = \frac{3}{4}$

9) p 223 rules of prob

1) The prob $P(A)$ of any event A satisfies $0 \leq P(A) \leq 1$

2) $P(S) = 1$ where S is the sample space

3) complement rule $P(\text{A does not occur}) = 1 - P(A) = P(\text{not } A)$

A and B are disjoint $P(A \text{ or } B) = P(A) + P(B)$ 20.

A or B mean A or B or both occur.

If A_1, A_2, \dots, A_k are disjoint, $P(A_1 \text{ or } A_2 \dots \text{ or } A_k) = P(A_1) + P(A_2) + \dots + P(A_k)$

ex) grades A B C D F addition rule for disjoint events
 .2 .3 .2 .15 .15

are disjoint $P(S) = .2 + .3 + .2 + .15 + .15 = 1.0$

$P(C \text{ or higher}) = .2 + .3 + .2 = .7$

$P(D \text{ or lower}) = .15 + .15 = 1 - .7 = .3$

ex) Die toss 2 fair die

sg given on p 221 2nd

	1	2	3	4	5	6
1st	(1/1)	(1/2)	(1/3)	(1/4)	(1/5)	(1/6)
						(6/6)

list all possibilities using order

eg 1st toss 2nd toss

etc

each outcome has prob $\frac{1}{36} =$

(10) p225 Finite $S = \{s_1, s_2, \dots, s_n\}$
 $P_i = P(s_i)$ $P_n = P(s_n)$
 i) $0 \leq P_k \leq 1$ and $\sum_{i=1}^n P_i = 1$
 each outcome in S has a prob between 0 and 1 and the sum of all the probs equals 1
 iii) $P(A) = \sum_{s \in A} P(\text{outcomes in } A)$

ex) die $A = \text{sum of 2 die} \Rightarrow .4 = P(11), .2 = P(12), .2 = P(21), .2 = P(22)$

ex) toss fair coin 3 times

0282 21

list all possibilities using order of toss

8 outcomes in S
each with prob $\frac{1}{8}$

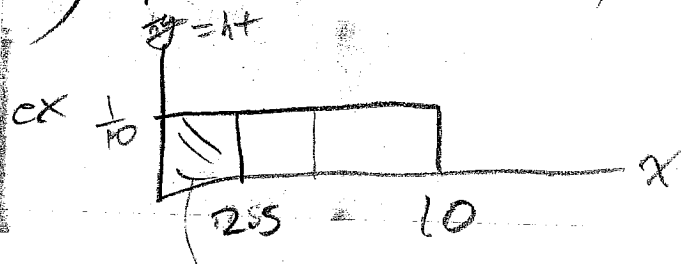
1st	2nd	3rd	1st	2nd	3rd
H	H	H	T	H	H
H	H	T	T	H	T
H	T	H	T	T	H
H	T	T	T	T	T

	A	B	C	D	E
Prob	.1	.2	.3	.1	\square

$$P(F) = 1 - (.1 + .2 + .3 + .1) = .3$$

- ii) If outcomes are equally likely, since the prob's sum to 1
- 12) Prob from density curve

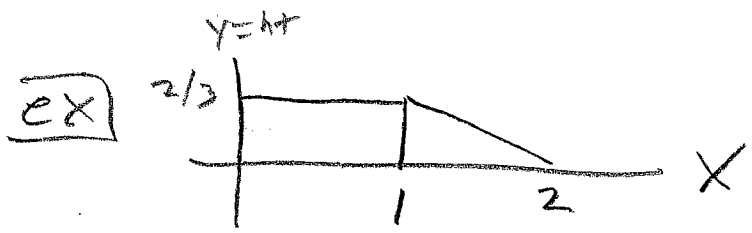
$$P(A) = \frac{\# \text{ outcomes in } A}{\# \text{ outcomes in } S}$$



rect area = base $(ht) = 1$
 $10 \left(\frac{1}{10}\right) = 1$

$$P(X < 2.5) = 2.5 \frac{1}{10} = .25$$

triangle area = $\frac{1}{2}$ base (ht)



$$1(ht) + \frac{1}{2} 1(ht) = 1$$

$$\text{so } \frac{3}{2} ht = 1 \quad ht = \frac{2}{3}$$

$$\text{Prob}(X > 1) = \frac{1}{2} 1\left(\frac{2}{3}\right) = \frac{1}{3}$$

$$\text{Prob}(X < 1) = 1 - \frac{1}{3} = 1\left(\frac{2}{3}\right) = \frac{2}{3}$$

ex) normal curve

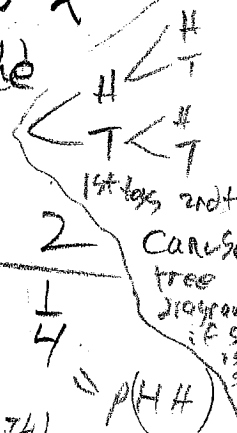
13) p231 a random variable X is a variable whose value is a numerical outcome of a random phenomenon

c9 toss coin 2 times $X = \# \text{ heads}$ $S = \{0, 1, 2\}$

(14) p231 The probability distribution of RV X tells the possible values of X along with the probabilities of those values.

→	HH	HT	TH	TT
	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$

X	0	1	2
Prob	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$
	"	"	"
	$P(TT)$	$P(HT) + P(TH)$	$P(HH)$

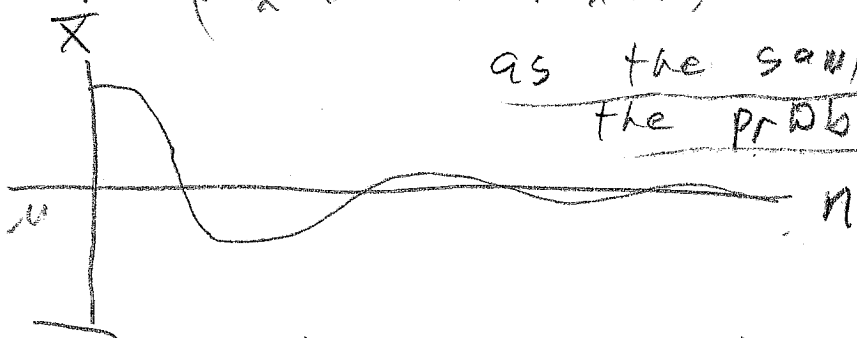


§4.3 (5) p237 Law of large numbers

Suppose you want to estimate μ_X , the mean of X with $\bar{X} = \frac{\sum x_i}{n}$, where the x_i are drawn at random. (Given any tolerance value, say $\delta > 0$, delete for large enough n)

$P(\mu_X - \delta < \bar{X} < \mu_X + \delta)$ is arbitrarily close to 1,

as the sample size increases the prob that \bar{X} is close to μ becomes closer and closer to one



ex You think μ_X is between 90 and 110 and want to estimate μ_X to within ± 1.001 with 99% probability. Then there is some integer N so that for any randomly drawn sample of size $n \geq N$

$$P(\mu_X - 1.001 \leq \bar{X} \leq \mu_X + 1.001) \geq 0.99$$

Implications | i) Casino gambling: If 0282 22

You make many bets that are as small as allowed, you will lose a lot of money.

Bet as much money as you can in as few bets as possible to maximize your chances of winning ^{a lot of money} (but you will still probably lose)

ii) Stock market stocks tend to increase

Buy small amounts of lots of shares and your portfolio should behave like the stock market average. If you put all your money in one stock, your gains or losses could be huge.

16) p241 The sampling distribution of

a statistic is the distribution of the statistic of the values taken by the statistic in all possible samples of the same size = pop of statistic

variable	population	mean	SD
ex X	pop of X	μ	σ
\bar{X}	sampling dist of \bar{X} = pop of \bar{X}	$\mu_{\bar{X}}$	$\sigma_{\bar{X}}$

ex pop net worths of 3 people Bill Gates 90 Billion
MJ 0.8 B you 0.0 B Find the sampling
distribution of \bar{X} for samples of size 2

BG	MJ	\bar{x} $90.8/2 = 45.4$
BG	Y	$90/2 = 45$
MJ	Y	$0.8/2 = 0.4$

223

\bar{x}	45.4	45	0.4
prob	$1/3$	$1/3$	$1/3$

17) p239 X comes from an underlying population with μ_x σ_x . A statistic, eg \bar{X} , comes from a pop = sampling distribution of the statistic which also has a center and spread eg $\mu_{\bar{X}}$, $\sigma_{\bar{X}}$.

18) p243 typically the center of the sampling distribution is close to what the statistic estimates (eg μ_x for \bar{X}) and the spread of the sampling distribution decreases as the sample size n increases. See Fig 4.10 p242.

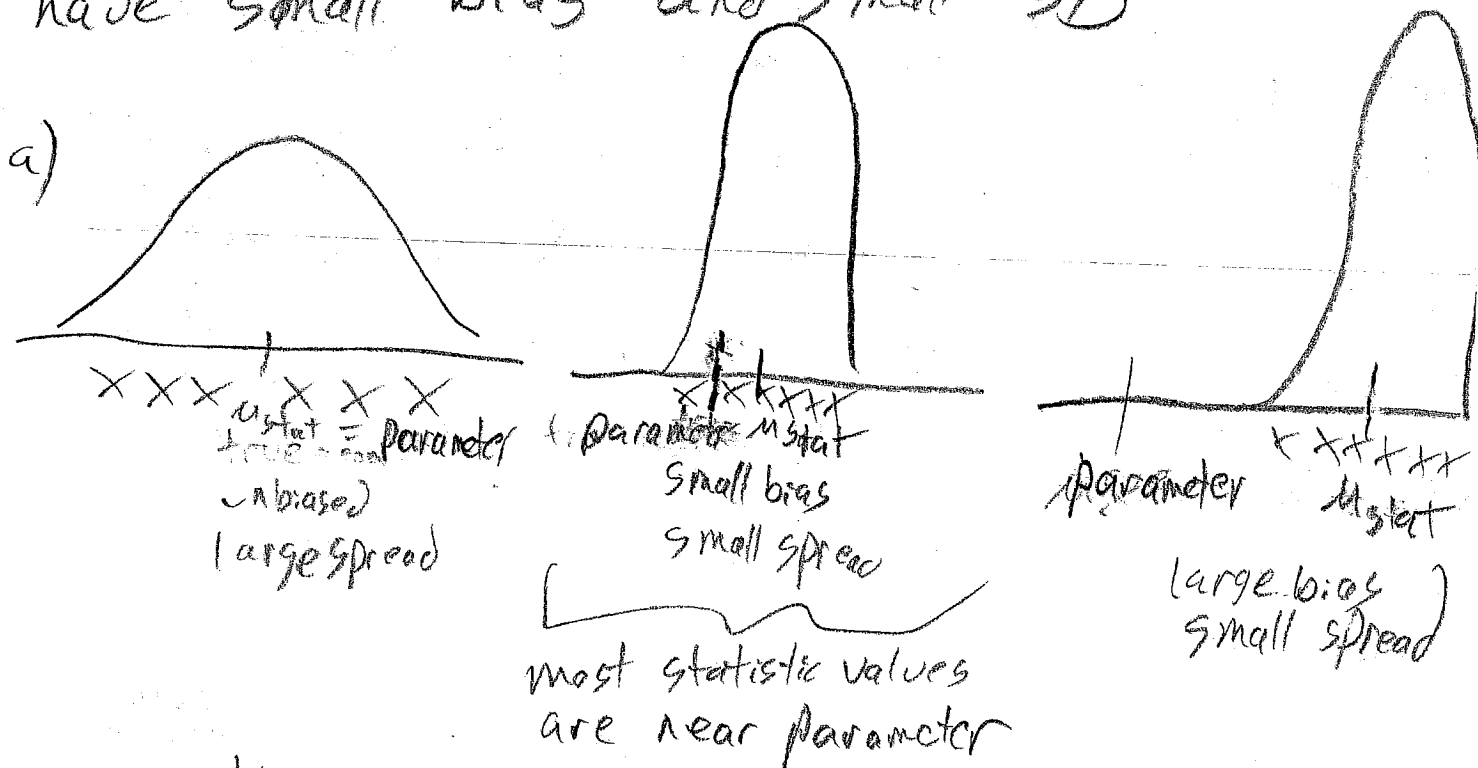
19) p242 A statistic is an unbiased estimator for a parameter if the mean of the sampling distribution of the statistic is equal to the parameter.

\bar{X} is unbiased for μ_x \hat{p} (or p)
 s^2 is not unbiased for σ^2

20) P239 The variability of a statistic ^{0282 23} is described by the SD of the sampling dist of the stat. For many statistics, larger samples have a smaller SD than smaller samples.

eg $\sigma_{\bar{X}_{1000}} < \sigma_{\bar{X}_{100}}$ Often the sampling dist SD does not depend on the pop size provided that the pop size is a lot larger than the sample size.

21) Want the sampling distribution to have small bias and small SD



22) ^{Know} P242, 244 The sampling distribution of \bar{X} obtained from a SRS of size n from the population of X having mean μ_X and SD σ_X has

i) mean $\mu_{\bar{x}} = \mu$

ii) SD $\sigma_{\bar{x}} = \sigma/\sqrt{n}$ if $n \in \text{pop size}$

iii) If X is $N(\mu_x, \sigma_x)$ then \bar{X} is $N(\mu_x, \frac{\sigma_x}{\sqrt{n}})$

iv) Central limit theorem (CLT) If n is large and X is from any distribution with finite σ_x and mean μ_x then

$\bar{X} \approx N(\mu_x, \frac{\sigma_x}{\sqrt{n}})$

Note: i) implies that \bar{X} is unbiased for μ_x

ii) implies that \bar{X} is less variable than X if $n \geq 2$

iii) implies \bar{X} is normal if X is normal

iv) implies that \bar{X} is approx normal if n is large

23) The results for \bar{X} also hold if the n observations come from n experiments done under identical conditions. This is like a SRS

from the pop of all possible measurements (expt results)

Luke Tierney
2000-02-14

24) Producing a sample with randomization : the impartial use of chance, "works out the effects of lurking variables, minimizes biases, and can make the sampling dist of the statistics easy to approach, eg normal

So \bar{X} is efficient for $\tau(\theta) = \frac{\theta - 1}{\theta} = \frac{\log \theta}{\theta}$ Done in class. 7.37

2.9) How large should n be to use CLT? ~~282~~ 28

i) $n \geq 1$ for X normal 24

ii) $n \geq 5$ for X close to normal

iii) IF X has highly skewed underlying POP, don't use the normal approx if $n \leq 30$.

iv) If $n \geq 100$ CLT usually holds in this class.

Usually I will say " X is normal"

(so \bar{X} is too) or "assume n is large enough so that the CLT holds"

ex) Common exam problem See fig 4.11 p 245 exception iii)

$n=100$ $X_i = i$ th IQ score
 \approx Normal with $\mu=100$ $\sigma=15$

i) Find $P(\bar{X} > 115)$

Key step =

Step 1 \bar{X} picture



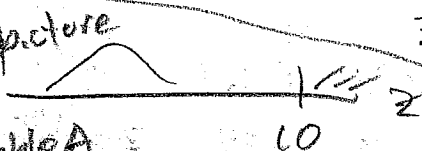
Step 0

$$\bar{X}, \mu_{\bar{X}} = 100 \quad \sigma_{\bar{X}} = \frac{15}{\sqrt{100}} = 1.5$$

Step 2 Standardize $Z = \frac{115 - 100}{1.5} = 10$

triplet, get from POP, not from sample

Step 3 Z picture



Z score is always from a triplet: $\frac{\text{value} - \text{mean}}{\text{SD value}}$

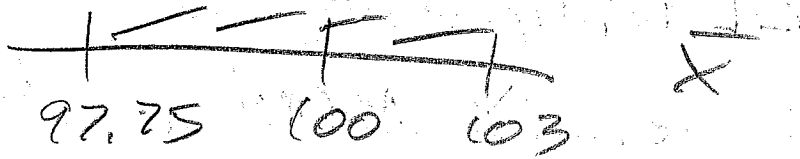
Step 4 table A

$$P(Z > 10) = 1 - P(Z < 10) = 1 - 1 = 0$$

ii) Find $P(97.75 < \bar{X} < 103)$

24.5

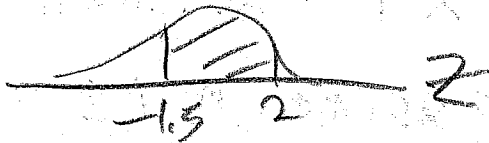
Step 1



Step 2 $z = \frac{97.75 - 100}{1.5} = -1.5$

$$\frac{103 - 100}{1.5} = 2$$

Step 3



Step 4 $= 0.9772 - 0.0669 = 0.9104$

$C P(-1.5 < z < 2) = P(97.75 < \bar{X} < 103)$

← likely exam, quiz question

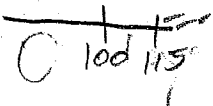
Ex) a) Suppose X comes from a highly skewed distribution $n=10$, $\mu=100$, $\sigma=15$.

