

Fr 1 June 2016
Colors

Black : possible student exercise (4)

SS = Synchronous Session Blue : reminders for me (1)

Q = Question Red & Teaching notes (3)

Green : calculations (2)

don't deal with samples - much less complicated

Part 1 (1-7) : JUD = Judgment, Uncertainty,
- and Decisions

Part 2 (8-14) : DAD = Data Analysis
and Decisions

Week

1	SS 1 intro	8
2	SS 2 probabilities	9
3	3 special distributions	10
4	office hours w/me?	11
5	4 regression analysis	12
6	ditto	13
7	SSS summary & review?	14 10 th summary & review

? front 15 ~ $\frac{1}{2}$ hr Qs, $\frac{1}{2}$ hr directed
end 15 ~ $\frac{1}{2}$ hr Qs, $\frac{1}{2}$ hr directed

- ingredients : a) problem discussion
b) lecture discussion
c) reading discussion → papers



Session 1 : Exploring Data ("in depth")

"characterize
efficiently"

- "essence of a body of data"
- "fanning Q" → data; but data → decision; with these data

- best way to view data

- frequency = # time
- "absolute" vs. "relative"

13:36 works 'cause calculating "proportion years"; \sum = cumulative

$$\text{THINK: what? how?} \quad \text{proportion} = \frac{\# \text{ people} \leq x}{\text{total people}} \quad \frac{\# \text{ people}}{\text{people}} \rightarrow \text{ave yrs per person}$$

use Anscombe's Quartet ? yes

- importance of plotting

- NB: can subtract years from average (in calc) 'cause they emphasize have the same dimensions, first, and units [years]

* exercise: demonstrate original mean calculation equivalent to professor's method:

$$B) \bar{x} = \frac{1}{N} \sum_{i=1}^N (\mu - \bar{x})_i^2 \quad \left. \begin{array}{l} \text{A)} \\ \frac{1}{N} \sum_{i=1}^N (\Delta t)_i = \mu \end{array} \right\}$$

$$\sum_{i=1}^{10} (\mu - \bar{x}_i)^2$$

$$10\mu - \sum_{i=1}^{10} (\Delta t)_i - \sum_{i=2}^{10} (\Delta t)_i$$

$$10\mu - \frac{n_1(\Delta t)_1}{N} -$$

$$\mu - \frac{n_1(\Delta t)_1}{N} + \mu - \frac{n_2(\Delta t)_2}{N} + \dots$$

$$\bar{x} = \frac{\sum \Delta t_i}{N}$$

$$= \frac{n_1(\Delta t)_1}{N} + \dots$$

$$= \sum_{i=1}^{10} \frac{(\Delta t)_i}{N}$$

$$\frac{1}{N} ((\Delta t) - (\Delta t)_1 + (\Delta t) - (\Delta t)_2) + \dots$$

* don't see how we can understand skewness (or kurtosis?)
w/o calculating? exercise?

* What is a "random variable"? e.g. in this example, \boxed{Q}
+ interquartile range - middle 50% truly random?

Qpww: around dots min at max?
[at box ~ whisker]

* IQ = normal?!

"policy implications" → decisions possible

* $\frac{\sigma}{\mu}$ = coefficient of variation; "relative dispersion"

1) key measures

2) policy implications



Session 2: Probability Concepts

- classical, empirical, subjective
 ↗ prior knowledge ↗ repeat "on the other hand"
 ↗ estimate "judgment" ↗ "all ok"
 ↗ 1) discussion: examples ↗ 2) use most? ↗ 3) best? "strengths & weaknesses"
 ↗ 4) evaluate ↗ SS: REQUIRE QUESTIONS

6 June

2) What is meant by "large"? (law of large numbers)

3) experiment: process separable

in which the results are uncertain

4) event vs simple event (outcome)
 (collection) of

what do we mean by this?

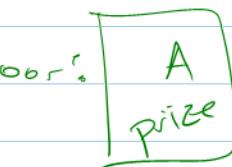
* distribution → what we expect: management
 ← neutrino experiment ← future

- or: "U" : $P(A \cup B) \rightarrow A \text{ or } B \text{ or } (A \text{ and } B)$
- and: $A \rightarrow$ both
- complement: $\overline{P(A)} = 1 - P(A)$
- $P(A|B) =$ conditional probability
 * \rightarrow independent if $P(A|B) = P(A)$
- * general addition law: $P(B \text{ or } C) = P(B) + P(C) - P(B \cap C)$
 ↗ zero if no e.
- Does everyone understand that writing as a decimal gives probability directly?
 that "per cent" means "per 100"?
- * (e.g., in smoking example) → explain chart numbers? (29% 20%)
- * Was "marginal probability" explained?
 Can also go through frequency, i.e., just numbers, first, before division by total
 → cause easier to think in terms of counting people (integers?)
- Pt: how you think about this can make a difference
 ↗ and in what order
- * doesn't really explain this conditional probability example
 ↗ $\frac{55 \text{ non smoke fundos}}{100 \text{ fundos}}$
 $= \frac{55/250}{100/250}$
- Exercise: create integer table first - then do probabilities
- * exercise? marginal, conditional, joint
- * weight lifting → have to do it, e.g., Monte Hall Problem
- * ↗ 5 problems for me!!

