### STAT6061/STAT5008 – Causal Inference

# Part 2-2. Completely Randomized Experiment

### **An-Shun Tai**

<sup>1</sup>Department of Statistics National Cheng Kung University <sup>2</sup>Institute of Statistics and Data Science National Tsing Hua University

# **Key questions for inference in randomized experiments**

How can statistical methods be applied to assess causal effects in randomized experiments?

> What is the causal hypothesis underlying inference in randomized experiments?

> Should prognostic covariates be adjusted in the analysis of randomized experiments?

# **Inference for completely randomized experiments**

*"Experiments should be analyzed as experiments, not as observational studies"* – Freedman (2006, p. 691)

> Methods of causal inference in randomized experiments

	Randomization-Based Inference (Design-Based Inference)	Sampling-Based Inference
Features	<ul> <li>✓ Assumes potential outcomes are fixed for each subject.</li> <li>✓ Treats treatment assignment as random.</li> <li>✓ Used for causal inference based on the experimental design.</li> </ul>	<ul> <li>✓ Assumes treatment assignments are fixed.</li> <li>✓ Treats outcomes as random, considering subjects as a random sample from a larger population.</li> <li>✓ Used for generalizing results beyond the study sample.</li> </ul>
Methods	<ol> <li>Fisher randomization test</li> <li>Neyman repeated sampling inference</li> </ol>	<ul><li>3. Regression methods</li><li>4. Model-based imputation</li></ul>

### **Fisher randomization test**

Randomized Experiment on Nocturnal Cough (Paul et al., 2007)

#### **Study Overview:**

- Investigates effects of buckwheat honey vs. no active treatment
- Focus on nocturnal cough and sleep difficulties in children with upper respiratory tract infections

### **Dataset Details:**

- 72 children participated
- Measured cough frequency and cough severity
- Scale: 0 (not at all) to 6 (extremely severe)

	Potential	outcomes	Observed variables		
Unit	Potential (A=1)	Control (A=0)	Treatment	Observed outcome <i>Y</i>	
1	$Y_1(1) = 3$	$Y_1(0) = ?$	A=1	$Y_1 = 3$	
2	$Y_2(1) = 5$	$Y_2(0) = ?$	A=1	$Y_2 = 5$	
3	$Y_{3}(1) = 0$	$Y_3(0) = ?$	A=1	$Y_3 = 0$	
4	$Y_4(1) = ?$	$Y_4(0) = 4$	A=0	$Y_4 = 4$	
5	$Y_5(1) = ?$	$Y_5(0)=0$	A=0	$Y_5=0$	
6	$Y_6(1) = ?$	$Y_6(0) = 1$	A=0	$Y_{6} = 1$	

#### Table. Cough frequency for the first six units from the honey study

#### **Question:**

What types of null hypotheses can be considered, and how can corresponding inferences be made?

$$H_0: \mathbb{E}(Y(1)) = \mathbb{E}(Y(0))$$

or

 $H_0:Y(1)=Y(0)$ 

# Fisher randomization test (1): Sharp null hypothesis

1. Specify the null hypothesis

Sharp null hypothesis of no treatment/causal effect (Rubin, 1980)

 $H_{0F}: Y_i(0) = Y_i(1)$  for all units i = 1, ..., N

- Under the sharp null hypothesis, it allows each unit to be assigned a hypothetical value (observed outcome) for its unobserved potential outcome.

	Potential	outcomes	Obser	rved variables	Potential outcomes under the sharp null hypnosis	
Unit	Potential (A=1)	Control (A=0)	Treatment	Observed outcome	Potential (A=1)	Control (A=0)
1	$Y_1(1) = 3$	$Y_1(0) = ?$	A=1	$Y_1 = 3$	$Y_1(1) = 3$	$Y_1(0) = 3$
2	$Y_2(1) = 5$	$Y_2(0) = ?$	A=1	$Y_2 = 5$	$Y_2(1) = 5$	$Y_2(0) = 5$
3	$Y_3(1) = 0$	$Y_3(0) = ?$	A=1	$Y_3=0$	$Y_{3}(1) = 0$	$Y_3(0) = 0$
4	$Y_4(1) = ?$	$Y_4(0) = 4$	A=0	$Y_4 = 4$	$Y_4(1) = 4$	$Y_4(0) = 4$
5	$Y_5(1) = ?$	$Y_5(0)=0$	A=0	$Y_5=0$	$Y_5(1)=0$	$Y_5(0) = 0$
6	$Y_6(1) = ?$	$Y_{6}(0) = 1$	A=0	$Y_{6} = 1$	$Y_{6}(1) = 1$	$Y_{6}(0) = 1$

