

The Benefits of Probability-Proportional-to-Size Sampling in Cluster-Randomized Experiments

Yeng Xiong*
Michael J. Higgins*

First Draft: 1/20/2014

This Draft: 2/20/2020 (02:48)

Abstract

In a cluster-randomized experiment, treatment is assigned to clusters of individual units of interest—households, classrooms, villages, etc.—instead of the units themselves. The number of clusters sampled and the number of units sampled within each cluster is typically restricted by a budget constraint. Previous analysis of cluster randomized experiments under the Neyman-Rubin potential outcomes model of response have assumed a simple random sample of clusters. Estimators of the population average treatment effect (PATE) under this assumption are often either biased or not invariant to location shifts of potential outcomes. We demonstrate that, by sampling clusters with probability proportional to the number of units within a cluster, the Horvitz-Thompson estimator (HT) is invariant to location shifts and unbiasedly estimates PATE. We derive standard errors of HT and discuss how to estimate these standard errors. We also show that results hold for stratified random samples when samples are drawn proportionally to cluster size within each stratum. We demonstrate the efficacy of this sampling scheme using a simulation based on data from an experiment measuring the efficacy of the National Solidarity Programme in Afghanistan.

Keywords— cluster randomized experiment, Neyman-Rubin model, probability-proportional-to-size sampling, Horvitz-Thompson estimator

*Department of Statistics, Kansas State University. We are grateful to Graeme Blair, Ben Fifield, Kosuke Imai, Ben Johnson, James Lo, the Imai Research Group, and the Higgins Research Group for helpful discussions and advice.

1 Introduction

Frequently in experiments, treatment is randomized across clusters, or groups, of units of interest instead of the units themselves. These are referred to as *cluster-randomized experiments* (CREs). Clusters of units are often formed *a priori* to the design of the experiment and without researcher intervention. Estimation of treatment effects is more precise when treatment is randomized across units [Cornfield, 1978]; hence, logistical issues (rather than increased precision of treatment effect estimates) motivate the randomization of treatment across clusters. Reasons for such randomization include addressing issues with the ethicality, legality, or feasibility of randomizing treatment across units, reducing risk of treatment contamination, and mimicking the implementation of a proposed program (e.g. an educational intervention) [Donner, 1998, Donner and Klar, 2004, Hayes and Moulton, 2009]. Common settings for cluster-randomized experiments include: testing an educational intervention that is implemented within classrooms [Raver et al., 2009]; evaluating efficacy of a health intervention that is implemented within clinics or medical practices [Bruce et al., 2004, King et al., 2007, Small et al., 2008, Imai et al., 2009]; measuring increases in compliance and turnout from mailers sent to households [Gerber and Green, 2000]; and identifying effects of interventions implemented within villages or other geographic regions [Wantchekon, 2003, Paluck, 2009, Beath et al., 2013].

To estimate and perform inference on the *population average treatment effect* (PATE), a CRE will require at least two stages of sampling: sampling clusters from a larger population of clusters (e.g. a sample of villages within a country) and sampling individual units from each of the sampled clusters—samples may be comprised of the entire sampling frame. After a sample of clusters is obtained, but before units are sampled within each cluster, treatment is allocated across sampled clusters. Researchers often improve the precision of treatment effect estimates by drawing a stratified sample of clusters and/or blocking sampled clusters before treatment assignment [Gail et al., 1996, Lewsey, 2004, Imai et al., 2009, Hayes and Moulton, 2009, Imbens, 2011, Hansen et al.,

2014]. When researchers are interested in heterogeneous treatment effects across subpopulations of interest, within-cluster samples may also be stratified (for an example, see Kerry et al. [2005]).

When clusters are sampled using simple random sampling (SRS) or stratified random sampling (StRS), current estimators of the PATE have undesirable properties. The unbiased Horvitz-Thompson (HT-SRS) estimator [Horvitz and Thompson, 1952] is not invariant to location shifts of responses, which inflates its variance. The location-invariant difference-in-means (DIM) estimator will be biased when treatment effects are correlated with cluster *sizes*—the number of units contained within each cluster [Middleton and Aronow, 2015]. Thus, this estimator is only unbiased in special cases such as under sharp null of no unit-level treatment effect [Hansen et al., 2014] or when clusters are blocked or stratified exactly on cluster sizes [Donner and Klar, 2004, Imai et al., 2009]. Moreover, as we will show, when within-cluster samples are not drawn proportional to the cluster size, DIM may estimate a quantity different from the PATE. In fact, the only current estimator of the PATE that is both unbiased and location-invariant is the Des Raj estimator (DR) [Middleton and Aronow, 2015], which requires the introduction of an additional parameter; however, estimating this parameter will induce bias in the estimator.

We propose an adjustment in the *design* of the experiment—as opposed to adjusting weights of estimators after the experiment—for differences in cluster sizes: to sample clusters with *probability proportional to size* (PPS) [Hansen and Hurwitz, 1943, Cochran, 1977, Lohr, 2010]. We show that, under this sampling scheme, the Horvitz-Thompson estimator (HT-PPS) is both unbiased and location invariant.

