

INTRODUCTION: STATISTICS IS THE

SCIENCE OF (1) OBTAINING DATA,

(2) SUMMARIZING DATA, AND

(3) USING DATA TO MAKE

IMPORTANT INFERENCE.

STATISTICAL/PROBABILISTIC

(1) OBTAINING DATA.

(a) DESIGNED EXPERIMENT: "A RANDOMIZED CONTROLLED" EXPERIMENT THAT IS SPECIFICALLY DESIGNED TO EXAMINE A SPECIFIC THING(S)

Exs: SALK VACCINE; BLOCK STUDIES (METHODS, COLA)

(b) OBSERVATIONAL STUDY: OBSERVATION OF RELATIONSHIPS WITH SUBSEQUENT ANALYSIS TO ADJUST FOR CONFOUNDING FACTORS.

Exs: EPA; left-handedness; human smoking

(c) SAMPLING: UTILIZING GOOD RANDOM SELECTION FROM POPULATION

Exs. Literary Digest; ESPN call-in-poll; Gallup poll.

(d) PUBLISHED REPORTS: WHICH ARE HOPEFULLY EITHER GOOD STUDIES USING (a), (b), OR (c), ABOVE, OR ARE A COMPLETE (CENSUS) OF USEFUL INFORMATION.

(2) SUMMARIZATION OF DATA.

- BASIC DESCRIPTIVE STATISTICS ^(D.S.) ON EACH VARIABLE
THESE D.S. INCLUDE MEASURES OF CENTER,
LIKE AVERAGE, AND MEDIAN;
MEASURES OF SPREAD, LIKE STANDARD DEVIATION
AND INTERQUARTILE RANGE
- PICTURES/GRAPHIC ON EACH VARIABLE,
LIKE HISTOGRAM, STEM/LEAF CHART, BOX PLOT, ETC.
- SUMMARY STATISTICS LOOKING AT RELATIONSHIP
BETWEEN VARIABLES, LIKE RELATIVE %'s.
- SUMMARY PICTURES/GRAPHICS
COMPARING VARIABLES, LIKE ^{VAR. VAR} X, Y PLOTS,
ETC.

(3) USE SCIENCE OF PROBABILITY AND RANDOM VARIABLES TO MAKE INFERENCES LIKE CONFIDENCE INTERVALS AND TESTS OF HYPOTHESES ABOUT POPULATION PARAMETERS.

DESCRIPTIVE STATISTICS

QUALITATIVE RESPONSE: (YES/NO); (ELECTRIC, GAS, OIL)

e.g. DID Child get Polio $\begin{cases} \text{Yes} \\ \text{No} \end{cases}$

HOUSE'S MAIN HEATING SOURCE $\begin{cases} \text{ELECTRIC} \\ \text{GAS} \\ \text{OIL} \end{cases}$

USE BAR CHARTS OR PIE CHARTS TO
SHOW PROPORTIONS (p's) OR $p \cdot 100\% = \%$'s.

QUANTITATIVE RESPONSE: HEIGHT, WEIGHT, SCORE, ETC.

• USE HISTOGRAM OR STEM/LEAF CHART
TO SHOW DISTRIBUTION OF VALUES FOR A
VARIABLE

• USE BOXPLOT TO VISUALLY SUMMARIZE
KEY SUMMARIES & LOOK FOR OUTLYING OBSERVATION
FOR A
VARIABLE

• USE PLOTS TO LOOK AT RELATIONSHIP
BETWEEN VARIABLES

(P): PROBABILITY : CHANCE ON A SCALE FROM [0.0 to 0.5 to 1.0]

↑ can't happen ↑ Equally likely to happen vs not happen ↑ will happen

Ex 1: A fair coin is tossed once.

$$P(\text{get a head}) = P(\text{don't get head}) \text{ so } P(\text{get head}) = 0.5 = \frac{1}{2}$$

[In very many repeats of Ex 1, about $\frac{1}{2}$ of time you'd get head] $1 - \frac{1}{2} = \frac{1}{2}$ time you'd NOT get a head

Ex 2. Roll A FAIR DIE ONCE.

$$\underbrace{P(\text{get a } \boxed{2})}_{P(\boxed{2})} = \frac{1}{6} = 1 - P(\text{don't get a } \boxed{2}) = 1 - \frac{5}{6} = \frac{1}{6}$$

GENERAL RULE 1: For ANY EVENT A, $P(A) = 1 - P(\text{not } A)$
same as: $P(\text{not } A) = 1 - P(A)$

In Ex 2, $P(\text{get an even number}) = \frac{3}{6}$, specific things that can occur

because there are a total of 6 (possible outcomes) of the roll, namely, $\{1, 2, 3, 4, 5, 6\}$

Each of these are equally likely, and.

For 3 of the 6, an even number occurs.

GENERAL RULE 2: IF ALL THE POSSIBLE OUTCOMES ARE EQUALLY LIKELY, THEN FOR ANY EVENT A,

$$P(A) = \frac{\# \text{ possible outcomes so that } A \text{ occurs}}{\text{total } \# \text{ possible outcomes}}$$

Box: R1 R2 B1 B2

Ex 3 A box has 2 red cards: R1, R2 and 2 blue cards: B1, B2. A

A hand of 2 cards are drawn (without replacement) from the box.

NOTE there are $4 \cdot 3 = 12$ possible outcomes of expt: namely,

$(R1, R2), (R2, R1), (R1, B1), (R1, B2), (B1, R1), (B2, R1), (R2, B1), (B1, R2), (R2, B2), (B2, R2), (B1, B2), (B2, B1)$

so, $P(\text{get red on 2nd draw}) = \frac{6}{12}$; $P(\text{red on 1^{st}} and red on 2nd}) = \frac{2}{12}$

NOW, IF YOU KNEW A RED WAS DRAWN ON 1st DRAW, THEN THERE IS

ONLY A 1 in 3 (or $\frac{1}{3}$ or $\frac{2}{6}$) chance that a RED IS DRAWN ON 2nd,

so, $\left\{ P(\text{RED ON 2<sup>ndst, this is a conditional probability,
 given that got</sup>$

A CONDITIONAL PROBABILITY IS ALSO A PROBABILITY, ITS JUST UNDER THE PARTICULAR CONDITION,

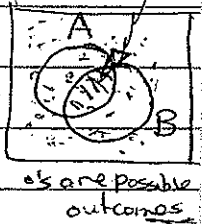
In last ex., e.g. $P(\text{red on 2nd} | \text{red on 1st}) = 1 - P(\text{NOT red on 2nd} | \text{red on 1st}) = 1 - \frac{2}{3} = \frac{1}{3}$.

Continuing last ex.,

we see $\frac{2}{12} = P(\text{red on 1st and red on 2nd}) = P(\text{red on 1st}) \cdot P(\text{red on 2nd} | \text{red on 1st}) = \frac{6}{12} \cdot \frac{1}{3} = \frac{2}{12}$

CALVIN
HOBBS

A and B



GENERAL RULE 3: FOR ANY 2 EVENTS A and B,
 $P(A \text{ and } B) = P(A) \cdot P(B|A)$

BOX: R1 R2 G1 B1

NEW Ex. 5 CONSIDER A BOX WITH 2 REDS: R1, R2; 1 GREEN: G1 and 1 Blue: B1

ONE DRAW IS MADE FROM BOX,

$P(\text{get GREEN OR get BLUE ON DRAW}) = \frac{2}{4} = P(\text{G1 or B1}) = P(\text{G1}) + P(\text{B1}) = \frac{1}{4} + \frac{1}{4} = \frac{2}{4}$

GENERAL RULE 4: IF 2 events D and E are mutually exclusive (one event occurring excludes other event from occurring), then $P(D \text{ or } E) = P(D) + P(E)$.



•'s are possible outcomes.

NOTE: GENERAL RULE 4 EASILY EXTENDS TO MORE

THAN 2 EVENTS, e.g. IF D1, D2, D3, ARE MUTUALLY EXCLUSIVE

THEN $P(D1 \text{ OR } D2 \text{ OR } D3) = P(D1) + P(D2) + P(D3)$

ALSO, GENERAL RULE 3 EASILY EXTENDS TO MORE THAN 2 EVENTS, e.g.

FOR ANY 3 EVENTS A1, A2, A3; $P(A1 \text{ and } A2 \text{ and } A3) = P(A1) \cdot P(A2|A1) \cdot P(A3|A1 \text{ and } A2)$

Ex. 6

Standard Deck of 52 cards (13 Hearts, 13 Spades, 13 Clubs, 13 Diamonds, & for each of these Suits have 2, 3, ..., 10, J, Q, K, A)

A hand of 3 cards is dealt (without repl.)

What's $P(\text{get exactly 2 Q's in hand})$

$$\begin{aligned} & \begin{array}{l} \text{Q and Q and Not Q} \\ \text{OR} \\ \text{Q Not Q Q} \\ \text{OR} \\ \text{Not Q Q Q} \end{array} \\ & \begin{array}{l} \frac{4}{52} \cdot \frac{3}{51} \cdot \frac{48}{50} \\ + \\ \frac{4}{52} \cdot \frac{48}{51} \cdot \frac{3}{50} \\ + \\ \frac{48}{52} \cdot \frac{4}{51} \cdot \frac{3}{50} \end{array} \leftarrow \begin{array}{l} \text{Rule 3} \\ \text{extended} \\ \text{Rule 4} \\ \text{extended} \end{array} \end{aligned}$$

Ex. 6 b

NOTE: If instead had dealt

3 cards WITH replacement.

$$\begin{aligned} \text{from deck then } P(\text{exactly 2 Q's}) &= \left\{ \begin{array}{l} \frac{4}{52} \cdot \frac{4}{52} \cdot \frac{48}{52} \\ + \\ \frac{4}{52} \cdot \frac{48}{52} \cdot \frac{4}{52} \\ + \\ \frac{48}{52} \cdot \frac{4}{52} \cdot \frac{4}{52} \end{array} \right\} = \left[3 \cdot \left(\frac{4}{52} \right)^2 \left(\frac{48}{52} \right) \right] \\ & \text{in 3 independent cards} \quad = \left[\frac{(3)}{(2)} \left(\frac{4}{52} \right)^2 \left(\frac{48}{52} \right) \right] \end{aligned}$$

In general two trials are independent, if no matter what you get on the 1st trial, ^{it} has no influence on probs for 2nd trial, so under WITH replacement, the 3 dealt cards are independent of one another. (under WITHOUT replacement, they are NOT indpt.)

* NOTE: If 2 trials are indpt & E_1 from 1st trial & E_2 from 2nd trial, then $P(E_1 \text{ and } E_2) = P(E_1) \cdot P(E_2 | E_1) \stackrel{\text{indpt}}{=} P(E_1)P(E_2)$

GENERAL RULE 5

for independent trials

"Binomial Rule"

$$P(\text{particular "thing" occurs on exactly } k \text{ out of } n \text{ (independent) trials}) = \binom{n}{k} \cdot [p^k (1-p)^{(n-k)}]$$

where p is chance the thing occurs on any given trial and

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} = \# \text{ of ways to choose the } k \text{ trials out of the } n \text{ total trials for which the thing occurs.}$$

NOTE: $n! = n(n-1)(n-2) \dots (2)(1)$

NOTE: our (ex. 6 b) above, is an illustration of General Rule 5, with $n=3$; "thing" is Q; $k=2$ and $p=4/52$

NOTE 2: If roll a fair die 10 (indpt) times, $P(\text{get } [2] \text{ on exactly } 4 \text{ of } 10 \text{ rolls}) = \binom{10}{4} \left[\left(\frac{1}{6} \right)^4 \left(\frac{5}{6} \right)^6 \right]$

Discrete Random Variables

Ex. 7 A fair coin is flipped two (indpt) times.

#

Let the random variable X be the number of heads in the 2 flips

Let's find the "probability distribution" for X .

Possible Outcomes of 2 Tosses			X		Also labelled $f(x)$
Prob.	1 st and 2 nd				
$\frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$	(T, T)	0	so, $P(X=0) = \frac{1}{4}$		
$\frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$	(T, H)	1	$\xrightarrow{\text{GR. 4}} P(X=1) = \frac{1}{4} + \frac{1}{4} = \frac{2}{4}$		
$\frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$	(H, T)	1			
$\frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$	(H, H)	2			$P(X=2) = \frac{1}{4}$

GR ↑ 3

GR. 3

Here, NOTE: $P(X \geq 1) = P(X=1) + P(X=2) = \frac{2}{4} + \frac{1}{4} = \frac{3}{4}$.