Table. Cough frequency for the first six units from the honey study

# Fisher randomization test (2): test statistic

- 2. Choose a test statistic
- Test statistic  $T = T(\widetilde{A}, \widetilde{Y})$ : a real-valued function of observed outcomes and treatment assignments
- Any test statistic for quantifying the contrast between the treatment and control groups can be used.

**Difference-in-means**  
$$T_1(\widetilde{A}, \widetilde{Y}) = \frac{\sum_{i=1}^N A_i Y_i}{N_1} - \frac{\sum_{i=1}^N (1 - A_i) Y_i}{N_0}$$

where  $N_1 = \sum_{i=1}^{N} A_i$  and  $N_0 = \sum_{i=1}^{N} (1 - A_i)$ 

#### Wilcoxon rank sum

$$T_2(\widetilde{A}, \widetilde{Y}) = \frac{\sum_{i=1}^N A_i R_i}{N_1} - \frac{\sum_{i=1}^N (1 - A_i) R_i}{N_0}$$
  
where  $R_i = \sum_{j=1}^N I(Y_j \le Y_i)$ 

Table. Cough frequency for the first six units from the honey study

	Potential	outcomes	Obser	rved variables	Potential outcomes under the sharp null hypnosis	
Unit	Potential (A=1)	Control (A=0)	Treatment	Observed outcome	Potential (A=1)	Control (A=0)
1	$Y_1(1) = 3$	$Y_1(0) = ?$	A=1	$Y_1 = 3$	$Y_1(1) = 3$	$Y_1(0) = 3$
2	$Y_2(1) = 5$	$Y_2(0) = ?$	A=1	$Y_2 = 5$	$Y_2(1) = 5$	$Y_2(0) = 5$
3	$Y_3(1)=0$	$Y_3(0) = ?$	A=1	$Y_3=0$	$Y_3(1) = 0$	$Y_3(0)=0$
4	$Y_4(1) = ?$	$Y_4(0)=4$	A=0	$Y_4 = 4$	$Y_4(1) = 4$	$Y_4(0) = 4$
5	$Y_5(1) = ?$	$Y_5(0)=0$	A=0	$Y_5=0$	$Y_5(1)=0$	$Y_5(0)=0$
6	$Y_6(1) = ?$	$Y_{6}(0) = 1$	A=0	$Y_{6} = 1$	$Y_{6}(1) = 1$	$Y_{6}(0) = 1$

# Fisher randomization test (3): randomization distribution

treatment assignments				nmei	nts	Statistics		
$A_1$	$A_2$	$A_3$	$A_4$	$A_5$	$A_6$	<b>Difference-in-means</b>	Wilcoxon rank sum	
0	0	0	1	1	1	-1.00	-0.67	
0	0	1	0	1	1	-3.67	-3.00	
0	0	1	1	0	1	-1.00	-0.67	
0	0	1	1	1	0	-1.67	-1.67	
0	1	0	0	1	1	-0.33	0.00	
0	1	0	1	0	1	2.33	2.33	
0	1	0	1	1	0	1.67	1.33	
0	1	1	0	0	1	-0.33	0.00	
0	1	1	0	1	0	-1.00	-1.00	
0	1	1	1	0	0	1.67	1.33	
1	0	0	0	1	1	-1.67	-1.33	
1	0	0	1	0	1	1.00	1.00	
1	0	0	1	1	0	0.33	0.00	
1	0	1	0	0	1	-1.67	-1.33	
1	0	1	0	1	0	-2.33	-2.33	
1	0	1	1	0	0	0.33	0.00	
1	1	0	0	0	1	1.67	1.67	
1	1	0	0	1	0	1.00	0.67	
1	1	0	1	0	0	3.67	3.00	
1	1	1	0	0	0	1.00	0.67	

#### 3. Generate the randomization distribution

#### **Definition 2.1 (Randomization Distribution)**

The randomization distribution is the distribution of a test statistic (such as the difference in means) obtained by

(1) permuting the treatment assignments

(2) while keeping the observed outcomes fixed.

