

TRANSLATING RESEARCH INTO ACTION

# Sampling and Sample Size

Benjamin Olken MIT

### Course Overview

- 1. What is evaluation?
- 2. Measuring impacts (outcomes, indicators)
- 3. Why randomize?
- 4. How to randomize
- 5. Threats and Analysis
- 6. Sampling and sample size
- 7. RCT: Start to Finish
- 8. Cost Effectiveness Analysis and Scaling Up

### Course Overview

- 1. What is evaluation?
- 2. Measuring impacts (outcomes, indicators)
- 3. Why randomize?
- 4. How to randomize
- 5. Threats and Analysis
- 6. Sampling and sample size
- 7. RCT: Start to Finish
- 8. Cost Effectiveness Analysis and Scaling Up

#### What's the average result?

• If you were to roll a die once, what is the "expected result"? (i.e. the average)



## Possible results & probability: 1 die



# Rolling 1 die: possible results & average



#### What's the average result?

• If you were to roll two dice once, what is the expected average of the two dice?



#### Rolling 2 dice: Possible totals & likelihood



# Rolling 2 dice: possible totals 12 possible totals, 36 permutations

		Die 1					
Die 2		2	3	4	5	6	7
		3	4	5	6	7	8
		4	5	6	7	8	9
		5	6	7	8	9	10
		6	7	8	9	10	11
		7	8	9	10	11	12

# Rolling 2 dice: Average score of dice & likelihood



### Outcomes and Permutations

- Putting together permutations, you get:
  - 1. All possible outcomes
  - 2. The likelihood of each of those outcomes
  - Each column represents one possible outcome (average result)
    - Each block within a column represents one possible permutation (to obtain that average)

# Rolling 3 dice: 16 results $3 \rightarrow 18$ , 216 permutations



# Rolling 4 dice: 21 results, 1296 permutations



# Rolling 5 dice: 26 results, 7776 permutations



# Rolling 10 dice: 50 results, >60 million permutations



Looks like a bell curve, or a normal distribution

# Rolling 30 dice: 150 results, $2 \ge 10^{23}$ permutations\*



>95% of all rolls will yield an average between 3 and 4

## Rolling 100 dice: 500 results, 6 x 10<sup>77</sup> permutations



>99% of all rolls will yield an average between 3 and 4

# Rolling dice: 2 lessons

- The more dice you roll, the closer most averages are to the <u>true</u> average (the distribution gets "tighter")
  THE LAW OF LARGE NUMBERS-
- 2. The more dice you roll, the more the distribution of possible averages (the *sampling distribution*) looks like a bell curve (a *normal* distribution) -THE CENTRAL LIMIT THEOREM-

### Which of these is more accurate?



- A. I.
- B. II.
- C. Don't know



#### Accuracy versus Precision



#### Accuracy versus Precision



#### THE basic questions in statistics

• How confident can you be in your results?

•  $\rightarrow$  How big does your sample need to be?

# THAT WAS JUST THE INTRODUCTION

### Outline

- Sampling distributions
  - population distribution
  - sampling distribution
  - law of large numbers/central limit theorem
  - standard deviation and standard error
- Detecting impact

### Outline

- Sampling distributions
  - population distribution
  - sampling distribution
  - law of large numbers/central limit theorem
  - standard deviation and standard error
- Detecting impact

#### Baseline test scores



#### Mean = 26



#### Standard Deviation = 20



#### Let's do an experiment

- Take 1 Random test score from the pile of 16,000 tests
- Write down the value
- Put the test back
- Do these three steps again
- And again
- 8,000 times
- This is like a random sample of 8,000 (*with replacement*)

### What can we say about this sample?



Good, the average of the sample is about 26...

#### But...

- ... I remember that as my sample goes, up, isn't the sampling distribution supposed to turn into a bell curve?
- (Central Limit Theorem)
- Is it that my sample isn't large enough?

# Population v. sampling distribution



This is the distribution of my sample of 8,000 students!

## Outline

- Sampling distributions
  - population distribution
  - sampling distribution
  - law of large numbers/central limit theorem
  - standard deviation and standard error
- Detecting impact

### How do we get from here...



# Draw 10 random students, take the average, plot it: Do this 5 & 10 times.

**Frequency of Means With 5 Samples** 



#### **Frequency of Means With 10 Samples**



## Draw 10 random students: 50 and 100 times

#### **Frequency of Means With 50 Samples**



#### **Frequency of Means with 100 Samples**


# <u>Draws 10 random students</u>: 500 and 1000 times

**Frequency of Means With 500 Samples** 



#### Frequency of Means With 1000 Samples



## Draw 10 Random students

- This is like a sample size of 10
- What happens if we take a sample size of 50?