The paper is organized as follows: Section 2 introduces notation. Section 2.7 demonstrates problems with HT-SRS, DIM, and DR estimators of PATE under SRS of clusters. Section 3 demonstrates that the HT-PPS estimator is both unbiased and location-invariant under PPS-without-replacement sampling of clusters, gives standard errors and estimates of standard errors for HT-PPS, and shows equivalence of HT-PPS and DIM (under PPS) estimators when within-cluster sample sizes are the same across clusters. Section 4 extends results to the case where the sample of

clusters and the within-cluster sample of units are stratified. Section 5 gives simulations on a data example, which shows that the HT-PPS estimator has the smallest mean squared error compared to the other estimators. This is due to the HT-PPS estimator being as efficient as the DIM estimator and being unbiased. It also shows that the estimated variance is conservative for the variability of HT-PPS estimator.

2 Notation, assumptions, and preliminaries

We consider a finite population of n units partitioned into ℓ clusters. Clusters are numbered 1 through ℓ . Let n_c denote the number of units within cluster c . Suppose units are ordered in some way within each cluster; let (k, c) denote the k^{th} unit in cluster c . We now introduce sampling and treatment assignment notation in the order in which they are performed in a CRE.

2.1 Sampling clusters

A total of s clusters are sampled; we assume s is fixed and chosen by the researcher. Let S_c denote a cluster sampling indicator; $S_c = 1$ if and only if cluster c is contained in the sample.

$$S_c = \begin{cases} 1, & \text{cluster } c \text{ is sampled,} \\ 0, & \text{otherwise.} \end{cases} \quad (1)$$

By definition, $\sum_{c=1}^{\ell} S_c = s$.

2.2 Treatment assignment

Each of the s sampled clusters is assigned to either treatment or control. Let T_{ct} denote a treatment indicator; $T_{ct} = 1$ if and only if cluster c receives treatment $t \in \{0, 1\}$.

$$T_{ct} = \begin{cases} 1, & \text{cluster } c \text{ receives treatment } t, \\ 0, & \text{otherwise.} \end{cases} \quad (2)$$

We define $T_{ct} = 0$ when $S_c = 0$. Let $\#T_t$ denote the number of clusters that receive treatment t .

We suppose that treatment assignment is *symmetric* across sampled clusters [Miratrix et al., 2013]. That is, conditioned on the number of treated clusters $\#T_t$, each of the $\binom{s}{\#T_t}$ possible treatment assignments is equally likely. Symmetric treatment assignment implies that, for any treatment $t \in \{0, 1\}$ and distinct clusters c, c' :

$$\mathbb{E}(T_{ct} | \mathbf{S}) = \frac{\#T_t}{s}, \quad (3)$$

$$\mathbb{E}(T_{ct}T_{c't} | \mathbf{S}) = \frac{\#T_t(\#T_t - 1)}{s(s - 1)}. \quad (4)$$

where $\mathbf{S} = (S_1, S_2, \dots, S_n)$ denote a random set of cluster sampling indicator variables under a sampling design. Complete randomization is a special case of symmetric treatment assignment. When the sample of clusters is stratified, symmetric treatment assignment also requires independence of treatment assignment across strata, which is discussed in Section 4.

2.3 Within-cluster sampling

After treatment is assigned across clusters, a SRS of s_c units is drawn within each sampled cluster c . This sample is drawn independently of treatment assignment and independently across clusters. We assume that these sample sizes are non-random and do not depend on the set of clusters sampled.

Let S_{kc} denote unit sampling indicator; $S_{kc} = 1$ if and only if the k^{th} unit in cluster c is sampled.

$$S_{kc} = \begin{cases} 1, & \text{unit } (k, c) \text{ is sampled,} \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

We define $S_{kc} = 0$ when $S_c = 0$. By definition, $\sum_{k=1}^{n_c} S_{kc} = s_c$.

2.4 Model of response: Neyman-Rubin Causal Model

Let y_{kct} denote the *potential outcome* of unit (k, c) given treatment t —the value of unit (k, c) we would have observed had that unit received treatment t . Note that y_{kct} is known if and only if that unit is sampled and receives treatment t (i.e., $S_c T_{ct} S_{kc} = 1$). Potential outcomes are assumed to be nonrandom. Let $\mathbf{y} = (y_{kct})_{k=1, c=1, t=0}^{n_c, \ell, 1}$ denote the vector of potential outcomes.

Let Y_{kc} denote the observed response of unit (k, c) had that unit been sampled. We assume responses follow the Neyman-Rubin Causal Model (NRCM) [Splawa-Neyman et al., 1923, Rubin, 1974, Holland, 1986]:

$$\begin{aligned} Y_{kc} &= y_{kc1} T_{c1} + y_{kc0} T_{c0} \\ &= y_{kc1} T_{c1} + y_{kc0} (1 - T_{c1}). \end{aligned} \quad (6)$$

Inherent in this model is the *stable-unit treatment value assumption* (SUTVA), which is often referred to as the *no-interference assumption*; the value of Y_{kc} only depends on the treatment assigned to cluster c and is not affected by the treatment assignment of any other cluster c' . Observe that, since each cluster receives a single treatment condition, this assumption only needs to hold across sampled clusters and does not need to hold for units within each cluster.