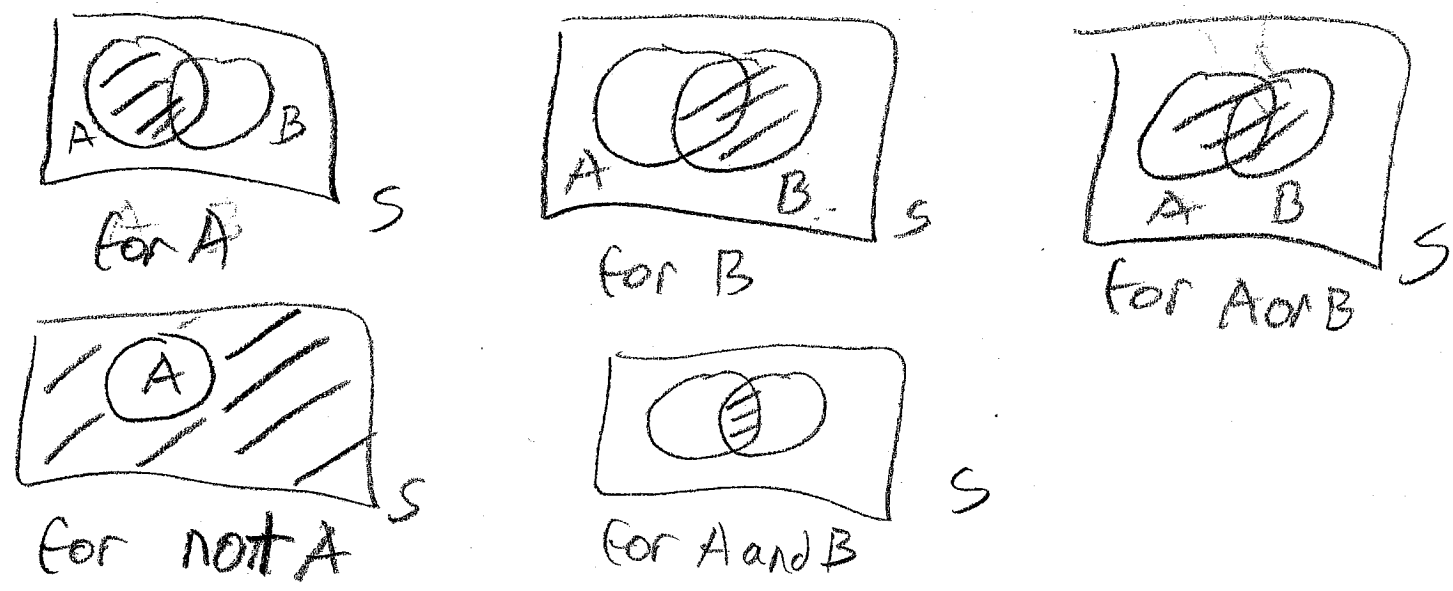
Find $P(\bar{X} \geq 115)$ if possible.

and X comes from a normal dist with

b) If $n=10$, find $P(\bar{X} \geq 115)$ if possible.

soln a) CLT does not apply, can't do it

b) C  \bar{X} , $\mu_{\bar{X}} = 100$, $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} = \frac{15}{\sqrt{10}} = 4.74$ $z = \frac{115 - 100}{4.74} = 3.16$
 $\Delta = 1 - .9992 = .0008$



2) p260 2 events A and B are independent
 if knowing that one occurs does not change
 the prob that the other occurs

n events A_1, \dots, A_n are independent if
 knowing that any subset of 1 to n-1 events
 occurs does not change the prob's of the
 others.

3) If events are not independent, then they
 are dependent.

ex If $0 < P(A), P(B) < 1$ and A and B are
 disjoint, then A and B are dependent
 since knowing A occurred means B did
 not occur (an extreme form of dependence)

ex Identical experiments done under identical
 conditions are independent eg coin toss 1

COMMON sense often works

height & IQ

hair color and height

Your final grade counts

25.9

4) Two events A and B are independent if $P(A \text{ and } B) = P(A)P(B)$ and vice versa

5) Multiplication rule p260

If A and B are independent $P(A \text{ and } B) = P(\text{both events A and B occur}) = P(A)P(B)$.

If events A_1, A_2, \dots, A_n are independent, then

$P(A_1 \text{ and } A_2 \text{ and } \dots \text{ and } A_n) = P(\text{all } n \text{ events occurred})$

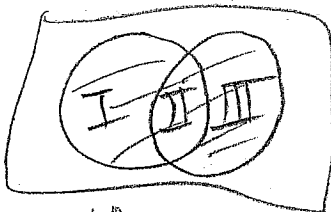
$= P(A_1)P(A_2)\dots P(A_n)$ (*) but if $n > 2$,

(*) does not imply that A_1, \dots, A_n are independent.

6) General P addition rule p264

~~$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$~~

~~$= P(A) + P(B) - P(A \text{ and } B)$~~



note $P(A) + P(B) = P(A \text{ or } B) + P(A \text{ and } B)$

I A and not B
II A and B
III B and not A

so $P(A \text{ or } B) = P(I) + P(II) + P(III) = P(A) + P(B) - P(II)$

7) Summary i) $0 \leq P(A) \leq 1$

ii) $P(S) = 1$

iii) complementary rule $P(\text{not } A) = 1 - P(A)$

iv) $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$ 0282 26
addition rule

v) If A and B are disjoint, then
 $P(A \text{ and } B) = 0$ while $P(A \text{ or } B) = P(A) + P(B)$ ←
add rule for disjoint events

vi) If A_1, A_2, \dots, A_n are disjoint, then
 $P(A_1 \text{ or } A_2 \text{ or } \dots \text{ or } A_n) = P(A_1) + \dots + P(A_n)$

vii) If A and B are independent, $P(A \text{ and } B) = P(A)P(B)$
and $P(A \text{ or } B) = P(A) + P(B) - P(A)P(B)$

viii) If A_1, A_2, \dots, A_n are independent then
 $P(A_1 \text{ and } A_2 \dots \text{ and } A_n) = P(A_1)P(A_2) \dots P(A_n)$

ix) $P(A) = \sum P(\text{outcomes in } A)$ if S is finite

x) List all possibilities using order of trials
eg roll 2 die, find sum toss coin 3 times, find # heads

B) P262-263 combining prob rules

Perform n identical experiments.

Interest is in a single outcome of the expt, say

Find i) $P(D \text{ occurred in none of the } n \text{ experiments})$

ii) $P(D \text{ occurred in at least one of the } n \text{ experiments})$

iii) $P(\text{all } n \text{ outcomes were } D)$

iv) $P(\text{not all outcomes were } D)$

i) $P(\text{none}) = \frac{1 - P(D)}{1} \cdot \frac{1 - P(D)}{2} \dots \frac{1 - P(D)}{n} = [1 - P(D)]^n$
comp rule mult rule

ii) $P(\text{at least one}) = P(\text{not none}) = 1 - P(\text{none})$

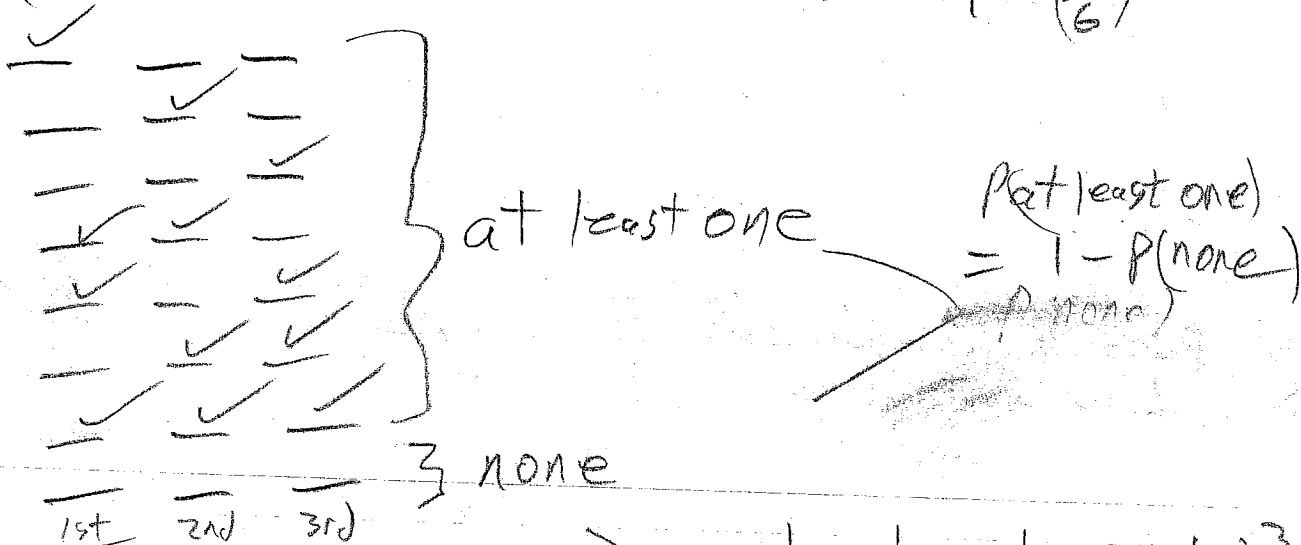
iii) $P(\text{all}) \stackrel{\text{mult rule}}{=} \frac{P(0)}{1st} \frac{P(0)}{2nd} \dots \frac{P(0)}{nth} = [P(0)]^n$

iv) $P(\text{not all}) = 1 - P(\text{all})$

ex) roll die 3 times \Rightarrow dice was a 5 in the next

i) $P(\text{none of the rolls are 5's}) = \frac{5}{6} \frac{5}{6} \frac{5}{6} = \left[1 - \frac{1}{6}\right]^3$

ii) $P(\text{at least one of the 3 rolls is a 5}) = 1 - \left(\frac{5}{6}\right)^3$



iii) $P(\text{all 3 rolls are 5's}) = \frac{1}{6} \frac{1}{6} \frac{1}{6} = \left(\frac{1}{6}\right)^3$

iv) $P(\text{not all 3 rolls are 5's}) = 1 - \left(\frac{1}{6}\right)^3$

ex) 70 rolls: $P(\text{none are 5's}) = \frac{5}{6} \frac{5}{6} \dots \frac{5}{6} = \left(\frac{5}{6}\right)^{70}$

Read ex 5.3 and 5.4 p 262-263

Carefully

ex: $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ 0.27

4 given only 3 - find the remaining

eg: A and B are ind. $P(A) = .2$ $P(B) = .6$
 Find $P(A \cup B)$.

Soln: $P(A \cap B) = .2(.6) = .12$

$P(A \cup B) = .2 + .6 - .12 = 0.68$

ex: $P(A) = .6$ $P(B) = .5$ are A and B disjoint or is more information needed?

Soln: if A and B are disjoint

$P(A \cup B) = P(A) + P(B) = .6 + .5 = 1.1$ impossible

so A and B are not disjoint

ex: 100 students 40 men 60 women
 5 men with beards randomly select a person

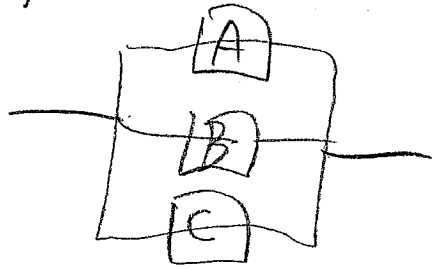
$P(A) = P(\text{man with beard is selected}) = \frac{5}{100} = .05$

$P(B) = P(\text{man selected}) = \frac{40}{100} = .4$

$P(A \cap B) = P(A) = \frac{5}{100} \neq P(A)P(B)$

so A and B are dependent

step ↓ ex i)



device fails if all components fail

27.5

exii) thief tries to break in house there are 3 triggers for alarm A B C

$$P(\text{device fails}) = P(\text{trigger goes off}) = \begin{matrix} A & B & C \\ 0.01 & 0.02 & 0.03 \end{matrix}$$

Prob

$$0.01 \cdot 0.02 \cdot 0.03 = \approx 0.000$$

$$0.99 \cdot 0.98 \cdot 0.97 = \approx 0.94$$

$$1 - 0.01 \cdot 0.02 \cdot 0.03 = \approx 1.000$$

$$1 - 0.99 \cdot 0.98 \cdot 0.97 = \approx 0.06$$

ex i)

all components fail
= device fails

none of the components failed

at least one component does not fail
= device runs

at least one component fails

ex ii)

thief set off all 3 triggers
(^{very} stupid thief)

thief does not set off any of the triggers
alarm does not go off

at least one trigger does not go off

thief fails to disable at least one trigger
= alarm goes off

↑ SHIP

§5.2 9) Let X be a count ^{of successes} (of 028)

eg take a SRS of size $n=100$ and
 $S>$ are F then $X=S>$ if F was counted

the proportion $\hat{p} = \frac{X}{n}$ so $X = n\hat{p}$

10) p269 The distribution of X can be derived

if i) there are n observations

ii) these are independent

iii) each observation is a 1 if it is
a "success" (what you want to count)

0 if not a success. Hence there
are 2 categories success and not a success

iv) $\text{prob}(\text{success}) = \text{Prob}(1) = p$

$\text{prob}(0) = 1-p$

11) p270 If i) - iv) hold, the distribution
of X is called a binomial distribution

with parameters n and p . X is $\text{bin}(n, p)$.

X	0	1	2	...	n	
prob	p_0	p_1	p_2	...	p_n	$p_i = P(X=i)$

12) p271 The count of X from a SRS
is approx $\text{bin}(n, p)$ if the sample size
 n is small compared to the pop size.

(3) p 272 The number of ways of arranging k successes among n obs's is

$$\binom{n}{k} = \frac{n!}{k! (n-k)!} = \text{binomial coefficient}$$

"n factorial"

(4) p 272 $n! = n(n-1)(n-2) \dots 3 \cdot 2 \cdot 1$

$= n(n-1)! = n(n-1)(n-2)! \text{ etc}$

$0! = 1$ by definition

$5! = 5 \cdot 4 \cdot 3 \cdot 2 \cdot 1 = 120$

$3! = 3 \cdot 2 \cdot 1 = 6 \text{ etc}$

ex $\binom{13}{3} = \frac{13!}{10! 3!} = \frac{13 \cdot 12 \cdot 11 \cdot (10!)}{(10!) 3!} = \frac{13 \cdot 12 \cdot 11}{3!}$

\downarrow
sum to 13

$= \frac{13 \cdot 12 \cdot 11}{6} = 286$

~~$\binom{n}{0} = \frac{n!}{0! n!} = 1$~~

~~$\binom{n}{n-k} = \binom{n}{k}$~~

~~$\binom{n}{1} = \frac{n!}{1! (n-1)!} = n$~~

~~$\binom{n}{n-1} = n$~~

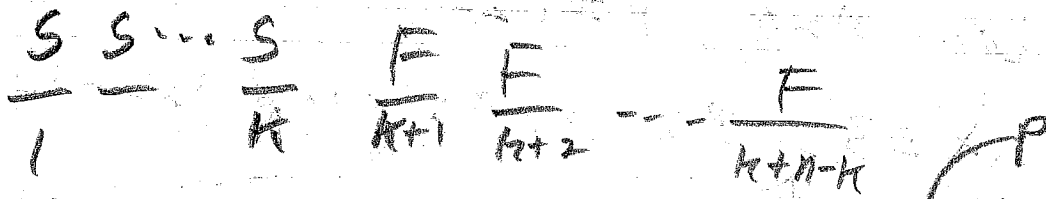
Next Up Previous
UNIVERSITY OF MINNESOTA
Twin Cities
Next: Week 8 Up: STATISTICS 8101 Homework Assignments Previous: Week 6
Luke Tierney
1999-12-13

School of Statistics

So $h(z) = \frac{1}{2} F(z)$

73) P 273 $P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$ F28226
 028229
 (for $k=0, 1, \dots, n$)

Idea: If $X=k$ then there were k success and $n-k$ failures.