## $\mathbf{N}=10$

# N = 50

#### 10 8 6 4 2 0 1 3 5 7 9 11 13 15 17 19 21 23 25 27 29 31 33 35 37 39 41 43 45 47 49 51

**Frequency of Means With 5 Samples** 

#### **Frequency of Means With 10 Samples**



**Frequency of Means With 5 Samples** 



#### **Frequency of Means With 10 Samples**



## Draws of 10

# Draws of 50

#### **Frequency of Means With 50 Samples**



#### Frequency of Means with 100 Samples



**Frequency of Means With 50 Samples** 



Frequency of Means With 100 Samples



### Draws of 10

# Draws of 50

Frequency of Means With 500 Samples



#### Frequency of Means With 500 Samples



#### Frequency of Means With 1000 Samples



#### Frequency of Means With 1000 Samples



# Outline

- Sampling distributions
  - population distribution
  - sampling distribution
  - law of large numbers/central limit theorem
  - standard deviation and standard error
- Detecting impact

## Population & sampling distribution: Draw 1 random student (from 8,000)



# Sampling Distribution: Draw 4 random students (N=4)



### Law of Large Numbers : N=9



### Law of Large Numbers: N = 100



### Central Limit Theorem: N=1



• The white line is a theoretical distribution

### Central Limit Theorem : N=4



### Central Limit Theorem : N=9



### Central Limit Theorem : N = 100



# Outline

- Sampling distributions
  - population distribution
  - sampling distribution
  - law of large numbers/central limit theorem
  - standard deviation and standard error
- Detecting impact

# Standard deviation/error

- What's the difference between the standard deviation and the standard error?
- The standard error = the standard deviation of the sampling distributions

### Variance and Standard Deviation

• Variance = 400  

$$\sigma^{2} = \frac{\sum (Observation \, Value \ - Average)^{2}}{N}$$

 Standard Deviation = 20
 σ = √Variance

 Standard Error = <sup>20</sup>/<sub>√N</sub>
 SE = <sup>σ</sup>/<sub>√N</sub>

## Standard Deviation/ Standard Error



## Sample size $\uparrow$ x4, SE $\downarrow$ $\frac{1}{2}$



# Sample size $\uparrow x9$ , SE $\downarrow$ ?



# Sample size $\uparrow$ x100, SE $\downarrow$ ?



# Outline

- Sampling distributions
- Detecting impact
  - significance
  - effect size
  - power
  - baseline and covariates
  - clustering
  - stratification

#### Baseline test scores



# We implement the Balsakhi Program



#### Endline test scores



After the balsakhi programs, these are the endline test scores

# The impact appears to be?

- A. Positive
- B. Negative
- C. No impact
- D. Don't know



### Post-test: control & treatment



• Stop! That was the control group. The treatment group is red.

# Is this impact statistically significant?



A. YesB. NoC. Don't know



# One experiment: 6 points



# One experiment



# Two experiments



### A few more...



### A few more...



# Many more...



### A whole lot more...





• • •
# Running the experiment thousands of times...



By the Central Limit Theorem, these are normally distributed

# Hypothesis testing

- In criminal law, most institutions follow the rule: "innocent until proven guilty"
- The presumption is that the accused is innocent and the burden is on the prosecutor to show guilt
  - The jury or judge starts with the "null hypothesis" that the accused person is innocent
  - The prosecutor has a hypothesis that the accused person is guilty

# Hypothesis testing

- In program evaluation, instead of "presumption of innocence," the rule is: <u>"presumption of insignificance</u>"
- The "<u>Null hypothesis</u>" (<u>H</u><sub>0</sub>) is that there was no (zero) impact of the program
- The burden of proof is on the evaluator to show a significant effect of the program

# Hypothesis testing: conclusions

• If it is very unlikely (less than a 5% probability) that the difference is solely due to chance:

- We "reject our null hypothesis"

• We may now say:

- "our program has a statistically significant impact"

# What is the significance level?

- Type I error: rejecting the null hypothesis even though it is true (false positive)
- Significance level: <u>*The probability*</u> that we will reject the null hypothesis even though it is true

# Hypothesis testing: 95% confidence

		YOU CONCLUDE	
		Effective	No Effect
THE TRUTH	Effective		Type II Error (low power)
	No Effect	Type I Error (5% of the time)	

## What is Power?

- Type II Error: Failing to reject the null hypothesis (concluding there is no difference), when indeed the null hypothesis is false.
- Power: If there is a measureable effect of our intervention (the null hypothesis is false), the probability that we will detect an effect (reject the null hypothesis)