GENERAL RULE 7: FOR ANY DISCRETE R.V. X AND ANY INTERVAL " I " OF REAL LINE

$$P(X \text{ is in } I) = \sum_{\substack{\text{all possible} \\ X \text{ in } I}} P(X=x) \text{ also NOTED as } \sum_{\substack{\text{all possible} \\ X \text{ in } I}} f(x)$$

GOING BACK TO EX. 7, NOTE that the long-run average value that X would take is:

$$0\left(\frac{1}{4}\right) + 1\left(\frac{2}{4}\right) + 2\left(\frac{1}{4}\right), \text{ here equals } 0 + \frac{2}{4} + \frac{2}{4} = 1, \text{ since probs are long-run fractions}$$

The long-run average value for X is also labelled as the mean of X , and also called the expectation for X and also called expected value for X and written as $E(X)$; also written as μ

General Rule 8: For any discrete r.v. X , $E(X) = \sum_{\substack{\text{all possible} \\ x}} x f(x) = \mu$

Note: $E(X)$ is a measure of the center of the distribution of X .

The Variance (X) $\equiv \text{Var}(X) = E((X-\mu)^2)$ = long-run average of the squared distance from μ . $\equiv \sigma^2$

It is a measure of the spread of the distribution of X .

Standard Deviation (X) $\equiv \text{SD}(X) = \sqrt{\text{Var}(X)} \equiv \sigma$

GOING BACK TO EX 7 (# of heads in 2 TOSSES: X).

Recall $\mu=1$, so

$$\sigma^2 = \text{Var}(X) = E((X-\mu)^2) = E((X-1)^2) = \sum_{\text{all poss. } x} (x-1)^2 P(X=x) = \sum_{\text{all poss. } x} (x-1)^2 f(x)$$

$$= (0-1)^2(\frac{1}{4}) + (1-1)^2(\frac{2}{4}) + (2-1)^2(\frac{1}{4}) = \frac{1}{4} + 0 + \frac{1}{4} = \frac{1}{2}$$

$$\sigma = \text{SD}(X) = \sqrt{\text{Var}(X)} = \sqrt{\frac{1}{2}} \approx .71$$

roughly speaking the $\text{SD}(X)$ is "typical" distance X is away from μ in long-run.

Ex. 8 Say Tom spends \$1 on a lottery ticket, he has a $\frac{1}{1000}$ chance of winning \$100 ^{jackpot}, otherwise he wins nothing. ^(no jackpot win) What's $E(\text{Tom's net gain})$, What's $\text{StdDev}(\text{Tom's net gain})$? Tom's net gain is a r.v., call it X .

Note: $P(X = -1 + 100) = \frac{1}{1000}$ while $P(X = -1) = \frac{999}{1000}$

$$\text{so } E(X) = 99 \left(\frac{1}{1000} \right) + (-1) \left(\frac{999}{1000} \right) = -0.9$$

$$\text{while } \text{SD}(X) = \sqrt{(99 - (-0.9))^2 \left(\frac{1}{1000} \right) + (-1 - (-0.9))^2 \left(\frac{999}{1000} \right)}$$

One important "type" of Discrete r.v. is Binomial (n, p).

It is as follows:

Let X be the # of times a particular "thing" occurs in n independent trials, where there is p chance it will occur on given trial.

by GENERAL RULE 5:

$$\text{for } k=0,1,2,\dots,n, \quad P(X=k) = \binom{n}{k} \cdot [p^k (1-p)^{(n-k)}]$$

Ex. 9

An illustration of this is: \rightarrow so indpt draws.

Make 3 draws WITH replacement from box that has 6 reds and 4 non-red tickets in it. 6 R's 4 N's

let r.v. X be # of reds drawn, then $X \sim \text{binomial}(n=3, p=\frac{6}{6+4})$
with for $x=0,1,2,3$, $P(X=x) = \binom{n}{x} p^x (1-p)^{n-x} = \binom{3}{x} \left(\frac{6}{10}\right)^x \left(\frac{4}{10}\right)^{(3-x)}$

NEW DISTRIBUTION.

Ex. 10 Now go back to Ex. 9, but have the 3 draws instead be drawn WITHOUT replacement from the: 6 R's 4 N's
some obs

Let r.v. X be # of reds drawn.

$$\text{then for } x=0,1,2,3, P(X=x) = \frac{\binom{6}{x} \binom{4}{3-x}}{\binom{10}{3}} = f(x)$$

$X \sim \text{hypergeometric}(3, 6, 4)$

"

"SELF-SIMILAR"

The # OF OCCURRENCES IN AN INTERVAL OF TIME
OFTEN \approx HAS A POISSON(λ) DISTRIBUTION,
WHERE λ IS THE MEAN VALUE.

Ex. 11 # OF TRAFFIC ACCIDENTS ON SILAS CREEK PARKWAY
FOR A MONTH'S RUSH HOURS. $\sim \text{Poisson}(\lambda)$

$$\text{for } x=0,1,2,3,\dots, P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!} = f(x)$$

Ex. 12 : # OF PEOPLE SERVED @ LOCAL MED'S FROM 2-5 PM $\sim \text{Poisson}$

CONTINUOUS RANDOM VARIABLES.

Continuous R.V's take values on a "continuum" (e.g. on an interval)

Ex. 12

(a) For ex. ^{adult} height of randomly selected ^{male} might take any value in the interval $(48.0", 86.5")$;

(b) randomly chosen ^(position) point of a yardstick could be any value in the interval $(0.0", 36.0")$;

(c) yield on a cotton plot might be any value in $(20.0, 90.0 \frac{\text{bu}}{\text{ac}})$.

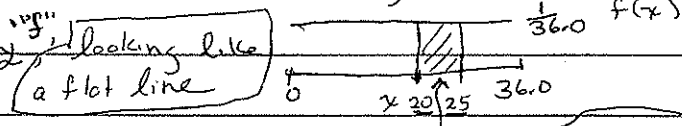
(d) amount of time until the next traffic accident occurs.

UNIFORM

One Major "Family" of Continuous R.V's is UNIFORM (a, b)

The ex¹² (b), above, randomly chosen ^(position) point on yardstick $\equiv X$.

Might have a Uniform $(0.0", 36.0")$ distribution, with relative likeli-hood, density, ^{"f"} looking like



To get $P(X \text{ is in } (20, 25)) = \text{Area under } f(x) \text{ between 20 and 25}$ (shaded above)

which here = $\frac{5}{36}$

NOTE: $P(X \text{ is in } (20, 25)) = \int_{20}^{25} f(x) dx$

NOTE: for continuous rv, f is a "density" not a probability

GENERAL RULE 9

FOR ANY CONTINUOUS R.V. X , WITH DENSITY, f ,

$$P(X \text{ is in interval } A) = \int_{x \text{ in } A} f(x) dx = \text{Area under } f \text{ over region } A.$$

$$\mu = E(X) = \int_{-\infty}^{\infty} x f(x) dx$$

$$\sigma^2 = \text{Var}(X) = E((X - \mu)^2) = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

$$\text{StdDev}(X) = \sqrt{\text{Var}(X)} = \sigma$$

ANOTHER MAJOR "FAMILY" IS NORMAL (μ, σ)

NORMAL

For ex 12(a) for male ht. might be \sim Normal $(\mu = 70", \sigma = 5")$, call it X

with $P(X \text{ is in } (70, 75))$ is the shaded Area $= \int_{70}^{75} f(x) dx$

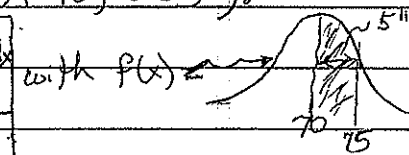
To approx. this area, 1st ^(standardize) convert the endpoints x 's

$$Z = \frac{X - \mu}{\sigma}$$

AND then use the standard normal table:

$$\text{here } P(X \text{ is in } (70, 75)) = P(Z \text{ is in } (\frac{70 - \mu}{\sigma}, \frac{75 - \mu}{\sigma})) = P(Z \text{ is in } (0, 1)) = .84 - .50 = .34$$

using (std) normal table



CONTINUOUS RANDOM VARIABLES (CONT.)

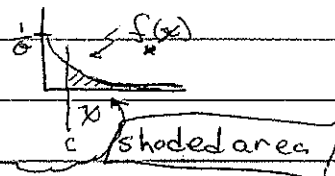
ANOTHER MAJOR FAMILY IS EXPONENTIAL(θ)

EXPONENTIAL For 12 d, time until next accident occurs, call it X
often has $\text{Exp}(\theta)$ distribution, where θ is mean of dist.

the density $f(x) = \frac{1}{\theta} e^{-x/\theta}$ for $x > 0$.

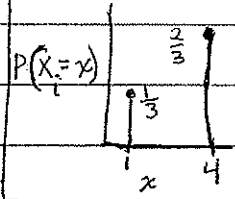
For any
constant
 $c > 0$

$$P(\text{time exceeded "c"}) = P(X > c) = \int_c^{\infty} \underbrace{\frac{1}{\theta} e^{-x/\theta}}_{f(x)} dx = e^{-c/\theta}$$



SAMPLING DISTRIBUTION

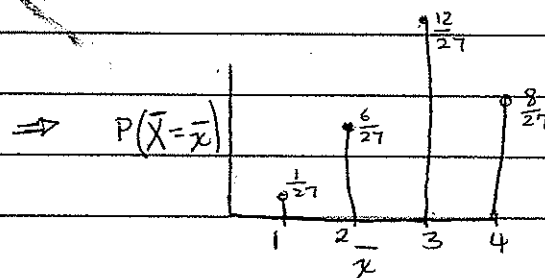
Ex. 13 CONSIDER A VERY SIMPLE DISCRETE DISTRIBUTION ("r.v. X "), where possible values for X are either "1" or "4" AND $P(X=1) = \frac{1}{3}$; $P(X=4) = \frac{2}{3}$. \Rightarrow



YOU CAN VIEW THIS AS A POPULATION OF VALUES, WHERE $\frac{1}{3}$ OF THE VALUES ARE "1" AND $\frac{2}{3}$ OF THE VALUES ARE "4" AND THINK OF X AS A RANDOM DRAW FROM POPULATION

NOW SUPPOSE THAT $n=3$ INDEPENDENT (WITH REPLACEMENT) RANDOM DRAWS (X_1, X_2, X_3) ARE MADE FROM POPULATION. Then Note that the possible VALUES for the (X_1, X_2, X_3) triples (samples), their probabilities, and "sample statistic \bar{X} " ARE:

prob.	possible (X_1, X_2, X_3)	Sample average $\bar{X} = \frac{1}{n} \sum X_i$	
$\frac{1}{27} = (\frac{1}{3})(\frac{1}{3})(\frac{1}{3})$	(1, 1, 1)	$\frac{1}{3}(1+1+1) = 1$	} NOTE: $P(\bar{X}=1) = \frac{1}{27}$
$\frac{2}{27} = \frac{2}{3} \cdot \frac{1}{3} \cdot \frac{1}{3}$	(4, 1, 1)	$\frac{1}{3}(4+1+1) = 2$	
$\frac{2}{27} = \frac{1}{3} \cdot \frac{2}{3} \cdot \frac{1}{3}$	(1, 4, 1)	$\frac{1}{3}(1+4+1) = 2$	
$\frac{2}{27} = \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{2}{3}$	(1, 1, 4)	$\frac{1}{3}(1+1+4) = 2$	} NOTE: $P(\bar{X}=2) = \frac{2}{27} + \frac{2}{27} + \frac{2}{27} = \frac{6}{27}$
$\frac{4}{27} = \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{1}{3}$	(4, 4, 1)	$\frac{1}{3}(4+4+1) = 3$	
$\frac{4}{27} = \frac{2}{3} \cdot \frac{1}{3} \cdot \frac{2}{3}$	(4, 1, 4)	$\frac{1}{3}(4+1+4) = 3$	
$\frac{4}{27} = \frac{1}{3} \cdot \frac{2}{3} \cdot \frac{2}{3}$	(1, 4, 4)	$\frac{1}{3}(1+4+4) = 3$	} NOTE: $P(\bar{X}=3) = \frac{4}{27} + \frac{4}{27} + \frac{4}{27} = \frac{12}{27}$
$\frac{8}{27} = \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{2}{3}$	(4, 4, 4)	$\frac{1}{3}(4+4+4) = 4$	
			} NOTE: $P(\bar{X}=4) = \frac{8}{27}$



Ex. 13 (continued).