1) Monty Hall

①

Door:



B C

switch

i) choose	A	$\xrightarrow{\text{open}}$	$\xrightarrow{\text{or}}$	$\xrightarrow{\text{open}}$	\rightarrow	lose
in } i.e. }	X	B		O	\rightarrow	win
	X	O		C	\rightarrow	win

$\frac{2}{3}$

2) Bayes Rule: $P(A|X)P(X) = P(X|A)P(A)$

$$\text{Door A is chosen} \left\{ \begin{array}{l} \text{in A} \\ \text{B has been opened} \\ \text{(with X)} \end{array} \right\} \quad P(A|X) = \frac{P(X|A)P(A)}{P(X)} = \frac{\frac{1}{2}}{\frac{2}{3}} = \frac{1}{2}$$

$$(1/2) : \text{opening}$$

prize is in A, B, or C; prob of B

$$\underbrace{\frac{1}{3} \cdot \frac{1}{2} + \frac{1}{3} \cdot 0 + \frac{1}{3} \cdot 1}_{= \frac{1}{2}}$$

Q numerical coincidence?

? $P(X|A) = P(X)$ do

or are they truly independent?

The results are very counterintuitive. For the three-door problem, if the contestant is correct on the first choice, then he or she will be incorrect after a switch. Or, if a contestant was incorrect on the first choice, he or she will be correct after a switch. Since the probability of being incorrect on the first choice is 0.667, then the probability of being correct after a switch is 0.6667.

onlinestatbook.com/2/probability/monty_hall_demo.html

$\frac{1}{2} + 0 + \frac{1}{3}$,
then divide
by 3

7 June 2016



Session 3: Special Probability Distributions

- factorials

- spreadsheet created for binary calc



\Rightarrow each trial has 1 of only 2 outcomes
- same probability; no hysteresis

tires on
car, truck

- point out: 0.9^n gives same result as calculation

did it the opposite way ("success" = problem)

$$P(y) = \frac{\frac{n!}{y!(n-y)!}}{s^y (1-s)^{n-y}} ; \text{ e.g. } n=3, s=\frac{1}{3}$$

$$E(y) = ns$$

$$\sigma^2(y) = ns(1-s)$$

an infinite pool of samples,
i.e., continual replacement

as $n \rightarrow \infty$, $P \rightarrow$ Gaussian (unless s extreme) unlike ↘

* Hypergeometric (Hypergeomdist)
↳ depends on history

finite, known
number of objects

(i.e., no longer have const prob.
across trials)

- "sampling without replacement"

$$P(y) = \frac{\binom{N_s}{y} \binom{N_f}{n-y}}{\binom{N}{n}}$$

$s = \text{success}$
 $f = \text{failure}$

$$E(y) = \frac{nN_s}{N_{\text{total}}}$$

* Geometric "very different"
special case of the Negative Binomial Dist
- for modelling the runs of consecutive failures
in repeated independent trials before getting first

$$P(y) = (1-s)^{y-1} s ; E(y) = 1/s ; \sigma^2 = \frac{1-s}{s^2}$$

Negbinomdist (F, k, s) in Excel

↳ discussion of # failures number of successes (more general)

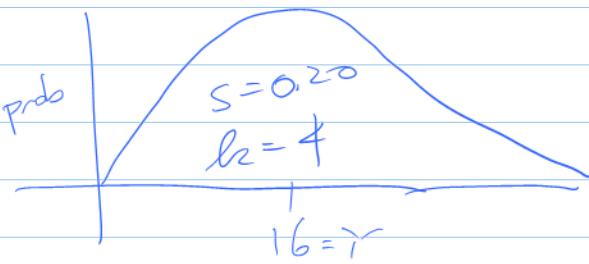
why it might look like that. — ; always 1 in geometric dist

$$P(y) = \frac{(y-1)!}{k! (1-s)^{y-k} s^k} \rightarrow \text{geo when } k=1$$

number of trials to obtain k successes

Part b: show ↗

$$E(Y) = \frac{h}{s} \quad ; \quad \sigma^2 = \frac{h(1-s)}{s^2}$$



* Poisson

$$P_Y(y) = \frac{e^{-\mu} \mu^y}{y!}$$

μ = mean; no upper limit
(unlike binomial)

* probability tree too complex to draw

→ Google it? discuss?? SS?

*

Poisson ($y|\mu$) cum or not

- three assumptions:
- 1) exactly one event → prob small & constant
 - 2) two events simultaneously ≈ 0
 - 3) events independent

ss
Discussion:
What if 2) is violated?

*

Exponential Distribution → distribution of waiting time before next event

- gives example of toll booth; isn't traffic sinusoidal?

makes sense

calculus → not covered

*

Normal Distribution

$$f(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{y-\mu}{\sigma}\right)^2}$$

what does this mean here? ?? $Y=f(y)$

NORMDIST ($x, \mu, \sigma, \text{true}$)

everyone get the difference calculation
to get 130 - 100 range?

*

7 distributions reviewed



Session 4: Decision Making Under Uncertainty

- 1) set of decisions
 - 2) set of outcomes, w/ probabilities
 - 3) "payoffs"
- key elements
in decision making
associated w/ decisions

three possible strategies:

- a) Maximin: minimizes downside risk
- b) Maximax: maximizes upside potential
- c) expected value

- point of "unfair coin": you don't know

sensitivity analysis: systematic altering of inputs

* p230f: Decision Trees

□ t, decision ^{decision known}

but "unfair" coin is one

○ t, probability ^{but}

result known ←

△ end: no uncertainty

conditional on previous events

~~specific~~

check if ~~diagram(s)~~
still used in book

no

"folding back procedure"

Q in SS?