2.5 Parameter of interest

Our quantity of interest is the *population average treatment effect* (PATE):

$$\delta = \delta(\mathbf{y}) \equiv \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \frac{y_{kc1} - y_{kc0}}{n} = \mu_1 - \mu_0, \quad (7)$$

where

$$\mu_t = \mu_t(\mathbf{y}) \equiv \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \frac{y_{kct}}{n}. \quad (8)$$

denotes the *population mean for treatment t*. Let

$$\mu_{ct} \equiv \sum_{k=1}^{n_c} \frac{y_{kct}}{n_c} \quad (9)$$

denote the *population mean of cluster c for treatment t*. We can write the population mean as:

$$\mu_t = \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \frac{y_{kct}}{n} = \sum_{c=1}^{\ell} \frac{n_c}{n} \sum_{k=1}^{n_c} \frac{y_{kct}}{n_c} = \sum_{c=1}^{\ell} \frac{n_c}{n} \mu_{ct}. \quad (10)$$

We then define

$$\sigma_{ct}^2 \equiv \sum_{k=1}^{n_c} \frac{(y_{kct} - \mu_{ct})^2}{n_c - 1}, \quad (11)$$

$$\sigma_{t,bet}^2 \equiv \sum_{c=1}^{\ell} \frac{n_c}{n} (\mu_{ct} - \mu_t)^2, \quad (12)$$

respectively, as the variance of potential outcomes within cluster c under treatment t and as the weighted across-cluster variance for treatment t .

2.6 Properties of estimators

A function of potential outcomes f is *monotonically increasing* if $f(\mathbf{y}^*) \geq f(\mathbf{y})$ whenever

$$y_{kct}^* \geq y_{kct}, \quad \text{for all } k \in \{1, \dots, n_c\}, c \in \{1, \dots, \ell\}, t \in \{0, 1\}. \quad (13)$$

A transformation of potential outcomes $\mathbf{y} \rightarrow \mathbf{y}^*$ is *linear* if, for constants a, b :

$$y_{kct}^* = a + by_{kct}, \quad \text{for all } k \in \{1, \dots, n_c\}, c \in \{1, \dots, \ell\}, t \in \{0, 1\}. \quad (14)$$

For simplicity, we may write this as $\mathbf{y}^* = a + b\mathbf{y}$. A *location transformation* or *shift* is a linear transformation in which $b = 1$.

Observe that the population mean is a monotone increasing function that is linear in potential outcomes,

$$\mu_t(a + b\mathbf{y}) = a + b\mu_t(\mathbf{y}), \quad (15)$$

whereas the PATE is *location-invariant*—that is, the value does not change given a location shift of potential outcomes,

$$\delta(a + \mathbf{y}) = \delta(\mathbf{y}). \quad (16)$$

2.7 Methods for estimating PATE under SRS of clusters

In CREs, clusters are typically sampled using SRS, and the common estimators under this sampling procedure include the Horvitz-Thompson (HT-SRS), the difference-in-means (DIM), and the Des Raj (DR) estimators. The HT-SRS estimator weights each unit's outcome with the inverse of the probability that the unit is treated and selected. Therefore, it is unbiased, which recommends it as an appropriate estimator of PATE, but Imai et al. [2009] shows that it can be criticized on two counts. The first is that the estimator is known to have huge variability since it does not account for

varying cluster size. Larger clusters will have greater sums of responses whereas smaller clusters will have smaller sums. The second being that it is not location-invariant. This poses a dilemma for variance calculation since the variance will change as a changes. The HT-SRS estimator will only be location-invariant when the number of treated *units* (not clusters) is equal to the number of units assigned to control, something of which researchers cannot control.

The DIM estimator is the difference between the sample means for treated and control units. The estimator, being elegantly simple, is favored among many researchers. Furthermore, contrary to the HT-SRS estimator, it is efficient and invariant to location shifts. However, Middleton and Aronow [2015] shows that it is biased in CREs. In actuality, the DIM estimator will be unbiased only when treatment effects are not correlated with cluster sizes and when within-cluster sample sizes are proportional to cluster sizes.

Middleton and Aronow [2015] instead advocate the Des Raj (DR) estimator, which adds a regression component on cluster size to the HT-SRS estimator. This helps alleviate the two criticisms on HT-SRS, but unfortunately, the solution itself poses a problem. Estimating the regression coefficient will bias the estimator. Having an estimate of the coefficient prior to the experiment will eliminate the bias, but this is often not feasible. Aronow and Middleton [2013] expands the DR estimator to allow for additional covariates, but the same issue still persists. In Appendix A, we prove the discussed shortcomings of these estimators.

3 Estimation of PATE under PPS sampling

Cluster size plays an important role in efficiently estimating the PATE in CREs. Both the DIM and DR estimators give each cluster an equal chance of being selected, regardless of cluster size, but account for it during the estimation stage. Staying true to the design-based philosophy of the Neyman-Rubin model, we advise instead to change the cluster sampling scheme to probability-proportional-to-size sampling (PPS), which can accommodate varying cluster sizes when sam-

pling. Under PPS, we derive the HT estimator, which is unbiased, location-invariant, and efficient.

3.1 PPS sampling of clusters

To be precise, we define a PPS sample with s draws as any sample in which the probability of any cluster c of being sampled is $n_c s/n$. While, generally, PPS samples can be drawn with or without replacement, we focus exclusively on PPS samples drawn without replacement (PPSWOR), where the number of unique clusters sampled are fixed. This allows researchers to have greater control in designing a CRE under a budget constraint. A PPSWOR sampling scheme requires each cluster to contain no more than n/s units.