NOTE: $E(\bar{X}) = \sum_{\text{all } \bar{x} \text{ possible}} \bar{x} P(\bar{X} = \bar{x}) = 1\left(\frac{1}{27}\right) + 2\left(\frac{6}{27}\right) + 3\left(\frac{12}{27}\right) + 4\left(\frac{8}{27}\right) = 3.00$

label as $\mu_{\bar{X}}$

Also NOTE that the mean μ of population:

$$\mu_{X_i} = E(X_i) = 1\left(\frac{1}{3}\right) + 4\left(\frac{2}{3}\right), \text{ which also } = 3.00$$

SINCE $E(\bar{X}) = \mu$, \bar{X} IS AN UNBIASED
ESTIMATOR OF μ .

NOTE: $\text{StdDev}(\bar{X}) = \sqrt{\text{Var}(\bar{X})} = \sqrt{E[(\bar{X} - \mu_{\bar{X}})^2]} \equiv \sigma_{\bar{X}}$

$$= \sqrt{\sum_{\bar{x}} (\bar{x} - 3)^2 P(\bar{X} = \bar{x})} = \frac{\sqrt{2}}{\sqrt{3}}$$

ALSO NOTE that the StdDev σ OF population

$$\sigma_{X_i} = \sqrt{\text{Var}(X_i)} = \sqrt{E[(X_i - \mu_{X_i})^2]} = \sqrt{2}$$

NOTE: here, $\sigma_{\bar{X}} = \frac{\sigma_{X_i}}{\sqrt{n}}$ AND $\mu_{\bar{X}} = \mu_{X_i}$

GENERAL RULE 10 FOR n indpt draws from popul with mean μ , SD σ ,
 $\bar{X} \equiv \left[\frac{1}{n} \sum X_i\right]$ ("sample average")

$$E(\bar{X}) = \mu \quad ; \quad \underbrace{\text{StdDev}(\bar{X})}_{\text{SD}} = \frac{\text{StdDev}(X_i)}{\sqrt{n}} = \frac{\sigma}{\sqrt{n}}$$

CENTRAL LIMIT THEOREM (CLT).

GOING BACK TO OUR LAST EX₆ (EX 13).

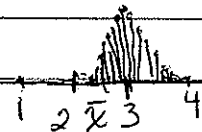
BUT NOW INSTEAD WE SAMPLED $n=100$

(SAME)
DRAWS FROM THE POPULATION

THEN (BY G.R. 10), $E(\bar{X}) = \mu = 3$ AND $SD(\bar{X}) = \frac{\sigma}{\sqrt{100}} = \frac{\sqrt{2}}{10} = .1414$

NOW, since n large, the distrib. of \bar{X} will be close to being symmetric and bell-shaped, like a normal curve and, for this ex., the distrib. of $\bar{X} \approx$

$P(\bar{X} = \bar{x})$



i.e. \bar{X} is Approx Distrib. Normal $(\mu = \mu = 3, \sigma = \frac{\sigma}{\sqrt{n}} = .1414)$

THIS IS AN ILLUSTRATION OF THE
CENTRAL LIMIT THEOREM.

GENERAL RULE 1)

CENTRAL LIMIT THM (CLT)

FOR LARGE n SAMPLE FROM POPULATION (μ, σ)

$\bar{X} \approx \text{normal}(\mu, \frac{\sigma}{\sqrt{n}})$, NO MATTER THE POPUL. DIST.

NOTE: WE CAN USE CLT, TO GET APPROX. PROBS RELATING TO \bar{X} WHEN LARGE n
in last ex, $n=100$; $\mu=3$; $\sigma=\sqrt{2} \approx 1.414$

$$P(\bar{X}_{100} \leq 3.1414) = P\left(\frac{\bar{X}_{100} - 3}{\frac{\sqrt{2}}{\sqrt{100}}} \leq \frac{3.1414 - 3}{.1414}\right) = P(Z \leq 1) = .84$$

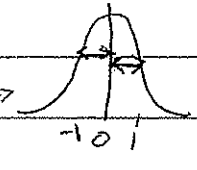
↑
std normal

from normal table

LARGE SAMPLE CONFIDENCE INTERVALS FOR μ

AGAIN, n large, X_1, X_2, \dots, X_n (indpt) drawn from population with mean μ & SD σ .

$$\bar{X}_n \approx \text{Normal}(\mu, \frac{\sigma}{\sqrt{n}})$$

$$\Rightarrow \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} \sim Z \sim \text{std normal}$$


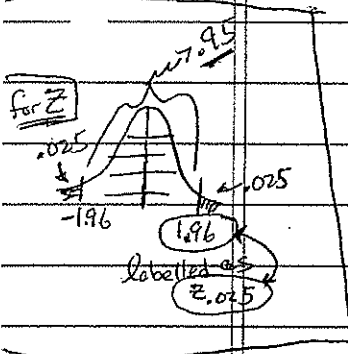
ALSO, if n large, then sample SD, $s = \sqrt{\frac{1}{n} \sum (x_i - \bar{x})^2}$ is "very close to σ ",

$$\text{So } \frac{\bar{X}_n - \mu}{\frac{s}{\sqrt{n}}} \sim Z$$

From Z table, $P(Z > 1.96) = 1 - P(Z \leq 1.96)$

$$= 1 - .975 = .025$$

$$\text{ALSO } P(Z < -1.96) = .025$$



$$\text{So, } P\left(-1.96 \leq \frac{\bar{X}_n - \mu}{\frac{s}{\sqrt{n}}} \leq 1.96\right) = P(-1.96 \leq Z \leq 1.96) = 1 - 2(.025) = .95$$

$$\Downarrow$$

$$P\left(-1.96 \frac{s}{\sqrt{n}} \leq \bar{X}_n - \mu \leq 1.96 \frac{s}{\sqrt{n}}\right) = .95$$

$$\Downarrow$$

$$P\left(\bar{X} - 1.96 \frac{s}{\sqrt{n}} \leq \mu \leq \bar{X} + 1.96 \frac{s}{\sqrt{n}}\right) = .95$$

$$\Downarrow$$

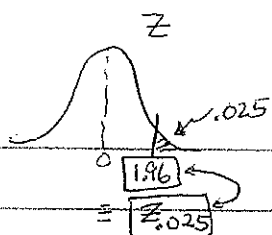
$$P(\mu \text{ is in interval } (\bar{X} - 1.96 \frac{s}{\sqrt{n}}, \bar{X} + 1.96 \frac{s}{\sqrt{n}})) = .95$$

from the sample (of size n), you get an observed sample avg \bar{x} and an observed $s = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2}$

and a $(.95) \times 100\% = 95\%$ confidence interval ^(C.I.) for μ is:

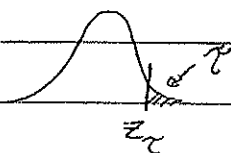
$$\left(\bar{x} - 1.96 \frac{s}{\sqrt{n}}, \bar{x} + 1.96 \frac{s}{\sqrt{n}}\right)$$

RECALL:



Since .025 of Area under Z curve is to right of 1.96
(NOTE: 1.96 is also the $1-.025$ quantile or 97.5th percentile of Z .)

IN GENERAL, z_{α} is value so that area under Z curve to the right of it is α .



GENERAL RULE 12. FOR LARGE n , A $(1-\alpha)*100\%$ C.I for μ is:

$$\left(\bar{X} - z_{\frac{\alpha}{2}} \frac{S}{\sqrt{n}}, \bar{X} + z_{\frac{\alpha}{2}} \frac{S}{\sqrt{n}} \right)$$

Also written as:

$$\bar{X} \pm z_{\frac{\alpha}{2}} \frac{S}{\sqrt{n}}$$

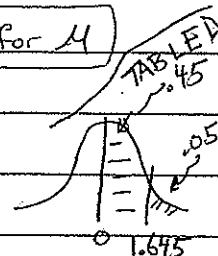
Check: A 95% CI for μ = $(1-\alpha)*100\%$ CI for $\mu \Rightarrow \alpha = .05$

$\Rightarrow \alpha/2 = .025$; To GET $z_{\alpha/2} = z_{.025} \Rightarrow 1.96$

\Rightarrow 95% CI for μ OF: $\bar{X} \pm 1.96 \frac{S}{\sqrt{n}}$, as before

Ex. 14 A sample of $n=400$ is taken from very large popul. the observed sample average is 23.2, while the observed SD is 8.0. Then, construct a 90% C.I. for μ

Since $90\% = (1-\alpha)*100\% \Rightarrow \alpha = .10 \Rightarrow \alpha/2 = .05$ AND
so $z_{\alpha/2} = z_{.05} = 1.645$



SHOW
DEMONSTRATIONS

so 95% CI for μ is: $23.2 \pm 1.645 \left(\frac{8.0}{\sqrt{400}} \right) \Rightarrow 23.2 \pm 0.65$

GENERAL RULE

$$\underline{X} \equiv (X_1, X_2, \dots, X_n)$$

(ANY)
GENERAL RULE FOR PARAMETER, SAY τ .

Let $u_1(\underline{X})$ and $u_2(\underline{X})$ be two functions of sample \underline{X} and constants so that:

$$P(u_1(\underline{X}) \leq \tau \leq u_2(\underline{X})) = 1 - \alpha,$$

THEN $(u_1(\underline{X}), u_2(\underline{X}))$ is a $(1 - \alpha) \times 100\%$ CI for τ .

Px. SAY X_1 is a draw from a $\text{Normal}(\mu, \sigma=3)$.

CONSTRUCT A 95% CI for μ , based only on X_1 & constants.

$$\text{NOTE: } \frac{X_1 - \mu}{3} \sim Z \text{ (std normal).}$$

$$\text{so, } P\left(-1.96 \leq \frac{X_1 - \mu}{3} \leq 1.96\right) = .95$$

$\underbrace{\qquad\qquad\qquad}_{Z_{1-.95/2}}$

$$\Rightarrow P\left(\underbrace{X_1 - 1.96(3)}_{u_2(\underline{X})} \leq \underbrace{\mu}_{\tau} \leq \underbrace{X_1 + 1.96(3)}_{u_1(\underline{X})}\right) = .95$$

so, $(X_1 - 1.96(3), X_1 + 1.96(3))$ is a 95% CI for μ

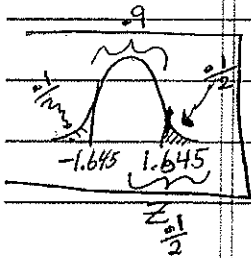
TEST OF HYPOTHESIS :

SAMPLE FROM SINGLE POPULATION.

N LARGE : TEST ON μ

FOR Ex. 14, WE HAD A SAMPLE OF $N=400$ WITH
SAMPLE AVG, $\bar{X}=23.2$ AND SAMPLE SD, $S=8.0$.

§ WE FORMED A 90% [i.e. $(1-\alpha)*100\%$] C.I FOR μ



$$\text{By } \bar{X} \pm Z_{\alpha/2} \cdot \frac{S}{\sqrt{n}} \Rightarrow \bar{X} \pm 1.645 \left(\frac{8.0}{\sqrt{400}} \right) \Rightarrow$$

$$\bar{X} \pm 0.65 \Rightarrow 23.2 \pm 0.65 \Rightarrow (22.55, 23.85)$$

so "90% confident" that μ is in interval: $(22.55, 23.85)$

Now, continuing with above data ($n=400, \bar{x}=23.2, s=8.0$)

So, if wanted to do a "two-sided" test of hypothesis

Null hypothesis: $H_0: \mu = 24.0$ vs.