↳ inconsistent (? too strong) with

* PWW Q: why 20, 25, 30 on diagram? ? = $\Delta(\text{win}-\text{lost})$?

SS Q: as Phil says, make sure understand all numbers
on diagrams

* "utility" ≈ "usefulness"; risk averse

→ utility maximizers vs EMV maximizers?

* mention Bazerman book



Session 5: Regression Analysis, Time Series Analysis, & Forecasting

"modeling is at the forefront of what managers do"
~~SSQ? T/F?? Should it be?~~

ind var (x) → "predictor"

↳ Q: bad terminology? I argue that it
contains data

dep var (y) → "criterion" even worse. and fit

"population regression function" → $E(Y) = \beta_0 + \beta_1 X$

straight line formula: $y = mx + b$

- graphing: good to put zero point

The slope has dimensions! SS definitely!

and so do x & y

- and label the graph → "sale price / \$ 1000"

or "sale price in thousands of \$"

\times Excel: Intercept ($\bar{Y}_{range}, \bar{X}_{range}$)
Google Slope ($\bar{Y}_{range}, \bar{X}_{range}$)

Point out Google S.S. work too
(from the beginning)

- method of "ordinary" L.S.

↳ hand

- \hat{y} means predicted

- outliers

"subsequent course"

→ part 2?

[OPWH]

\times correlation vs. regression Q (PWW)

↳ no distinction b/w dependent + ind. variable

Excel: Correl ($\bar{X}_{range}, \bar{Y}_{range}$)

$$\sum (x_i - \bar{x})(y_i - \bar{y}) / (n-1)$$

$$r_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{s_x s_y}}$$

\times multiple linear r: $\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$

SS: △ quantitative + qualitative; e.g., near a school can
 $\left[PWW \right]$ (vs numeric/non-numeric) be qualitative easily!

- additive vs. interactive (non-additive)
 - ↳ slope changes, e.g., w/quadrant of city
 - ↳ between variables (not people)

time series: $x \rightarrow t$ and so $y(x) \rightarrow y(t)$

log y vs. t ↳

$$\log y = mt + b$$

$\log y$ |
t

"change in $y = f(t)$ "

$$\log y_2 - \log y_1 \rightarrow \text{ratio}$$

$$= mt_2 - mt_1 \rightarrow m \Delta t$$

$$10 = \frac{y_2}{y_1}$$

≈ 37:30

? SS →

(x slide sequence off? this time or always?)

means it's good in forecasting models; explore SS??
should be on last slide, Alternative Models
#19

⇒ switch last two slides



SECTION 6: Optimization

QSS: What's the difference between "optimization" and numerical solution?

- can use Google Sheets, too
- Simplex Method: one of many to do the search
- constraints; QSS: LP used in your organization?

PC Tech Optimization problem still in textbook?

- how to avoid VoPP??

- "Key Elements": a) input variables
b) decision variables (changing)
c) objective function (maximized)
d) other { \$80 n_{Basics} + \$129 n_{XPs}
e) constraints
- but Google helps

- algebraic model is an intermediate step
- graphical solution i) not clear why plot XP_{Basics}?
ii) what if 3 dimensions?
Do in 3D?

solution:
560 Basics
1200 XPs

Q:

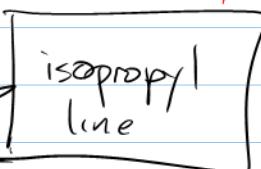
derive straight form for assemble hours constraint

$$5x_1 + 6x_2 \leq 10,000 \quad \left\{ \begin{array}{l} \text{and testing hour} \\ \text{constraints} \\ \text{+ answer} \end{array} \right.$$

x_1 (basis) x_2 (RP)

x_1 on plot \rightarrow max is $\frac{10,000}{6} = 1666.6$

j + units are ?

define  isopropyl line

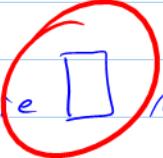
$x_2 \leq \frac{10000 - 5x_1}{6}$ why is

$= 1666.6 - \frac{5}{6}x_1$ THAT THE ANSWER?

$\Omega x = 2000$ $\frac{10,000}{6} - \frac{10,000}{6} = 0$

(QWQ) want to sell more XPs } want to maximize assembly hours, too } want area under line maximized?
 "higher"

- Simplex Method for LP

→ says book is excellent here  make one there?!

+ Google SS

- "quite a lot of time" SS discussion?

- postal problem ← should be easy to change to doctor model
 - blending models
 - logistics models
 - aggregate planning models; workforce levels → production schedules
 - financial models

"comparatively little wiggle room" → inputs specified highly

- mentions 2 more, too hard: integer & non-linear

but these at least

same principle (I-wellague)

probably doesn't use Simplex (but who cares?)



Session 7: Simulation

→ "brief"; "convincing me"; not generally Excel

→ FCS
sim 3D ^{SN-dominated}
FSM AND jet sim.