Drawing a PPSWOR sample is a deceptively unintuitive task and quite a bit of work has been devoted to efficient and/or exact selection of PPSWOR samples [Hanurav, 1967, Vijayan, 1968, Sinha, 1973, Brewer and Hanif, 1982, Berger and Till, 2009]. Unlike SRS or sampling with replacement, PPSWOR sampling schemes are not uniquely defined solely by the property that the marginal probability of sampling a cluster is $n_c s/n$. Instead, for each pair of unique clusters c, c' a PPSWOR sampling scheme requires knowing the joint probability $\pi_{cc'}$ of having both of these clusters included in the sample. To reduce variance in estimators, it is useful to choose a sampling scheme such that

$$\pi_{cc'} \geq P(S_c = 1)P(S_{c'} = 1) = n_c n_{c'} s^2 / n^2 > 0. \quad (17)$$

Sunter [1977, 1986] provides list-sequential methods for drawing an approximate PPSWOR sample of general size n satisfying (17).

3.2 Horvitz-Thompson estimator under PPS sampling

We define the *Horvitz-Thompson estimator under PPS sampling (HT-PPS)* for the population mean under treatment t as:

$$\hat{\mu}_{t,\text{HT,PPS}} = \hat{\mu}_{t,\text{HT,PPS}}(\mathbf{y}) \equiv \sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c}. \quad (18)$$

In words, this estimate is obtained by finding each cluster c that receives treatment t , computing the average response within each of these clusters, and then taking the average of these within-cluster averages. The *HT-PPS estimator for the PATE* is the difference of the HT-PPS estimator for the population mean under treatment and under control:

$$\hat{\delta}_{\text{HT-PPS}} = \hat{\mu}_{1,\text{HT-PPS}} - \hat{\mu}_{0,\text{HT-PPS}}. \quad (19)$$

Note that if a mean estimator is linear in potential outcomes, then the PATE estimator consisting of the mean estimators will be location-invariant. The HT-PPS estimator for the population mean is linear in potential outcomes, which is formally stated in the following lemma:

Lemma 1 *Suppose that clusters are sampled according to PPSWOR sampling, and suppose that treatment is symmetric across clusters. Then:*

$$\hat{\mu}_{t,\text{HT,PPS}}(a + \mathbf{y}) = a + \hat{\mu}_{t,\text{HT,PPS}}(\mathbf{y}). \quad (20)$$

The location invariance, unbiasedness, and variance of the HT-PPS estimator for PATE is then provided in the following theorem:

Theorem 2 *Suppose that clusters are sampled according to PPSWOR sampling, and suppose that*

treatment is symmetric across clusters. Then:

$$\hat{\delta}_{HT,PPS}(a + \mathbf{y}) = \hat{\delta}_{HT,PPS}(\mathbf{y}) \quad (21)$$

$$\mathbb{E} \left(\hat{\delta}_{HT,PPS} \right) = \delta, \quad (22)$$

$$\begin{aligned} \text{Var} \left(\hat{\delta}_{HT,PPS} \right) &= \sum_{t=0}^1 \left[\mathbb{E} \left(\frac{1}{\#T_t} \right) \sigma_{t,bet}^2 + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \left(\sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2 \right) \right] \\ &+ \sum_{t=0}^1 \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \left(1 - \frac{s_c}{n_c} \right) \frac{\sigma_{ct}^2}{s_c} \\ &- 2 \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{\pi_{cc'}}{s(s-1)} - \frac{n_c n_{c'}}{n^2} \right] \mu_{c1} \mu_{c'0} + 2 \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mu_{c1} \mu_{c0}. \end{aligned} \quad (23)$$

The standard error for the HT-PPS estimator of PATE is then the square root of eq. (23). A proof of the lemma and theorem is given in Appendix B.

A PPS sample naturally gives larger clusters a greater probability of being selected. Hence, the sample will be biased towards larger clusters. However, the HT-PPS estimator takes this into consideration as weights when estimating the PATE, thereby, eliminating the bias. Moreover, if the same number of units are sampled from each cluster (say, s_u), this will give each treated (controlled) unit in the population an equal probability of being sampled, which does not hold for a SRS of clusters:

$$P(S_c T_{ct} S_{kc} = 1 | \text{PPS}) = \frac{\#T_t s_u}{n} \quad (24)$$

$$P(S_c T_{ct} S_{kc} = 1 | \text{SRS}) = \frac{\#T_t s_u}{\ell n_c}. \quad (25)$$

Under this condition, then, the HT-PPS estimator and the DIM estimator (given a PPS sample of clusters) will be the same.

3.3 Variance estimator for HT-PPS estimator

Since

$$\text{Var}(\hat{\delta}) = \text{Var}(\hat{\mu}_1) + \text{Var}(\hat{\mu}_0) - 2\text{Cov}(\hat{\mu}_1, \hat{\mu}_0), \quad (26)$$

estimating each of the three components will give an estimator for the variance of the HT-PPS estimator for the PATE. The Sen-Yates-Grundy (SYG) variance estimator is an unbiased estimator for the first two parts involving the sampling variance of $\hat{\mu}_t$ [Lohr, 2010]. On the contrary, since the last term of eq. (23) requires clusters being both treated and controlled, there is no unbiased estimator for the covariance between $\hat{\mu}_1$ and $\hat{\mu}_0$. Consequently, the variance of the HT-PPS estimator cannot be unbiasedly estimated, but a conservative bound is instead provided:

$$\begin{aligned} \widehat{\text{Var}}_C(\hat{\delta}_{\text{HT,PPS}}) &= \sum_{t=0}^1 \frac{1}{2} \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{s(s-1)}{\pi_{cc'} \#T_t (\#T_t - 1)} \frac{n_c n_{c'}}{n^2} - \frac{1}{\#T_t^2} \right] S_c T_{ct} S_{c'} T_{c't} (\hat{\mu}_{ct} - \hat{\mu}_{c't})^2 \\ &\quad - 2 \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{\pi_{cc'}}{s(s-1)} - \frac{n_c n_{c'}}{n^2} \right] \frac{s(s-1)}{\pi_{cc'}} \frac{S_c S_{c'} T_{c1} T_{c'0}}{\#T_1 \#T_0} \hat{\mu}_{c1} \hat{\mu}_{c'0} \\ &\quad + \sum_{t=0}^1 \sum_{c=1}^{\ell} \frac{S_c T_{ct} n_c^2}{\#T_t n^2} \hat{\mu}_{ct}^2 \end{aligned} \quad (27)$$

where $\hat{\mu}_{ct} = \sum_{k=1}^{n_c} y_{kct} S_{kc} / s_c$ is the within-cluster sample mean for t . Note that the first term is SYG variance estimator for μ_t , and the last two terms make up the covariance bound. In appendix B.3.4, eq. (27) is shown to be positively biased for eq. (23). Taking the square root of eq. (27) will give the estimated standard error of the PATE estimator.

Estimating the variance requires knowledge about the $\pi_{cc'}$ under the specific sampling procedure used to obtain a PPS of clusters, but this is rarely given in practice. Therefore, the $\pi_{cc'}$ needs to be estimated too. This can be achieved using analytical approximations [Lohr, 2010, Berger and Till, 2009] or Monte Carlo simulations [Fattorini, 2009].

4 Allowing for stratification

Current literature recommends stratifying and/or blocking on cluster size to further reduce sampling variability. Since PPS sampling already incorporate size variation, stratifying on other prognostic cluster covariates, rather than cluster size, can drastically improve estimation. For example, villages may be stratified based on whether they are in a rural/urban environment or based on the villages' geographic region. Suppose that the population of ℓ clusters are partitioned into m strata based on a categorical cluster characteristic (or a discretized numerical one). Cluster sampling and treatment assignment is done within each stratum and independently across strata. The cluster-stratified HT-PPS estimator is then defined (without the use of indicator variables) as

$$\hat{\delta}_{\text{CS,HT,PPS}} = \sum_{u=1}^m \frac{n_u}{n} \left[\sum_{\substack{c \in u, \\ t=1}}^{\#T_1} \frac{1}{\#T_1} \sum_{k=1}^{s_c} \frac{y_{kcu1}}{s_c} - \sum_{\substack{c' \in u, \\ t=0}}^{\#T_0} \frac{1}{\#T_0} \sum_{k^*=1}^{s_{c'}} \frac{y_{k^*c'u0}}{s_{c'}} \right] \quad (28)$$

$$= \sum_{u=1}^m \frac{n_u}{n} \hat{\delta}_{u,\text{HT,PPS}} \quad (29)$$

where n_u is the population of units in stratum u . The statistical properties for the cluster-stratified HT-PPS estimator can be easily derived from Theorem 2.

Theorem 3 *Suppose clusters are first stratified. Suppose also that clusters are sampled with PPSWOR and treatments are randomized within stratum and independently across strata. Then:*

$$\mathbb{E} \left(\hat{\delta}_{\text{CS,HT,PPS}} \right) = \delta \quad (30)$$

$$\text{Var} \left(\hat{\delta}_{\text{CS,HT,PPS}} \right) = \sum_{u=1}^m \frac{n_u^2}{n^2} \text{Var} \left(\hat{\delta}_{u,\text{HT,PPS}} \right) \quad (31)$$

$$\hat{\delta}_{\text{CS,HT,PPS}}(a + \mathbf{y}) = \hat{\delta}_{\text{CS,HT,PPS}}(\mathbf{y}). \quad (32)$$

Plugging eq. (27) into eq. (31) will give a conservative estimate of the sampling variability for the cluster-stratified HT-PPS estimator.

Stratification may be applied to units within clusters instead of on clusters. In this setting, the n_c units in cluster c are divided into q_c strata with n_v units in each stratum. A SRS sample of s_v units is taken. The unit-stratified HT-PPS estimator is

$$\hat{\delta}_{\text{US,HT,PPS}} = \sum_{\substack{c=1, \\ t=1}}^{\#T_1} \frac{1}{\#T_1} \sum_{v \in c} \frac{q_c}{n_c} \sum_{k \in v} \frac{s_v}{n_v} \frac{y_{kvc1}}{s_v} - \sum_{\substack{c=1, \\ t=0}}^{\#T_0} \frac{1}{\#T_0} \sum_{v \in c} \frac{q_c}{n_c} \sum_{k \in v} \frac{s_v}{n_v} \frac{y_{kvc0}}{s_v}. \quad (33)$$

Since the stratification is on the units within a cluster, we need to only adjust the within-cluster variance in Theorem 2 to get the statistical properties for the unit-stratified HT-PPS estimator. Hence,

$$\text{Var}(\hat{\mu}_{ct}) = \frac{1}{\#T_t} \sum_{c=1}^{\ell} \frac{n_c}{n} \left(1 - \frac{s_c}{n_c}\right) \frac{\sigma_{ct}^2}{s_c} \quad (34)$$

will instead be

$$\text{Var}(\hat{\mu}_{ct}) = \frac{1}{\#T_t} \sum_{c=1}^{\ell} \frac{n_c}{n} \sum_{v \in c} \frac{q_c}{n_c} \frac{n_v^2}{n_c^2} \left(1 - \frac{s_v}{n_v}\right) \frac{\sigma_{ct}^2}{s_v}. \quad (35)$$

Similarly, eq. (27) is still a conservative estimate of the sampling variability, but $\hat{\mu}_{ct}$ will instead be the stratified estimator of the within-cluster sample mean for t . Naturally, if stratification is desired at both the cluster- and unit-levels, combining eq. (28) and eq. (33) will give an unbiased estimator for the PATE.