Alternative hypothesis: $H_A: \mu \neq 24.0$ AT $\alpha = 0.1$ LEVEL,

WE CONSTRUCT THE TEST STATISTIC:
$$\frac{\bar{X} - 24.0}{\frac{S}{\sqrt{n}}} \sim Z_{H_0} \quad \mu = 24.0$$

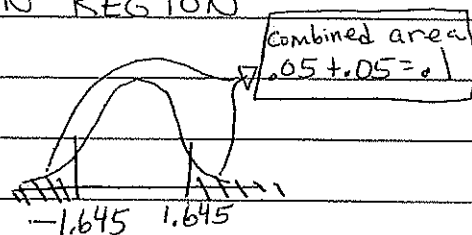
WE'LL REJECT H_0 IF OUR DATA'S VALUE OF TEST
STATISTIC IS "EXTREME" VS WHAT WOULD BE
SUGGESTED BY H_0 .

DATA'S VALUE OF TEST STATISTIC IS:

$$\frac{23.2 - 24.0}{\frac{8.0}{\sqrt{400}}} = \boxed{2.0}$$

Now, A (2-SIDED) $\alpha = .1$ REJECTION REGION

FOR A STD. NORMAL (Z) IS.



A VALUE OF TEST STATISTIC

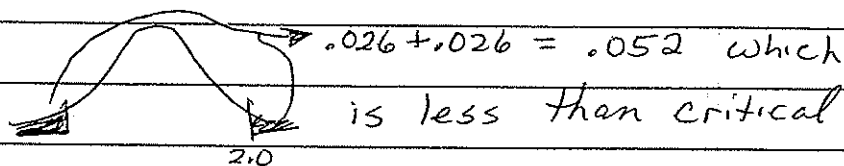
THAT'S GREATER THAN 1.645 OR LESS THAN -1.645,

HERE OUR VALUE $\boxed{2.0}$ IS IN THE REJECTION REG.,
SO WE REJECT H_0 AND CLAIM H_A ($\mu \neq 24.0$)

AN EQUIVALENT WAY TO CONDUCT THIS TEST

IS WITH P-VALUES (OBSERVED SIGNIFICANCE LEVEL).

HERE PVALUE =



is less than critical

set α level of .10, so reject H_0

DEFINITION:

PVALUE $\equiv P_{H_0}$ (GET A VALUE OF TEST STATISTIC THAT
IS AT LEAST AS SUPPORTIVE OF H_A
AS THE DATA'S VALUE OF TEST STATISTIC)

FOR THIS DATA,

(C.I.)

NOTE: THE CONFIDENCE INTERVAL/RESULT AND THE ^(2-sided) TEST OF HYPOTHESIS ^(T.H.) RESULT ARE COMPLEMENTARY, IN SENSE THAT THE 90% i.e. $(1-\alpha) \times 100\%$ C.I. FOR μ DOES NOT CONTAIN 24.0, WHILE THE $\alpha=0.1$ T.H. THAT $H_0: \mu=24.0$ IS REJECTED AND IT IS CLAIMED THAT $\mu \neq 24.0$

n LARGE: TEST ON POPULATION PROPORTION: P

RECALL \hat{p} (SAMPLE PROPORTION) IS LIKE A SAMPLE AVG OF 0's AND 1's; $SD(\hat{p}) = \frac{\sqrt{p(1-p)}}{\sqrt{n}}$ $E(\hat{p}) = p$, so. FOR

$$H_0: p = p_0$$

$$\frac{\hat{p} - p_0}{\frac{\sqrt{p_0(1-p_0)}}{\sqrt{n}}}$$

$\sim Z$
 $p=p_0$
 n large

so. e.g. to test $H_0: p=0.5$
($\alpha=0.1$ level)

$$H_A: p > 0.5$$

one-sided here

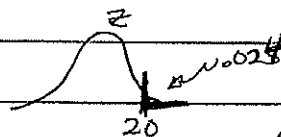
would reject H_0 if data value of test stat is 'sufficiently large'.

Ex. $n=100$ flips of coin by Tom; he gets #1 for each head; we get #1 for each tail.

In the 100 flips, 60 heads (i.e. thus 40 tails) were obtained.

DATA'S VALUE OF TEST STAT:

$$\frac{\hat{p} - p_0}{\frac{\sqrt{p_0(1-p_0)}}{\sqrt{n}}} = \frac{\frac{60}{100} - 0.5}{\frac{\sqrt{0.5(1-0.5)}}{\sqrt{100}}} = +2.0 \Rightarrow P\text{VALUE} =$$

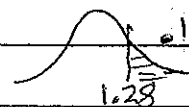


H_0 : "coin fair" ($p=0.5$ head)
 H_A : coin biased for heads
 $p > 0.5$

since PVALUE OF 0.024 $< 0.1 \Rightarrow$ REJECT H_0 , CLAIM H_A

ALTERNATIVELY $\alpha=0.1$ ^(one-sided) REJECTION REGION IS:

DATA VALUE OF TEST STAT > 1.28 ,
WHICH IT IS.

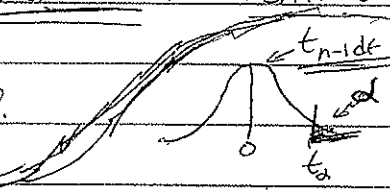


(TEST OF HYPOTHESIS)

n small T.O.H. ON μ UNDER NORMAL POPULATION:

RECALL HERE

$$\frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} \sim t \text{ WITH } n-1 \text{ d.f.}$$



So, e.g. to test
(at critical level)

$$H_0: \mu = \mu_0$$

$$H_A: \mu > \mu_0$$

THEN REJECT H_0 IFF $\frac{\bar{X} - \mu_0}{\frac{s}{\sqrt{n}}} > t_{\alpha, n-1}$

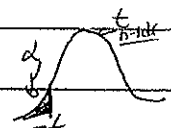
if instead test

$$H_0: \mu = \mu_0$$

$$H_A: \mu < \mu_0$$

THEN REJECT H_0 IFF

$$\frac{\bar{X} - \mu_0}{\frac{s}{\sqrt{n}}} < -t_{\alpha, n-1}$$



To test

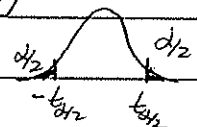
$$H_0: \mu = \mu_0$$

$$H_A: \mu \neq \mu_0$$

THEN REJECT H_0 IFF

$$\left(\frac{\bar{X} - \mu_0}{\frac{s}{\sqrt{n}}} > t_{\alpha/2} \right) \text{ OR } \left(\frac{\bar{X} - \mu_0}{\frac{s}{\sqrt{n}}} < -t_{\alpha/2} \right)$$

2-sided test



THE 2-sided T.H. result is complementary to the corresponding confidence interval (C.I) result of $\bar{x} \pm t_{\alpha/2, n-1} \frac{s}{\sqrt{n}}$ IN SENSE THAT IF μ_0 IS NOT IN C.I., THEN WOULD REJECT $H_0: \mu = \mu_0$

ACTUAL DATA EX. SAY THE "ALLOWED AVERAGE CONTAMINATION LEVEL IN A LAKE IS 30.0 $\frac{g}{m^3}$, SAY TOOK $n=16$ SAMPLES TAKEN

& GOT AVG RESULT OF 56.5, WITH SD OF RESULTS = 25.0, CAN ONE "PROVE"

AT $\alpha=0.05$ LEVEL, THAT LAKE CONTAMINATION LEVEL IS ABOVE ALLOWABLE (EXCEEDS)

$$\text{TEST } H_0: \mu = 30.0$$

$$H_A: \mu > 30.0$$

$$\text{Test Stat: } \frac{\bar{X} - 30.0}{\frac{s}{\sqrt{n}}} \sim t \text{ WITH } \frac{16-1}{15} \text{ d.f.}$$

NORMAL INDIV. LEVELS

WILL REJECT H_0 , CLAIM H_A & THUS "PROVE" (@ $\alpha=0.05$ LEVEL) EXCESS

IF $\left(\text{OBSERVED } \frac{\bar{X} - 30.0}{\frac{s}{\sqrt{n}}} > t_{\alpha, n-1} \right)$ here $\frac{\bar{X} - 30.0}{\frac{s}{\sqrt{n}}} = \frac{56.5 - 30.0}{\frac{25}{\sqrt{16}}} = 4.24 > t_{0.05, 15}$

⇒ PROVED EXCESS

TEST OF HYPOTHESIS ON σ^2 UNDER NORMAL POPULATION

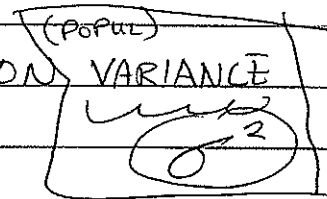
RECALL IF X_1, X_2, \dots, X_n RANDOM (INDPT) SAMPLE FROM

THEN $\boxed{\frac{(n-1)S^2}{\sigma^2}} \sim \chi^2_{\text{WITH } n-1 \text{ df.}} \text{ NORMAL}(\mu, \sigma).$

SO, FOR NORMAL POPULATION, CAN USE THIS TO PERFORM A TEST ON σ^2 (POP'L VARIANCE)

Ex. GO BACK TO OUR LAKE CONTAMINATION STUDY, WHERE "ASSUMED" THAT INDIV. CONTAMINATION LEVELS IN LAKE ARE NORMALLY DISTRIBUTED,

BUT NOW WANT TO DO A TEST ON VARIANCE OF INDIV. LEVELS IN LAKE (POPUL.)



RECALL SAMPLE SIZE: $n = 16$.

RECALL SAMPLE SD, $S = 25.0$, SO, SAMPLE VARIANCE $S^2 = 25^2 = 625$

$\alpha = 0.05$ LEVEL

SAY WANT TO DO TEST: $H_0: \sigma^2 = 500$ (i.e. $\sigma = \sqrt{500}$)

vs $H_A: \sigma^2 > 500$ (i.e. $\sigma > \sqrt{500}$)

THEN TEST STATISTIC IS: $\frac{(16-1)(S^2)}{500} \sim \chi^2_{\text{WITH } 15 \text{ df.}}$

H_0
NORMAL INDIV LEVELS

THEN REJECT H_0 & CLAIM H_A IFF

$$\frac{(16-1) \cdot 625}{500} > \chi^2_{0.05, 15 \text{ d.f.}}$$

INTERPRETATION:

n small ^{CAN} ~~Do~~ NOT ASSUME NORMAL POPULATION.

YOU CAN STILL DO A TEST ON
THE POPULATION MEDIAN, η , AS LONG
AS ^{ANY} CONTINUOUS POPULATION.

RECALL FOR ANY CONTINUOUS POPUL., $P(X_i < \eta) = 0.5 = P(X_i > \eta)$

FOR EX. LET'S LOOK AT SIMPLER VERSION OF "CONTAMINATED"

LAKE, WHERE ONLY $n = 7$ SAMPLES ARE TAKEN,

WE DO NOT ASSUME NORMAL (BUT ONLY CONTINUOUS)

DISTRIB. OF LEVELS IN LAKE; SAY NOW 'CAN TEST (@ $\alpha = 0.05$)
(LEVEL)

Popl MEDIAN

$$H_0: \eta = 30.0$$

$$H_A: \eta > 30.0$$

§ SUPPOSE OUR $n = 7$ LEVELS ARE:

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6 \quad x_7$
35.8, (24.2), 33.4, 45.8, 39.2, 36.1, 38.7

NOTE THAT 6 OF THE 7 VALUES EXCEED 30.0 (ALL BUT THE 24.2)

HERE,

TEST STATISTIC IS: # OF X_i 'S THAT EXCEED 30.0 \sim BINOMIAL ($n=7$, $p=0.5$)
 $H_0: \eta = 30.0$

$$P\text{-VALUE} = P(\text{Binomial}(n=7, p=0.5) \geq 6) = 1 - P(\text{BIN}(7, 0.5) \leq 5) < 0.05$$

SO, REJECT H_0 &
SO, CLAIM $\eta > 30.0$

OBSERVED VALUE
TEST STAT

CAN USE
BINOMIAL TABLE

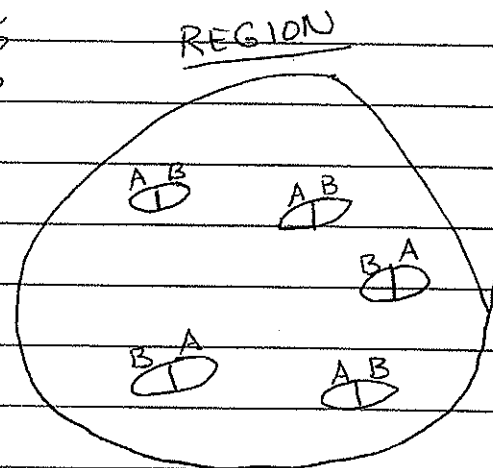
critical α

TWO POPULATIONS : PAIRED DATA

EXAMPLE: SUPPOSE $n=5$ PLOTS ARE RANDOMLY SELECTED FROM A REGION. EACH PLOT IS "SPLIT" IN TWO, WITH ONE HALF RANDOMLY SELECTED TO GET TRT A, WHILE THE OTHER HALF GETS TRT B. THEN @ END OF GROWING SEASON, THE CORN YIELD IS MEASURED ON EACH "HALF PLOT".

