# 36-303: Sampling, Surveys and Society

Statistics of Sampling I Brian Junker 132E Baker Hall brian@stat.cmu.edu

27 January 2010

#### Outline / Announcements

- Project Proposals (Team Project Part I.1)
- Team Project Part I.2 Due Next Tues
  - □ I will email detailed feedback tonight or tomorrow
  - Revise A,B,C, and add D,E,F,G for <u>each</u> project proposal you made
- Statistics of Surveys
  - Part I of an occasional series in the class!
  - Partial review of basic tools
  - Examples related to surveys
  - □ Foreshadowing: Survey Statistics is Different!
- Review: Lohr's Appendix B

#### Handouts

- Appendix B of Lohr (review of probability)
- Lecture Notes

27 January 2010

#### Project Proposals

- <u>I looked at them all</u>; I will email feedback to each team later today.
- **<u>Grades</u>** (50/project x 2 projects = 100 pts)
  - □ A: Is this interesting? 20 pts
  - B: General questions/research questions 20 pts
  - C: One article with description from each team member
     10 pts
- <u>Revise everything</u> especially the parts where you got less than full credit!
- Each team proposed at least one doable project!
  - The project we decide for your team may or may not be the high-scoring project! Depends on <u>feasibility</u>, <u>my interest</u>, etc.

1

#### Team Project Part I.2 Due Next Tues

- The projects should to be interesting enough to make an impact (what can someone do about it?)
- I will email detailed feedback tonight or tomorrow
- For <u>each</u> project you proposed:
  - Revise A, B, C: Interesting topic? General research questions? Articles about past research in the area?
  - Add D, E, F, G: Target population? Sampling Frame? Mode of Data Collection? Major Variables?

27 January 2010

#### Pointers for I.2

- F. Mode of Data Collection How will you get the data?
- Invite people to website with online SAQ, using email, postcards, etc.
- Approach people on the street/sidewalk/etc. and use P&P SAQ, CAPI, etc.
- □ Go to a certain intersection at a certain time and observe buses, people, accidents, or other events of interest.
- Go to a school and interview some/all students
- Give a sense of how many intersections, times, schools, students, etc. might be needed to "represent" the population.
- <u>G. Variables to Measure</u> List (and define) two to five variables that you must measure to have a successful survey.

#### Pointers for I.2

- <u>E. Target population</u> What are the individual units that give you information?
  - students? buses? faculty members? times of day? locations? events ("the bus is late" or "10 students walked by", etc.)
- <u>D. Sampling Frame</u> In most (but not all) cases there will be a real or hypothetical list of units that you could sample from. E.g.:
  - Numbers in the phone book (which one? or maybe random digit dialling? which exchanges? etc)
  - Email addresses in C-Book
  - In some cases there will be no natural sampling frame. E.g.:
  - Interview people as they pass by the fence
  - Wait for instances of late buses
  - In these cases give a very specific description of <u>what</u> <u>kinds of units</u> you will be looking for, and <u>how you</u> <u>will find them.</u>

27 January 2010

#### Statistics of Surveys

- Survey Statistics is different from other kinds of Statistics
  - □ Sampling from a finite population is *different*
  - Design features (stratification, clustering, weights) increase information at the cost of more complex analysis
- We will get there, in occasional smallish steps
  - Today:
    - Partial Review of Probability Tools
    - Application: Sample Size Calculations
    - Application: Randomized Response
  - Future:
    - Urn models
    - What is random about finite population sampling?
    - Accounting for complex survey designs

#### Partial Review of Probability Tools

- Discrete Random Variables
- Expected Value, Mean, Variance
- More than One Random Variable
  - Covariances, Independence, Linear Combinations, Normal Approximation (CLT)
  - Application: Sample Size Calculations
- Conditioning
  - Conditional Probability, Conditional Distribution, Conditional Expectation
  - Application: Randomized Response

27 January 2010

#### Discrete Random Variable

- A <u>discrete</u> random variable X has a sample space that you can "count" (1, 2, 3, …)
  - □ Toss a die, let *X* be the side that comes "up"
  - Toss a coin until "heads" comes up, let X be the number of "Tails" until first "Heads"
  - Spin a spinner, let *X* be the exact angle in degrees at which the spinner comes to rest.
- A <u>continuous</u> random variable X has a sample space that includes a continuous interval (so there are uncountably many outcomes)
  - □ Which of the above X's is discrete, which is continuous?