By independence $P(S \dots S F \dots F) = [P(S)]^k (P(F))^{n-k}$

(A and B ind $\Rightarrow P(A \text{ and } B) = P(A)P(B)$)

But there are $\binom{n}{k}$ ways to choose the

k S's among the n outcomes

So $P(X=k) = \binom{n}{k} (p)^k (1-p)^{n-k}$ Fact $p^0 = 1$
 $(1-p)^0 = 1$

76) P274 IF X is bin n p then

$\mu_x = np$

$\sigma_x = \sqrt{np(1-p)}$

ex) a machine produces light bulbs, 80% meet specifications. A sample of 6 bulbs is taken from the machine's production. If 3 or more fail to meet specifications, the

Machine will be recalibrated.

Find the probability that the machine will be recalibrated.

Soln want $P(\text{at least 3 (are defective)})$

Let X count the number of defectives

X is bin $n=6$ $p=.2$

want $P(X \geq 3) = P(X=3) + P(X=4) + P(X=5) + P(X=6)$

add by $p=4$, the 4 events have no outcomes in common

$$= \binom{6}{3} (.2)^3 (.8)^3 + \binom{6}{4} (.2)^4 (.8)^2 + \binom{6}{5} (.2)^5 (.8)^1 + \binom{6}{6} (.2)^6 (.8)^0$$

$$= 20 (.008) (.512) + 15 (.0016) (.64) + 6 (.00032) (.8) + 1 (.000064)$$

$$= 0.08192 + 0.01536 + 0.001536 + 0.000064$$

$$= 0.09888 \approx 10\%$$

$$\mu_x = 6(.2) = 1.2$$

$$\sigma_x = \sqrt{np(1-p)} = \sqrt{6(.2)(.8)} = \sqrt{.96} = .9798$$

~~4.5-36) The sampling distribution of X obtained from a SRS of size n from the population of X having mean~~

17) Relationship with p 261-263

X counts number of successes

$$P(X=0) = P(\text{none were successes}) = \binom{n}{0} p^0 (1-p)^n = (1-p)^n = \frac{1-p}{1} \dots \frac{1-p}{n}$$

$$P(X=n) = P(\text{all were successes}) = \binom{n}{n} p^n (1-p)^0 = p^n = \frac{p}{1st} \dots \frac{p}{nth}$$

$$P(\text{at least one success}) = P(X \geq 1) = 1 - P(X=0) = 1 - (1-p)^n$$

$$P(\text{not all}) = P(X \leq n-1) = 1 - P(X=n) = 1 - p^n$$

18) at most and at least

$$P(\text{at least } j \text{ successes}) = P_j + P_{j+1} + \dots + P_n = P(X \geq j)$$

$$P(\text{at most } j) = 1 - P_0 - P_1 - \dots - P_{j-1}$$

$$P(\text{at most } j \text{ successes}) = P(0) + P(1) + \dots + P(j) = P(X \leq j) = 1 - P_n - P_{n-1} - \dots - P_{j+1} \quad \text{where } P_k = P(X=k)$$

ex] 20 multiple choice questions

for each, student randomly guesses

Let X = # questions answered correctly

i) $P(X=16) = \binom{20}{16} (.2)^{16} (.8)^4$ $n=20$ $p = \frac{1}{5} = .2$
 Not surprising always adds to n

$$= \frac{20!}{16! 4!} (.2)^6 (.8)^4 = 4845 (.2)^6 (.8)^4 = 1.3 \times 10^{-4}$$

$$P(\text{student gets at most 2 correct}) = P(X \leq 2)$$

$$= P_0 + P_1 + P_2 = \binom{20}{0} (.2)^0 (.8)^{20} + \binom{20}{1} (.2)^1 (.8)^{19} + \binom{20}{2} (.2)^2 (.8)^{18}$$

$$= (.8)^{20} + 20(.2)(.8)^{19} + 190(.2)^2 (.8)^{18}$$

$$= .015292 + .05776461 + .1369094$$

$$= \boxed{0.2062}$$

iii) $P(\text{student gets at least 3 correct})$

$$= P(X \geq 3) = 1 - P_0 - P_1 - P_2$$

$$= 1 - .2062 = \boxed{.7938}$$

$$= 1 - \binom{20}{0} (.2)^0 (.8)^{20} - \binom{20}{1} (.2)^1 (.8)^{19} - \binom{20}{2} (.2)^2 (.8)^{18}$$

18b) p 278 Suppose X is $\text{bin}(n, p)$
 (and $n \leq \frac{\text{pop size}}{10}$ usually true).

IF $np \geq 10$ AND $n(1-p) \geq 10$

then $X \approx N(\mu_x = np, \sigma_x = \sqrt{np(1-p)})$

20) If $np < 10$ OR $n(1-p) < 10$ DON'T
 use the normal approx. KNOW

ex) 20 multiple choice questions

0M282 3 /

$$P = \text{prob student gets correct answer} = \frac{4}{5}$$

knew 75
guessed
on 25
if took
many
over and
over again

i) If $n=20$ can normal approx be used?

$$nP = 20 \frac{4}{5} = 16 \geq 10 \quad n(1-p) = 20 \frac{1}{5} = 4 < 10$$

NO

ii) If $n=100$ find the prob student gets at least 85 questions correct.

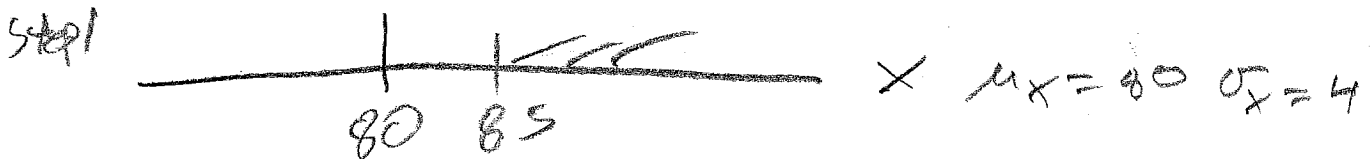
$$\text{Soln } nP \geq n(1-p) = 100 \frac{4}{5} = 80 \geq 10$$

use normal approx

$$\text{step 0) } \mu_x = np = 100 \frac{4}{5} = 80$$

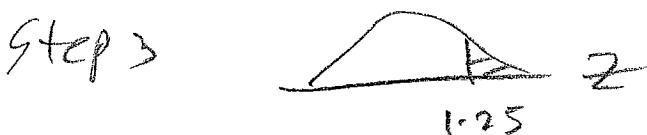
$$\sigma_x = \sqrt{np(1-p)} = \sqrt{100 \frac{4}{5} \frac{1}{5}} = \sqrt{16} = 4$$

sums to 1



Step 2 Standardize $z = \frac{x \text{ value} - \mu_x}{\sigma_x} = \frac{85 - 80}{4}$

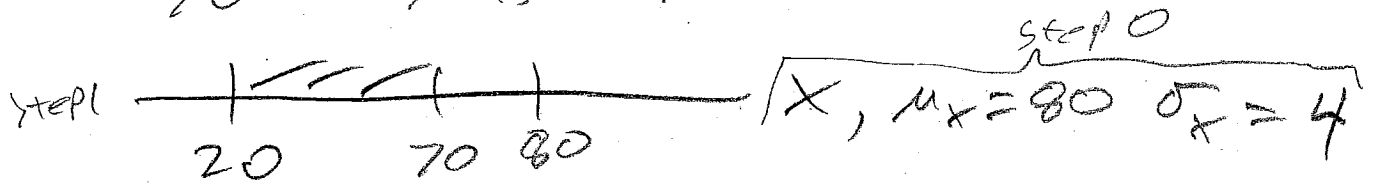
$$= \frac{5}{4} = 1.25$$



Step 4 Table A $P(z > 1.25) = 1 - P(z \leq 1.25)$

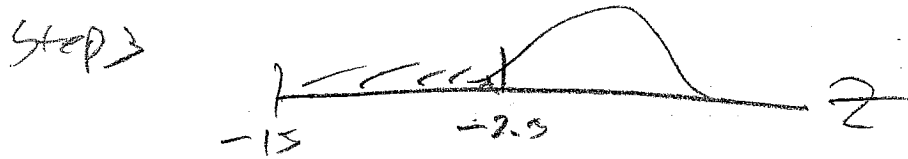
$$= 1 - .8944 = .1056$$

iii) Find Prob Student gets between 20 and 70 questions correct.



step 2

$$\frac{70-80}{4} = -2.5 \qquad \frac{20-80}{4} = -15$$

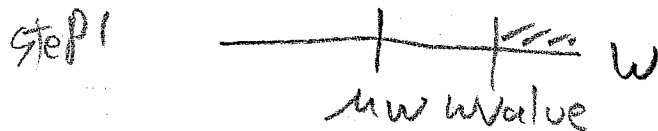


$$= P(Z \leq -2.5) - P(Z \leq -15)$$

$$= .0062 - 0 = .0062$$

21) 3 calculations with table A ← later

step 0) Find w, μ_w, σ_w



step 2) Find $Z = \frac{w \text{ value} - \mu_w}{\sigma_w}$



step 4) table A

key words

step 0)

X	μ_x	σ_x
\bar{X}	$\mu_{\bar{X}} = \mu$	$\sigma_{\bar{X}} = \sigma/\sqrt{n}$
X	$\mu_x = NP$	$\sigma_x = \sqrt{NP(1-P)}$

individual $n=1$
 mean $n>1$
 count of categorical

decide which to use

\$6 | 1) P 298

Statistical inference assumes

#0282

32

that data comes from a random sample or randomized experiment. If this assumption is false, your conclusions may be way off.

ex vol resp sample Ann Landers

2) After an experiment, all probabilities are 0 or 1.

ex] Flip coin. Before the exp't (Flip)

$P(H) = P(T) = \frac{1}{2}$. Suppose that after the flip, the coin is H. The $P(\text{it was H}) = 1$

$P(\text{it was a T}) = 0$.

3) P 306

Large sample confidence interval (CI)

for μ when σ is known.

Take a SRS of size n (large enough for CLT to hold, eg pop normal) from a pop with unknown mean μ and known s.d.

Then a $100 - C\%$ CI = $(1 - \alpha)100\%$ CI for

μ is $\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}}$ where $P(Z \leq z^*) = 1 - \frac{\alpha}{2}$

so $P(-z^* \leq Z \leq z^*) = C$, when Z is $N(0,1)$.

4) P^{306} z^*	1.645	1.96	2.576	(325)
C 100%	90%	95%	99%	

eg 90% CI $P(Z \leq z^*) = 1 - \frac{\alpha}{2}$

$1 - \alpha = .9$ so $\alpha = .1$ $1 - \frac{\alpha}{2} = .95$

$P(Z \leq z^*) = .95$ so $z^* = 1.645$

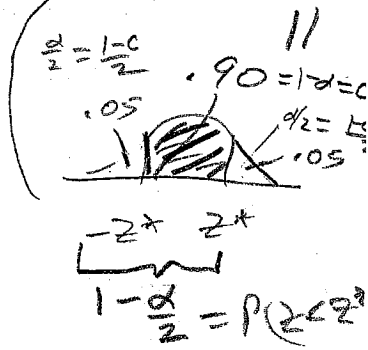


Table C row z^* has many values

ex 99% $z^* = 2.326$

5) Idea: i) \bar{X} is approx $N(\mu, \frac{\sigma}{\sqrt{n}})$ if CLT holds

Let $A = P\left(\mu - 2\frac{\sigma}{\sqrt{n}} \leq \bar{X} \leq \mu + 2\frac{\sigma}{\sqrt{n}}\right) \approx .95$

Since $P(\bar{X}$ is within $2\sigma_{\bar{X}}$ of $\mu_{\bar{X}}) \approx .95$ by the 68-95-99.7 rule.

So $A \approx P\left(-\frac{2\sigma}{\sqrt{n}} \leq \bar{X} - \mu \leq \frac{2\sigma}{\sqrt{n}}\right) =$

$P\left(-\bar{X} - \frac{2\sigma}{\sqrt{n}} \leq -\mu \leq -\bar{X} + \frac{2\sigma}{\sqrt{n}}\right)$

$= P\left(\bar{X} - \frac{2\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + \frac{2\sigma}{\sqrt{n}}\right) \approx .95$

ii) Suppose $\sigma = 10, n = 100$ 0282 33
 $\bar{x} = 90$

and CLT holds.

$$\text{Then } \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = 1$$

and $2\sigma_{\bar{x}} = 2$. Want to know for

what values of μ is it "plausible"

that \bar{X} came from a $N(\mu, \sigma = 10)$ distribution.

Fix μ and take "plausible" to mean

that the chance that \bar{X} is within $2\sigma_{\bar{x}}$ of μ is at least 95%.

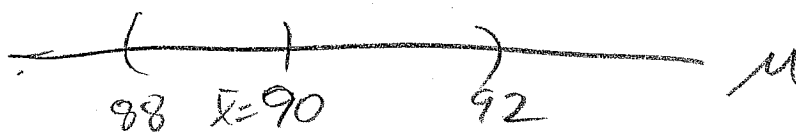
eg $\mu_1 = 90$ $\mu_2 = 400$ $\mu_3 = 88$ $\mu_4 = 87$
 $\mu_5 = 92$.

μ_3 is the smallest "plausible" value

(by the 68 - 95 - 99.7 rule) and

$\mu_5 = 92$ is the largest plausible value,
1.96 is better

So the 95% CI for μ is $\approx \bar{x} \pm 2 = (88, 92)$



plausible
values
of μ

6) ^{p301} why 95% confidence instead of 95% prob?