\*Observed values are shown in red.

# Fisher randomization test (3): randomization distribution

treatment assignments				nmei	nts	Stati	stics
$A_1$	$A_2$	$A_3$	$A_4$	$A_5$	$A_6$	Difference-in-means	Wilcoxon rank sum
0	0	0	1	1	1	-1.00	-0.67
0	0	1	0	1	1	-3.67	-3.00
0	0	1	1	0	1	-1.00	-0.67
0	0	1	1	1	0	-1.67	-1.67
0	1	0	0	1	1	-0.33	0.00
0	1	0	1	0	1	2.33	2.33
0	1	0	1	1	0	1.67	1.33
0	1	1	0	0	1	-0.33	0.00
0	1	1	0	1	0	-1.00	-1.00
0	1	1	1	0	0	1.67	1.33
1	0	0	0	1	1	-1.67	-1.33
1	0	0	1	0	1	1.00	1.00
1	0	0	1	1	0	0.33	0.00
1	0	1	0	0	1	-1.67	-1.33
1	0	1	0	1	0	-2.33	-2.33
1	0	1	1	0	0	0.33	0.00
1	1	0	0	0	1	1.67	1.67
1	1	0	0	1	0	1.00	0.67
1	1	0	1	0	0	3.67	3.00
1	1	1	0	0	0	1.00	0.67

3. Generate the randomization distribution

#### **Procedure:**

(i) Conduct all possible treatment assignments:  $\sim$  (i)  $\sim$  (ii)

 $\widetilde{A}^{(1)}, \widetilde{A}^{(2)}, \dots, \widetilde{A}^{(M)}$ 

where  $M = C_{N_1}^N$ 

(For example, in the randomized experiment on nocturnal cough, there are  $C_3^6 = 20$  possible treatments assignments)

(ii) Calculate the values of the test statistic for each possible treatment assignment:

 $\mathbb{F} = \{T(\widetilde{A}^{(1)}, \widetilde{Y}), T(\widetilde{A}^{(2)}, \widetilde{Y}), \dots, T(\widetilde{A}^{(M)}, \widetilde{Y})\}$ 

(Note: In a CRE, each possible treatment assignment to units has an *equal probability*.)

(iii)  $\mathbb{F}$  forms the randomization distribution

\*Observed values are shown in red.

# Fisher randomization test (4): Fisher's exact p-value

treatment assignments			nmer	nts	Statistics		
$A_1$	$A_2$	$A_3$	$A_4$	$A_5$	$A_6$	<b>Difference-in-means</b>	Wilcoxon rank sum
0	0	0	1	1	1	-1.00	-0.67
0	0	1	0	1	1	-3.67	-3.00
0	0	1	1	0	1	-1.00	-0.67
0	0	1	1	1	0	-1.67	-1.67
0	1	0	0	1	1	-0.33	0.00
0	1	0	1	0	1	2.33	2.33
0	1	0	1	1	0	1.67	1.33
0	1	1	0	0	1	-0.33	0.00
0	1	1	0	1	0	-1.00	-1.00
0	1	1	1	0	0	1.67	1.33
1	0	0	0	1	1	-1.67	-1.33
1	0	0	1	0	1	1.00	1.00
1	0	0	1	1	0	0.33	0.00
1	0	1	0	0	1	-1.67	-1.33
1	0	1	0	1	0	-2.33	-2.33
1	0	1	1	0	0	0.33	0.00
1	1	0	0	0	1	1.67	1.67
1	1	0	0	1	0	1.00	0.67
1	1	0	1	0	0	3.67	3.00
1	1	1	0	0	0	1.00	0.67

4. Compute Fisher's exact p-value

**One-sided** 

$$P_F = \frac{1}{M} \sum_{i=1}^{M} I\left(T\left(\widetilde{A}^{obs}, \widetilde{Y}\right) \ge T\left(\widetilde{A}^{(i)}, \widetilde{Y}\right)\right)$$

or

$$P_F = \frac{1}{M} \sum_{i=1}^{M} I(T(\widetilde{A}^{obs}, \widetilde{Y}) \leq T(\widetilde{A}^{(i)}, \widetilde{Y}))$$