"specialized simulation language" ↳ Python?
wasn't listed (later)

QSS: A "simulation" + "model"? ("general purpose")
Wikipedia definition: → (disagreed) ↳ optimization ↳ programming

- objectives a) model

b) test vs. empirical data but can also produce values ↳ STELLA
c) → future

d) $f(\Delta \text{parameters}) \rightarrow \text{"what if"}$ } use of probabilities here, too

- bank/teller example; business, gov, manufacturing
w/in e.g.s, ⇒ health-care delivery

- Four advantages: a) no disruption
b) test before resources
c) i.d. bottlenecks
d) insight into key factors of system performance

- Discrete vs. Continuous

↪ but hydro units are a hybrid,
 $\Delta t \neq 0$

- "Model of a System"

yes

now agreeing w/ my disagreement above

* "look at the book"

(DATA ANALYSIS & DECISIONS)



Session 8: Sampling + Sampling Distributions [1]

refers to test
for more info
details

vs., say Hi-stage

"mu"

learn about population from sample data
probability sampling vs. judgment sampling
↪ systematic but then?!

4 types: simple,
mostly systematic,
in course stratified, cluster

Sampling error ↳ CONCEPT SSQ

\Rightarrow Capture, at the least cost, the heterogeneity of the population
What P.W. calls "characteristics" I would call "goals."
 \rightarrow efficiency: smallest w/ least possible error

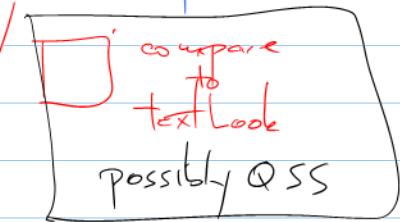
Simple Random Sampling (6)

simplicity, enumeration, representativeness, cost,
geographical dispersion, efficiency / less sampling error
(\leftarrow after \rightarrow heterogeneity) (after low in simple)

Jr June 2016

B * Gallup, "multi-stage"; but it doesn't use "simple"!

- in course \rightarrow random sample assumed QSS What about not in course?
- "empirical rule" (from where?): QSS?
 16^{moo} can change fast food to E.R.?
- pop parameter vs. sample statistics \rightarrow sampling error
(typically not known \rightarrow estimate \Rightarrow estimation process)
 - df taken from "dictionary reference to com" / "the fed dictionary"
 - 18^{m14} "the fed dictionary"
 - 18^{m51} : wikipedia definitions!! (oh my)
- point vs. interval



QSS [vs.] confidence intervals (frequentist)
credible intervals (Bayesian)

[sampling error vs. non-sampling error] QSS

QSS: what is truly random?

\rightarrow "estimation error" \rightarrow distribution (prob. dist. of a statistic
from random samp.)

- Central Limit Theorem; remember mean
(\bar{x} : roulette problem \rightarrow linear normal estimates $\Rightarrow \Delta$ (ii))
 - get a better definition! (the book?!) as n increases.
 - $n=30$ "fairly obvious" ?! 1920s?!
- if pop = normal, sample mean = normal, ind of n

$s = \sqrt{\frac{s_{\text{pop}}}{n}}$ "The Empirical Rule" (Applied to Sample Means)
standard error of sample mean $\approx 36^{moo}$ nice 5 step process (SS); if done 100 times, μ right 75% of the time

- but then not telling people how numerical answers obtained
 → make sure no questions

$$\sqrt{\frac{N-n}{N-1}} = \text{finite population correction}$$

$\rightarrow \left(\frac{1 - \frac{n}{N}}{1 - \frac{n}{N}} \right)^{\frac{1}{2}} = \frac{(1-x)^{\frac{1}{2}}}{\sqrt{A}}, A = (1 - \frac{n}{N})^{\frac{1}{2}}, x = \frac{n}{N}$

- needed when n/N is not small, i.e., $\geq 5\%$

★ Session 9 [2] Confidence Interval Estimation

- ch 8 omissions

Announcements

• Sections from Chapter 8 which you will not be expected to calculate the

confidence interval of but will still be expected to know about: confidence interval of, but will still be expected to know about:

- Confidence interval for a total
- Confidence interval for a standard deviation
- Paired samples case of confidence interval for difference in means.
- You will be held responsible on quizzes and the Final Exam for knowing how to calculate the sample size only for the estimate of a single population mean and of a single population proportion

check sections t & x^2

- using rule of α : $\alpha = 1 - \text{level of confidence}$

$$95\%: \bar{x} \pm 2 \frac{s_{\text{sampled}}}{\sqrt{n}}$$

$$= n-1 \quad (\text{for mean})$$

$$\text{degrees of freedom}$$

t distribution: "The sampling distribution of the standardized sample mean when the sample standard deviation is used in place of the population standard deviation"

} Text glossary Book defines/explains t!
 equation introduce P here too?

So if (in flight example) adding 2-digit integers, you would $\mu = 18.6$ (Correct - e.g. else silly, e.g. upper and lower bounds out to 5 decimals)

19th: nice "generic approach to confidence interval (at least for a mean)"

$$\text{sample proportion } \hat{p} \rightarrow \text{statistic } \hat{p} \pm z \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = s$$

standard error
 "certainty",
 ; stressing similarity of form

mug problem: 2 mugs sampled out of 106 → $\hat{p} = 0.02$
 $z \rightarrow z = \text{NORMINV}(1 - \frac{\alpha}{2}, 0, 1)$ QSE? What is \hat{p} ? 0, 1?