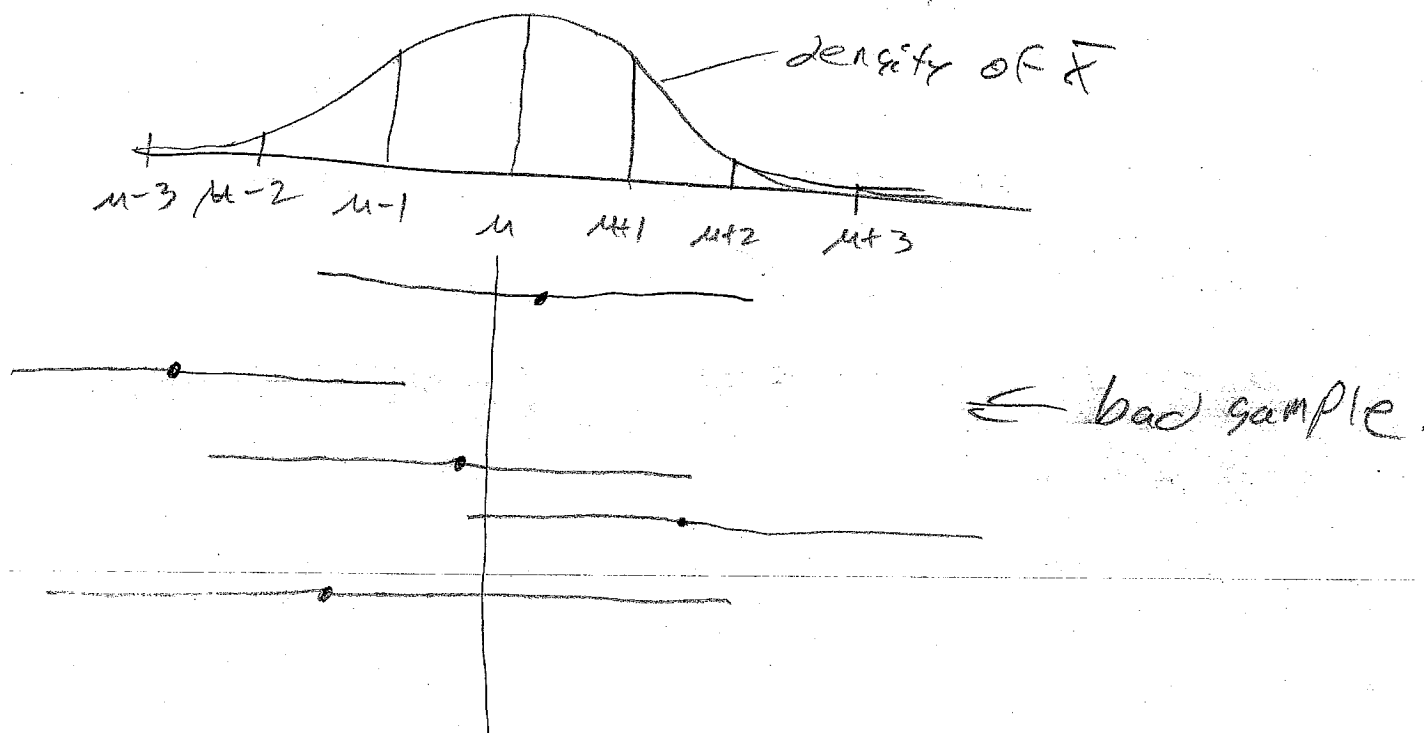
A CI is computed after the experiment.

Hence $P(\text{CI contains } \mu) = \begin{cases} 0 & \text{if it does not} \\ 1 & \text{if it does} \end{cases}$

If you make 100 95% CI's (100 different SR's), approx 95 will contain μ and 5 will not.

see Figure 6.4 p 302

suppose \bar{X} is $N(\mu_{\bar{X}} = \mu, \sigma_{\bar{X}} = 1)$



Each CI is $(\bar{x} - 1.96, \bar{x} + 1.96)$ and has length $2(1.96) = 3.92$. If \bar{x} is more than 1.96 away from μ , then the CI will NOT contain μ .

If you get a "bad sample", the CI does not contain μ . Chance of bad sample = 5% for a 95% CI.

2) p308 margin of error $(= z^* \frac{\sigma}{\sqrt{n}})$ ~~0.282~~ 34

$= \frac{1}{2}$ CI length is a measure of how accurate we believe our estimate of the parameter (μ) is.

ex) $n = 91$ $\bar{x} = 124510$ $\sigma = 180000$
MBA salaries in 1994 if got MBA in 70's and are sole source of income

Find a 90% CI for μ .

$z^* = 1.645$ eg. table C

$$\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}} = 124510 \pm 1.645 \frac{180000}{\sqrt{91}}$$

$$= 124510 \pm 3103.97$$

$$= (121406.03, 127613.97)$$

8) p309 $\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}}$ CI length increases if c do
if σ do,
decreases if n increases

9) ^{know} p341 For a $c(100\%) = (L, U)$ 100% CI to estimate μ within a margin of error m , take a SRS of size $n = \lceil \frac{z^* \sigma}{m} \rceil^2$ round up.

Want to estimate ave weight

34.5

to within 0.1 oz with 95% confidence, and $\sigma = 1.02$.

Find n .

Soln $m = 0.1$ $1 - \alpha = .95$

$\alpha = .05$ $P(Z < 1.96) = 1 - \frac{\alpha}{2} = .975$

So $z^* = 1.96$ (or table \leftarrow).

$$n = \left(\frac{z^* \sigma}{m} \right)^2 = \left(\frac{(1.96)(1)}{0.1} \right)^2 =$$

$$(19.6)^2 = 384.16$$

So take $(n = 385)$
(round up so level is at least 95%)

10) $P = 392$ Use CI if

i) data is SRS from pop

or ii) from a randomized experiment

ii) n large enough for CLT to hold
from all random experiments

ii) Sampling methods other than SRS's

such as multistage samples give

distributions that do not have $\sigma_x = \frac{\sigma}{\sqrt{n}}$

90 $\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}}$ is not correct for probability methods more complex than simple random samples.

ex CI from voluntary response sample is not trustworthy

§6.2 12) p 321 The null hypothesis = H_0

is a claim about a pop characteristic (parameter)

such as $H_0: \mu = \mu_0$ (eg $\mu_0 = 125 \text{ lbs}$)

13) The alternative hypothesis H_A is a competing

claim eg $H_A: \mu \neq \mu_0$

In ch 6 and 7,

ex H_0 is of the form "parameter = value"
 H_A " " " " " "

which is correct?

a) $H_0: \mu > 9$ $H_A: \mu = 9$

b) $H_0: \bar{x} = 9$ $H_A: \bar{x} > 9$

c) $H_0: \mu = 9$ $H_A: \mu > 9$

14) A test statistic is used to decide whether to reject H_0 or fail to reject H_0

15) p 321 The p value is the probability, assuming H_0 is true, of obtaining a test statistic at least as extreme or

"contradictory" to H_0 as what was (39.9) actually observed, "contradictory" is determined by the form of H_A ,

H_0 says: the difference (of the test statistic) from the H_0 model is due to chance

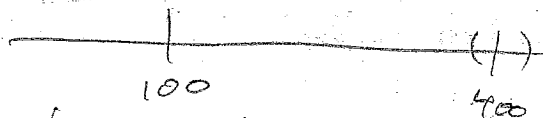
H_A : the difference of the test statistic from the H_0 model is real

P-value near 0 means test statistic was very unlikely to occur if

ex) $H_0 \mu = 100$ (σ known to be 1) H_0 was true
 $H_A \mu > 100$

Sample of size 100 $\bar{X} = 400$

$$\sigma_{\bar{X}} = \frac{1}{\sqrt{100}} = .10$$



would expect difference to be real,

If $\bar{X} = 90$, H_0 is less contradictory than H_A although neither seems likely,

16) A test has 4 steps p. 328 and examples

- i) State H_0 and H_A
- ii) calculate the test statistic
- iii) Find the P value.
- iv) State a conclusion!

If p value $\leq \alpha$, reject H_0
otherwise fail to reject H_0 .

Say in words what rejecting or failing to reject H_0 means.

17) P327 IF α is not given, use $\alpha = 0.05$.
 Test is statistically significant if $p\text{-value} \leq \alpha$.
 Note: Small p value is evidence against H_0

0282 (36)

18) P329 SRS of size n , unknown μ , known σ
 CLT holds if H_0 is true.

$H_0 \mu = \mu_0$ vs $H_A \mu \neq \mu_0$

test statistic $Z_0^* = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \approx N(0,1)$ if H_0 is true
 (obtained from H_0)

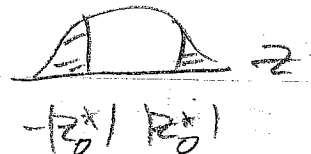
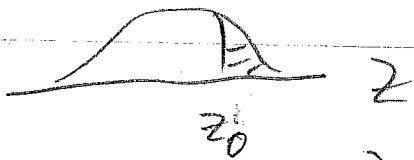
p values

major typo in text

right tail test
 $H_A \mu > \mu_0$

left tail test
 $H_A \mu < \mu_0$

2 tail
 $H_0 \mu \neq \mu_0$



$p\text{val} = P(Z > z_0^*)$
 $= 1 - P(Z < z_0^*)$ table A
 recall absolute value

$p\text{val} = P(Z < z_0^*)$

$p\text{val} = 2P(Z > |z_0^*|)$
 $= 2P(Z < -|z_0^*|)$ table A

so $| -1 | = 1, | 7 | = 7, | 18 | = 18$ etc

19) P337 For a level α 2 tail test,

$H_0 \mu = \mu_0$ is rejected if μ_0 is outside of a $(1-\alpha)100\%$ CI for μ , otherwise fail to reject H_0 .

(eg $\alpha = 5\%$ test is rejected if μ_0 is outside a 95% CI)

4 p313

20) To decide the form of H_A , choose the form that goes with strong evidence.

$\bar{x} = 35 \quad \sigma = 5 \quad n = 25$

ex) Ford says Escort gets at least 37 MPG

i) $H_0 \mu = 37 \quad H_A \mu > 37$ ii) $z_0^* = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} = \frac{35 - 37}{5/\sqrt{25}} = -2$

iii) $pval = P(Z > -2) = 1 - 0.0228 = 0.9772$

iv) Fail to reject H_0 Escort does not get at least 37 MPG

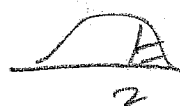
ex) Consumer group says Escort gets less than 37 MPG

i) $H_0 \mu = 37 \quad H_A \mu < 37$ ii) $z_0^* = \frac{35 - 37}{5/\sqrt{25}} = -2$

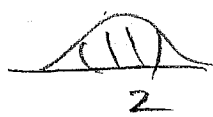
iii) $pval = P(Z < -2) = 0.0228$

iv) $pval < \alpha = 0.05$ so reject H_0 . There is strong evidence that Escort does not get 37 MPG.

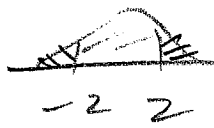
ex) $H_0 \mu = 100 \quad \alpha = 0.05 \quad z_0^* = 2.0$

right tail $H_A \mu > 100$  $pval = P(Z > 2)$

$= 1 - 0.9772 = 0.0228 < 0.05$ reject H_0

left tail $H_A \mu < 100$ 

$pval = 0.9772 > 0.05$ fail to reject H_0

2-tail $H_A \mu \neq 100$ 

$pval = 2 P(Z < -2) = 2(0.0228) = 0.0456 < 0.05$
reject H_0

ex) $H_0 \mu = 100 \quad \alpha = 0.05 \quad z_0^* = -2.0$

right tail $H_A \mu > 100$ $P_{val} = \frac{1 - 0.0228}{2} = 0.4886$

$= 1 - P(Z < -2) \Rightarrow 1 - 0.0228 = 0.9772 > 0.05$
fail to reject H_0

left tail $H_A \mu < 100$ $\frac{0.0228}{2} = P_{val} = 0.0228 < 0.05$
reject H_0

2 tail $H_A \mu \neq 100$ $\frac{0.0228}{2} = P_{val} = 2 P(Z < -2)$

$0.0456 < \alpha = 0.05$ reject H_0

2) P_{345} Assume σ (skip $P_{337} - 335$)
To use test need SRS

or independent observations from pop.

Need n large enough for CLT to hold

see ex's in book I will do ex's when

both μ and σ are unknown,
(need $n \approx \frac{\text{pop size}}{10}$ & almost always true)

532) P_{343} $P_{value} = 0.049$ $P_{value} = 0.051$

are not practically different. ~~stat~~

~~a practical~~. Always using $\alpha = 0.05$ makes little sense. State the p value and let your audience decide if the result is significant.

23) P_{344} Statistical significance \neq

Practical significance.

$H_0 \mu = 475$

$\mu = 478$

= mean math score

475 and 478 ... hardly different

but the p value could be 0.000000001 / 3
if the sample size is large enough,

24) p 345 ~~Statistical significance~~ says that something other than chance caused a small p value but does not say what the cause is. Surgery reduced pain but so did giving patients a gel capsule placebo effect.

25) p 346 Suppose $\alpha = .05$. If you test many effects, about 5% will be declared significant (reject H_0) even when they are not. Alternatively, if someone unethical wants to state an effect is significant, they could take 20 samples. About 1 of them should lead to rejecting H_0 even if H_0 is true.

EX) There is a huge on going study of nurses. Every nurse fills out a big questionnaire. Every few months some association between eg diet weight smoking and say cancer obesity etc is found and published.

The study can suggest associations, but much more work needs to be done to declare that the behavior is causing bad health. Need controlled experiments (Section 3.2)

§ 6.4 26) p 350

decision	reject H_0	+ truth H_0 true	H_0 false
		Type I error	correct decision
Fail to reject H_0	Fail to reject H_0	correct decision	Type II error

Type I error: reject H_0 when H_0 is true 02-82 -38

Type II error: fail to reject H_0 when H_0 is false

27) p 350 $\alpha = P(\text{type I error if } H_0 \text{ is true})$

So for $\alpha = .05$, H_0 will be rejected

5% of the time when H_0 is true; (100 samples about 5 will reject H_0).

Step power p 354

Ch 7 If you like beer read p 364.

1) For ch 7 assume $n \leq \frac{\text{pop size}}{10}$ unless pop size is small.

2) p 367 An SD estimated from the data is called a standard error SE.

$$\text{ex) } \sigma_{\bar{x}} = SD(\bar{x}) = \frac{\sigma}{\sqrt{n}} \quad SE(\bar{x}) = \frac{s}{\sqrt{n}}$$

2) p 367 If a SRS of size n is drawn from a $N(\mu, \sigma)$ population, then the sampling distribution of

$$t_0 = \frac{\bar{x} - \mu}{s/\sqrt{n}} \text{ has a } \underline{t \text{ distribution with}}$$

$n-1$ degrees of freedom, denoted by $t_{(n-1)}$.

3) P370 Suppose a SRS of size n is drawn from a POP with unknown mean μ and SD σ . If n is large enough for the CLT to hold, then a $100(1-\alpha)\% = 100c\%$ CI for μ is

$$\bar{x} \pm t^* \frac{s}{\sqrt{n}} \quad \text{where}$$

$$P(t_{n-1} \leq t^*) = 1 - \frac{\alpha}{2} \quad \text{or} \quad P(-t^* \leq t_{n-1} \leq t^*) = c$$

4) Get t^* from table C. If $df = n-1 > 30$, use z^* rather than t^* for $df = 40-1000$.

ex) $\bar{x} = 100 \quad n = 16 \quad s = 40$

$df = 16 - 1 = 15$ Find a 95% CI
assume CLT holds

df	
15	2.131
	95% conf level C

$$95\% \text{ CI is } 100 \pm 2.131 \frac{40}{\sqrt{16}} = 100 \pm 21.31$$

$$= (78.69, 121.31)$$

ex) $\bar{x} = 100 \quad n = 32 \quad s = 40$ Find 95% CI
 $df = 31 > 30$ so go to $z^* = 1.96$

$$100 \pm 1.96 \frac{40}{\sqrt{32}} = 100 \pm 13.81 = (86.19, 113.81)$$

Using $df=30$ would give $100 \pm \dots$ U-82 39

But rather than interpolating or taking closest, just use z^* . A $t_{(n-1)}$ distr is almost a $N(0,1)$ dist for $n-1 \geq 30$ or $n \geq 31$

5) Key assumption is n large enough for CLT to hold, eg pop is normal

5) p-370 One sample t-test
GR 4 of size $n \geq 2$ large enough for CLT to hold from a pop with unknown

SD: σ

Test $H_0: \mu = \mu_0$

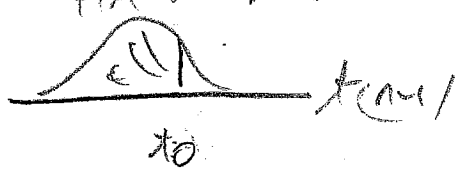
test stat $t_0 = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \sim \begin{cases} t_{(n-1)}, & n-1 \leq 30 \\ N(0,1), & n-1 > 30 \end{cases}$

use table A if $n-1 \geq 30$ or $n \geq 31$
 E if $n-1 \leq 30$ or $n \leq 31$ to find p-values.

left tail
 $H_A: \mu < \mu_0$

right tail
 $H_A: \mu > \mu_0$

2 tail
 $H_A: \mu \neq \mu_0$



- same 4 step procedure
- i) State H_0 and H_A
 - ii) calculate test statistic t_0
 - iii) Find p-value
- $P_{val} = P(t_{(n-1)} > t_0), n-1 \leq 30$
 $P(z \leq t_0), n-1 > 30$
 $P_{val} = P(t_{(n-1)} \geq t_0), n-1 \leq 30$
 $1 - P(z \leq t_0), n-1 > 30$
 $P_{val} = 2P(t_{(n-1)} \geq |t_0|), n-1 \leq 30$
 $2P(z > |t_0|), n-1 > 30$

iv) state a conclusion
 if $p\text{val} < \alpha$ reject H_0 39.5
 otherwise fail to reject H_0
 use $\alpha = .05$ if not given,
 state in words what rejection or
 failing to reject H_0 means.

ex $H_0: \mu = 0$ $H_a: \mu > 0$ $n = 15$
 normal pop $t_0 = 1.97$

- a) Find df $df = n - 1 = 14$
 b) Find the 2 t^* values from table
 C that bracket t_0 .

df	upper tail prob	
	.05	.025
14	1.761	2.145

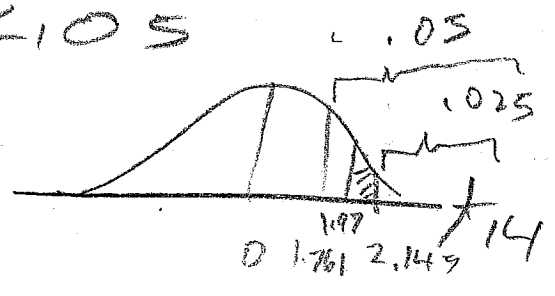
c) what are the right tail probs of these
 t^* values?

