

Lecture 7

Stratified randomized experiments

Outline

- Stratified randomized experiment
 - Fisher's exact p-value
 - Neyman's repeated sampling approach
 - Regression analysis
- Post stratification
- Suggested reading: Imbens and Rubin Chapter 9.1-9.6, Peng's book Chapter 5

STAR (Student-Teacher Achievement Ratio) Project in Tennessee

(Mosteller. 1997. Bull. Am. Acad. Arts Sci.)

- What is STAR? (1985-1989)

- A large-scale, four-year, longitudinal, experimental study of reduced class size
- One the historically most important educational investigations
- Cost of about \$12 million
- Conclusion: small classes have an advantage over larger classes in reading and math in the early primary grades

- Why was STAR needed?

- Legislators and school administrators doubted the significance of smaller classes
- Conducted at the elementary-school level as this is where the foundation is laid for children's success in school.
- The most credible study of class size



STAR (Student-Teacher Achievement Ratio) Project in Tennessee

(Mosteller. 1997. Bull. Am. Acad. Arts Sci.)

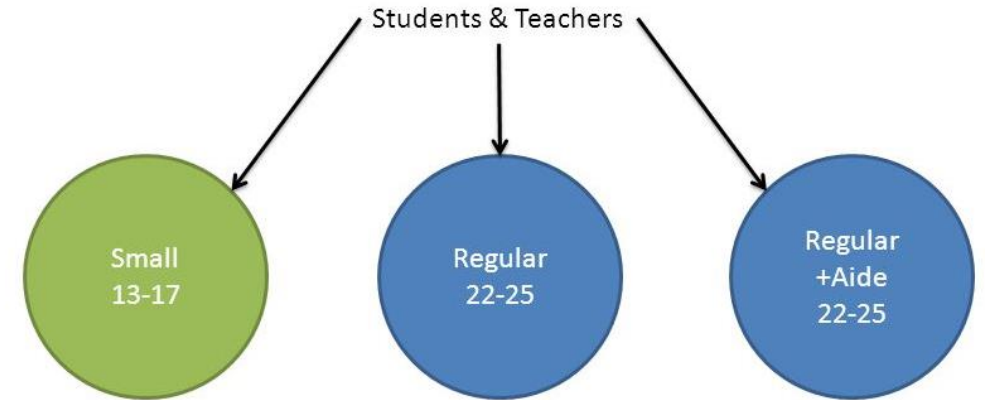
- How is the experiment designed?

- Three levels of “treatment”: three types of classes

- Regular class + Aide:
One teacher plus a full-time teacher’s aide.
- Difference between Class Size and Pupil/Teacher Ratio

- What need to be randomized, students or teacher?

- Both students and teachers were randomly assigned to the one of the 3 arms
- What is a unit, student or teacher?
 - The unit is a teacher in a class, instead of a student to avoid violation of no interference assumption



The project STAR example

(Mosteller. 1997. Bull. Am. Acad. Arts Sci.)

- Two randomizations happen in the experiment
 - Randomization of teachers
 - Randomization of students
- Our causal analysis only relies on the randomization of teachers
 - The treatment effect on a particular teacher in a particular school is comparing the test score of being randomly assigned to a type of class and the test score of being randomly assigned to another type of class
- The randomization of students helps interpreting our results
 - Treatment effect between two arms can be explained by the classroom size difference instead of the systematic differences of students

The project STAR example

(Mosteller. 1997. Bull. Am. Acad. Arts Sci.)

- How to conduct a large experiments across schools?
 - The study included 79 schools resulting in over 6,000 students per grade
 - potentially large differences in resources, teachers and students between schools
 - How to deal with that?
- Stratified randomization procedure
 - A school need to have a minimum of 57 students in kindergarden (at least one for each type of class)
 - Once a school is admitted, a decision was made on the number of classes per arm
 - Randomization within each school
 - Students and teachers within the school were randomly assigned to the one of the 3 arms

The project STAR example

(Mosteller. 1997. Bull. Am. Acad. Arts Sci.)

- The interventions were initiated as the students entered school in kindergarten and continued through third grade.