# Before the experiment



• Assume two effects: no effect and treatment effect  $\beta$ 

# Impose significance level of 5%



Anything between lines cannot be distinguished from 0

# Can we distinguish H<sup>\beta</sup> from H0 ?



Shaded area shows % of time we would find H<sup>β</sup> true if it was

# What influences power?

- What are the factors that change the proportion of the research hypothesis that is shaded—i.e. the proportion that falls to the right (or left) of the null hypothesis curve?
- Understanding this helps us design more powerful experiments

# Power: main ingredients

- 1. Effect Size
- 2. Sample Size
- 3. Variance
- 4. Proportion of sample in T vs. C
- 5. Clustering

# Power: main ingredients

- 1. Effect Size
- 2. Sample Size
- 3. Variance
- 4. Proportion of sample in T vs. C
- 5. Clustering

#### Effect Size: 1\*SE



#### Effect Size = $1 \times SE$



# Power: 26% If the true impact was 1\*SE...



The Null Hypothesis would be rejected only 26% of the time

#### Effect Size: 3\*SE



Bigger hypothesized effect size  $\rightarrow$  distributions farther apart

#### Effect size 3\*SE: Power= 91%



Bigger Effect size means more power

What effect size should you use when designing your experiment?

- A. Smallest effect size that is still cost effective
- B. Largest effect sizeyou estimate yourprogram to produce
- C. Both
- D. Neither



### Effect size and take-up

- Let's say we believe the impact on our participants is "3"
- What happens if take up is 1/3?
- Let's show this graphically

#### Effect Size: 3\*SE



Let's say we believe the impact on our participants is "3"

# Take up is 33%. Effect size is 1/3rd



#### Back to: Power = 26%



Take-up is reflected in the effect size

# Power: main ingredients

- 1. Effect Size
- 2. Sample Size
- 3. Variance
- 4. Proportion of sample in T vs. C
- 5. Clustering

# By increasing sample size you increase...



- A. Accuracy
- B. Precision
- C. Both
- D. Neither
- E. Don't know



# Power: Effect size = 1SD, Sample size = N



# Power: Sample size = 4N



#### $\overline{P}$ ower: $64^{\circ}/_{\circ}$



# Power: Sample size = 9



#### Power: 91%



# Power: main ingredients

- 1. Effect Size
- 2. Sample Size
- 3. Variance
- 4. Proportion of sample in T vs. C
- 5. Clustering

# What are typical ways to reduce the underlying variance

- A. Include covariates
- B. Increase the sample
- C. Do a baseline survey
- D. All of the above
- E. A and B
- F. A and C



#### Variance

- There is sometimes very little we can do to reduce the noise
- The underlying variance is what it is
- We can try to "absorb" variance:
  - using a baseline
  - controlling for other variables
    - In practice, controlling for other variables (besides the baseline outcome) buys you very little

# Power: main ingredients

- 1. Effect Size
- 2. Sample Size
- 3. Variance
- 4. Proportion of sample in T vs. C
- 5. Clustering

# Sample split: 50% C, 50% T



Equal split gives distributions that are the same "fatness"

#### Power: 91%


## If it's not 50-50 split?

- What happens to the relative fatness if the split is not 50-50.
- Say 25-75?

## Sample split: 25% C, 75% T



Uneven distributions, not efficient, i.e. less power

#### Power: 83%



#### Allocation to T v C



## Power: main ingredients

- 1. Effect Size
- 2. Sample Size
- 3. Variance
- 4. Proportion of sample in T vs. C
- 5. Clustering

## Clustered design: intuition

- You want to know how close the upcoming national elections will be
- Method 1: Randomly select 50 people from entire Indian population
- Method 2: Randomly select 5 families, and ask ten members of each family their opinion

## Low intra-cluster correlation (Rho)



## HIGH intra-cluster correlation (rho)



All uneducated people live in one village. People with only primary education live in another. College grads live in a third, etc. Rho on education will be..

- A. High
- B. Low
- C. No effect on rho
- D. Don't know



If rho is high, what is a more efficient way of increasing power?

A. Include more clusters in the sample

- B. Include more people in clusters
- C. Both
- D. Don't know



# Testing multiple treatments

5	50 <u>←0.15 SD</u> →	260
Control Group	<b>N N</b>	Balsakhi
0.15 SD	0.0 <b>8.29</b> SD	0.10 SD
$\checkmark$	К Л	$\mathbf{V}$
<b>3</b> 0	00 ←0.10 SD→	100
CAL program		Balsakhi + CAL