What about difference between 2 population means? (\bar{x})

$$\left\{ \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2} \left[\frac{1}{n_1} + \frac{1}{n_2} \right] \right\}^{1/2}$$

assuming equal s^2 .

actually say "plug + chug"
= standard error of the difference

$\rightarrow (n_1-1) + (n_2-2) = n_1+n_2 - 2$
again use t_{inv} , but now called "critical value"

then, 2 pop proportions ($\hat{\pi}$) difference

$$\left\{ \frac{\hat{\pi}_1(1-\hat{\pi}_1)}{n_1} + \frac{\hat{\pi}_2(1-\hat{\pi}_2)}{n_2} \right\}^{1/2} = \text{standard error of difference in proportion}$$

- same $Z = norminv(1 - \alpha/2, 0, 1)$

x

Sample size for estimating population mean:

$$Z \frac{\sigma}{\sqrt{n}} = B$$
$$Z \frac{s}{\sqrt{n}} = B$$

QSS
Make student again from } desired
derive these
two formulas?

$$n = \left(\frac{Z \cdot \text{test}}{B} \right)^2$$

estimate of standard dev = s

margin of error

For proportion, same Z :

and $Z \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}}$ (from above) set to B

uses $\hat{\pi} = 0.5$

as best guess

$$\rightarrow B = Z \frac{\sqrt{\hat{\pi}(1-\hat{\pi})}}{n}$$

$$\Rightarrow n = \left(\frac{Z}{B} \right)^2 \sqrt{\hat{\pi}(1-\hat{\pi})}$$

'cause you choose about
a specific value about
which to calculate



Session 10 Hypothesis Testing, Part 1

H_0 must always contain our = sign { hypotheses refer to the population (μ) not the sample}

H_0 = null \rightarrow disprove, vs. H_A = alternative hypothesis

try to see if it is disproven { not true that finding alternative true is goal

\Rightarrow not guilty \neq innocent } only confidence, not proof

("reasonable doubt" = α (= 1 - confidence))

- don't "accept" hypothesis
"frames of the null hypothesis are never accepted. We either reject them or fail to reject them." jerry dallas

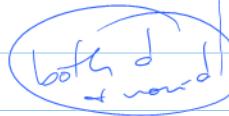
"constitution implicitly understood it." ERIC CIOFFI 17 June 2016

- directional alternative hypotheses vs non-directional

1 tail

2 tail

→ "test of significance approach"
(p value)



"confidence interval"

vs only for non-d.

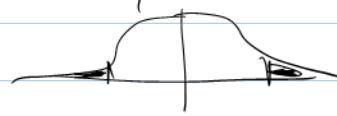
$$t_{crit} = TINV(2\alpha, v) \text{ for directional}$$

→ significance level changing? (no)

function obscures

$$(\alpha, v) \text{ for non directional}$$

so why 2α vs α ?



$$t = \frac{\bar{x} - \mu}{S.E. \sim} = \frac{s}{\sqrt{n}}$$

$|t| > t_{crit}$ and direction consistent

vs. $t > t_{crit}$? w/o α / H_A

"if you can justify directional" QSS did he do that?

↳ i.e., more info \Rightarrow more "statistical power"

says $H_0: \mu \leq 1.4$ "fully equivalent" to $H_0: \mu = 1.4$

(30") $H_a: \mu > 1.4 \rightarrow$ directional

(referring to textbooks)

$H_a: \mu > 1.4$

Eric likes

"operating on that value (1.4)"
and justification for directionality
is apparent & consistent

now proportion $\rightarrow \pi$ (text uses p)

and
 π

$$z_{crit} = NORMINV(1 - \alpha, 0, 1) \text{ directional}$$

$\left\{ \begin{array}{l} (1 - \alpha) \\ \alpha \end{array} \right\}$ non directional

$$S.E. = \sqrt{\frac{\pi(1-\pi)}{n}}$$

π is number you want

$$z = \frac{(\hat{\pi} - \pi_0)}{S.E.}; \hat{\pi} \text{ you measured (sample)}$$

again, if $|z| > z_{crit}$ and direction is consistent w/ H_A

- h. testing, difference in means ; again t
- $$s.e. = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1+n_2-2)} \left[\frac{1}{n_1} + \frac{1}{n_2} \right]}$$
- $$t = \frac{\bar{x}_1 - \bar{x}_2 - \mu_0}{s.e.}$$
- the hypothesized difference
- same as before, $|t| > t_{\text{crit}}$, direction $\rightarrow H_A$
- "LAB REPORT"
- Today: 53:50
- 58:50
- hypothesis testing approach
- 1) t-test of a single population mean; proportion
- 2) independent sample test, or
- 3) a t-test of two population means
- "things You Should Know How To Do"

"in general the approach is the same"

give em a page or two??

summary
each week



Session II: Advanced Hypothesis Testing

[4]