	Kindergarten	Grade 1	Grade 2	Grade 3
Inner City	17	15	15	15
Suburban	16	15	15	15
Rural	38	38	38	38
Urban	8	8	7	7
Total	79	76	75	75

- Practical issues faced in real experiment
 - Longitudinal experiment
 - Schools may drop out of the project
 - Classes may gain/lose students so that can become too small or too big
 - Selection bias in students' involvement
 - Students' parents were informed so may want their children to be in the smaller class

Table 9.1. Class Average Mathematics Scores from Project Star

School/ Stratum	No. of Classes	Regular Classes ($W_i = 0$)	Small Classes ($W_i = 1$)
1	4	−0.197, 0.236	0.165, 0.321
2	4	0.117, 1.190	0.918, −0.202
3	5	−0.496, 0.225	0.341, 0.561, −0.059
4	4	−1.104, −0.956	−0.024, −0.450
5	4	−0.126, 0.106	−0.258, −0.083
6	4	−0.597, −0.495	1.151, 0.707
7	4	0.685, 0.270	0.077, 0.371
8	6	−0.934, −0.633	−0.870, −0.496, −0.444, 0.392
9	4	−0.891, −0.856	−0.568, −1.189
10	4	−0.473, −0.807	−0.727, −0.580
11	4	−0.383, 0.313	−0.533, 0.458
12	5	0.474, 0.140	1.001, 0.102, 0.484
13	4	0.205, 0.296	0.855, 0.509
14	4	0.742, 0.175	0.618, 0.978
15	4	−0.434, −0.293	−0.545, 0.234
16	4	0.355, −0.130	−0.240, −0.150
Average (S.D.)		−0.13 (0.56)	0.09 (0.61)

- We focus on two arms (regular classes v.s. small classes) and 16 schools that have at least two classes per arm

Stratified randomized experiment

- Basic procedure:
 1. Blocking (Stratification): create groups of similar units based on pre-treatment covariates, let $B_i \in \{1, \dots, J\}$ be the block indicator
 2. Block (Stratified) randomization: completely randomize treatment assignment within each group
- Blocking can improve the efficiency by minimizing the variance of the potential outcomes within each strata

“Block what you can and randomize what you cannot”

Box, et al. (2005). Statistics for Experimenters. 2nd eds. Wiley

- Assignment probability

$$P(\mathbf{W} = \mathbf{w} | \mathbf{X}) = \begin{cases} \prod_{j=1}^J \binom{N(j)}{N_t(j)}^{-1} & \text{if } \sum_{i: B_i=j} w_i = N_t(j) \text{ for } j = 1, \dots, J \\ 0 & \text{otherwise} \end{cases}$$

Compare treated v.s. control? Simpson's paradox

- Compare the success rates of two treatment of kidney stones
- Treatment A: open surgery; treatment B: small pictures

	Treatment A	Treatment B
Small stones	93% (81/87)	87% (234/270)
Large stones	73% (192/263)	69% (55/80)
Both	78% (273/350)	83% (289/350)

- Large difference in treatment assignment probability across strata
 - Small stone: assignment probability $\frac{87}{87+270} = 0.24$
 - Large stone: assignment probability is $\frac{263}{263+80} = 0.77$
- Compare within each strata and take a weighted average:
 - True average causal effect: $83.2\% - 78.2\% : (93\% - 87\%) \times 0.51 - (73\% - 69\%) \times 0.49$

Fisher's exact p-value

- We still focus on the **Sharp null**: $H_0: Y_i(0) \equiv Y_i(1)$ for all $i = 1, \dots, N$
- **Choice of test statistics:**

Denote sample means for every strata / block

$$\bar{Y}_c^{\text{obs}}(j) = \frac{1}{N_c(j)} \sum_{i: G_i=j} (1 - W_i) \cdot Y_i^{\text{obs}}, \quad \bar{Y}_t^{\text{obs}}(j) = \frac{1}{N_t(j)} \sum_{i: G_i=j} W_i \cdot Y_i^{\text{obs}}$$

- Weighted combination of group mean differences across blocks

$$T^{\text{dif}, \lambda} = \left| \sum_{j=1}^J \lambda(j) \cdot (\bar{Y}_t^{\text{obs}}(j) - \bar{Y}_c^{\text{obs}}(j)) \right|$$

- Weights based on relative sample size $\lambda(j) = \frac{N(j)}{N}$
sample difference is more accurate in larger strata
- **“inverse-variance-weighting”**: assume that per-strata potential outcomes sample variances $S_c^2(j) \equiv S_t^2(j) \equiv S^2$ for all j , then under stratified randomization

$$\mathbb{V}_W[\bar{Y}_t^{\text{obs}}(j) - \bar{Y}_c^{\text{obs}}(j) | \mathbf{Y}(0), \mathbf{Y}(1)] = S^2 \left(\frac{1}{N_c(j)} + \frac{1}{N_t(j)} \right)$$