5 Data example

Beath et al. [2013] perform an experiment in Afghanistan to investigate whether development programs with mandatory women contribution can change villagers' perspectives on women's political participation. Five hundred villages, ranging from sizes 60 to 9000, were sampled and matched into pairs. Within each pair, one village is randomly assigned to receive the National Solidarity Programme (NSP), and the other village serves as a control to receive the NSP after the experiment. The NSP creates a community development council and provides grants for village development

projects. The council is then responsible for distributing the grants among the projects. However, the NSP stipulates that half of the council must be women and at least one of the projects must be a priority for the women. After two years, ten head-of-household men and ten head-of-household women from each village are selected for a follow-up survey. Respondents are asked whether women should have equal decision making in the village council.

We perform Monte Carlo simulations to compare the HT-PPS estimator to its SRS counterparts. We generate the potential outcomes from a fitted LOWESS line of the NSP data and then mimic a simplified experiment in which clusters are randomly sampled using either PPS or SRS. The R package *TeachingSampling* is used to perform Sunter’s PPSWOR sampling. We vary the number of sampled clusters from 20 to 200. Treatments are assigned completely at random to the sampled clusters. For ease, we fix the number of treated clusters to be half of the sampled clusters, but this will not drastically change the theoretical results.

The PATE is then estimated with the HT-PPS, DIM, HT-SRS, biased DR, and Hájek estimators. For the DR estimator, θ is optimally estimated as described in Middleton and Aronow [2015] using the simulated sample data. The Hájek estimator (see [Hajek, 1971]) is a ratio estimator similar to the DIM. It estimates the population mean for treatment t as a ratio of the estimated treated (controlled) cluster total over the total number of treated (controlled) units in the sampled clusters. The other estimators are as described in section 2.7.

Figure 1 compares the MSE of the estimators as the number of sampled clusters are varied. To get an exact PPSWOR sample of the NSP data, the number of sampled clusters can be at most 45, and thus, the samples of sizes 60 and up are only approximately PPS. Even so, the HT-PPS estimator performs best out of all the estimators, including the omnipresent DIM, across all sample sizes of clusters. Figure 2 gives a more thorough comparison of the sampling distributions for the PATE estimators when 40 clusters are sampled.

In addition, Monte Carlo simulations are done to examine the performance of estimating the sampling variance of the HT-PPS estimator. The $\pi_{cc'}$ are estimated using analytical approxima-