Tough one \rightarrow make sure have SS

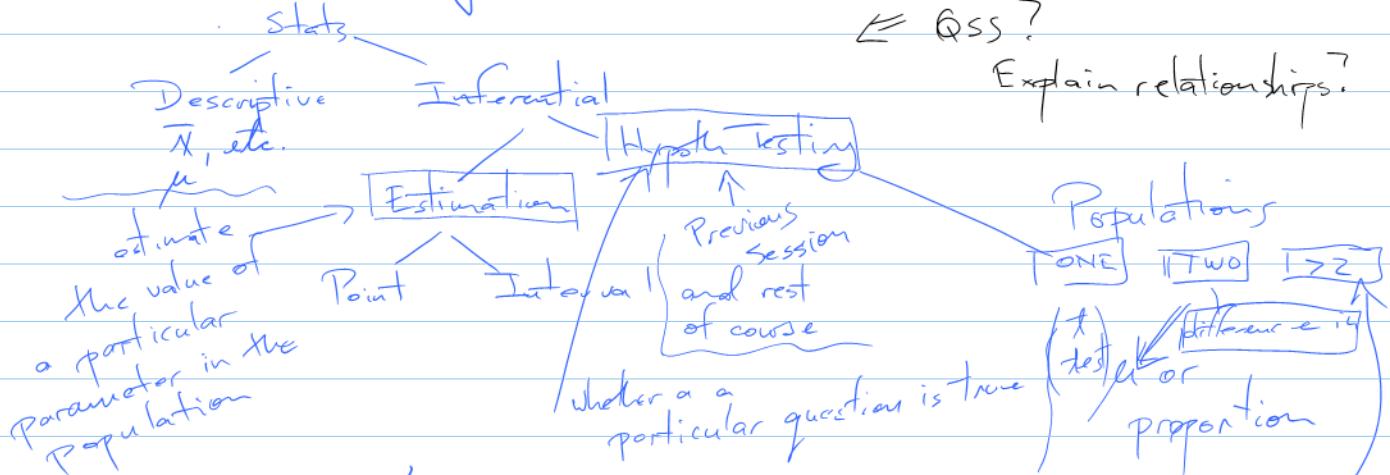
You will be held responsible for the substance, but not the computation of, hypothesis tests of

- equal population variances
- tests of normality

"understanding"

Excel Analysis ToolPak : Analysis Group on Data Tab
 \rightarrow used in most of course

Statistical Roadmap



Hyp. Test for $\Delta \text{ in } \mu$ / $H_0: \mu_1 = \mu_2$ $\xrightarrow{\text{pt. out?}}$ "analysis of variance"

$H_A: \mu_1 \neq \mu_2$

DO IT \rightarrow t-test: two sample assuming variance

- e.g., TINV in EXCEL \rightarrow critical value: "considerable disadvantage" $\Rightarrow p\text{ value} \equiv \text{conclusion}$

$\text{if } p < \alpha \rightarrow \text{reject } H_0$

otherwise, "fail to reject H_0 " = fail to accept H_a "

p : prob making Type I error
if reject H_0 (false positive: incorrect rejection)

23 June

QSS 1.53 vs 1.52
25^m → argue there's reality and there's statistics! discuss

overhead of an α to conference { "one chance in x statistically..."
"one chance in a million it's real!"

* - To make it directional, divide two-tail P by 2 (\equiv one-tail [ct])

Ind. Samples vs. Paired / Matched Samples

(2 pops (or more) vs ↑ 1 pop w/ one sample \rightarrow measure 2 variables
Pepsi vs. Coke example
 \rightarrow each row represents a different person

"we conclude w/ reasonable certainty"

extending: μ_y same across all populations] null $H_0: \mu_a = \mu_b = \mu_c$
↳ dependent, continuous variable] alternative

μ_y different in at least one pop \Rightarrow ANOVA = analysis of variance { $H_a: \mu_a \neq \mu_b = \mu_c$
 $\mu_a = \mu_b \neq \mu_c$

can F test too
various sum of squares - "group mean"

here,
3 groups

1) Total: $\sum_{ij} (y_{ij} - \bar{y})^2$

\Rightarrow not $(\mu_a = \mu_b = \mu_c)$

2) Between: $\sum_i n_i (\bar{y}_i - \bar{y})^2$

(SS) over ranges (notation)
go over ranges make them do it?

3) Within: $\sum_{ij} (y_{ij} - \bar{y}_i)^2$

y	:	:	:
\downarrow	A	B	C

"test statistic" $\Rightarrow F = \frac{\text{Between} / (I-1)}{\text{Within} / (n-I)}$ # of sets
 "a measure of departure from the null hypothesis" \rightarrow total # of points

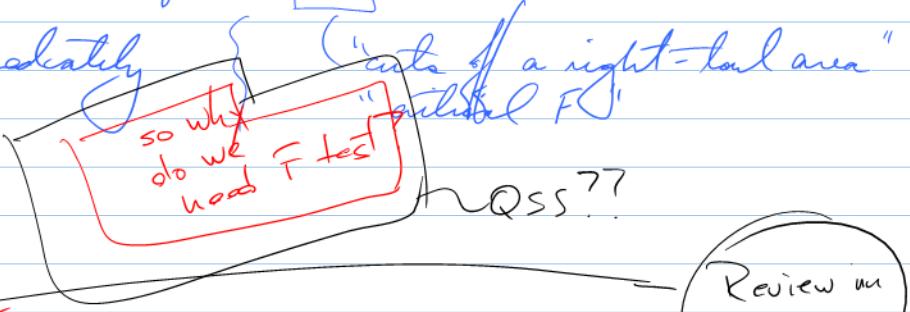
"rejection rule": reject H_0 if $F > F_{\alpha}$

- but get answer immediately from P-value

* "Failing to Reject"

(the null hypothesis)

~~go back earlier when first referenced~~ when e.g., "we will never be able to talk how true the null H is"



Review in SS

maybe write them up for earlier SS??