Fisher's exact p-value

- We still focus on the **Sharp null**: $H_0: Y_i(0) \equiv Y_i(1)$ for all $i = 1, \dots, N$
- **Choice of test statistics:**

Denote sample means for every strata / block

$$\bar{Y}_c^{\text{obs}}(j) = \frac{1}{N_c(j)} \sum_{i: G_i=j} (1 - W_i) \cdot Y_i^{\text{obs}}, \quad \bar{Y}_t^{\text{obs}}(j) = \frac{1}{N_t(j)} \sum_{i: G_i=j} W_i \cdot Y_i^{\text{obs}}$$

- Weighted combination of group mean differences across blocks

$$T^{\text{dif}, \lambda} = \left| \sum_{j=1}^J \lambda(j) \cdot (\bar{Y}_t^{\text{obs}}(j) - \bar{Y}_c^{\text{obs}}(j)) \right|$$

- Weights based on relative sample size $\lambda(j) = \frac{N(j)}{N}$
sample difference is more accurate in larger strata
- **“inverse-variance-weighting”**: weights

$$\lambda(j) = \frac{1}{\left(\frac{1}{N_c(j)} + \frac{1}{N_t(j)} \right)} / \sum_{k=1}^J \frac{1}{\left(\frac{1}{N_c(k)} + \frac{1}{N_t(k)} \right)}$$

Fisher's exact p-value

- Can we simply use the two-sample mean difference statistic $T = |\bar{Y}_t^{\text{obs}} - \bar{Y}_c^{\text{obs}}|$?
 - This is still one test statistic and we will still get valid Fisher's exact p-value if we follow the stratified randomization procedure to generate the reference distribution

Simpson's paradox:

- We may not always get small value of T even when the sharp null is true
 - Example:
 $Y_i(0) \equiv Y_i(1) = 1$ for strata 1 and $Y_i(0) \equiv Y_i(1) = 2$ for strata 2,
 $N_c(1) = N_t(1) = 5$, $N_c(2) = 15$ and $N_t(2) = 5$
Then $\bar{Y}_t^{\text{obs}} = 1.5$ and $\bar{Y}_c^{\text{obs}} = 1.75$
- Power of the Fisher's test is affected

Fisher's exact p-value and the project STAR

- Choice of test statistics:
 - Rank-based statistics
 - Get R_i^{strat} as the within-strata rank of each individual i (definition page 196 of Imbens and Rubin's book)
 - Average difference of within-strata ranks between treatment and control

$$|\bar{R}_t^{\text{strat}} - \bar{R}_c^{\text{strat}}|$$

- Calculate the null distribution of test statistics
 - Randomly simulate treatment assignments following the same stratified randomization

- Project STAR results
 - P-values for the first 3 are similar as most schools have 4 classes
 - Large p-value for rank-based statistics as # classes too few in most schools

Test statistics	P-value
Weights	
$\lambda(j) = \frac{N(j)}{N}$	0.034
"inverse-variance-weighting"	0.023
$ \bar{Y}_t^{\text{obs}} - \bar{Y}_c^{\text{obs}} $	0.025
Rank-based statistics	0.15

Neyman's repeated sampling approach

- **Target:** PATE or SATE $\tau = \sum_j \frac{N(j)}{N} \tau(j)$ where $\tau(j)$ is the PATE or SATE for strata j

- **Analysis procedure**

1. Apply Neyman's analysis to each strata / block

$$\hat{\tau}^{\text{dif}}(j) = \bar{Y}_t^{\text{obs}}(j) - \bar{Y}_c^{\text{obs}}(j), \quad \text{and} \quad \hat{V}^{\text{neyman}}(j) = \frac{s_c(j)^2}{N_c(j)} + \frac{s_t(j)^2}{N_t(j)}$$

- Variance estimator is conservative within each strata as discussed before

2. Aggregate block-specific estimates and variances

$$\hat{\tau}^{\text{strat}} = \sum_j \frac{N(j)}{N} \hat{\tau}^{\text{dif}}(j), \quad \hat{V}(\hat{\tau}^{\text{strat}}) = \sum_j \left(\frac{N(j)}{N} \right)^2 \hat{V}^{\text{neyman}}(j)$$

- Both treatment assignments and potential outcomes are independent across strata