.05 and .025

d) what is the p value?

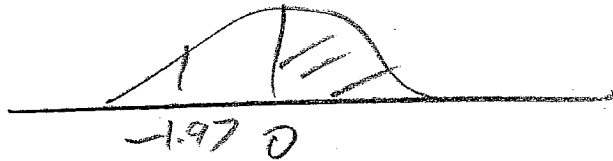
$.025 < p\text{value} < .05$

e) Is the test
 significant at 5% level?
 1% level?



$p\text{val} < 5\%$ so significant at 5%
 $p\text{val} > 1\%$ so not sig at 1%.

i) if $t_0 = -1.97$, what is pval = 0.282 48

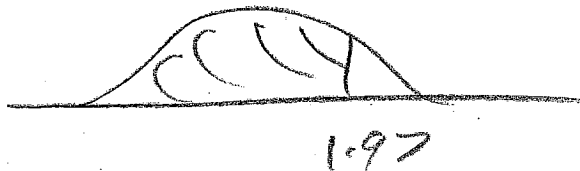


pval > 0.5 fail to reject H_0

in fact $.95 < pval < .975$

\uparrow \uparrow
 $1-.05$ $1-.025$

ii) $t_0 = 1.97$ $H_A: \mu < 0$ what is pval



pval > .5

in fact $.95 < pval < .975$, fail to reject H_0

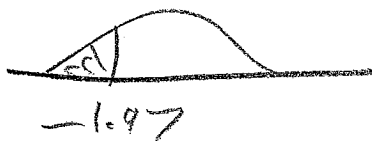
iii) $t_0 = 1.97$ $H_A: \mu \neq 0$

$2(.025) < pval < 2(.05)$

or $.05 < pval < .1$ fail to reject H_0

iv) $t_0 = -1.97$ $H_A: \mu < 0$

by symmetry
 \downarrow
 \equiv



so $.025 < pval < .05$

v) $t_0 = 1.97$ $n = 32$ $H_A: \mu > 0$.

df = 31 > 30 so use table A

$= 1 - .9796 = 0.0204$



ex] $\alpha = .01$ Golpher tries out 40.5

club to see if new club drives balls further than old club. with old club he believes he can hit ball 200 yards. Will buy club if it is better. Shop lets him take 10 drives. $n=10$ $\bar{x} = 204.60$ $s = 9.03$

soln $\mu =$ mean drive length with new club

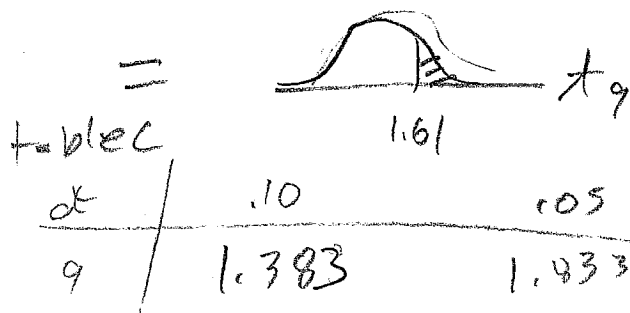
step 1) $H_0: \mu = 200$
 $H_A: \mu > 200$

μ_0 is never from sample

step 2) $t_0 = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{204.60 - 200}{9.03/\sqrt{10}} =$

$\frac{4.6}{2.8555} = 1.611$

step 3) $df = n - 1 = 9$ $pval = P(t_9 > 1.61)$



$.05 < pval < .1$

step 4) Fail to reject H_0 since $pval > \alpha = .01$.

Not enough evidence to conclude that the new club is better than the old club.

Common final problem

Test from

computer output

0282

amount of relief from sugar pill

48

ex Test of $\mu = 0.000$ vs $\mu \neq 0.00$

	N	MEAN	STDEV	SEMEAN	T	P
sugar pill	8	4.000	4.3094	1.5236	2.63	0.034

Step 1) $H_0: \mu = 0$ $H_A: \mu \neq 0$

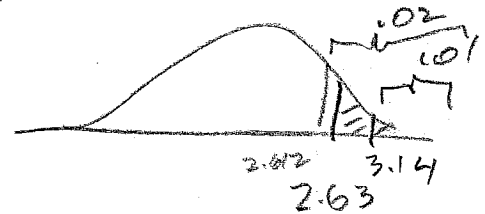
Step 2) $t_0^* = 2.63$

Step 3) $p\text{-value} = 0.0341$

Step 4) $p\text{-val} = .0341 < .05$ so reject H_0

sugar pill brings some relief

Note computer pvalue is better than values in table C.



df	.02	.01
7	2.612	3.143

2 tail test so

$$2(.01) = .02 \leq p\text{-val} \leq .04 =$$

7) Matched Pairs + IP 375 - 1122

Data comes in pairs $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

x_i and y_i are dependent, that is knowledge of x_i should give some information about y_i . The pairs (x_i, y_i) and (x_j, y_j) are independent for $i \neq j$.

That is knowledge of the 1st pair gives no info about the 2nd pair etc. uncorrelated

Often 2 measurements are taken from the same subject
eg X_i = cholesterol level for ^{the person} i th person before the
 Y_i after

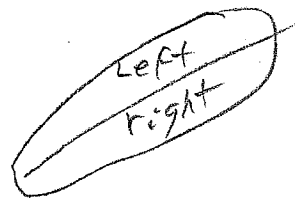
ex) correlated so use matched pairs
age of husband age of wife

compute D_1, \dots, D_n where $D_i = X_i - Y_i$
to get 1 sample of D's.

- A) Perform a t CI on the D's
- B) perform a t test on the D's

ex Matched paired t test

Use 2 different viral preparations on a tobacco leaf. Divide leaf into left half right. Flip coin, put virus 1 on left if heads otherwise on



Count number of lesions due to each
Dependent since strong leaf should resist both viruses weak plant or will have less resistance.

Suppose researchers expect virus 1 to cause more lesions than virus 2.

$$\bar{D} = \frac{\sum D_i}{n} = \frac{\sum (X_i - Y_i)}{n} = 1.450 \quad S_D = \sqrt{\frac{\sum (D_i - \bar{D})^2}{n-1}}$$

$n=20$ μ_D = mean of $X_i - Y_i$

4 step test

$x_i > \mu_0$ so $x_i - \mu_0 > 0$ #0282/42

Step 1 $H_0 \mu_D = 0$ $H_A \mu_D > 0$

Step 2 $t_0 = \frac{\bar{D} - \mu_0}{s_D / \sqrt{n}} = \frac{1.450 - 0}{3.203 / \sqrt{20}} = 2.02$

Step 3 $df = 19$ table C

	0.05	0.025
19	1.729	2.093

so $0.025 < p\text{-val} < 0.05$

Step 4) reject H_0 since $p\text{-val} < 0.05$.

Virus 1 causes more legions than Virus 2.

Find a 90% CI

$\bar{D} = 1.450$ $s_D = 3.203$ $n = 20$ $df = 19$

table C $t^* = 1.729$

90% CI for μ_D is

$1.450 \pm 1.729 \frac{3.203}{\sqrt{20}} = 1.450 \pm 1.238 = (0.212, 2.688)$

ex) Suppose $n = 36$ instead of 20

Step 1 $H_0 \mu_D = 0$ $H_A \mu_D > 0$

Step 2 $t_0 = \frac{\bar{D} - \mu_0}{s_D / \sqrt{36}} = \frac{1.450 - 0}{3.203 / \sqrt{36}} = \frac{1.450}{0.5338} = 2.72$

Step 3 $df = 36 - 1 = 35 > 30$ so use table A.

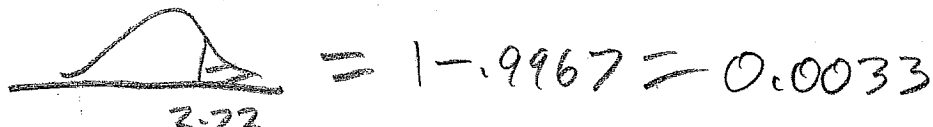


table C w
say 0.002
pval
table A is more

Step 4. $p\text{-val} < \alpha = .05$ so reject H_0

ex) p377 output) i) $H_0: \mu_D = 0$ $H_A: \mu_D > 0$ ii) $T^* = .35$ iii) $p\text{-val} = .365$ iv) fail to reject H_0

8) p380-5 when to use t test + CI

SRS from population (with finite SD) or data measurements from an experiment

a) $n < 15$ Use if data is close to normal. If data is skewed or outliers are present don't use t procedures.

b) $n \geq 15$ Use t test unless outliers are present or the data is strongly skewed.

c) If n is large enough, t procedures can be used even on strongly skewed data.

Text suggests $n \geq 40$ for clearly skewed data.

d) Don't use a t procedure if the pop mean is known

see ex 7.5 p380

9) 2 sample procedures (other than matched pairs) need 2 independent samples

or an experiment that used randomization to form a treatment and control group ($\frac{1}{2}$ get medic, $\frac{1}{2}$ get plac)

	size	mean	SD
10) * p392	n_1	\bar{x}_1	s_1
SRS			
	n_2	\bar{x}_2	s_2
SRS			

p395

$df = k = \text{smaller of } n_1 - 1, n_2 - 1$

if output is not being used,

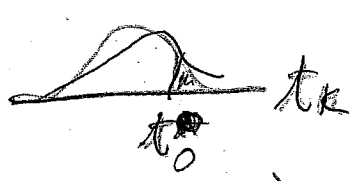
$H_0 \mu_1 = \mu_2$ or $H_0 \mu_1 - \mu_2 = 0$

(or $H_0 \mu_1 - \mu_2 = a$).

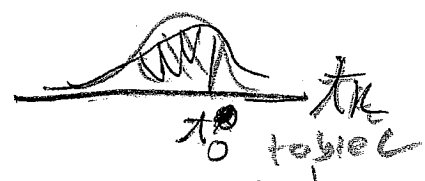
test statistic $t_0 = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$

p values

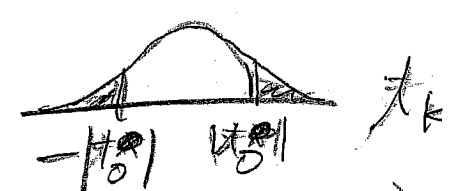
right tail	left tail	2 tail
$H_A \mu_1 > \mu_2$	$H_A \mu_1 < \mu_2$	$H_A \mu_1 \neq \mu_2$



$P(t_k > t_0^*)$ $k \leq 30$
 $1 - P(z \leq t_0)$ $k > 30$



$P(t_k < t_0^*)$ $k \leq 30$
 $P(z < t_0)$ $k > 30$



$2P(t_k > t_0^*)$ $k \leq 30$
 $2P(z < -|t_0|)$ $k > 30$

100% CI = $100(1 - \alpha)\%$ CI for $\mu_1 - \mu_2$

is $(\bar{x}_1 - \bar{x}_2) \pm t^* \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$

where $P(-t^* < t_k < t^*) = C = 1 - \alpha$

or $P(t_k < t^*) = 1 - \alpha/2$, $k = df = \text{smaller of } n_1 - 1, n_2 - 1$

12) *P408 Use 2 sample t if $\mu_1, \mu_2, \sigma_1, \sigma_2$ are unknown, NO reason to believe that $\sigma_1 = \sigma_2$.

pop	mean	SD
x_1	μ_1	σ_1
x_2	μ_2	σ_2

If $n_1 = n_2$ and the shapes of the 2 histograms are similar, can use t test even for $n_1 = n_2$ small ($n_1 = n_2 \geq 5$)
 If $n_1 \neq n_2$ want CLT to hold for both samples.

13) P405 \checkmark pooled 2 sample t tests and CI's are used if $\mu_1, \mu_2, \sigma_1, \sigma_2$ are unknown and it is believed (before the experiment) that $\sigma_1 = \sigma_2$.

You only need to know how to do the test from output.

ex P405 variances T DF Prob > |T|
 assume SD's are equal equal 2.9912 10.0 0.0135

step 1 $H_0: \mu_1 = \mu_2$ $H_A: \mu_1 \neq \mu_2$

step 2 $\bar{T}_0 = 2.9912$

df	10	10
	2.264	3.16

step 3 pval = .0135

pval = .0135 < $\alpha = .05$ table C $2(.005) < pval$ or $.01 < pval < .$

step 4) ~~reject H0~~ reject H_0
 DDT mean is not equal to control group mean C.O.L.

Warning Do not look at sample variances to decide whether to use 2 sample t or pooled 2 sample t. Use pooled 2 sample t if the story problem says the population SD's or variances are approximately

equal.

ex) Bank has 2 proposals to increase the amount of credit charged on its credit cards. Each proposal is offered to a SRS of 150 of its customers.

proposal	n	\bar{x}	s
A	150	1987	392
B	150	2056	413

\bar{x}_A estimates μ_A = mean (total) annual amount charged by customers who used proposal A

Test whether there is a difference in μ_A and μ_B .

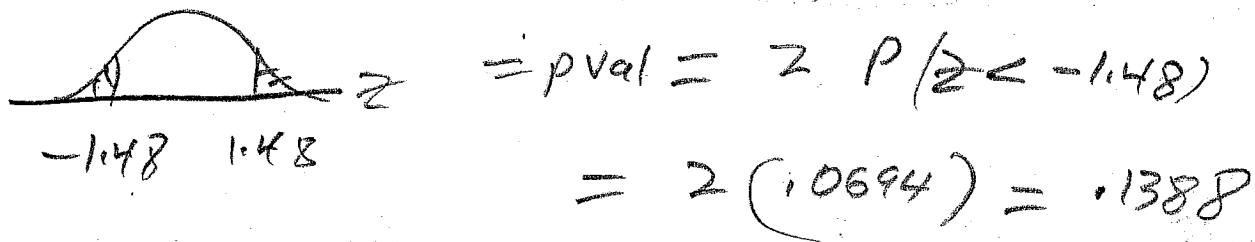
step 1 $H_0: \mu_A - \mu_B = 0$ $H_A: \mu_A - \mu_B \neq 0$

$$2) t_0 = \frac{\bar{x}_A - \bar{x}_B}{\sqrt{\frac{s_A^2}{n_1} + \frac{s_B^2}{n_2}}} = \frac{1987 - 2056}{\sqrt{\frac{(392)^2}{150} + \frac{(413)^2}{150}}} =$$

$$\frac{-69}{\sqrt{1024.43 + 1137.13}} = \frac{-69}{\sqrt{2161.55}} = -1.484$$

44.5

Step 3 $df = n - 1 = 149$ so use table A



iv) $pval > 0.05$ fail to reject H_0 .