Session 12: Regression Analysis [5]



- already in JUD, but only w/ total pop; here sample (w/ errors)

see expanded notes,
next page

that we generalize to a population

POSS

x: ind variable or "predictor" I don't like this term
y: dep "criterion" or this one

Auscomes have
Quoted quickly!
again

and yet values on line are

revind's
+ not
intercept

- maybe in 5-min video

"predicted" values or "expected" values

$E(y) = \beta_0 + \beta_1 x$

OLS = ordinary least squares

$\Delta b_0 + \beta_0, b_1 + \beta_1 \rightarrow$ sampling error

now $b_0 + \beta_0$, are the estimates, from the sample, of the population

"the Q of interest [often] focus on

relationships in the pop, "not the sample"

- "specification error" refers to not having all info about what we need to ask to get a valid answer; fun before [?]
- and sampling error

X Belongs above

26 June

* PWW: if only specification error, how well can we predict?
add sampling error \rightarrow can we predict "at all"?

QSS discuss: I argue that in reality, should always be asking "at all"
These QSS in bear question → "significance" meant "to signify,"
not "important"
X = I have X question
should discuss X
discuss X
X of any "hypothesis" analysis { perhaps check OED}

X

- new Q: can we predict dep variable "at all" Q I've been asking all along!
 - can we be sure that $\beta_1 \neq 0$, i.e., is there any relationship??
- notation: $b_0 = \hat{\beta}_0$ and $b_1 = \hat{\beta}_1$
- $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ error term
- but $y = \hat{\beta}_0 + \hat{\beta}_1 + \epsilon$
- Hypothesis Testing: $\beta_1 = 0 \equiv H_0$
- "permits an inference about a pop based on a sample" \Rightarrow "never if the alternative is true" *only if null can be rejected*
- $\Rightarrow b_0, \hat{\beta}_0 \}$ "parameter estimates" vs β_0, β_1 , which indeed are "parameters"
- $b_1, \hat{\beta}_1 \}$ but he's droppin' later "true" says
- H_0 research hypothesis
- H_1 alternative like this!
- "standard error of the slope" $\equiv \sqrt{s_{b_1}}$ from many slopes; but one works well!
- "under certain assumptions," remember, distribution of $\frac{b_1}{s_{b_1}} \rightarrow t\text{-dist}$ with $n-p$ deg of freedom) \{ (=2 parameters number of d.f. in the model: slope + intercept)
- Critical values of $t = \frac{b_1}{s_{b_1}}$ is a value "so extreme that it would be considered very unlikely to occur if the null hypothesis is true." * compare to earlier hypothesis testing again $\alpha = 0.05 \rightarrow$ Type I Error Rate $\Rightarrow 2.5\%$ on each end
- "classical" CONFIDENCIAL But can I see it better from TINV when I think of the dist and critical t value. get st. to draw and verify numbers
- is the p-value less than α ? \rightarrow reject H_0
- [ct] possibly refer st. here first
- p value a percentage vs t value which b_1 measured in std. fractional areas under the tails but probably earlier, when first encountering p values

Two Step Rule for directional alternative hypotheses

$$H_0: \beta_1 \leq 0, H_a: \beta_1 > 0$$

- 1) divide p-value in half
- 2) directionally consistent?

[d]

P vs. t
5 minutes somewhere?
explained earlier?



Session 13: Multiple Regression

25 June 12 topics in intro! in 63 min $\rightarrow \approx 5.25$ min each;
 now (vs. in JSD) samples differ
 we don't specify perfectly
 Sampling + specification errors
 $e_i = y_i - \hat{y}_i$
 was this ever defined previously?
 - residual value

* CONSIDER Tool Pak
from the get-go *

QSS

is that enough

What else is going on?
Look at the data

$$\text{RMS} = \left[\frac{\sum e_i^2}{n-2} \right]^{\frac{1}{2}} = \text{root mean square error}$$

standard error of the estimate
 can go to ∞

$$\left\{ \begin{array}{l} \text{standard deviation of the residuals} \end{array} \right.$$

Coefficient of Determination (R^2) = the proportion of reduction in the error

between y and \hat{y}

$$1 - \frac{\sum e_i^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{\sum y_i^2 - 2\bar{y}\hat{y} + \hat{y}^2}{\sum y_i^2 - 2\bar{y}y + \bar{y}^2}$$

26 Jun

measuring "goodness of fit"

$$= 1 - \frac{\sum 1 - 2\hat{y}_i/y_i + (\hat{y}_i)^2}{\sum 1 - 2\bar{y}_i + (\bar{y}_i)^2}$$

if all $\hat{y}_i = y_i$ numerator goes to zero

Multiple Regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{p-1} x_{p-1} + e$$

partial slope coefficients no. of parameters

β_i = slope 'tween y and x_i "holding constant"

or "controlling for"

plus "additivity" assumption \Rightarrow slope doesn't change all other x_j if x_i

e.g., $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$, β_1 holds "at all levels of all the other ind. variables"
 i.e., find β_1 ; same as x_2 varies

$\Rightarrow \beta_1$ shows how y changes when x_1 changes given that x_2 doesn't change QSS?