SUPPOSE RESULTING DATA IS:

(PAIR) PLOT	Trt A yield	Trt B yield	A-B yield difference
1	24	23	+1
2	32	33	-1
3	37	37	0
4	36	34	2
5	40	37	3



SO, IN ANALYSIS, VIEW THIS AS A SAMPLE OF $n=5$ DIFFERENCES OUT OF A "SINGLE" POPULATION OF ALL POSSIBLE PAIRED (A-B) DIFFERENCES.

SO, E.G. IF ASSUME POPL OF DIFFERENCES IS ABOUT NORMAL THEN CAN DO A t-test ON μ_{DIFF} (THE POPL MEAN DIFF)
~~(Treat A & B equally effective)~~

$$H_0: \mu_{DIFF} = 0$$

$$H_A: \mu_{DIFF} \neq 0$$

$$\text{Test stat: } \frac{\bar{X}_{diff} - 0}{S_{diff}/\sqrt{5}} \sim t \text{ with } 5-1=4 \text{ df.}$$

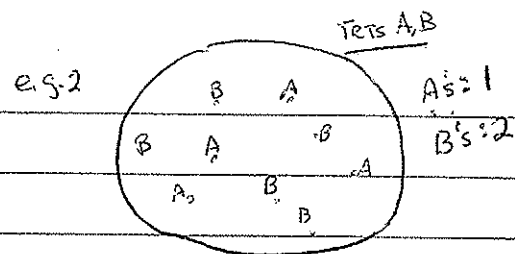
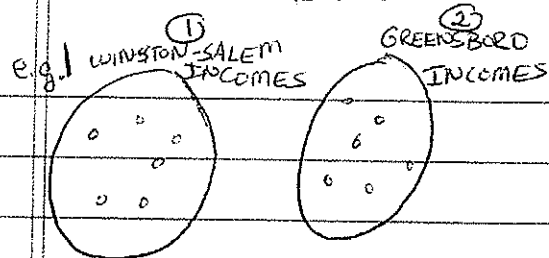
NOTE: $\bar{X}_{diff} = \frac{+1 + (-1) + 0 + 2 + 3}{5} = +1$

$$S = \sqrt{\frac{\sum (d_i - \bar{X}_{diff})^2}{n-1}} = 1.58$$

$\Rightarrow \text{obs test stat} = 1.41 \neq t_{0.05, 4 \text{ df}}$

so do NOT rej. H_0

TWO INDEPENDENTLY SAMPLED POPULATIONS.



TESTS ON DIFFERENCES IN POPULATION MEANS μ_1, μ_2 .

$$H_0: \mu_1 = \mu_2 \quad (\mu_1 - \mu_2 = 0)$$

LARGE SAMPLES SETTING:

IF n_1 AND n_2 LARGE, THEN BY CLT,

$$\frac{(\bar{X}_1 - \bar{X}_2) - E(\bar{X}_1 - \bar{X}_2)}{SD(\bar{X}_1 - \bar{X}_2)} \sim Z$$

$\sigma_{(\bar{X}_1 - \bar{X}_2)}$

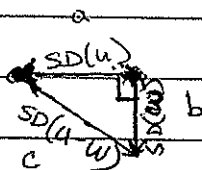
NOTE: $E(\bar{X}_1 - \bar{X}_2) = E(\bar{X}_1) - E(\bar{X}_2)$
 $= \mu_1 - \mu_2$

RECALL $SD(\bar{X}_1) = \frac{\sigma_1}{\sqrt{n_1}} = \frac{\sigma_1}{\sqrt{n_1}}$

GENERAL RULE: IF 2 THINGS U & W ARE INDPT,
 (PYTHAGOREAN THEOREM) THEN $SD(U - W) = \sqrt{[SD(U)]^2 + [SD(W)]^2}$

$$c^2 = a^2 + b^2$$

so $c = \sqrt{a^2 + b^2}$



so ABOVE BECOMES

$$\frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\left(\frac{\sigma_1}{\sqrt{n_1}}\right)^2 + \left(\frac{\sigma_2}{\sqrt{n_2}}\right)^2}} \sim Z$$

Recall n_i large that $s_i \approx \sigma_i$.

2 INDEPLY SAMPLED POPULATIONS ; n_i 's LARGE (CONTINUED)

so to test $H_0: (\mu_1 - \mu_2) = 0$

use Test statistic

$$\frac{(\bar{X}_1 - \bar{X}_2) - 0}{\sqrt{\left(\frac{S_1}{\sqrt{n_1}}\right)^2 + \left(\frac{S_2}{\sqrt{n_2}}\right)^2}} \quad \text{which is } \overset{H_0}{\sim} Z$$

For $H_A: (\mu_1 - \mu_2) \neq 0$ [i.e. for $\overset{H_A}{\mu_1 \neq \mu_2}$], Rej H_0 & Claim H_A

iff observed value of test statistic is $\begin{cases} > Z_{\alpha/2} \\ \text{or} \\ < -Z_{\alpha/2} \end{cases}$

NOTE: (FOR LARGE n_1, n_2), THE ABOVE TEST CAN BE "INVERTED" TO OBTAIN THE $(1-\alpha)*100\%$ C.I FOR $(\mu_1 - \mu_2)$

IT IS
$$(\bar{X}_1 - \bar{X}_2) \pm Z_{\alpha/2} \cdot \sqrt{\left(\frac{S_1}{\sqrt{n_1}}\right)^2 + \left(\frac{S_2}{\sqrt{n_2}}\right)^2}$$

IT IS COMPLIMENTARY TO THE ABOVE TEST IN SENSE THAT IF $[(\mu_1 - \mu_2) \text{ VALUE OF}] 0$ IS NOT IN THE ABOVE C.I, THEN WOULD REJECT H_0

2 INDEPLY SAMPLED POPULATIONS.
 n_1 AND n_2 SMALL

IF Assume both populations ARE NORMAL

& HAVE A COMMON VARIANCE

(i.e. $X_{11}, X_{12}, \dots, X_{1n_1} \sim \text{Normal}(\mu_1, \sigma^2)$ the same variance,
 $X_{21}, X_{22}, \dots, X_{2n_2} \sim \text{Normal}(\mu_2, \sigma^2)$ but unknown.

THEN to do A TEST ON $H_0: \mu_1 - \mu_2 = 0$

you first get a pooled estimate for σ^2 ,

namely
$$S_p^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1-1 + n_2-1} = \frac{\sum_{i=1}^2 \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2}{n_1-1 + n_2-1}$$

Now, test statistic looks like last one, except
 put in s_p for both s_1 & s_2

AND NOW (since n 's small, s_p is NOT great est of σ^2)

the test statistic has a t distribution (with n_1-1+n_2-1 df)

Specifically $\frac{\bar{X}_1 - \bar{X}_2 - (0)}{H_0} \sim t$ with n_1-1+n_2-1 d.f.

$$\sqrt{\left(\frac{S_p}{\sqrt{n_1}}\right)^2 + \left(\frac{S_p}{\sqrt{n_2}}\right)^2}$$
, so for $H_A: \mu_1 - \mu_2 \neq 0$, Rej H_0 iff
 test stat $\} > t_{\alpha/2}$
 value $\} \text{ or } < -t_{\alpha/2}$

corresponding $(1-\alpha) \times 100\%$ CI for $(\mu_1 - \mu_2)$ is:

$$(\bar{x}_1 - \bar{x}_2) \pm t_{\alpha/2}^{df=n_1-1+n_2-1} \cdot \sqrt{\left(\frac{S_p}{\sqrt{n_1}}\right)^2 + \left(\frac{S_p}{\sqrt{n_2}}\right)^2}$$

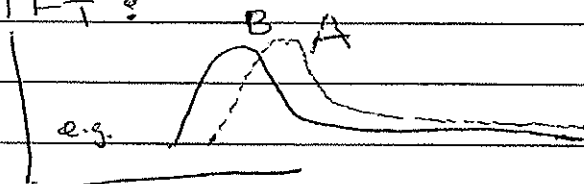
2 INDEPTLY SAMPLED POPULATIONS n_1 AND n_2 SMALL

CAN/DO NOT NECESSARILY ASSUME
NORMAL DISTRIBUTIONS,
ONLY ASSUME CONTINUOUS
DISTRIBUTIONS WITH POSSIBLE
LOCATION SHIFT:

To Test H_0 : A & B have same dist

H_A : A is shifted to right of B

(i.e. $P(X_{A_j} > X_{B_k}) > \frac{1}{2}$) ← more likely that obs from A
is larger than obs from B



Use "Sum-RANK Test" (Wilcoxon) - SUPPLEMENTAL Chpt 14.
DISK

Ex.
Emphasizing
Concept

SAY GOT $n_A = 5$ VALUES FROM A OF: 17.1, 22.3, 9.2, 20.5, 15.9

INDPTLY

GET $n_B = 4$ VALUES FROM B OF: 8.1, 7.4, 12.2, 6.9, 8.5

NOW IF RANK ALL $5+4=9$ VALUES (IN ORDER) GET:

6.9, 7.4, 8.1, 8.5, 9.2, 12.2, 15.9, 17.1, 20.5, 22.3
(B B B B A A A A A)

NOTE: under H_0 OF A & B having SAME dist, that the
5 positions for the A's in the 9 total "ordered" slots
are equally likely, there are $\binom{9}{5} = \frac{9!}{5!(9-5)!} = 126$
total (equally likely under H_0) sets of 5 positions for the A's

only 2 of these are at least as supportive of H_A
as our observed positions (namely
(B B B A A A A A A) and (B B B A A A A A A))
(B B B A A A A A A)
4 6 → 9

so P-value for our test is $\frac{2}{\binom{9}{5}} = \frac{2}{126} = \frac{1}{63} \approx 0.016$ so e.s.
would reject
 H_0 @ $\alpha = .05$ level

TWO INDEP POPULATIONS

ASSUME BOTH POPULATIONS NORMAL

TEST ON POPULATION VARIANCES

SPECIFICALLY $X_{11}, X_{12}, \dots, X_{1n_1} \overset{\text{indpt.}}{\sim} \text{Normal}(\mu_1, \sigma_1^2)$
 $X_{21}, X_{22}, \dots, X_{2n_2} \overset{\text{indpt.}}{\sim} \text{Normal}(\mu_2, \sigma_2^2)$

Recall we know: $\frac{(n_i - 1) S_i^2}{\sigma_i^2} \sim \chi^2$ with d.f. $n_i - 1$
 For either popl i

Now

$$\frac{\frac{(n_1 - 1) S_1^2}{\sigma_1^2}}{\frac{(n_2 - 1) S_2^2}{\sigma_2^2}} \overset{\text{OF FORM}}{=} \frac{\frac{\chi_{n_1-1}^2}{\text{df for } \chi_{n_1-1}^2}}{\frac{\chi_{n_2-1}^2}{\text{df for } \chi_{n_2-1}^2}} = F_{\substack{\text{num df} = \text{df for "1"} \\ \text{den df} = \text{df for "2"}}}$$