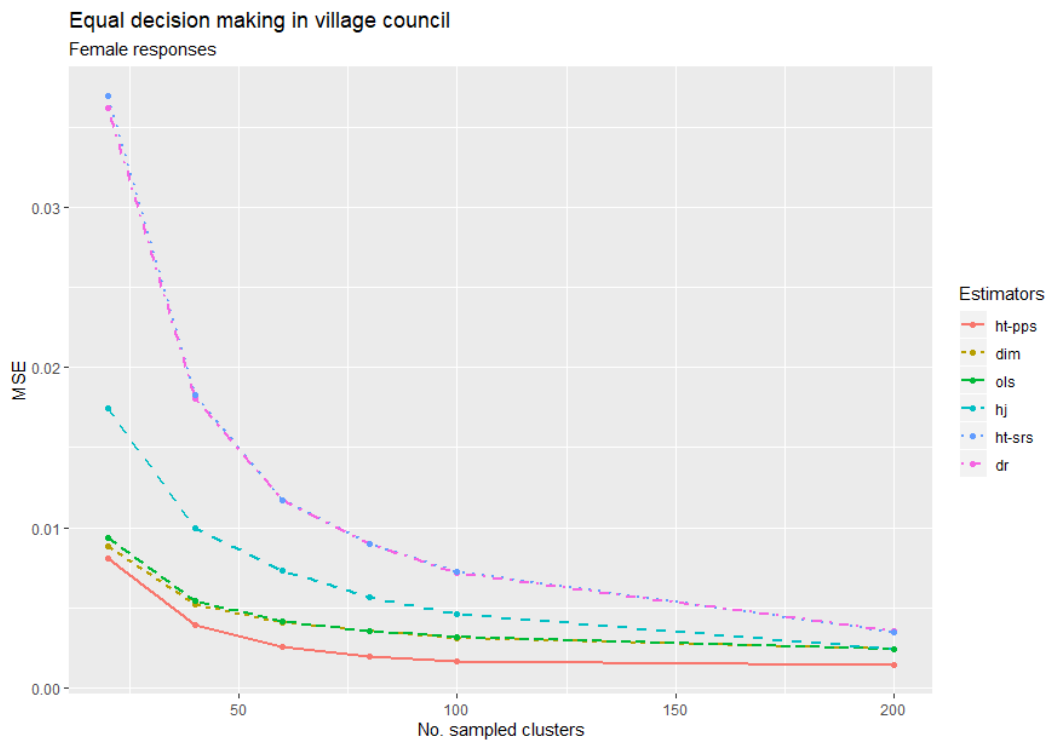
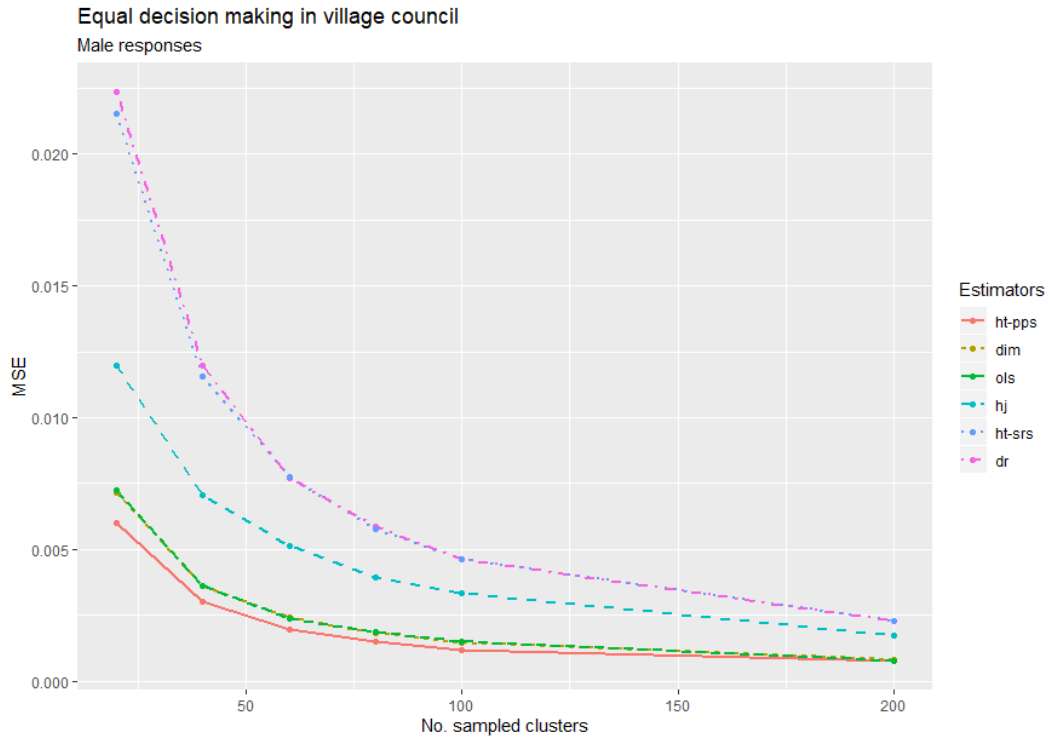


Figure 1: The number of sampled clusters are 20, 40, 60, 80, 100, and 200. Results are based on 10,000 simulations. Our estimator, HT-PPS (red and solid), performs better than the SRS estimators.