Unique vs Shared Reduction in Error (Venn Diagrams)

R^2 : Model Coefficient of Determination \rightarrow shared plus unique gives reduction in error by all x 's

Squared Semipartial Correlation Coefficient: unique std. by each x_i .
[not calculated in this course]

N.B. if $b_i (\beta_i) \approx 0$, small unique relationship w/ y

$$\text{Root MSE} = \sqrt{\frac{\sum e_i^2}{n-p}}$$

1 entire shaded area 2 number of parameters in model ($\stackrel{?}{=}$ single regression)
(in Venn Diagrams)

Standard Error of the Slope from our sample is s_{bk}
 \leq estimates s.d. of b_k on repeated sampling

$\frac{b_k}{s_{bk}}$ \approx normal \rightarrow actually t w/ n-p d.o.f.

again, critical value of t; typically $\alpha = 0.05$
again, p values and t w/ no good explanation
of why they are equivalent

- make sure to do first time t introduced?

again, if $p < \alpha$ (unidirectional), reject H_0
- if not, offers no unique contribution

Global F Test - tests null that all $\beta_s = 0$ ($\stackrel{?}{=} R^2 = 0$)
larger F, more likely null to be rejected

"significance F" is p-value associated w/ F

* if you want single best β , run multiple simple regressions,
'cause multiple regression tells you about unique contribution
of each β

* Two Step Conditions, same as before

* all above assumed ADDITIVITY



Session 14: Regression Assumptions [7]

↑ n.b., not multiple, but any

QSS -
on all?

1) Normally Distributed Errors → average error is zero
assumed for the population AND the sample
if not, t calc no good

2) Homoskedasticity / previously: equal variances in two populations
variance of observations around values are same
now, in regression (i.e., diff x values)

3) Linearity! "frequently unwanted"

says can test, but difficult

⇒ can't look!! (Anscombe again)

Today:

"more substantial 4) Uncorrelated Error Terms ; Serial Uncorrelation Assumption
discussion"
(among different observations)

5) Additivity - slope for x, doesn't depend on x
(only in multiple regression)

→ Data Analysis Course looks at these in detail

here, relationship between y and x, say, is the slope
→ that doesn't change (parallel lines)

- pretty much same answer
as last week

but now addition ⇒ only 1 p value!

asks: "is it reasonable?" → "in many cases, no!"

⇒ we assume it because it is "convenient" ⇒ "what a silly reason!"

⇒ go to an Interactive Model

⇒ add $x_1 \cdot x_2$, = "the interaction term"

- b₃ relationships "much more complicated"

→ now illustrates relationship of y when all other x_i = zero!

- include only interaction terms that "make sense" if i

$$y = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 \cdot x_2$$

$$= (\beta_1 + \beta_3 x_2) x_1 + \beta_2 x_2$$

(slope when holding x₂ fixed)

QSS - derive

Why not teach this from the start

back to #4, serially correlated error: $y = 2x + u \equiv \text{Actual}$
 if $u = u(x)$: great example at 30° 50'
 can get completely mis-leading line } $y = -3 + 3x + u$ always ≥ 0 !

\Rightarrow inflated t or F statistic $\Rightarrow s_u^2 = \text{zero!}$

\rightarrow greater probability of Type I error

? serial correlation of data \equiv auto correlation

another example (u_t correlated w/ u_{t+1}): defense spending!

First Order Auto Regressive Process

$$y(t) = b_0 + b_1 t + u_t \quad u_t = \underbrace{\sum_{j=1}^{p-1} u_{t-j}}_{\text{first-order correlation}} + v_t \quad \begin{cases} \text{"white noise"} \\ \text{time stochastic error} \end{cases}$$

$\Rightarrow \langle b_0 \rangle, \langle b_1 \rangle$ unaffected by first order correlation

but variance is ("sometimes greatly")

$\sigma^2 > 0$ more serious (magnification of effect)

$\rightarrow t$ value too large \rightarrow reject H_0
 P " " small $\left\{ \begin{array}{l} \text{too other} \\ \text{(Type I)} \end{array} \right.$

- assuming first-order: go back one time period
 "auto regressive process"
 to capture auto correlation (all of the effect)

- Conventional TESTS for Autocorrelation ? only with time series?

- two basic types: 1. distribution-free tests

a) Geary Test: count sign changes
 AGAIN, look! in residuals - not too many, not too few

b) Durbin-Watson (d-statistic)
 Test

$$d = \frac{\sum_{t=1}^n (\hat{e}_t - \hat{e}_{t-1})^2}{\sum_{t=1}^n \hat{e}_t^2} \quad \text{where } \hat{e}_t \text{ is OLS regression residual}$$

\rightarrow positive autocorrelation, ratio small
 negative

larger

$$\rho = 0 \rightarrow d = 2$$

no

$$\rho = -1 \rightarrow d = 4$$

negative

$$\rho = 1 \rightarrow d = 0 : \text{strong positive correlation}$$

"REGIONS OF SIGNIFICANCE" from 0 to 4

"reject" vs. "uncertainty"

$d \leq 1.5$ positive
 $d \geq 2.5$ negative

\triangleleft correlation suspected
 "suspected"

Six-Step Process in Excel!

last
page

$$\boxed{X \quad Y} = \text{sumXYZ } (C2:C50, C1:C49) / \text{sumSA } (C1:C50)$$

check Google Sheets! ✓

- if n.g., Excel n.g.! → SAEs

✗ "Retrospective Overview of the Course"