NOTE: THAT IF $\sigma_1^2 = \sigma_2^2$ [i.e. if $\frac{\sigma_1^2}{\sigma_2^2} = 1$], then left-hand side *

simplifies to just $\frac{S_1^2}{S_2^2}$

so to test $H_0: \sigma_1^2 = \sigma_2^2$ (i.e. $\frac{\sigma_1^2}{\sigma_2^2} = 1$) vs $H_A: \sigma_1^2 > \sigma_2^2$ (i.e. $\frac{\sigma_1^2}{\sigma_2^2} > 1$)

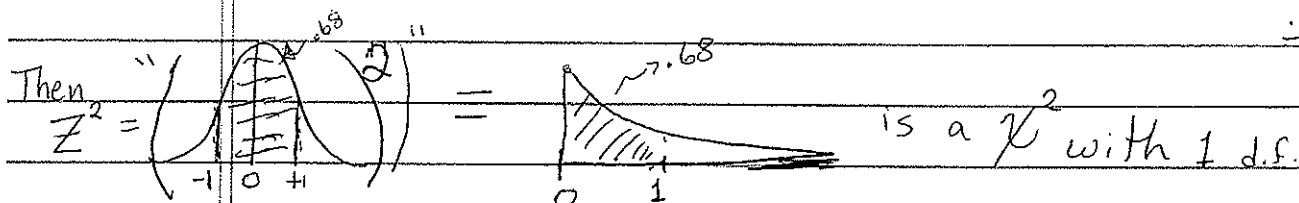
use test stat. of $\frac{S_1^2}{S_2^2} \sim F_{\substack{\text{num df} = n_1 - 1 \\ \text{den df} = n_2 - 1}}$ & rej H_0 iff $\boxed{\text{obs. } \frac{S_1^2}{S_2^2}} > F_{\alpha, \substack{\text{num df} = n_1 - 1 \\ \text{den df} = n_2 - 1}}$

[ex] sample of size 16 from popl 1 got $S_1^2 = 6.2$
 (indpt) ... 2-1 2 got $S_2^2 = 2.1$

Then $\text{obs. } \frac{S_1^2}{S_2^2} = \frac{6.2}{2.1} = 2.95$ while $F_{\alpha, \substack{\text{num df} = 15 \\ \text{den df} = 20}} =$, so

χ^2 : Chi-Square DISTRIBUTIONS

Start with a standard normal (" Z ") distrib.



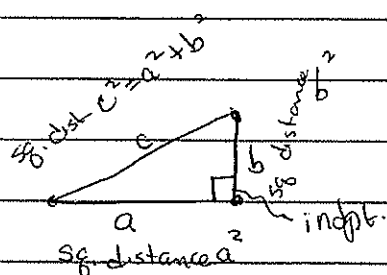
A χ^2 distrib. can also be thought of as a "random" "distance"

GENERAL RESULT.

The sum of independent χ^2 's is a χ^2

whose d.f. is the sum of the individual χ^2 's d.f.'s.

This is an "extension of Pythagorean Thm"



For $X_1, X_2, \dots, X_n \sim \text{Normal}(\mu, \sigma)$

NOTE: "t" distrib. with $n-1$ df

$$\frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} = \frac{\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}}{\frac{s/\sigma}{\sqrt{n}}} = \frac{\frac{Z}{\sigma/\sqrt{n}}}{\frac{s/\sigma}{\sqrt{n}}}$$

$$\frac{Z}{\sqrt{\frac{\chi^2(n-1)}{n-1}}} \leftarrow \frac{\frac{Z}{\sigma/\sqrt{n}}}{\sqrt{\frac{(n-1)s^2/\sigma^2}{n-1}}}$$

If $X_i \sim \text{Normal}(\mu, \sigma)$: Then Note: $\sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2 = \sum_{i=1}^n \left(\frac{Z_i}{\sigma/\sqrt{n}} \right)^2 = \text{sum of } n \text{ indep } \chi^2(1\text{'s}) \Rightarrow \chi^2(n \text{ d.f.})$

Now, $\sum_{i=1}^n \left(\frac{X_i - \bar{X}}{\sigma} \right)^2 = \sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2 - \left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right)^2 = \chi^2(n) - \chi^2(1) = \chi^2(n-1 \text{ d.f.})$

$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1) \text{ d.f.}$

MANY PROBABILISTIC MODELS ARE OF FORM:

$$Y_i = \underbrace{E(Y_i)}_{\text{deterministic part}} + \underbrace{\epsilon_i}_{\text{random part}}$$

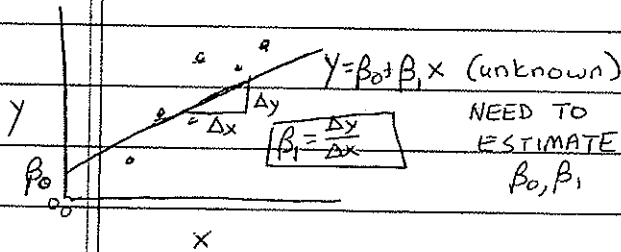
For person $i=1, 2, \dots, n$

e.g. Simple Linear Regression of Weight of i (Y_i) WITH HEIGHT of i (X_i)

is: $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$ where $E(\epsilon_i) = 0$
 $\text{Var}(\epsilon_i) = \sigma^2$
 ϵ_i are normal } $\epsilon_i \sim \text{normal}(0, \sigma)$
 ϵ_i are normal
 AND ϵ_i 's are indpt.
 simple linear deterministic part

here, β_0, β_1 and σ^2 are (unknown) parameters

Observe $n (X_i, Y_i)$ pairs:



$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

choose $\hat{\beta}_0, \hat{\beta}_1$

so as to

$$\text{minimize } \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

choose $\hat{\beta}_0, \hat{\beta}_1$

$$\text{to minimize } \sum_{i=1}^n (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i))^2 = f(\hat{\beta}_0, \hat{\beta}_1)$$

Take $\frac{\partial f}{\partial \beta_0}$ set $= 0$
 Take $\frac{\partial f}{\partial \beta_1}$ set $= 0$
 to get $\hat{\beta}_0, \hat{\beta}_1$

Estimate σ^2 by: $\frac{\sum (Y_i - \hat{Y}_i)^2}{n-1}$
 now label this as S^2

Now, test if X has a signif. cont (linear) "effect" on Y by:

$$H_0: \beta_1 = 0 \text{ (no effect of } x \text{ on } y)$$

$$H_1: \beta_1 \neq 0$$

$$\text{Test statistic } = \frac{\hat{\beta}_1 - 0}{\frac{S}{\sqrt{n}}}$$

$\sim t$ with $n-2$ d.f.

Rej H_0 if data value test stat

$$n-2 \text{ d.f.}$$

$$> t_{\alpha/2} \text{ or } < -t_{\alpha/2}$$

SAMPLE CORRELATION (r) BETWEEN X and y (here height) (here weight)

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{[\sum (x_i - \bar{x})^2][\sum (y_i - \bar{y})^2]}}$$

OR in book's notation

$$= \frac{SS_{xy}}{\sqrt{SS_{xx} SS_{yy}}}$$

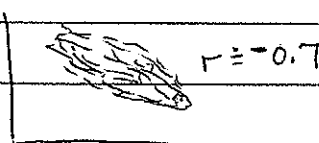
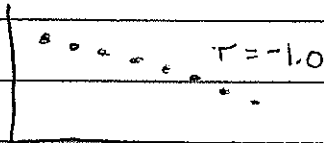
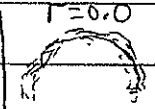
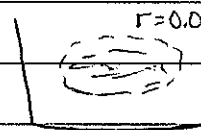
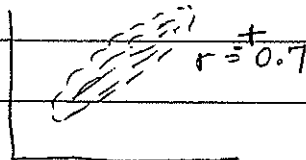
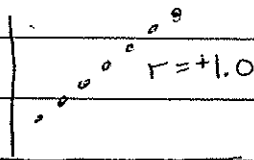
r can range from -1.0 to 0.0 to +1.0

perfect
negative linear
(x and y) association
(for sample)

no net
positive nor
negative
linear assoc.

perfect
positive linear
association

Some pictures:



long-run avg value for y_i for given value of x_i

CONFIDENCE INTERVALS FOR $E(y_i)$ FOR A GIVEN VALUE OF x_i

$$\hat{\beta}_0 + \hat{\beta}_1 x_i$$

$$\frac{\hat{y}_i - E(y_i)}{\hat{\sigma}_{\hat{y}_i}}$$

$$\hat{\beta}_0 + \hat{\beta}_1 x_i$$

FOR GIVEN x_i

$\sim t$ with $n-2$ d.f.

$\Rightarrow (1-\alpha) \times 100\%$ CI FOR $E(y_i)$ is:

$$\hat{y}_i \pm t_{\frac{\alpha}{2}} \left(\hat{\sigma}_{\hat{y}_i} \right) \quad df=n-2$$

PREDICTION INTERVALS FOR A FUTURE VALUE OF y_F SAY FOR A GIVEN VALUE OF x_F

$$\frac{y_F - \hat{y}_{\text{BASED ON } x_F}}{\hat{\sigma}_{(y_F - \hat{y}_{\text{BASED ON } x_F})}}$$

$\sim t$ WITH $n-2$ d.f.

OBTAINED USING PYTHAG. THM

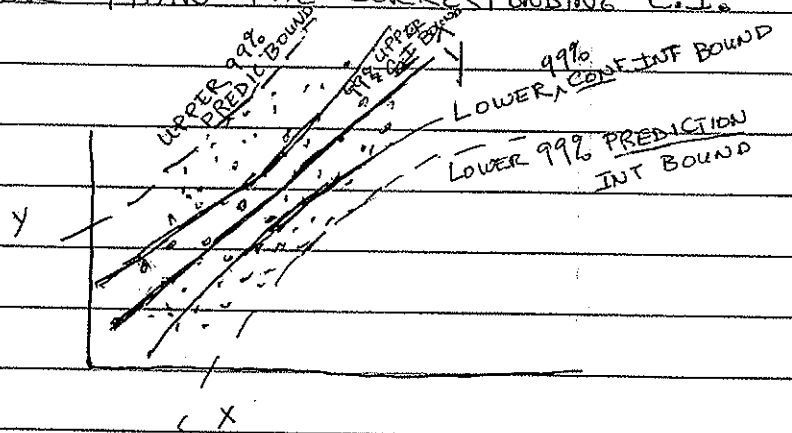
$$\hat{y}_i \pm t_{\frac{\alpha}{2}} \left(\hat{\sigma}_{(y_F - \hat{y}_{\text{BASED ON } x_F})} \right)$$

THE PREDICTION INTERVAL WILL ALWAYS BE WIDER THAN THE CORRESPONDING C.I.

"ROUGHLY"

GET RESULTS

LIKE THE FOLLOWING:



MULTIPLE REGRESSION

NOW INSTEAD OF A SINGLE PREDICTOR (X)

WE HAVE MULTIPLE "TERMS", CALL THEM Z_1, Z_2, \dots, Z_K

THAT CAN INFLUENCE Y .

THE MODEL IS NOW:

(FOR $i=1, 2, \dots, n$)

$$Y_i = \underbrace{\beta_0 + \beta_1 Z_{1i} + \beta_2 Z_{2i} + \dots + \beta_K Z_{Ki}}_{\substack{E(Y_i) \\ \text{deterministic part}}} + \underbrace{\epsilon_i}_{\substack{\text{random} \\ \text{part.}}} \quad \begin{array}{l} \epsilon_i \sim \text{Normal}(0, \sigma^2) \\ \epsilon_i \text{'s indpt.} \end{array}$$

(unknown) parameters are $\beta_0, \beta_1, \dots, \beta_K$ & σ^2 which must be estimated

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 Z_{1i} + \hat{\beta}_2 Z_{2i} + \dots + \hat{\beta}_K Z_{Ki}$$

Choose the estimates $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_K$ so as to minimize $\sum_{i=1}^n (Y_i - \hat{Y}_i)^2$

i.e. Choose ests. so as to minimize $\sum_{i=1}^n (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 Z_{1i} + \hat{\beta}_2 Z_{2i} + \dots + \hat{\beta}_K Z_{Ki}))^2$

conceivably could take partial derivative with respect to each β_j & set all = 0, we'll see later that doing this with vector/matrix notation greatly simplifies the concepts/math.