27 January 2010

#### Discrete Random Variable

- For us, X usually has a <u>finite sample space</u>
  - $\hfill X$  can take on only the values  $x_1,\,x_2,\,...,\,x_K$  , with probability  $p_1,\,p_2,\,...,\,p_K$
- Examples:
  - □ Biased coin, X=1 for "Heads", 0 for "Tails"
    - (this is a \_\_\_\_\_\_ random variable!)
    - *P*[X=1] = *p*, *P*[X=0] = 1-*p*
  - □ Flip a coin n times, let *X* be the number of "Heads"
    - (this is a \_\_\_\_\_\_ random variable!)
    - P[X=k] = \_\_\_\_\_, k=0, 1, 2, ..., n
  - Consider a population of 1,000 adults, and let x<sub>k</sub> be each adult's annual income, k=1, ..., 1000. Pick one adult at random and let X be that person's income.
    - *P*[*X*=*x<sub>k</sub>*] = \_\_\_\_\_, *k*=1, 2, ..., 1000

#### Expected Value, Mean, Variance

Let X be a discrete random variable taking on the values x<sub>1</sub>, ..., x<sub>K</sub> with probabilities p<sub>1</sub>, ..., p<sub>K</sub>:
 The probabilities *must* add to 1:

$$\sum_{i=1}^{K} p_i = 1,$$

• The <u>mean</u> of X is defined to be

$$\mu_X = E[X] = \sum_{i=1}^{K} x_i P(X = x_i) = \sum_{i=1}^{K} x_i p_i$$

• The <u>variance</u> of X is defined to be

$$\sigma_X^2 = Var[X] = E[(X - \mu_X)^2] = \sum_{i=1}^K (x_i - E[X])^2 P(X = x_i) = \sum_{i=1}^K (x_i - \mu_X)^2 p_i.$$

• More generally, for any function g(x), the <u>expected value</u> of g(X) is

$$E[g(X)] = \sum_{x} g(x)P(X = x).$$

9

27 January 2010

#### Expected Value Example

 Let X be a Bernoulli random variable, P[X=1]=.2, and suppose I pay you \$50 if X=1 and you pay me \$10 if X=0. What is the expected value of your income?

g(x) = 50 if x = 1, and g(x) = -10 if x = 0.

$$E[g(X)] = 50 \times p - 10 \times (1 - p)$$
  
= 50(0.2) - 10(0.8)  
= 2  
$$Var(g(X)) = (50 - 2)^{2}(0.2) + (-10 - 2)^{2}(0.8)$$
  
= 2304(0.2) + 144(0.8)  
= 576  
$$SD(g(X)) = \sqrt{576} = 24$$

13

15

27 January 2010

#### Covariance & Independence

- Recall that  $Var(X) = E[(X-\mu_X)^2]$
- Similarly,  $Cov(X,Y) = E[(X-\mu_X)(Y-\mu_Y)]$

$$Cov(X,Y) = \frac{1}{4} \left\{ (1-2.5)(2-6) + (2-2.5)(8-6) + (3-2.5)(6-6) + (4-2.5)(8-6) \right\}$$
  
= 2

If X and Y are independent, Cov(X,Y) = 0

$$Cov(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$
  
=  $E[(X - \mu_X)]E[(Y - \mu_Y)] = 0 \cdot 0 = 0$ 

#### More Than One Random Variable

<i>x</i>	y	xy  P[X = x, Y = y]	Note that
1	2	$2 \frac{1}{4}$	
2	8	$16 \frac{1}{4}$	$E[X]E[Y] = (6)(2.5) = 15 \neq 17 = E[XY]$
4	8	$\frac{32}{10}$ $\frac{1}{4}$	
3	6	$18 \frac{1}{4}$	thus X and Y cannot be independent.
E[X] $E[Y]$ $E[XY]$	=	$\frac{1}{4}(1+2+4+3) = 2.5$ $\frac{1}{4}(2+8+8+6) = 6$ $\frac{1}{4}(2+16+32+18) = 17$	More generally <i>X</i> and <i>Y</i> are <i>independent</i> if and only if P[X = x, Y = y] = P[X = x]P[Y = y] for all <i>x</i> and <i>y</i> .