The Δ means charge amounts from the mounts (2 proposals) are the same to chance, see ex 7.9 p 398-04

ex Is the angular velocity of the knee higher for skilled rowers than for novice rowers?

Output	N	mean	STD DEV	STD error	DF	P
skilled	10	4.1828	0.47906	0.15149		
novice	8	3.0100	0.95895	0.3390		

variances	T	DF	Prob > T
unequal	3.1583	9.6	0.0104
equal	3.3918	16.0	0.0037

Step 1 $H_0 \mu_S - \mu_N = 0$ $H_A \mu_S - \mu_N > 0$

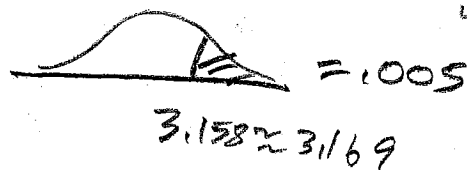
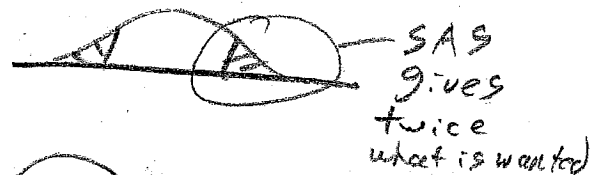
Step 2 $T_0 = 3.1583$

$pval = P(T > 3.1583) = 0.0104$

SAS $P(t_{150} > |t_0|)$

FOM 282 45

or check $df = 9.6 \approx 10$



Step 4) $pval < .05$ so reject H_0

Skilled rowers have higher knee velocities

ex Same but suppose we are told that the pop SD's of the 2 groups are approx equal.

Then use pooled 2 sample T .

Step 1 $H_0: \mu_S - \mu_N = 0$ $H_A: \mu_S - \mu_N > 0$

Step 2 $T_0 = 3.3918$

Ship

Step 3 $pval = \frac{.0037}{2} = .00185$

(from table C $.001 < pval < .0025$)

Step 4 reject H_0

skilled rowers have higher knee velocities

Ship P403, §7.3

Ch 8 CI's and Tests for proportions

1) Recall $\hat{p} = \frac{\text{count of successes}}{n}$

where $n =$ sample size

2) * p 437

A $100(1-\alpha)\% = c\%$ CI for p is

$\hat{p} \pm z^* \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$ where

$P(z^* \leq z \leq z^*) = 1-\alpha = c$ and

$P(z \leq z^*) = 1 - \frac{\alpha}{2}$

3) * p 437 test $H_0: p = p_0$

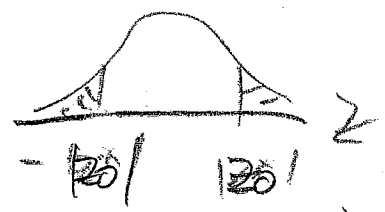
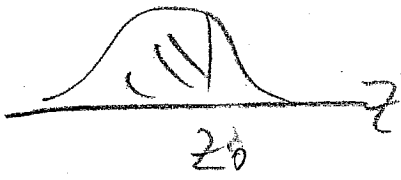
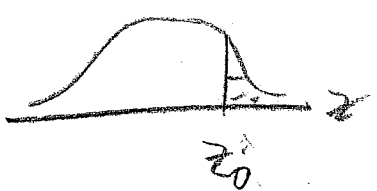
Statistic $z_0 = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \approx N(0,1)$ if H_0 is true.

P values

right tail
 $H_A: p > p_0$

left tail
 $H_A: p < p_0$

2 tail
 $H_A: p \neq p_0$



$P_{val} = P(z > z_0^*)$
 $1 - P(z < z_0)$

$P(z < z_0^*)$

$2 P(z < -|z_0|)$
 $= 2 P(z > |z_0|)$

4) * Assumptions

- a) Data are SRS from population.
- b) Pop is at least 10 times the sample size
- c) For test $n p_0 > 10$ and $n(1-p_0) > 10$

almost always true

For $2LI$ $n\hat{p} \geq 10$ and $n(1-\hat{p}) \geq 10$

(048)

5) p431-432 (similar to p278 $\hat{p} = \frac{x}{n}$)
 Sampling distribution of \hat{p}

If $n\hat{p} \geq 10$ $n(1-\hat{p}) \geq 10$ and \hat{p} is from a SRS of size $n = \frac{\text{pop size}}{10}$,

then $\hat{p} \approx N(\mu_{\hat{p}} = p, \sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}})$

ex) ^{Forward calculation for \hat{p}}
 Suppose a politician decides to run for office if 40% or more of a SRS of 2500 voters say they would vote for her. If the true pop percentage was 38%, what is the chance the politician will run?

Soln step 1 ~~_____~~ $\hat{p}, \mu_{\hat{p}}, \sigma_{\hat{p}}$
 .38 .4

Step 0) $\mu_{\hat{p}} = p = .38$ (always from pop never from sample)

$$\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}} = \sqrt{\frac{.38(1-.38)}{2500}} = \sqrt{.0000942}$$

$$= .00971 \text{ step 1) } \begin{array}{c} \text{_____} \\ .38 \quad .4 \end{array} \hat{p} \mu_{\hat{p}} = .38 \sigma_{\hat{p}} = .00971$$

Step 2) get z score

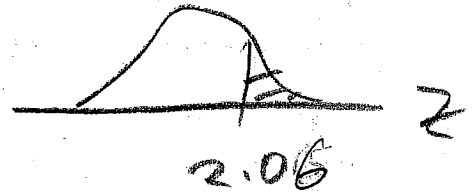
$$z = \frac{\hat{p} \text{ value} - \mu_{\hat{p}}}{\sigma_{\hat{p}}}$$

(triplet $\hat{p}, \mu_{\hat{p}}, \sigma_{\hat{p}}$)

$$\frac{140 - .38}{.00971} = 2.06$$

46.9

Step 3) z picture



Step 4) table A

$$1 - \frac{\Phi(2.06)}{2.06} = 1 - .9803 = .0197$$

That is there is about a 2% chance she will run.

on p 278 we had $X = \text{count}$

$$X \approx N(\mu_X = np, \sigma_X = \sqrt{np(1-p)})$$

Note $\hat{p} = \frac{X}{n}$ $\mu_{\hat{p}} = \frac{\mu_X}{n} = p$ $\sigma_{\hat{p}} = \frac{\sigma_X}{n} = \sqrt{\frac{p(1-p)}{n}}$

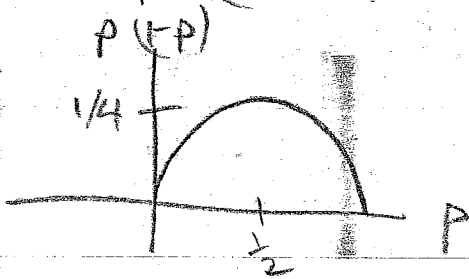
For CI $n\hat{p} \geq 10$ and $n(1-\hat{p}) \geq 10$ 0M4824

b) p44B Sample size for \hat{p} to be within m of p with $100(1-\alpha)\% = 100c\%$ confidence if

$$n = \left(\frac{z^*}{m}\right)^2 p^* (1-p^*) \quad \text{where}$$

p^* is a good guess for p and

$$P(Z < z^*) = 1 - \frac{\alpha}{2} \quad \text{or} \quad P(-z^* < Z < z^*) = 1 - \alpha =$$



If no good guess for p is available, use

$$p^* = \frac{1}{2} \quad \left(\text{and } n = \left(\frac{z^*}{2m}\right)^2\right)$$

Round up to make n an integer

(If p is far from 0.5, then using $p^* = \frac{1}{2}$ will result in using a much larger sample size than actually needed.)

ex) SRS of 1711 from records of people who died in bicycle crashes between 1987 and 1991. 386 of the 1711 had blood alcohol levels above 0.10% (were drunk).

a) Find a 95% CI for p .

47%

$$\hat{p} = \frac{386}{1711} = 0.2256$$

$$SE(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \sqrt{\frac{0.2256(1-0.2256)}{1711}} = 0.0101$$

$$z^* = 1.96$$

$$95\% \text{ CI} = \hat{p} \pm z^* SE(\hat{p})$$

$$= 0.2256 \pm 1.96 (0.0101)$$

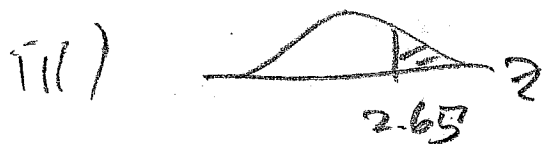
$$= 0.2256 \pm 0.0198 = (0.2058, 0.2454)$$

b) Test whether p is greater than 0.20.
Same 4 step procedure

i) $H_0: p = 0.20$ $H_A: p > 0.20$

ii) $z_0^* = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{0.2256 - 0.2}{\sqrt{\frac{0.2(1-0.2)}{1711}}} = \frac{0.0256}{0.0097}$

$$= 2.65$$



$$= p_{\text{val}} = 1 - 0.9960 = 0.004$$

iv) Reject H_0 since $p_{\text{val}} = 0.004 \leq \alpha = 0.05$.

The proportion of fatally injured cyclists that were drunk > 0.20 .

c) Test whether p is less than .25, ^{U. 100}

i) $H_0: p = .25$ vs $H_A: p < .25$

ii) $z_0 = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{.2256 - .25}{\sqrt{\frac{.25(.75)}{1711}}} = \frac{-0.0244}{.0105}$

$= -2.33$

iii) $p_{\text{value}} = \text{Area to the left of } z = -2.33 = .0099$

iv) Reject H_0 since $p_{\text{value}} \leq .05$

The proportion of fatally injured cyclists that were drunk < 0.25 .

(Big sample so $\hat{p} = .2256$ is fairly accurate)

D) Find n to estimate p to within 0.01 with 90% confidence if no good guess for p is available.

$n = \left(\frac{z^*}{E} \right)^2 \frac{1}{4} = \left(\frac{1.645}{.01} \right)^2 \frac{1}{4} = \frac{164.5^2}{4}$

$= \frac{27060.25}{4} = 6765.06$ Use $n = 6766$

E) Find n if $p^* = .2$

$n = \left(\frac{1.645}{.01} \right)^2 .2(1-.2) = 4329.64$

Use $n = 4330$

7) 2 proportions P_{447} SR & POP P_i n_i \hat{P}_i (48)

independent $\begin{cases} 1 & P_1 & n_1 & \hat{P}_1 \\ 2 & P_2 & n_2 & \hat{P}_2 \end{cases}$

8) P_{449} 100% $(1-\alpha)\%$ CI for $P_1 - P_2$ is

$$(\hat{P}_1 - \hat{P}_2) \pm z^* SE(\hat{P}_1 - \hat{P}_2) =$$

$$(\hat{P}_1 - \hat{P}_2) \pm z^* \sqrt{\frac{\hat{P}_1(1-\hat{P}_1)}{n_1} + \frac{\hat{P}_2(1-\hat{P}_2)}{n_2}}$$

where $P(Z < z^*) = 1 - \frac{\alpha}{2}$ or $P(-z^* < Z < z^*) = c = 1 - \alpha$

9) P_{453} test $H_0: P_1 = P_2$

Let $\hat{p} =$ count of "successes" in both samples

= pooled sample proportion, $\frac{n_1 + n_2}{n_1 + n_2}$

Test statistic $Z = \frac{\hat{P}_1 - \hat{P}_2}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \approx N(0,1)$ under H_0

P values

right tail

$H_A: P_1 > P_2$



$P(Z > z^*) = 1 - P(Z \leq z^*)$

left tail

$H_A: P_1 < P_2$



$P(Z < z^*)$

2 tail

$H_A: P_1 \neq P_2$



$2 P(Z > |z^*|) = 2 P(Z < -|z^*|)$

10) ^{PM 453} ASSUMPTIONS For 2 sample proportions procedures DOM 492 49

a) 2 independent SRS's
 or trt group vs ^{or another trt} control group created with random chance mechanism

b) $n_1 \hat{p}_1 \geq 5, n_1(1-\hat{p}_1) \geq 5$
 $n_2 \hat{p}_2 \geq 5, n_2(1-\hat{p}_2) \geq 5$

c) $\left. \begin{array}{l} \text{pop 1 size} \geq 10 \quad n_1 \\ \text{pop 2 size} \geq 10 \quad n_2 \end{array} \right\}$ almost always true

ex	100 men	31	# who say women are safest driver
	100 women	67	

$$\hat{p} = \frac{31 + 67}{200} = 0.49$$

(Source WWW.usatoday.com snapshots who is the safest driver)

A) Test i) $H_0: p_1 = p_2$ vs $H_A: p_1 < p_2$

$$ii) z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} = \frac{.31 - .67}{\sqrt{.49(.51)\left(\frac{1}{100} + \frac{1}{100}\right)}}$$

$$= \frac{-0.36}{\sqrt{0.0499}} = \frac{-0.36}{0.0707} = -5.09$$

iii)  $p\text{-val} = 0.0$

iv) reject H_0 there is strong evidence that the

Proportion of men who say that women are safer drivers than men is smaller than the proportion of women who say so. (49.3)

B) Find a 90% CI for $p_1 - p_2$

$$z^* = 1.645$$

$$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

$$= \sqrt{\frac{(0.31)(0.69)}{100} + \frac{(0.67)(0.33)}{100}}$$

$$= \sqrt{0.002139 + 0.002211} = \sqrt{0.00435}$$

$$= 0.06595 \quad 90\% \text{ CI is } (\hat{p}_1 - \hat{p}_2) \pm z^* SE$$

$$= (0.31 - 0.67) \pm 1.645(0.06595) = -0.36 \pm 0.108$$

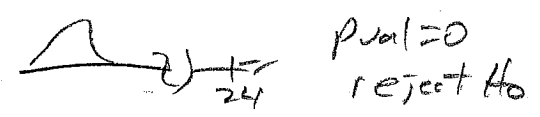
$$= (-0.468, -0.252)$$

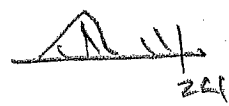
ch 9 For hypotheses tests so far, ^{not on review} we get a statistic \bar{x} , \hat{p} , $\bar{x}_1 - \bar{x}_2$, or $\hat{p}_1 - \hat{p}_2$ and find the appropriate z value (if $df > 30$) we reject H_0 if " z " is far from zero in the direction of H_A . If " z " is small or far away from zero in the direction away from H_A we fail to reject H_0 .

ex) $\mu = \text{mean IQ of } 1282 \text{ students}$

$H_0 \mu = 60 \quad H_A \mu > 60 \quad \bar{x} = 120, s = 15$

$z \approx \frac{120 - 60}{15/\sqrt{36}} = 24.0$
mean IQ is higher than 60



ex) $H_A \mu < 60$  $pval = 1$

no evidence that mean IQ of 282 students is less than 60

In this ex, H_0 is absurd, but H_A is more absurd.
pval = 1 means \bar{x} could not have occurred if H_A is true.
Ho

1) p483 Suppose there are 2 categorical variables and we want to test
 H_0 there is no relationship between 2 categorical variables
 H_A there is a relationship

2) p472 usually the 2 cat. variables can be displayed in a (two way) $r \times c$ table.
The row variable has r categories and the column variable has c categories.

3) p472 Each entry of the table is a cell. There are rc cells.