**Two-sided** 

$$P_F = \frac{1}{M} \sum_{i=1}^{M} I(|T(\widetilde{A}^{obs}, \widetilde{Y})| \le |T(\widetilde{A}^{(i)}, \widetilde{Y})|)$$

#### **Experiment on nocturnal cough**

Difference-in-means:  $P_F = 16/20 = 0.8$ Wilcoxon rank sum:  $P_F = 16/20 = 0.8$ 

\*Observed values are shown in red.

# Neyman vs. Fisher

- > The limitations of Fisher randomization test
- 1. Do not account for treatment effect heterogeneity
- 2. Do not support inference at the population level

During the development of statistical inference methods, Fisher and Neyman introduced distinct frameworks that have significantly influenced causal inference.

- > Neyman's two questions:
- 1. What would the average outcome be if all units were exposed to the treatment?
- 2. How did that compare to the average outcome if all units were exposed to the control?

# Neyman repeated sampling inference

Fisher emphasizes hypothesis testing, while Neyman focuses on parameter estimation.

#### Fisher randomization test

- Sharp null hypothesis (No individual treatment effect within the analytic sample)
- Considers only randomness from treatment assignment
- Derives exact p-values

#### Neyman repeated sampling inference

- Focuses on population/sample average treatment effect
- Accounts for two sources of randomness:
  1. Random sampling from a (super-) population
  - 2. Random treatment assignment
- Derives unbiased estimators and confidence intervals (considers variability across both repeated treatment assignments and repeated sampling)

### > The causal estimands of interest in the Neyman inference:

- Sample average treatment effect (SATE)
- Population average treatment effect (PATE)

# **SATE and PATE**

Sample average treatment effect (SATE)

$$\tau_s \equiv \frac{1}{N} \sum_{i=1}^{N} (Y_i(1) - Y_i(0)) = \overline{Y}(1) - \overline{Y}(0)$$

where N represents the total number of units in the experiment and  $\overline{Y}(a) = N^{-1} \sum_{i=1}^{N} Y_i(a)$ .

- > Population average treatment effect (PATE)  $\tau_p \equiv \mathbb{E}(Y_i(1) - Y_i(0))$
- Note 1: ITE is  $\tau_i \equiv Y_i(1) Y_i(0)$
- Note 2: SATE is precisely the ATE as previously defined.
- Questions for SATE
- 1. Is SATE a **parameter** or an **estimator**?
- 2. In SATE, are  $\{Y_i(1), Y_i(0)\}_{i=1}^N$  random or fixed? How does this compare to PATE
- 3. Who cares about SATE, and why it is important?

# Misunderstandings between experimentalists and observationalists

(Imai, King, and Stuart, 2008)

	<b>SATE</b> $(\tau_s)$	<b>PATE</b> $(\tau_p)$
Scope	Applies to study sample	Applies to entire population
Estimability	Directly computed from study data	Requires extrapolation beyond sample
Generalizability	Limited to sample	Requires assumptions for external validity
Bias risks	Susceptible to selection bias if sample is not random	More representative but harder to estimate

> SATE is useful for internal validity (causal inference within a study).

- > PATE is necessary for external validity (generalizing findings to a broader population).
- > **Bridging the gap** between SATE and PATE requires either:

1. Random sampling,

- 2. Statistical adjustments (e.g., weighting, regression models),
- 3. Replicating studies in diverse settings.

### **Estimator of SATE**

- ➤ Reminder: In SATE inference under CRE,  $\{\tilde{Y}(1), \tilde{Y}(0)\} = \{Y_i(1), Y_i(0)\}_{i=1}^N$  are fixed, and the only source of randomness comes from the random treatment assignments.
- > The difference-in-means estimator of SATE

$$\hat{\tau}_s = \frac{1}{N_1} \sum_{i=1}^{N} A_i Y_i - \frac{1}{N_0} \sum_{i=1}^{N} (1 - A_i) Y_i$$

where  $N_1$  and  $N_0$  represent the number of units in the treatment and control groups, respectively.