27 January 2010

#### Linear Combinations

• Exercise: Use the definitions so far to show

E[aX + bY + c] = aE[X] + bE[Y] + c

 <u>Exercise</u>: Use this fact to show that for any set of random variables X<sub>1</sub>, X<sub>2</sub>, ... X<sub>n</sub> that all have the same mean μ,

$$E\left[\overline{X}\right] = E\left[\frac{1}{n}\sum_{i=1}^{n}X_{i}\right] = \mu_{\text{(Definition of }\overline{X})}$$

This is the part to show!)

#### Mean and Variance of Sample Average

• Let  $X_1, ..., X_n$  all have the same mean  $\mu$ , and let  $\overline{\mathbf{x}} = 1 \sum_{n=1}^{n} \mathbf{x}$ 

 $\overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$ 

• We know  $E[\overline{X}] = \mu$ , what about  $Var(\overline{X})$ ?

Use the definitions to show:

 $Var(aX + bY + c) = a^2Var(X) + 2abCov(X,Y) + b^2Var(Y)$ 

We use this on the next page to work out  $Var(\overline{X})$ .

27 January 2010

#### Central Limit Theorem

 We have shown: If X<sub>1</sub>, ..., X<sub>n</sub> are independent, identically distributed (iid) with E[X<sub>i</sub>]=μ and Var(X<sub>i</sub>)=σ<sup>2</sup>, then

$$E[\overline{X}] = \mu$$
,  $Var(\overline{X}) = \frac{\sigma^2}{n}$ 

The Central Limit Theorem then tells us

$$\frac{X-\mu}{\sigma/\sqrt{n}} \sim N(0,1)$$

•  $\sigma$  is the SD of X<sub>i</sub>;  $\sigma/\sqrt{n}$  is the SE of  $\overline{X}$ 

### Mean and Variance of Sample Average

From

$$ar(aX + bY + c) = a^{2}Var(X) + 2abCov(X, Y) + b^{2}Var(Y)$$

we can calulate

$$\sqrt{ar}\left[\frac{1}{n}(X_1 + X_2)\right] = \frac{1}{n^2} \left( Var(X_1) + 2Cov(X_1, X_2) + Var(X_2) \right)$$

and applying this to n terms instead of 2 terms (induction!), we get the following mess

$$Var\left[\frac{1}{n}\sum_{i=1}^{n}X_{i}\right] = \frac{1}{n^{2}}\left\{\sum_{i=1}^{n}Var(X_{i}) + 2\sum_{i=1}^{n}\sum_{j=1}^{i-1}Cov(X_{i},X_{j})\right\}$$

We now assume  $X_1, X_2, ..., X_n$  have the same mean  $\mu$ , the same variance  $\sigma^2$ , and covariance  $Cov(X_i, X_j) = 0$  whenever  $i \neq j$ . Then the "mess" reduces to the more familiar:

$$Var(\overline{X}) = \frac{1}{n^2} \left\{ n\sigma^2 + 2 \cdot \binom{n}{2} \cdot 0 \right\} = \frac{1}{n}\sigma^2$$

27 January 2010

17

19

#### Application: Sample Size Calculation

- Let X<sub>1</sub>, ..., X<sub>n</sub> be an iid sample of people's heights, with a common mean μ=5.75 ft and SD σ=0.5ft.
- Then  $E[\overline{X}]$  = 5.75, with SE 0.5/ $\sqrt{n}$
- CLT: Approx 95% confidence interval for  $\mu$ :  $\left(\overline{X} - (1.96)(0.5)/\sqrt{n}, \overline{X} + (1.96)(0.5)/\sqrt{n}\right)$
- How large n to have 95% confidence that X is within 0.1 of µ?

• Roughly, need 0.1 >  $1/\sqrt{n}$  or n > 100.

#### Foreshadowing: Survey Statistics is Different!

- In real <u>Survey Sampling</u> work, Cov(X<sub>i</sub>,X<sub>j</sub>) is usually not zero!
- Hence

 $E[\overline{X}] = \mu$ 

but

27 January 2010

$$Var(\overline{X}) \neq \sigma^2/r$$

- The CLT is not quite true, as stated, either!
- But the basic CLT calculation is often a reasonable "crude guess"...