4) p473 A row marginal total is the sum of the row can
A column column

The table total is the sum of all of the counts = n

ex 1

50.9

r	gender		row tot	3x2 table 6 cells
	MF	MF		
civil eng	10	20	30	
nuclear eng	30	40	70	
industrial eng	40	60	100	
col tot	80	120	200 = N	

s) p 483 The test of H_0 vs H_A can be used for 2 situations

A) There are r independent SRS's from r populations. Each person is classified into one category of column categorical variable.

B) There is a single SRS of individuals. Each individual is classified according to both of two categorical variables.

ex In the last ex, take a SRS of 30 civil eng's 70 nuc eng's 100 ind eng's and classify each individual by gender. This gives situation A.

For situation B, take a SRS of 200 civil nuclear and ind engineers. Classify each individual by type of engineer and by gender.

6) In situation A, r SRS's, OM282 51
 suppose the categorical variable has only two values success, failure so $C=2$

eg ~~male~~ female "count females"

eg cured not cured "count cured"

Then the sample proportion from each SRS can be obtained and plotted in a bar graph, consider

H_0 there is NO relationship between the two categorical variables

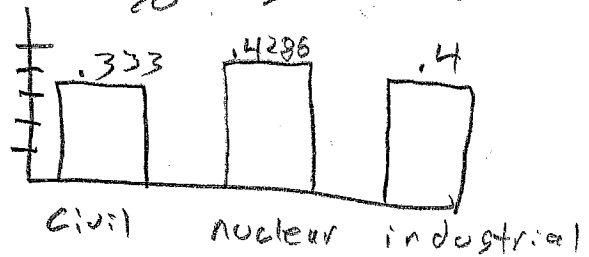
H_A there is a relationship

If H_0 is true then there is NO difference between the r SRS's ("treatments") and $p_1 = p_2 = \dots = p_r$. If H_0 is not true then there is a difference so not all of the p_1, p_2, \dots, p_r are equal.

ex) eng vs gender $\hat{p}_i = \text{prop of women from ith major}$

$$\hat{p}_1 = \frac{10}{30} = \frac{1}{3} \approx .333 \quad \hat{p}_2 = \frac{30}{70} = \frac{3}{7} \approx .4286$$

$$\hat{p}_3 = \frac{40}{100} = .4$$



Is there a difference in bar chart heights

or is the difference due to chance?

7) p 472 If H_0 is true, then

$$\text{the expected cell count} = \frac{(\text{row total})(\text{column total})}{\text{table total}}$$

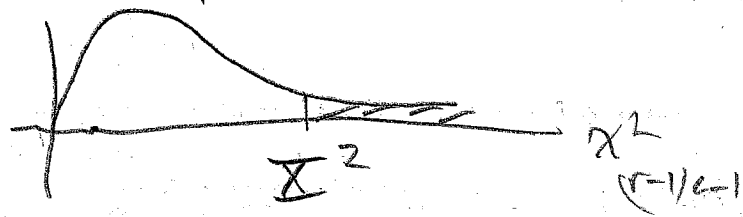
8) p476 The test statistic for H_0 is the chi-square statistic (5/5)

$$\chi^2 = \sum_{\text{sum over } r \times c \text{ cells}} \frac{(\text{observed count} - \text{expected count})^2}{\text{expected count}}$$

$r = \# \text{ rows}$
 $c = \# \text{ columns}$

Cell components of χ^2 p478

9) p480 If H_0 is true $\chi^2 \approx \chi^2_{(r-1)(c-1)}$ and the p value is from table E



$$p\text{val} = P(\chi^2_{(r-1)(c-1)} \geq \chi^2)$$

If $p\text{value} \leq \alpha$, reject H_0
 conclude there is a relationship.

Same 4 step test.

10) p485 Use the χ^2 test if no more than 20% of the expected counts are less than 5 and if all of the individual expected counts are 1 or greater. (Also see part 5.)

ex) A 2×2 table has 4 cells and $1.2(4) = 4.8 < 5$. So all 4 expected cell counts should be 5 or greater.

ex) eng

0 m 282 52

	Gender			proportions by major	
	F	M			
Fr	80	120	200	.4	.6
So	80	120	200	.4	.6
Jr	60	90	150	.4	.6
Sr	80	90	150	.4	.6

So knowing row proportions tells nothing about year, attend gender and year are unrelated.

ex)

	uses nail polish	does not use nail polish	row Prop's
Finger M	0	200	0.0 1.0
does not use F	120	80	1.0 .4

Now there is a relationship

Note

Expected count for a cell i, j in i th row j th column is $\frac{i\text{th row total} \cdot j\text{th column total}}{n}$

$$= n \left(\frac{i\text{th row total}}{n} \right) \left(\frac{j\text{th col tot}}{n} \right) \approx n (\text{prob in } i\text{th row}) (\text{prob in } j\text{th col})$$

$$= n p (\text{in } i\text{th row and } j\text{th column})$$

if col and row are independent (mult rule)

EX | p472-474

429

Drug	relapse		row total
	no	yes	
Doxip	14	10	24
Lith	6	18	24
Placebo	4	20	24
col tot	24	48	72

Test: i) H_0 no relationship between relapse & treatment
 H_A there is a relationship

ii) work

table of

obs

exp

χ^2 contrib ()

$$\frac{\text{row tot} \cdot \text{col tot}}{\text{table tot}} = \frac{(O-E)^2}{E}$$

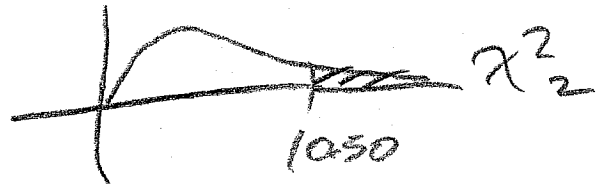
24	24	obs	14	10	row tot	24
$\frac{(14-8)^2}{8}$	exp	[8]	[16]	$\frac{(10-16)^2}{16}$	col tot	24
	cell χ^2	(4.5)	(2.25)		table tot	72
$\frac{(6-8)^2}{8}$		6	18	24		
		[8]	[16]	$\frac{(18-16)^2}{16}$		
		(0.5)	(0.25)			
$\frac{(4-8)^2}{8}$		4	20	24		
		[8]	[16]	$\frac{24(48)}{72} = 16$		
		(2.0)	(1.0)	$\frac{(20-16)^2}{16}$		
col tot		24	48	72		

$\chi^2 = \text{sum of } \chi^2 \text{ cell contributions}$
 $= 4.5 + 2.25$
 $+ 0.5 + 0.25$
 $+ 2.0 + 1.0 = 10.50$

$$(iii) df = (r-1)(c-1) = (3-1)(2-1) = 2 \quad \text{OM282 93}$$

table F

df	.001	.005
2	9.21	10.60



$$.005 < p\text{val} < .01$$

(iv) reject H_0 there is a relationship between treatment and relapse

ex Minitab output makes test much easier
I am likely to give table with χ^2 cell contributions added, but I will leave some expected counts and χ^2 cell contributions blank.

P477

MCY line

$$\text{chisq} = \dots$$

$$= 10.5$$

$$df = 2 \quad p\text{val} = .005$$

i) H_0 there is no relationship between trt and relapse
 H_a a relationship

ii) $\chi^2 = 10.5$

iii) $p\text{val} = .005$

iv) reject H_0 there is a relationship between trt and relapse

		days per week exercising				
		0-1	2-3	4-5	6-7	row tot
M	Obs	40	53	26	6	125
	exp	33.636	59.0	28.636	7.727	
	cell chisq	1.204	.073	.243	.386	
F	Obs	34	68	37	11	150
	exp	40.363	66.0	34.364	9.273	
	cell chisq	1.003	.061	.202	.322	
col tot		74	121	63	17	275

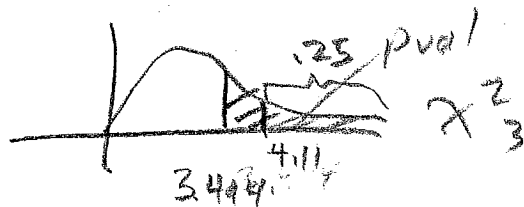
4 step test $\chi^2 = 1.204 + .073 + .243 + .386 + 1.003 + .061 + .202 + .322 = 3.494$

i) H_0 there is no relationship between gender and # days exercising
 H_A " " "a" "

ii) $\chi^2 = 3.494$

iii) $df = (r-1)(c-1) = (2-1)(4-1) = 3$

df	table F
3	4.11 4.61



$.25 \leq pval \leq 1$

iv) Fail to reject H_0
 there is no relationship between gender and # of days exercising

ex SRS of 2453 fatally injured pedestrians ^{0.5%}
 blood alcohol conc

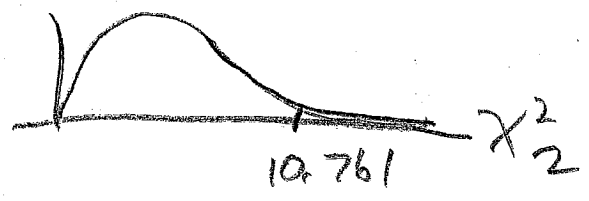
		0.0	0.01-0.09	0.10-	tot
age	16-34 obs	439	85	696	1220
	exp	473.477	91.015	655.508	
	cell chisq	2.510	0.398	2.501	
35-	obs	513	98	622	1233
	exp	478.523	91.985	662.442	
	cell chisq	2.484	0.393	2.475	
tot		952	183	1318	2453

$$\chi^2 = 2.510 + 0.398 + 2.501 + 2.484 + 0.393 + 2.475 = 10.761$$

- i) H_0 age at death and blood alcohol conc are not related
 H_A are related

ii) $\chi^2 = 10.761$

iii) $df = (2-1)(3-1) = 2$



df	.005	.0025
2	10.60	11.98

$.0025 < pval < .005$

- iv) reject H_0 age at death and blood alcohol conc are related.

ii) A chi square test can also be used for a GDS from a single population where each individual is classified according to both of two categorical variables.

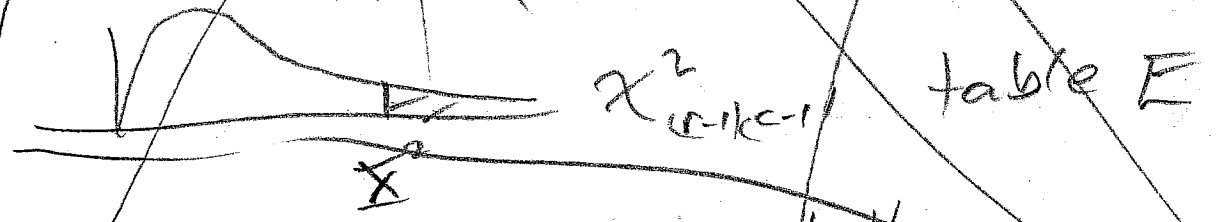
4 step test

- i) H_0 : no relationship between the 2 categorical variables
- ii) H_A : there is a relationship

iii)
$$\chi^2 = \sum \frac{(obs - exp)^2}{exp} \sim \chi^2_{(r-1)(c-1)}$$

sum over the rc cells.

iii)
$$P\text{value} = P(\chi^2_{(r-1)(c-1)} > \chi^2)$$



iv) If $p\text{val} \leq \alpha$, reject H_0
 there is a relationship.

12) p 540

EOM 282

55

Use χ^2 test if no more than 20% of the expected counts are less than 5 and if all expected counts are ≥ 1 or ≥ 2 .

ex) $\chi^2 = 1.2 < 1.2$ so all 4 expected counts are ≥ 5 or ≥ 2 .

$\chi^2 = 1.2 < 1.2$ so all

4 expected counts are ≥ 5 or ≥ 2 .

Ch 10 1) p 531-532 recall regression $\hat{y} = a + bx$

The population mean regression line is

$\mu_y = \alpha + \beta x = \text{mean of } y \text{ given that}$

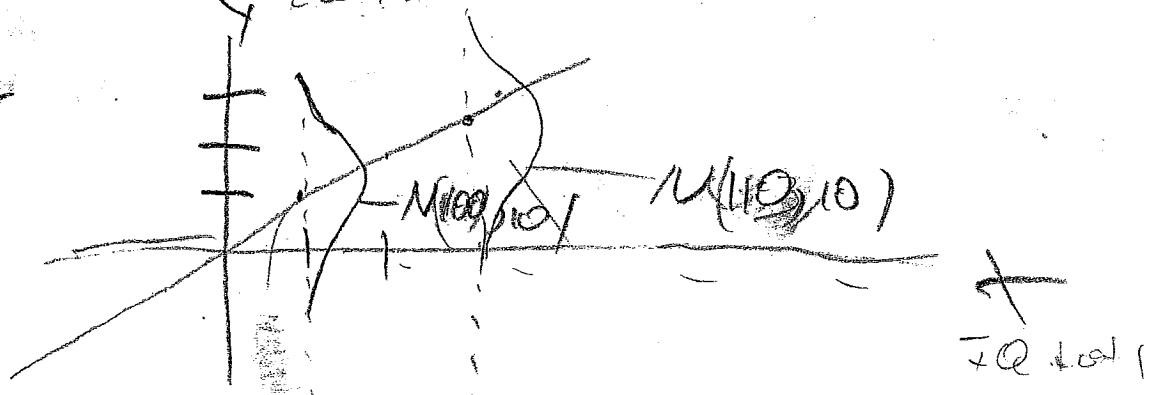
the value of the explanatory variable is x .

2) p 532 $y_i = \alpha + \beta x + \epsilon_i$

The ϵ_i ^{independent} $\sim N(0, \sigma)$. α, β, σ

are unknown, χ^2 test 2

3) p 532



Interpretation! For every value x the distribution of Y at x is $N(\mu_y, \sigma)$
 $= N(\alpha + \beta x, \sigma)$.

Tutorial

4) P533 $S = \sqrt{\frac{1}{n-2} \sum (y_i - \hat{y}_i)^2} = \sqrt{\frac{1}{n-2} \sum (\text{residual})^2} = \sqrt{\text{MSE}}$
 estimates σ

prog $\hat{y} = a + b x$ $b = r \frac{s_y}{s_x}$ $a = \bar{y} - b \bar{x}$

5) with residual $e_i = y_i - \hat{y}_i$
 = observed response - predicted response

6) The normality assumption of the ϵ_i is important. The e_i estimate the ϵ_i .

7) *P537 100% CI for slope β is $b \pm t^* SE_b$ where

$P(-t^* < t_{(n-2)} < t^*) = \alpha$

and SE_b would be given.

$df = n - 2$
 Get t^* from table α if $df =$
 use $t^* = z^*$ from table α if $df =$
 or use output

8) *P540 test $H_0: \beta = 0$
 test statistic $t_0^* = \frac{b}{SE_b} \approx t_{(n-2)}$

$df = n - 2$

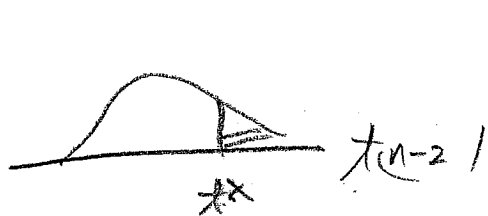
P values

right tail

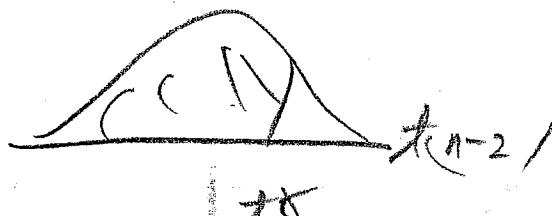
left tail

two tail

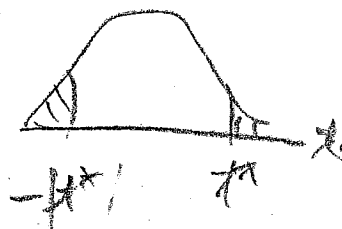
$H_A: \beta > 0$



$P(t_{n-2} > t^*)$



$P(t_{n-2} \leq t^*)$



$2 P(t_{n-2} \geq |t^*|)$

9) For inference we will always use output, same 4 step test.