3. Statistical inference

- Use normal approximation of the distribution of $\hat{\tau}^{\text{strat}}$

- Normal approximation works as long as N is large enough

- Either small strata size with many strata or large strata size with few strata

Power gain in Neyman's approach after stratification

- Variance decomposition

$$\underbrace{\mathbb{V}(X)}_{\text{total variance}} = \underbrace{\mathbb{E}\{\mathbb{V}(X | Y)\}}_{\text{within-block variance}} + \underbrace{\mathbb{V}\{\mathbb{E}(X | Y)\}}_{\text{across-block variance}}$$

- Assume that the treatment proportion $\frac{N(j)}{N}$ is the same across all strata
 - Then $\hat{\tau}^{\text{dif}} = \hat{\tau}^{\text{strat}}$
- $\mathbb{V}_{\text{complete}}(\hat{\tau}^{\text{dif}}) - \mathbb{V}_{\text{stratified}}(\hat{\tau}^{\text{strat}}) \geq 0$
 - Intuitively, we do not need to consider noise due to heterogeneity across blocks
 - For a rigorous proof, see Peng's book section 5.3.3
- Result in the project STAR
 - $\hat{\tau}^{\text{strat}} = 0.241, \widehat{\mathbb{V}}(\hat{\tau}^{\text{strat}}) = 0.092^2$
 - (Incorrect) if we analyze as if it is a completely randomized experiment
 - $\hat{\tau}^{\text{dif}} = \bar{Y}_t^{\text{obs}} - \bar{Y}_c^{\text{obs}} = 0.224$ can be a biased estimate for τ
 - $\widehat{\mathbb{V}}(\hat{\tau}^{\text{dif}}) = 0.141^2$ larger standard deviation

Linear regression

- Run separate linear regressions within each strata
 - Does not work if each strata size is too small
- Denote $B_i(j)$ as the indicator variable of whether sample i belong to strata j
- If there are no covariates, equivalently, we can write separate linear regression models into a joint regression model

$$Y_i^{\text{obs}} = \alpha_j + \tau(j)W_i + \varepsilon_i$$

- The underlying model for the potential outcomes

$$\mathbb{E}[Y_i(w) | \{B_i(j), j = 1, \dots, J\}] = \alpha_j + \tau(j)w$$

- Average causal effect for strata j is $\tau(j)$
- The strata indicators $B_i(j)$ are treated as pre-treatment covariates
- We need to adjust for the strata indicators as we only have conditional independence

$$(Y(0), Y(1)) \perp W \mid B(j)$$

- The homoscedastic error assumption for the joint model is assuming that

$$\mathbb{V}[Y_i(0) | \{B_i(j), j = 1, \dots, J\}] = \mathbb{V}[Y_i(1) | \{B_i(j), j = 1, \dots, J\}] = \sigma^2$$

Post-stratification

- In a completely randomized experiment, each assignment vector has the sample probability ($P(\mathbf{W} = \mathbf{w})$) if $\sum_{i=1}^N w_i = N_t$
- If we focus on a subgroup S , conditional on $N_{t,S} = \sum_{i \in S} W_i$, the assignment vector for the individuals in the subgroup also has the same probability ($P(\mathbf{W}_S = \mathbf{w}_S)$) if $\sum_{i \in S} w_i = N_{t,S}$
- So conditional on $N_{t,S}$, we can treat the treatment assignment as from a completely randomized experiment also for the subgroup
- **Post-stratification** (Miratrix. et al. 1971. J. Royal Stat. Soc. B.)
 - Blocking after the experiment is conducted
 - Analyze the experiment as from a stratified randomized experiment by conditioning on $N_{t,S}$ for each strata S
 - By post-stratification, we can stratify individuals into relatively homogenous subpopulations
 - Post-stratification is nearly as efficient as pre-randomization blocking

Meinert et. al. (1970)'s example

- A completely randomized experiment.
- Treatment is tolbutamide ($Z = 1$) and control is a placebo ($Z = 0$)
- Causal effect: difference in the survival probability

Age < 55			Age \geq 55		
	Surviving	Dead		Surviving	Dead
$Z = 1$	98	8	$Z = 1$	76	22
$Z = 0$	115	5	$Z = 0$	69	16
Total					
	Surviving			Dead	
$Z = 1$	174			30	
$Z = 0$	184			21	

Peng's book Section 5.4.1

- Subgroup and sample average estimates with post-stratification

	stratum 1	stratum 2	post-stratification	crude
est	−0.034	−0.036	−0.035	−0.045
se	0.031	0.060	0.032	0.033