

6/6/11

Lecture #9

Announcements

- Tue 8, 9 no class
talk about holiday
- Midterm review session
8-10pm here
- office hrs today & tomorrow
- 2e, & unassigned

P188-188 Log transform (do demo) same as p184!

Plan

- Log transform + demo
- Sampling ch 13
- Cont. r.v.'s *Adams goes*

Height	$\ln(H)$
64"	4.159
68"	4.229
70"	4.249
71"	4.263
81"	4.394



$\bar{X} = 70.8$	$\bar{Y} = 4.257$	$e^{\bar{Y}} = 70.6$	X
Med = 70"	Med = 4.249	$e^{\text{Med}} = 70$	✓
Boxile = 71"	Boxile = 4.263	$e^{\text{Boxile}} = 71$	✓

log is a "quantile-preserving" transformation

Why log transform? Residuals will get to zero.

Ch 13: Inference / Summary

Data: realizations of r.v.'s. Where do these r.v.'s come from? How should we "harvest" the realizations? This is called "Sampling".



$N = \infty$ "N"

for all examples in this class, we're going to forget about finite N

Temp scores show up

Population: all possible items, people, cars, etc.
Sample: the proper subset we collect data about.

In general, a "representative" sample reflects the population. A non-representative sample is "biased".

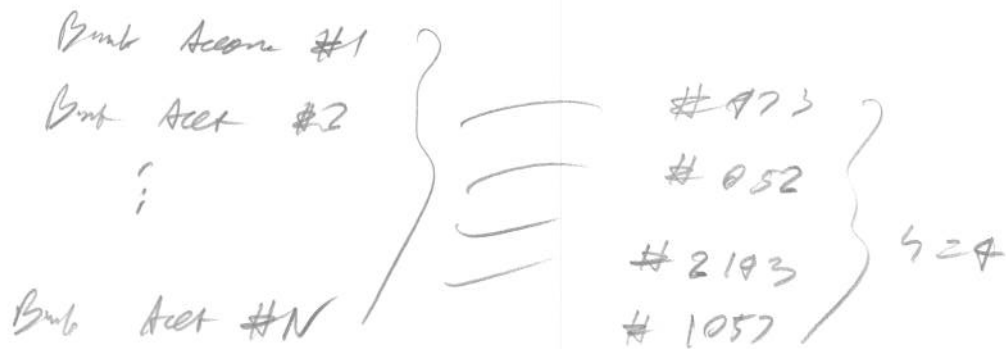
representative:

biased:

Randomization: pick members from population at random: simple means

No matter size of pop., sample does not need to change.

Ensures it's every.



members on avg the population.

Harder: Sample from all unemployed people (macro)

Who is unemployed? Definition... Need to study

Randomization: representative, no bias... we can use sample statistics from the sample to estimate the pop parameters of the population, e.g. \rightarrow book says \bar{Y} ... not sure why...

\bar{X} estimate in
 S^2 estimate of

Inference

p311
figure r, p, b, β

most powerful tool is STATISTICS.

⇒ much better to have small unbiased sample than large biased sample
It's not the "n" that matters... it's how you got the n.
once you did it right
Given the more n, the more accurate

Randomness forcibly:

$$\text{Pop} = \{X_1, X_2, \dots, X_N\}$$

$$\text{Smp} \subset \text{Pop} \quad \text{s.t.} \quad |\text{Smp}| = n$$

If each sample is equally likely, then it is considered a
"Simple random sample" (SRS), the gold standard.

How to get an SRS? Sampling frame: the list of items in the Pop.
Randomization: take random #'s and scan

①	Pop	Rand #		②
	X_1	0.85	← # born 0 and 1	Order ..., take <u>first</u> n
	X_2	0.87		
	X_3	0.11		
	\vdots			
	X_N	0.59		

Systemic sampling: get every 10^{th} person on an telephone list (if you don't have a computer or it's too poor)
You have to make sure there is not correlation with whatever you're measuring.

Simple Random

We flip a coin 10 times. We expect 5 to be H, but it could have been 4 H or 6 H or 3 H or 7 H...

Why? Data is random because it comes from r.v.'s.
Once we get to Hypothesis testing, we have to know we are making decisions from random data - think!

→ Every SRS will have different sample statistics.

\bar{X} : the r.v., \bar{x} : the realization for the random variable
huge variance

16

From now on, for everything else we use in the class,
we assume SRS.... If we don't have an SRS,
Everything we do is WRONG.

More Complex methods

① Stratified Random sample: Split sample frame into
homogenous groups called strata. p.314
has rather good example.

Imagine hotel: 99% tourists 1% business. If you
do an SRS, you won't learn anything about the
business guests, so you take a larger sample from
that strata. This is not an SRS but can be
adjusted if you want to account for this...
Overrepresented? Who?

A type of Stratified sample is Cluster sample:

The geographical units e.g. census tracts are the strata.

Ex 13.2 is a good example: inference estimation

Census: Sample all N ... costly, impractical, expensive

How to select a sample?

- a) Voluntary response: send a bunch of surveys out...
where sends them back call as SRS. Problem?
- b) Convenience Sampling: ask whom is around...
- c) Selection bias: selection people in some groups be too big or
due to early or make judgments about a parameter

What's on the modern?

Set up: random explicitly...

Counting: perm, comb methods

prob: definition, axiom, rules, application,
trees, ind / dep

Ch 1-4 (no p 69)

Ch 7.2 - ch 9

Ch 10.1, 3, 5, 11

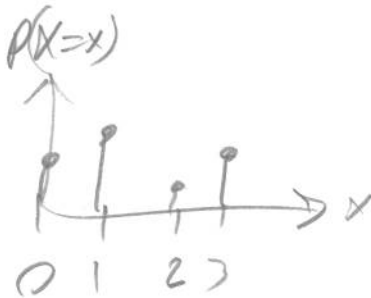
Ch 12

Conditional prob, Bayes Rule

r.v.'s: exp, un, sd, application, sum, constant, binomial, poisson,
data: basic stats, multivariate distribution, data types, special

Discrete r.v.'s: outcome: finite or countably finite # of choices, $\{1, 2, 3, \dots\}$
Cont. r.v.'s: $\{0, 1\}$, $\{0, 1, 2, 3, 4, 5\}$, $\{0, 1, 2, \dots\}$, $\{-1, 1, 1000, 10000\}$ etc.
 all interval / discrete / fractional values acceptable.

→ here PMF is:



→ here density function:



to get probs out of PMF's you sum:

$$P(1 \leq X \leq 4) = \sum_{i=1}^4 P(X=i)$$

to get probs out of PDF's you integrate

$$P(1 \leq X \leq 4) = \int_1^4 f_X(x) dx$$

For example... how wait 10 min



so-called uniform density $X \sim U(0, 10)$

What is prob you wait between 2 and 7 min?

$$P(2 \leq X \leq 7) = \int_2^7 (1) dx = \left[x \right]_2^7 = 7 - 2 = \boxed{\frac{1}{2}}$$

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

More special const. density:

$$f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \quad \text{has mean}$$