10) p 544 100% CI for μ_y when $x = x^*$

is $\hat{y} \pm t^* SE_{\hat{\mu}_y}$

where $\hat{y} = a + b x^*$, $SE_{\hat{\mu}_y}$ is given in output

and $P(-t^* < t_{n-2} < t^*) = C$

11) p 544 A prediction interval is for a new observation y_{new} rather than a mean. Analogy 68-95-99.7 rule.

$y_i - \hat{y}_i = e_i \approx \epsilon_i \approx N(0, \sigma)$

so there is about a 95% chance

that Y_{new} is in $(\mu_y - 2\sigma, \mu_y + 2\sigma)$ §6.9

if Y_{new} is from a $N(\mu_y, \sigma)$ pop.

12) p1544 A $100(1-\alpha)\% = 100C\%$ prediction interval for a future observation Y_{new} at an explanatory value of $X = x^*$ is

$$\hat{y} \pm t^* SE_{\hat{y}} \text{ where}$$

$$P(-t^* \leq t_{(n-2)} \leq t^*) = 1-\alpha = C.$$

Note Y_{new} is a random variable, not a parameter.

13) p533 $SE_b = \frac{s}{\sqrt{\sum (x_i - \bar{x})^2}} \rightarrow 0 \text{ as } n \uparrow$

p544 $SE_{\hat{\mu}_y} = s \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum (x_i - \bar{x})^2}} \rightarrow 0 \text{ as } n \uparrow$

p545 $SE_{\hat{y}} = s \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum (x_i - \bar{x})^2}} \rightarrow 0 \text{ as } n \uparrow$

14) ex

	df	coef.	Standard error	T	p-value
intercept	1	1.1630	0.0066	177.3	0.0000
log skin	1	-0.0631	0.0041	-15.2	0.0000

Want to predict $Y =$ body density from

§6.9

$x = \log$ skinfold

$\neq 0M28259$

a 95% CI for β is

$$b \pm t^* SE_b =$$

$$df = 92 - 2 \\ = 90 > 30$$

$$= -0.0631 \pm 1.96 (.0041)$$

$$= -0.0631 \pm .0080 = (-0.071, -0.055)$$

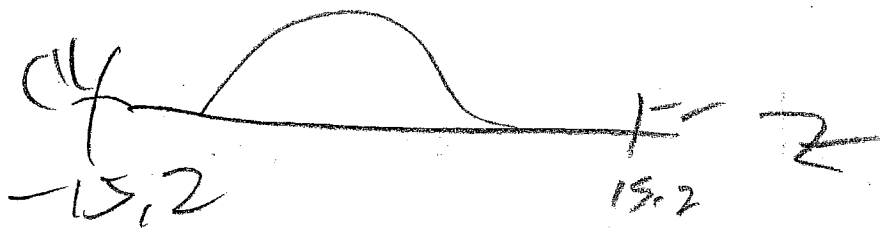
4 step test

1) $H_0: \beta = 0$ $H_A: \beta \neq 0$

2) $T = -15.2$

3) $P_{val} = 2P(t(90) > 15.2)$

$df = 90 > 30$
use table
A



$$= 2(0) = 0.0$$

4) or use output
Reject H_0 knowing the value
of \log skinfold helps predict
body density.

ex) Same example but $n=9$ 47.9

95% CI $df = n-2 = 7$ $t^* = 2.365$

$$b \pm t^* SE_b =$$

$$-0.0631 \pm 2.365 (0.0041) =$$

$$-0.0631 \pm .0097 = (-0.073, -0.053)$$

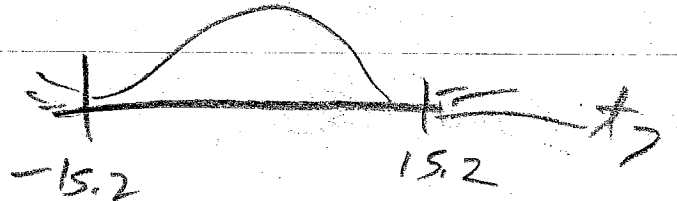
4 step test

s1) $H_0: \beta_1 = 0$ $H_A: \beta_1 \neq 0$ $\alpha = .005$

s2) $t_0^* = -15.2$

table C

s3) $df = 7$ $pval =$



≈ 0

$$pval < 2(.005) = .001$$

s4) reject H_0 knowing the value of log skinfold helps predict the body density,

ex p. 282

$$\hat{Y} = 1.0892 + 1.1890 X$$

OM 282 5

Find 95% CI for μ_Y and 95% PI

for Y_{new} if $X^* = 20$,

$$n = 16$$

$$SE_{\hat{\mu}_Y} = 0.0855$$

$$SE_{\hat{Y}} = 0.3496$$

Output

Fit	stdev fit	95% CI	95% PI
4.8692	0.0855	(4.6858, 5.0526)	(4.1193, 5.6191)

$$\hat{Y} = 1.0892 + 1.1890(20) = 4.8692$$

$$\text{stdev fit} = SE_{\hat{\mu}_Y}$$

So 95% CI for μ_Y is

$$\hat{Y} \pm t^* SE_{\hat{\mu}_Y} \quad df = n - 2 = 14$$

$$t^* = 2.145$$

$$4.8692 \pm 2.145 (0.0855) = 4.8692 \pm 0.183$$

$$= (4.6858, 5.0526) = 95\% \text{ CI for } \mu_Y$$

$$SE_{\hat{Y}} = 0.3496$$

$$4.8692 \pm 2.145 (0.3496)$$

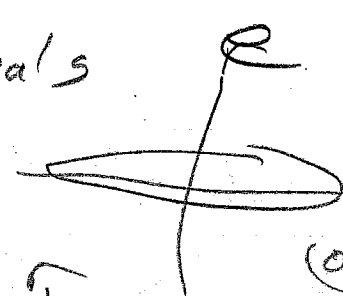
$$= 4.8692 \pm 0.7499 = (4.1193, 5.6191) = 95\% \text{ for}$$

(15)

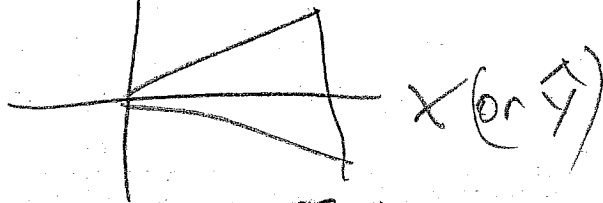
p546 i) scatter plot of y vs x should be 'football shaped' No curvature in a plot of e_i vs x

ii) $\epsilon_i \sim N(0, \sigma)$ σ does not depend on x

The plot of the residuals vs x should look like



bad



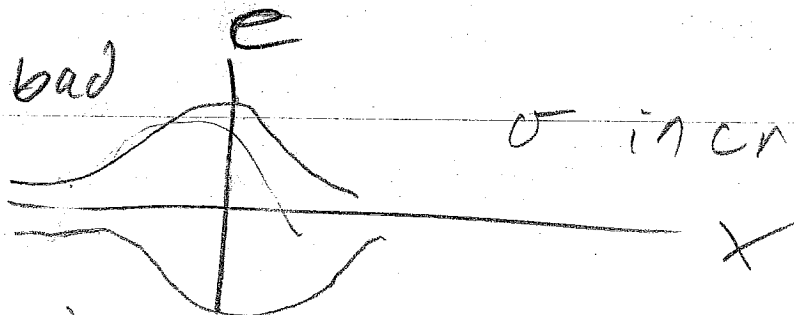
bad



σ increases with x

σ decreases with x

bad



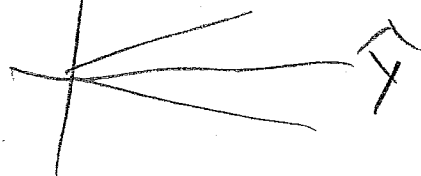
σ increases then decreases

iii) $\epsilon_i \sim N(0, \sigma)$ Normality

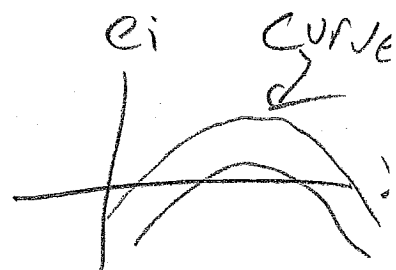
make a histogram of the residuals e_i ; check for outliers.

ex] Fig 11.8 p549

a)



b)



- 1) P166 In an experiment the investigators decide which individuals get a "treatment" and which do not, then observe a response.
- 2) P166 In an observational study (the individuals decide if they take the treatment) investigators simply observe a response.
- 3) Know the difference. Clue! the study is observational if the "treatment" eg smoking, illegal drugs, marriage divorce, etc is not something a person would do at a flip of a coin.

A study of a new medicine before FDA approval should always be an experiment, *Actual or* when the "individuals" are not people and the treatment is cheap, the study will often be an experiment.

A study of 2 established medicines could be observational or an experiment.

- ex Effect of smoking on cancer obs
- ex effect of Viagra on impotence exp
when the drug was being made
- ex study of people with 3 or more marriages obs
inter racial marriage obs
heavy alcohol consumption obs

ex) effect of irrigation fertilizer and herbicide on crop yields exp

ex) effects of smoking on rats exp

(49.3)

4) ^{P186} "The individuals" on which an experiment is done are called experimental units. They could be people, rats, farm's fields etc.

Each unit receives a treatment.

5) ^{P187} The purpose of an experiment is to reveal the response of one variable as explanatory variables change.

ex) sick people give medicine or sugar pill
response was person cured or not
explanatory medicine vs sugar pill

6) we compare the treatment with other treatments or a placebo (sham treatment ^{P189})

7) ^{P189, 193} The most important step in an experiment is randomization, ^{Use know} of subjects into 2 or more groups, each group gets a different ^{treatment or placebo}.
Use random numbers (impersonal chance) to decide which units get the treatment and which units get the placebo. ^{see point 17)}

ex] Salt vaccine for Polio (60.5)
2 studies NFIP ^{controlled} not randomized
and a double blinded randomized controlled
experiment (saline solution as placebo),
Think of factors that might make
contracting Polio more likely

treatment group control group

- income
- race
- gender
- hygiene
- hereditary
- # of siblings

For all the lurking variables you can think of, there are many more that you did not think of. Randomization causes the treatment and control groups to have approximately the same proportion of any lurking variable. Eg if 15% of the units are white female with brown hair, then about 15% of the treatment units and about 15% of the control units are too. However chance variation means you could get 14.7% or 15.3% say. As more units are used, the chance variation decreases.

The salt vaccine could only be ^{at 202-61} given to children if their parents gave consent

DBCRC Expt			NFIP				
	Size	rate per 100000		Size	rate		
Expt } Trt	200000	28	grade 2 tot	225000	25		
	control	200000	71	grade 1 & 3 control	725000	54	
NO consent			350000	46	grade 2 & NO consent	125000	44

not part of expt, but recorded anyway

The NFIP study was badly done,
 NO randomization
 Treatment group was all grade 2 consenters,
 The control group was all grade 1 & 3 students
 and contained many kids whose parents
 would not have given consent. Consenters
 are different than nonconsenters, Polio is
 one of the few diseases that affects the rich
 more than the poor, but a greater percentage
 of poor did not give consent, children in
 poorer families tend to get earlier milder
 cases of polio.

ii) ^{p203} A study is biased if it systematically
 favors certain outcomes. In the
 NFIP study the control group was less

likely to get polio than the treatment group because the treatment group only had 2nd grade consenters while the control group had all 1st and 3rd grades (would be consenters and nonconsenters), NEIP was biased against vaccine

(12) placebo effect Idea of treatment gives relief

25 days
54 days
not
looks as good as 29/71

ex) gel capsule of medicine for headache relief

If one group is given medicine and the other nothing, can't tell how much relief is due to idea of treatment and how much to the medicine (confounding)

To defeat confounding give one group medicine, one group placebo then compare.

ex) cirrhosis of liver portacaval shunt

SI studies 4 randomized controlled degree of enthusiasm

expt	marked	moderate	none
i) no controls	24	7	1
ii) controls not randomized	10	3	2
iii) randomized controlled	0	1	3

Some patients were too sick for surgery. The randomized controlled experiment divided patients into eligible and ineligible patients.

Half the eligible patients got the ^{DE 282} treatment
half did not, 62

In ii) doctors decided who got treatment
and tended not to give the treatment
to the sicker patients. There
study is biased in favor of the
treatment

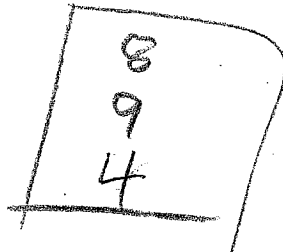
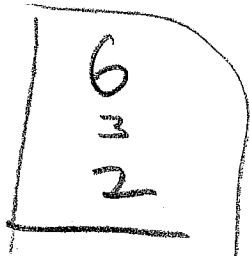
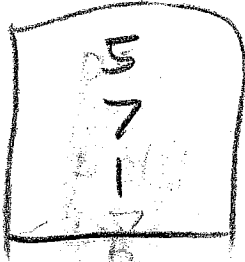
13) good experiments	bad
best double blinded ^{control} randomized controlled	not controlled
bad single blinded ^{control} randomized controlled	not randomized controlled,
good randomized controlled	

14) Often can't use placebo or blinding
ex surgery doctor and patient know
if patient got surgery, can't pull
out an appendix as a placebo
for liver surgery,

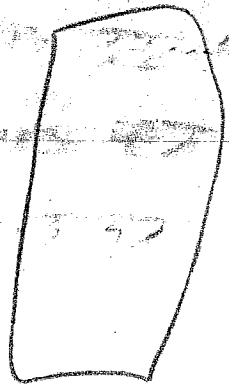
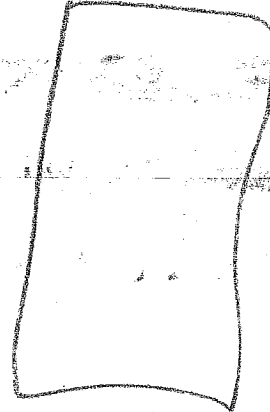
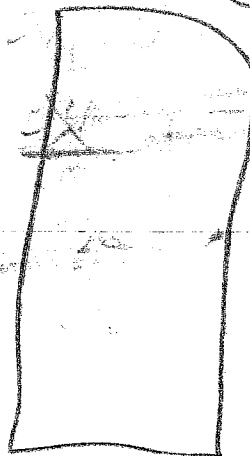
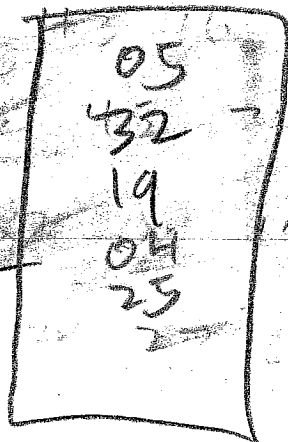
15) P206 There can be more than 2
treatments. When all experimental units
are allocated at random among all
treatments, the design is called
completely randomized

16) ex 3 treatments $\underbrace{19}_{1 \text{ digit}}$ units 625

line 131 05 00 71 66 32 | 81 19 41 48 73
~~04~~ 19 78 55 76 45 19 59 65 65
 13 73 25 52 59 84 29 20 87 96



ex 4 treatment $\underbrace{40}_{2 \text{ digits}}$ units



1st 10
1-10

next 10
11-20

next 10
21-30

last 10
31-40

line 131 05 00 71 66 32 81 19 41
 48 73 04 ~~19~~ 78 55 76 45 ~~19~~ 59 65 65

line 132 68 73 25 92 etc

17) Randomization minimizes the effects of lurking variables
 makes the groups as alike as possible except for
 different treatments, helps prevent biases of the experimenter
 from influencing the outcome of the expt, and causes \bar{X}_0 or \bar{P}_0
 to be approximately normal