#### Conditioning

• The conditional probability of event A, given event B, is

$$P[A|B] = \frac{P[A \cap B]}{P[B]}$$

It is often useful to write this as a formula for  $P[A \cap B]$ :

$$P[A \cap B] = P[A|B]P[B]$$

• The <u>conditional distribution</u> of X given Y = y is

$$P[X = x|Y = y] = \frac{P[X = x, Y = y]}{P[Y = y]} \quad [comma means "and"!]$$

• The *conditional expected value* of X given Y = y is the expected value with respect to the conditional distribution:

$$E[X|Y = y] = \sum_{x} xP[X = x|Y = y]$$

27 January 2010

#### Conditioning

				E[X] = 2.5
<u>x</u>	у	ху	P[X=x,Y=y]	$V_{at}(X) = \frac{1}{-1}[(1-25)^2 + (2-25)^2]$
1	2	2	$\frac{1}{4}$	$4^{1(1-2.5)+(2-2.5)}$
2	8	1 <b>6</b>	<u>1</u> 4	$+(3-2.5)^2+(4-2.5)^2]$
4	8	32	$\frac{1}{4}$	- 1.25
3	6	18	1	- 1,23
-	•	10	4	
			P[Y - 7 V -	$E[X Y=8] = \frac{1}{2}(2+4) = 3$
<i>P</i> [ <i>X</i> =	= 2 }	? = 8]	$= \frac{P[X - 2, Y - 2]}{P[Y = 8]}$	$\frac{3}{2} Var(X Y=8) = \frac{1}{2}[(2-3)^2 + (4-3)^2]$
			$= \frac{1/4}{1/2} = \frac{1}{2}$	= 1
<i>P</i> [X =	= 4 ]	7 = 8]	$= \cdots = \frac{1}{2}$	<b>Exercise:</b> Show that if X and Y are independent, then $E[X Y = y] = E[X]$ , for any y.

#### Application: Randomized Response

- "Flip a coin, but don't tell me whether it's heads or tails.
  - "If heads, answer truthfully: have you ever cheated in a CMU class?
  - "If tails, answer truthfully: is the last digit of your SSN odd?"
- Let p=P[Heads],  $\pi$ =P[Cheat],  $\lambda$ =P[Yes]. Then
  - $\lambda = P[Yes \cap Heads] + P[Yes \cap Tails]$ 
    - = P[Yes|Heads]P[Heads] + P[Yes|Tails]P[Tails]
    - $= \pi \cdot p + (1/2) \cdot (1-p)$

Therefore

$$\pi = \frac{\lambda - (1/2) \cdot (1-p)}{p}$$

21

27 January 2010

#### Application: Randomized Response

 $\pi = \frac{\lambda - \frac{1}{2}(1-p)}{p}$ 

Suppose the coin is fair  $(p = \frac{1}{2})$  and in our survey we get a fraction  $\hat{\lambda}$  of people answering "yes". Then

$$\hat{\pi} = 2(\hat{\lambda} - 1/4)$$

$$E[\hat{\pi}] = 2(E[\hat{\lambda}] - 1/4)$$

$$= 2(\lambda - 1/4) = \pi (Exercise!)$$

So  $\hat{\pi}$  is an <u>unbiased</u> estimator of  $\pi$ ; and

$$Var(\hat{\pi}) = Var[2(\hat{\lambda} - 1/4) \\ = 4Var(\hat{\lambda})$$

so  $Var(\hat{\pi})$  is inflated, relative to  $Var(\hat{\lambda})$ :  $\hat{\pi}$  is statistically inefficient. <u>Exercise</u>: The closer p = P[Answer Cheating Question] is to 1, the closer  $Var(\hat{\pi})$  is to  $Var(\hat{\lambda})$ .

27 January 2010

#### Review

- Feedback on Project Proposals
- Team Project part I.2 (target pop, frame, mode of data collection) Due Next Tuesday
   HW02 due next Tues also!
- Statistics of Surveys (Part I of Occasional Series)
- Read Lohr Appx B (handout today)

## Foreshadowing: Survey Statistics is Different!

 In a regular statistics course we would go on to say \$\overline{\lambda} = Y/n\$ where Y is the number of "Yes"'s among a sample of size n.

#### Therefore

- $\Box E[\hat{\lambda}] = \lambda$ , the true proportion of "Yes"'s.
- $SE(\hat{\lambda}) = \sqrt{\lambda(1-\lambda)/n}$ , because Y is a binomial random variable.
- In Survey Sampling
  - □ The expected value part is OK
  - The variance will be different; Y is not quite binomial!

26

27 January 2010