- Is  $\hat{\tau}_s$  an unbiased estimator for SATE?

- How can you conduct inference for SATE by using  $\hat{\tau}_s$ , such as deriving a confidence interval?

### **Estimator of SATE**

### **Theorem 2.1 (Unbiasedness and sampling variance in estimating SATE)**

Under a CRE, the difference-in-means estimator  $\hat{\tau}_s$ 

- (1) is unbiased for  $\tau_s$  (SATE) and
- (2) has sampling variance

$$Var\left(\hat{\tau}_{s}|\widetilde{\boldsymbol{Y}}(1),\widetilde{\boldsymbol{Y}}(0)\right) = \frac{S_{1}^{2}}{N_{1}} + \frac{S_{0}^{2}}{N_{0}} - \frac{S_{\tau}^{2}}{N}$$

where

$$S_a^2 = \frac{1}{N-1} \sum_{i=1}^{N} (Y_i(a) - \bar{Y}(a))^2 \text{ and } S_{\tau}^2 = \frac{1}{N-1} \sum_{i=1}^{N} (\tau_i - \tau_s)^2.$$

- $\blacktriangleright$  It demonstrates that  $\hat{\tau}_s$  is a *reasonable estimator* for SATE under a CRE.
- >  $S_{\tau}^2$  includes both potential outcomes  $Y_i(1)$  and  $Y_i(0)$ , making it non-estimable from observed data.
- > Neyman suggested that, in practice,  $S_{\tau}^2$  is often ignored when estimating the sampling variance of  $\hat{\tau}_s$ .

# **Estimator of the sampling variance**

### **Theorem 2.2**

Under the constant additive treatment effect assumption (i.e., the individual treatment effect  $\tau_i (\equiv Y_i(1) - Y_i(0))$  is constant), we have

(1)  $S_{\tau}^2 = 0$  and

(2) an unbiased estimator for the sampling variance

$$\widehat{\mathbb{V}} = \frac{\widehat{S}_1^2}{N_1} + \frac{\widehat{S}_0^2}{N_0}$$

where

$$\hat{S}_{1}^{2} = \frac{1}{N_{1} - 1} \sum_{i=1}^{N} A_{i} (Y_{i} - \bar{Y}_{A=1})^{2},$$
  

$$\hat{S}_{0}^{2} = \frac{1}{N_{0} - 1} \sum_{i=1}^{N} (1 - A_{i}) (Y_{i} - \bar{Y}_{A=0})^{2},$$
  

$$\bar{Y}_{A=1} = \frac{1}{N_{1}} \sum_{i=1}^{N} A_{i} Y_{i}, \text{ and } \bar{Y}_{A=0} = \frac{1}{N_{0}} \sum_{i=1}^{N} (1 - A_{i}) Y_{i}$$

# Further insights on $\widehat{\mathbb{V}}$

➢ It is important to note that if there is heterogeneity in treatment effects—meaning the constant additive treatment effect assumption does not hold—then  $\widehat{V}$  is a *conservative estimator* of the sampling variance:

$$\mathbb{E}\left(\widehat{\mathbb{V}}|\widetilde{\mathbf{Y}}(1),\widetilde{\mathbf{Y}}(0)\right) - Var\left(\widehat{\tau}_{s}|\widetilde{\mathbf{Y}}(1),\widetilde{\mathbf{Y}}(0)\right) = \frac{S_{\tau}^{2}}{N} \ge 0.$$

- ➢ Even when the assumption of an additive treatment effect may be known to be inaccurate, 𝒱 remains widely used (in two-sample testing).
- 1. The confidence interval for SATE constructed using  $\widehat{\mathbb{V}}$  is conservative.
- 2. For PATE, in a (super-)population,  $\widehat{\mathbb{V}}$  is an unbiased estimator for the sampling variance of  $\hat{\tau}_s$ .