Then σ^2 is estimated by $\hat{\sigma}^2 = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n - (K+1)} \equiv S^2$

To test $H_0: \beta_j = 0$

$H_A: \beta_j \neq 0$

Use Test Stat:

$$\frac{\hat{\beta}_j - 0}{\hat{\sigma}_{\hat{\beta}_j}} \bigg|_{H_0} \sim t \text{ with } n - (K+1) \text{ d.f.}$$

VECTOR/MATRIX EXPRESSION OF MULTIPLE REGR.

→ REVIEW SIMPLE LINEAR ALGEBRA ARGUMENTS

→ REVIEW MATRIX LOOK @ MULTIPLE REGRESSION

$$\underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} \beta_0 \cdot 1 + \beta_1 z_{11} + \beta_2 z_{21} + \dots + \beta_k z_{k1} + \epsilon_1 \\ \beta_0 \cdot 1 + \beta_1 z_{12} + \beta_2 z_{22} + \dots + \beta_k z_{k2} + \epsilon_2 \\ \vdots \\ \beta_0 \cdot 1 + \beta_1 z_{1n} + \beta_2 z_{2n} + \dots + \beta_k z_{kn} + \epsilon_n \end{pmatrix} = \begin{pmatrix} 1 & z_{11} & z_{21} & \dots & z_{k1} \\ 1 & z_{12} & z_{22} & \dots & z_{k2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & z_{1n} & z_{2n} & \dots & z_{kn} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

$\xrightarrow{\text{label as "Z"}}$ $\xrightarrow{\text{label as } \underline{\beta}}$ $\xrightarrow{\text{label as } \underline{\epsilon}}$

NOT NORMAL, IT'S A MATRIX OF "PREDICTORS"

$$\Rightarrow \underline{y} = \underline{Z} \cdot \underline{\beta} + \underline{\epsilon}$$

$$\Rightarrow \hat{\underline{y}} = \underline{Z} \hat{\underline{\beta}} = \underline{Z} \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_k \end{pmatrix}$$

Choose $\hat{\underline{\beta}}$ so as to minimize $\sum_{i=1}^n (y_i - \hat{y}_i)^2$ in vector/matrix notation

$$= (\underline{y} - \hat{\underline{y}})^T (\underline{y} - \hat{\underline{y}}) = (\underline{y} - \underline{Z} \hat{\underline{\beta}})^T (\underline{y} - \underline{Z} \hat{\underline{\beta}})$$

$$= (\underline{y}^T - \hat{\underline{\beta}}^T \underline{Z}^T) (\underline{y} - \underline{Z} \hat{\underline{\beta}}) = \underline{y}^T \underline{y} - \underbrace{\hat{\underline{\beta}}^T \underline{Z}^T \underline{y} - \underline{y}^T \underline{Z} \hat{\underline{\beta}}}_{-2 \hat{\underline{\beta}}^T \underline{Z}^T \underline{y}} + \hat{\underline{\beta}}^T \underline{Z}^T \underline{Z} \hat{\underline{\beta}}$$

Now Take $\frac{\partial}{\partial \underline{\beta}} \Rightarrow 0 - 2 \underline{Z}^T \underline{y} + 2 (\underline{Z}^T \underline{Z}) \hat{\underline{\beta}}$ set = 0 $\Rightarrow \underline{\hat{\beta}} = (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \underline{y}$

NOTE: $E(\hat{\underline{\beta}}) = E[(\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \underline{y}] = (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T E(\underline{y}) = (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T (\underline{Z} \underline{\beta}) = (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \underline{Z} \underline{\beta} = \underline{\beta}$ so $\hat{\underline{\beta}}$ is unbiased est. of $\underline{\beta}$

Using $E(cX) = cE(X)$

Variance Cov MATRIX FOR $\hat{\underline{\beta}}$

$$= (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \begin{pmatrix} \text{Var} \\ \text{Cov} \\ \text{FOR } \underline{y} \end{pmatrix} [(\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T]^T = (\underline{Z}^T \underline{Z})^{-1} \sigma^2$$

Using $\text{Var}(cX) = c \cdot \text{Var}(X) \cdot c^T$

general

$$\begin{pmatrix} \text{Var} \hat{\beta}_0 & \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) \\ \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) & \text{Var} \hat{\beta}_1 \\ \vdots & \vdots \\ \text{Cov}(\hat{\beta}_m, \hat{\beta}_k) & \text{Cov}(\hat{\beta}_k, \hat{\beta}_m) \end{pmatrix}$$

so $\text{Var} \hat{\beta}_j = [j, j] \text{ element of } (\underline{Z}^T \underline{Z})^{-1} \cdot \sigma^2$

use this in t-test on β_j

FIRST-ORDER MULTIPLE REGRESSION

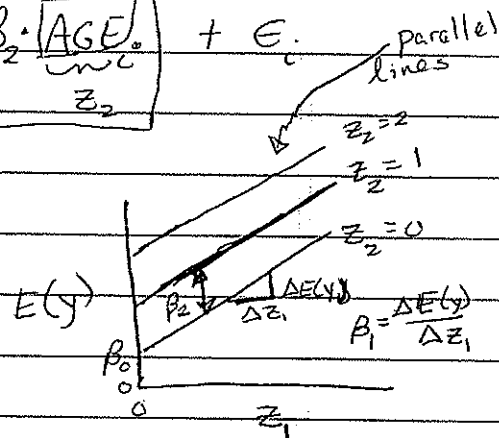
- ALL OF THE PREDICTING "TERMS" (z_1, z_2, \dots, z_k) ARE SEPARATE LINEAR (1ST-ORDER) VARIABLES

Ex. $\underbrace{\text{WEIGHT}_i}_y = \beta_0 + \beta_1 \cdot \underbrace{\text{HEIGHT}_i}_{z_1} + \beta_2 \cdot \underbrace{\text{AGE}_i}_{z_2} + \epsilon_i$

$E(y_i)$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 z_{1i} + \hat{\beta}_2 z_{2i}$$

$$\underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$



$Z =$ matrix $\begin{pmatrix} 1 & z_{11} & z_{21} \\ 1 & z_{12} & z_{22} \\ \vdots & \vdots & \vdots \\ 1 & z_{1n} & z_{2n} \end{pmatrix}$

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = (Z^T Z)^{-1} Z^T \underline{y}$$

$$\begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{pmatrix} = \underline{\hat{y}} = Z \hat{\beta} ; \hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-3} = \frac{(y - \hat{y})(y - \hat{y})^T}{n-3}$$

To test $H_0: \beta_j = 0$ vs $H_A: \beta_j \neq 0$.

Test Stat: $\frac{\hat{\beta}_j - 0}{\sqrt{\hat{\text{Var}}(\hat{\beta}_j)}} \underset{H_0}{\sim} t \text{ with } n-3 \text{ df}$

(where $\hat{\text{Var}}(\hat{\beta}_j) = [(j,j) \text{ element } (Z^T Z)^{-1}] \cdot \hat{\sigma}^2$)

Rej H_0 iff $\begin{matrix} \text{observed value} \\ \text{of test stat} \end{matrix} \begin{matrix} > t_{\alpha/2} \\ \text{OR} \\ < -t_{\alpha/2} \end{matrix}$

ALLOWING INTERACTIONS BETWEEN VARIABLES

EX.

EXPAND THE WEIGHT: HEIGHT, AGE 1ST ORDER

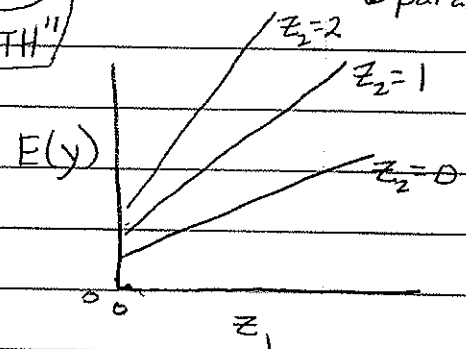
MODEL TO ALLOW THE LINEAR EFFECT

OF HEIGHT (z_1) TO DEPEND ON AGE (z_2) "INTERACT WITH" NOT parallel

FOR $i=1, 2, \dots, n$

$$\text{WEIGHT}_i = \beta_0 + \beta_1 z_{1i} + \beta_2 z_{2i} + \beta_3 \underbrace{z_{1i} z_{2i}}_{z_{3i}} + \epsilon_i$$

$E(y_i)$



Now the Z matrix has an extra 4th column & β has 4th value

Still $\hat{\beta}$ of form $(Z^T Z)^{-1} Z^T y$ $\hat{y} = Z \hat{\beta}$ Now $\hat{\sigma}^2 = \frac{(y - \hat{y})^T (y - \hat{y})}{n-4}$

& NOTE can test if there is an interaction effect by testing $H_0: \beta_3 = 0$ in a similar t-test manner.

ALLOWING QUADRATIC (2ND-ORDER) AND HIGHER ORDER TERMS

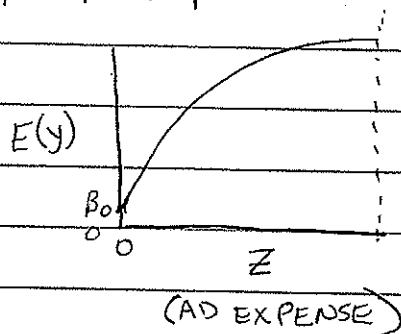
EX: EXPAND THE $\overbrace{\text{WEIGHT}}^Y$ & $\overbrace{\text{HEIGHT}}^{Z_1}$ $\overbrace{1^{\text{ST}} \text{ ORDER}}^{\text{LINEAR}}$

MODEL TO ALLOW QUADRATIC (2ND-ORDER) & HIGHER ORDERS

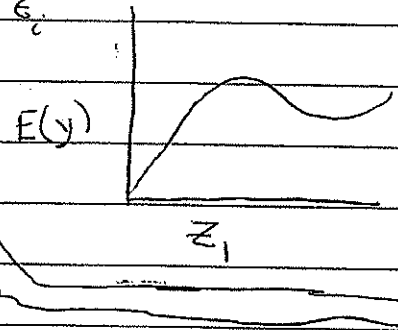
$$\text{SAY } Y_i = \beta_0 + \beta_1 Z_{1i} + \beta_2 Z_{1i}^2 + \beta_3 Z_{1i}^3 + \epsilon_i$$

KOHLS SALES VS. ADVERTISING EXPENSE: Z
Months $i=1, 2, \dots$

EX: $\underbrace{\text{SALES}}_{Y_i} = \beta_0 + \beta_1 Z + \beta_2 Z^2 + \epsilon_i$



HERE: β_2 is negative
 β_1 is positive.



ALLOWING MULTIPLE VARIABLES OF NON-LINEAR
ORDER AND INTERACTION(S) BETWEEN VARIABLES

BACK TO WEIGHT : HEIGHT, AGE BUT NOW ALLOW
 y z_1 z_2

2ND-ORDER FOR HEIGHT^(z₁) & AGE^(z₂) AND ALSO

ALLOW INTERACTION BETWEEN HEIGHT &

AGE @ LINEAR LEVEL,

LINEAR LEVEL
INTERACTION

$$\Rightarrow V_i = \beta_0 + \beta_1 z_{1i} + \beta_2 z_{1i}^2 + \beta_3 z_{2i} + \beta_4 z_{2i}^2 + \beta_5 z_{1i} z_{2i} + \epsilon_i$$

(NOTE: IF ALSO ALLOWED INTERACTION BETWEEN z_1 & z_2
AT QUADRATIC (2ND-ORDER) LEVEL,

THEN WOULD ALSO ADD A $\beta_6 z_{1i}^2 z_{2i}^2$ TERM.

INCORPORATING QUALITATIVE VARIABLES INTO MULTIPLE REGRESSION

BY USE OF "DUMMY" VARIABLES.

Ex. ONE WANTS TO COMPARE THE EFFECTS
OF 3 DIFFERENT TEACHING TECHNIQUES:
(TUTORING) A, B, C
ON STUDENT PERFORMANCE
TUTEE

"WE USE DUMMY" z_j VARIABLES SO WE CAN
DO THIS IN A MULTIPLE REGRESSION

SETTING. SET ONE LEVEL, SAY "A" AS BASE LEVEL.

THEN SET UP (0,1) DUMMY VARIABLES FOR EACH OTHER LEVEL.