### **Proof of Theorem 2.1**

### Unbiasedness

First, the difference-in-means estimator  $\hat{\tau}_s$  can be expressed as

$$\hat{\tau}_s = \frac{1}{N_1} \sum_{i=1}^{N} A_i Y_i(1) - \frac{1}{N_0} \sum_{i=1}^{N} (1 - A_i) Y_i(0)$$

Then we have

$$\begin{split} \mathbb{E}(\hat{\tau}_{s} | \widetilde{Y}(1), \widetilde{Y}(0)) &= \frac{1}{N_{1}} \sum_{i=1}^{N} \mathbb{E}(A_{i} | \widetilde{Y}(1), \widetilde{Y}(0)) Y_{i}(1) - \frac{1}{N_{0}} \sum_{i=1}^{N} \{1 - \mathbb{E}(A_{i} | \widetilde{Y}(1), \widetilde{Y}(0))\} Y_{i}(0) \\ &= \frac{1}{N_{1}} \sum_{i=1}^{N} \mathbb{E}(A_{i}) Y_{i}(1) - \frac{1}{N_{0}} \sum_{i=1}^{N} \{1 - \mathbb{E}(A_{i})\} Y_{i}(0) \\ &= \frac{1}{N_{1}} \sum_{i=1}^{N} \frac{N_{1}}{N} Y_{i}(1) - \frac{1}{N_{0}} \sum_{i=1}^{N} \frac{N_{0}}{N} Y_{i}(0) \\ &= \tau_{s} \end{split}$$

### **Estimator of PATE**

### **Theorem 2.3 (Unbiasedness and sampling variance in estimating PATE)**

Under a CRE, the difference-in-means estimator  $\hat{\tau}_s$ 

- (1) is unbiased for  $\tau_p$  (PATE) and
- (2) has sampling variance

$$Var(\hat{\tau}_s) = \frac{\sigma_1^2}{N_1} + \frac{\sigma_0^2}{N_0}$$

where  $\sigma_a^2$  is the population variance of Y(a) for a = 1,2

- Unbiasedness is straightforward to demonstrate, as the units/samples are randomly drawn from the population.
- The derivation of sampling variance (over repeated sampling and repeated treatment assignments):
  By the law of total variance,

$$Var(\hat{\tau}_{s}) = Var\left(\mathbb{E}\left(\hat{\tau}_{s} \middle| \widetilde{Y}(1), \widetilde{Y}(0)\right)\right) + \mathbb{E}\left(Var\left(\hat{\tau}_{s} \middle| \widetilde{Y}(1), \widetilde{Y}(0)\right)\right)$$
$$= \frac{1}{N} Var\left(Y_{i}(1) - Y_{i}(1)\right) + \frac{\sigma_{1}^{2}}{N_{1}} + \frac{\sigma_{0}^{2}}{N_{0}} - \frac{1}{N} Var\left(Y_{i}(1) - Y_{i}(1)\right) = \frac{\sigma_{1}^{2}}{N_{1}} + \frac{\sigma_{0}^{2}}{N_{0}}$$
$$Causal Inference, Part 1-3. An-Shun Tai$$

### **Confidence interval**

→ Li and Ding (2017) demonstrated the asymptotic Normality of  $\hat{\tau}_s$  based on the finite population central limit theorem.

➤ Wald-type large-sample confidence interval is given by

$$\hat{\tau}_s \pm z_{1-\alpha/2} \sqrt{\widehat{\mathbb{V}}}$$

where  $z_{1-\alpha/2}$  is the  $1-\alpha/2$  upper quantile of the standard Normal distribution.

### References

Imai, K., King, G., & Stuart, E. A. (2008). Misunderstandings between experimentalists and observationalists about causal inference. *Journal of the Royal Statistical Society Series A: Statistics in Society*, *171*(2), 481-502.

Li, X., & Ding, P. (2017). General forms of finite population central limit theorems with applications to causal inference. *Journal of the American Statistical Association*, *112*(520), 1759-1769.

Paul, I. M., Beiler, J., McMonagle, A., Shaffer, M. L., Duda, L., & Berlin, C. M. (2007). Effect of honey, dextromethorphan, and no treatment on nocturnal cough and sleep quality for coughing children and their parents. *Archives of pediatrics & adolescent medicine*, *161*(12), 1140-1146.

Rubin, D. B. (1980). Comment on "Randomization analysis of experimental data: The Fisher randomization test". *Journal of the American statistical association*, *75*(371), 591-593.