HERE: $z_{1i} = \begin{cases} 1 & \text{if } i \text{ GOT LEVEL B} \\ 0 & \text{if NOT} \end{cases}$ $z_{2i} = \begin{cases} 1 & \text{if } i \text{ GOT LEVEL C} \\ 0 & \text{if NOT} \end{cases}$

TUTEE:
PERFORMANCE

$$E(y_i) = \beta_0 + \beta_1 z_{1i} + \beta_2 z_{2i} + \epsilon_i$$

NOTE: if i got A, then $E(y_i) = \beta_0 + \beta_1(0) + \beta_2(0) = \beta_0$, so $\mu_A = \beta_0$
if i got B, then $E(y_i) = \beta_0 + \beta_1(1) + \beta_2(0)$ so $\mu_B = \beta_0 + \beta_1 \Rightarrow \beta_1 = \mu_B - \mu_A$
if i got C, then $E(y_i) = \beta_0 + \beta_1(0) + \beta_2(1)$ so $\mu_C = \beta_0 + \beta_2 \Rightarrow \beta_2 = \mu_C - \mu_A$

FOR DUMMY SETTING:

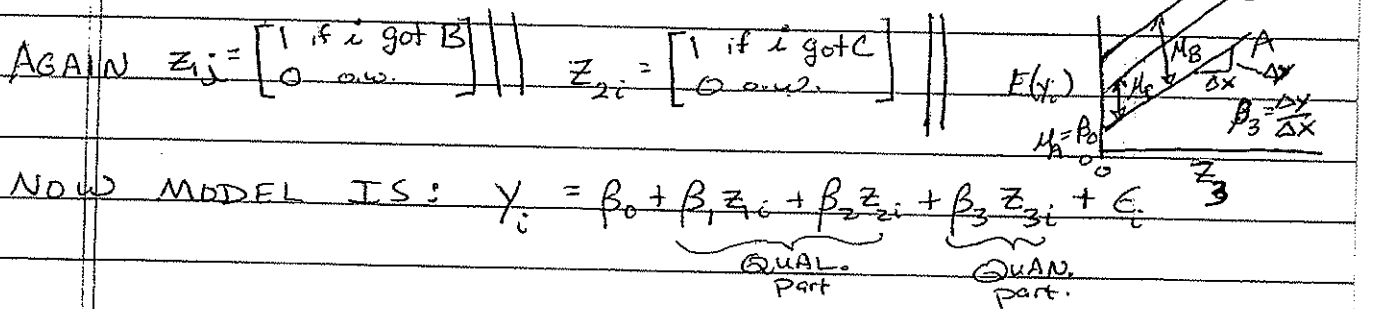
$$\beta_0 = \mu_{\text{BASE LEVEL}} ; \beta_j = \mu_{z_j \text{ level}} - \mu_{\text{BASE LEVEL}}$$

INCLUDING BOTH QUALITATIVE AND QUANTITATIVE VARIABLES IN MULTIPLE REGR.

Ex. GO BACK TO "TUTEE/TECHNIQUE" ^{QUALITATIVE A,B,C} EXAMPLE, BUT
 NOW, ^{ALSO} ^{LINEAR} ALLOW EFFECT DUE TO LENGTH z_{3i} , OF

TUTEE i 's TUTORING SESSION. ϵ_i DO NOT ALLOW

INTERACTION BETWEEN QUAL & QUANT. THEN

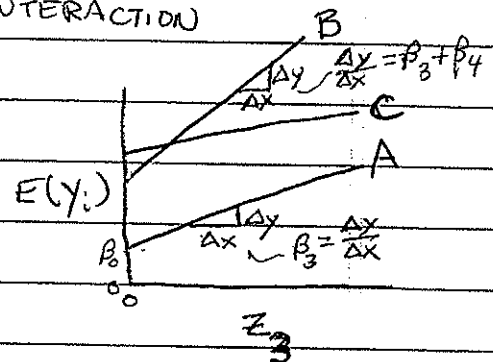


Ex. ABOVE ^{ALSO} BUT ALLOW INTERACTION BETWEEN QUAL & QUANT

QL,QT INTERACTION

Now, $y_i = \beta_0 + \beta_1 z_{1i} + \beta_2 z_{2i} + \beta_3 z_{3i} + \beta_4 z_{1i} z_{3i} + \beta_5 z_{2i} z_{3i} + \epsilon_i$

QUAL. QUANT. QUAL/QUANT. INTERACTION



TEST OF HYPOTHESIS WHEN ONE MODEL (LABELLED "REDUCED") IS NESTED WITHIN A LARGER MODEL (LABELLED "COMPLETE")

GENERAL SETTING

VIEW "COMPLETE" (LARGER) MODEL AS:

$$Y_i = \beta_0 + \beta_1 z_{1i} + \beta_2 z_{2i} + \dots + \beta_g z_{gi} + \beta_{g+1} z_{g+1i} + \dots + \beta_K z_{Ki} + \epsilon_i$$

AND "REDUCED" (SMALLER) MODEL AS:

$$Y_i = \beta_0 + \beta_1 z_{1i} + \beta_2 z_{2i} + \dots + \beta_g z_{gi} + \epsilon_i$$

$$H_0: 0 = \beta_{g+1} = \beta_{g+2} = \dots = \beta_K = 0 \text{ (i.e. REDUCED MODEL "SUFFICES")}$$

vs H_A : Not H_0

(At least one of

$z_{g+1}, z_{g+2}, \dots, z_K$ is important in predicting y)

(i.e. none of $z_{g+1}, z_{g+2}, \dots, z_K$ are important in predicting y)

1st FOR COMPLETE (C) MODEL, GET $\hat{\beta}_0^C, \hat{\beta}_1^C, \dots, \hat{\beta}_K^C$ to minimize $\sum (y_i - \hat{y}_i^C)^2 \equiv SSE_C$
(UNDER)

NEXT UNDER REDUCED (R) MODEL, GET $\hat{\beta}_0^R, \hat{\beta}_1^R, \dots, \hat{\beta}_g^R$ to minimize $\sum (y_i - \hat{y}_i^R)^2 \equiv SSE_R$

Test Statistic:

$$\frac{SSE_R - SSE_C}{K - g} \div \frac{SSE_C}{n - (K + 1)}$$

$\sim F_{\text{num df} = K - g, \text{den df} = n - (K + 1)}$
 H_0

Rej H_0 iff observed value of test stat. $> F_{\alpha}$

ILLUSTRATION OF NESTED MODEL.

GO BACK TO THE

EX,
GL, GT
INTERACTION

* WHERE
$$Y_i = \beta_0 + \underbrace{\beta_1 Z_{1i} + \beta_2 Z_{2i}}_{\text{QUALT}} + \underbrace{\beta_3 Z_{3i}}_{\text{QUANT}} + \underbrace{\beta_4 (Z_{1i} Z_{3i}) + \beta_5 (Z_{2i} Z_{3i})}_{\text{INTERACTION}} + \epsilon_i$$

like Z_{4i}^* like Z_{5i}^*

NOW SUPPOSE ONE WANTED TO TEST THE NULL HYPOTHESIS

THAT NEITHER THE QUANT VARIABLE NOR INTERACTION

ARE IMPORTANT IN PREDICTING y (i.e. $H_0: 0 = \beta_3 = \beta_4 = \beta_5 = 0$)

VERSUS H_A : NOT H_0 .

SO, COMPLETE MODEL IS GIVEN BY *, ABOVE, WITH $K=5$.

THE REDUCED MODEL IS: $Y_i = \beta_0 + \beta_1 Z_{1i} + \beta_2 Z_{2i} + \tilde{\epsilon}_i$, WITH $g=2$.

TEST STATISTIC:

$$\frac{SSE_R - SSE_c}{K - g} \bigg/ \frac{SSE_c}{N - (K+1)}$$

\sim_{H_0}

$$F_{\text{num d.f.} = K - g, \text{den d.f.} = n - (K+1)}$$

NOTE: $SSE_c = \sum (y_i - \hat{y}_i^c)^2$; $SSE_R = \sum (y_i - \hat{y}_i^R)^2$

Rej H_0 if observed

value of test stat $> F_{\alpha}$

(for num d.f. = $K - g$
den d.f. = $n - (K+1)$)

GEOMETRIC VECTOR / MATRIX INTERPRETATION OF NESTED MODEL EX.

RECALL IN GL, QT INTERACTION EX OF NESTED MODEL, THAT

FOR $i=1,2,\dots,n$,

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 z_{1i} + \hat{\beta}_2 z_{2i} + \hat{\beta}_3 z_{3i} + \hat{\beta}_4 z_{4i} + \hat{\beta}_5 z_{5i}$$

$$\Rightarrow \hat{y} = \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \hat{\beta}_0 + \begin{pmatrix} z_{11} \\ z_{12} \\ \vdots \\ z_{1n} \end{pmatrix} \hat{\beta}_1 + \begin{pmatrix} z_{21} \\ z_{22} \\ \vdots \\ z_{2n} \end{pmatrix} \hat{\beta}_2 + \begin{pmatrix} z_{31} \\ z_{32} \\ \vdots \\ z_{3n} \end{pmatrix} \hat{\beta}_3 + \begin{pmatrix} z_{41} \\ z_{42} \\ \vdots \\ z_{4n} \end{pmatrix} \hat{\beta}_4 + \begin{pmatrix} z_{51} \\ z_{52} \\ \vdots \\ z_{5n} \end{pmatrix} \hat{\beta}_5$$

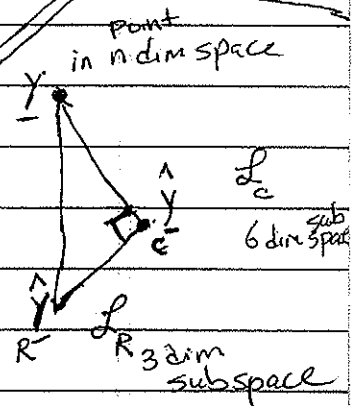
$\underbrace{\quad}_{\hat{z}_0} \quad \underbrace{\quad}_{\hat{z}_1} \quad \underbrace{\quad}_{\hat{z}_2} \quad \underbrace{\quad}_{\hat{z}_3} \quad \underbrace{\quad}_{\hat{z}_4} \quad \underbrace{\quad}_{\hat{z}_5}$

$$\hat{y} = \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{pmatrix} = \begin{pmatrix} 1 & z_{11} & z_{12} & z_{13} & z_{14} & z_{15} \\ 1 & z_{21} & z_{22} & z_{23} & z_{24} & z_{25} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & z_{n1} & z_{n2} & z_{n3} & z_{n4} & z_{n5} \end{pmatrix} \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \\ \hat{\beta}_4 \\ \hat{\beta}_5 \end{pmatrix}$$

with linear independence of $\hat{z}_0, \hat{z}_1, \dots, \hat{z}_5$,

\hat{y} is the point in the $5+1=6$ dimensional subspace of n dimensional space (generated by $\hat{z}_0, \hat{z}_1, \hat{z}_2, \hat{z}_3, \hat{z}_4, \hat{z}_5$) that

has the smallest squared distance from y



$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = y$$

Similarly

$$\hat{y}_R = \begin{pmatrix} \hat{y}_{R1} \\ \hat{y}_{R2} \\ \vdots \\ \hat{y}_{Rn} \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \hat{\beta}_0 + \begin{pmatrix} z_{11} \\ z_{12} \\ \vdots \\ z_{1n} \end{pmatrix} \hat{\beta}_1 + \begin{pmatrix} z_{21} \\ z_{22} \\ \vdots \\ z_{2n} \end{pmatrix} \hat{\beta}_2$$

$\underbrace{\quad}_{\hat{z}_0} \quad \underbrace{\quad}_{\hat{z}_1} \quad \underbrace{\quad}_{\hat{z}_2}$

with linear indep. of $\hat{z}_0, \hat{z}_1, \hat{z}_2$,

\hat{y}_R is the point in the $2+1=3$ dim (generated by $\hat{z}_0, \hat{z}_1, \hat{z}_2$) that has the smallest squared distance from y

SSE_C is the squared distance between y and \hat{y}_C
 SSE_R " " " " " " " " y and \hat{y}_R .