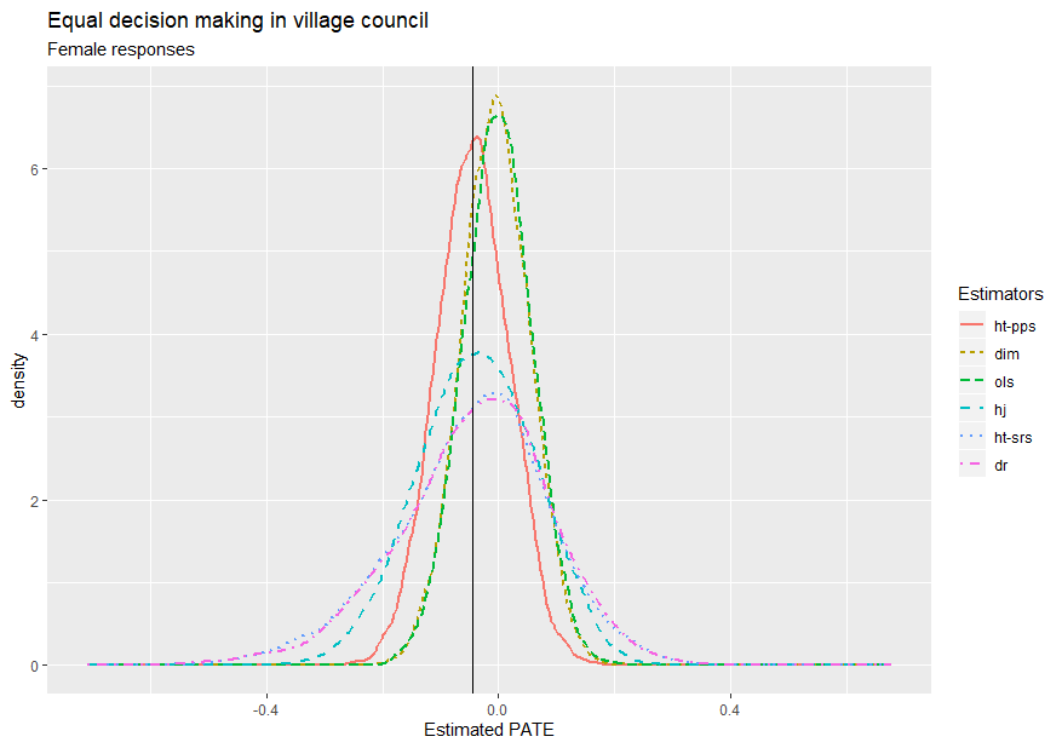
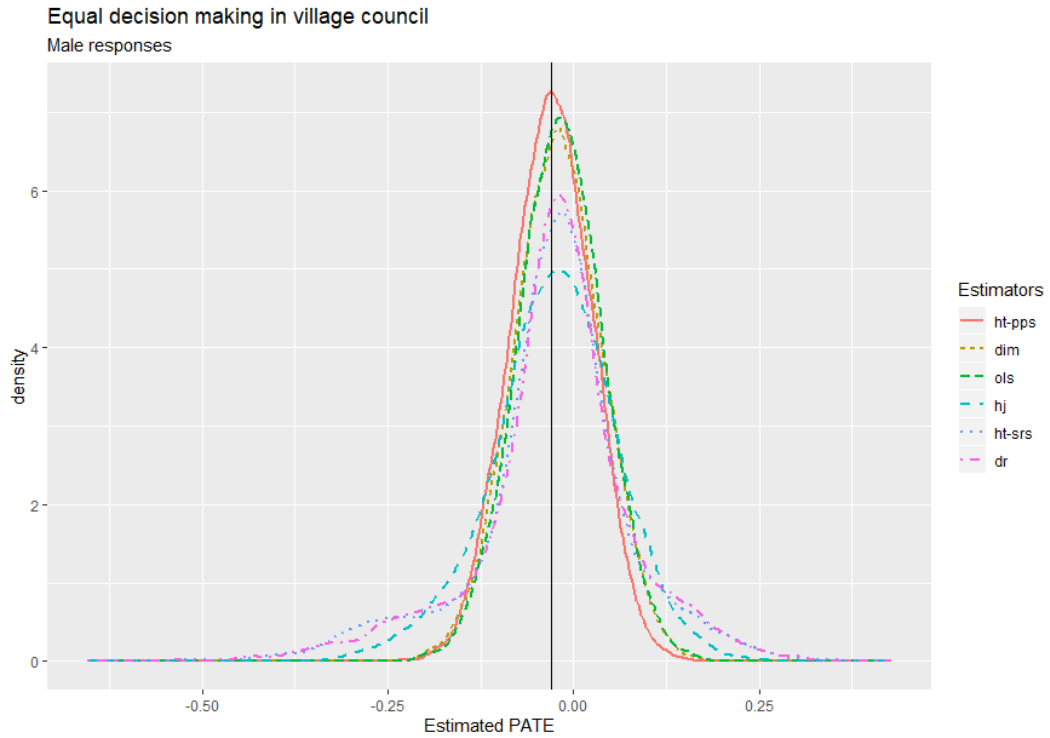


Figure 2: Results based on 10,000 simulations of sampling 40 clusters. The solid vertical line is the PATE (-0.0302 for male, -0.0448 for female). Our estimator, HT-PPS (red and solid), is unbiased and as efficient as the DIM.

tions. Table 1 provides statistics (estimated variance, bias, and true sampling variability) on the variance estimation as the number of sampled clusters are varied. Note that the bias is all positive so estimates are conservative. Estimates are also close to the true variance.

No. sampled clusters	Male responses			Female responses		
	Est. var.	Bias	Samp. var.	Est. var.	Bias	Samp. var.
20	0.0061	6.47E-05	6.37E-06	0.0082	2.48E-04	4.76E-06
40	0.0031	7.15E-05	7.39E-07	0.0041	1.9E-04	5.23E-07
60 ¹	0.0034	1.45E-03	4.83E-07	0.0069	4.29E-03	6.68E-07
80 ¹	0.0026	1.13E-03	2.06E-07	0.0055	3.56E-03	2.84E-07
100 ¹	0.0021	9.44E-04	1.07E-07	0.0046	3.07E-03	1.47E-07
200 ¹	0.0012	6.08E-04	1.31E-08	0.0026	1.92E-03	1.79E-08

Table 1: Results based on 10,000 simulations. Estimated variances are upwardly biased, but the bias is marginally small.

6 Conclusion

Experiments are the “gold standard” for investigating causal relationships, but traditionally, the causal inference is limited to the convenience sample recruited for the experiment. Often, though, researchers prefer to generalize to individuals beyond those in the sample. This then requires a random sample from the population of interest. Since populations are naturally structured in groups, it is easier to sampled groups, rather than individuals, to be in experiments; thus, cluster randomized experiments are a fitting design choice.

On the other hand, the multi-level constitution of CREs poses analytical adversities. Much of the difficulties arises from unequal cluster sizes. If clusters contain the same number of units, all estimators discussed would be the same, and the idea of choosing the “best” estimator would be nonexistent. However, varying cluster sizes are intrinsic to CREs. Hence, in this paper, we account for cluster sizes by sampling clusters with probability proportional to size. Estimating PATE can then be done with the HT-PPS estimator. The HT-PPS estimator is an attractive alternative to SRS-

¹Estimated variances uses the with-replacement variance estimator since samples are not exactly PPS.

based estimators since it is intuitive, unbiased, efficient, and location-invariant. We also derive a conservative variance estimator for the sampling variability of the HT-PPS estimator.

Stratification and blocking can still be used to further reduce the sampling variability, but with PPS sampling, other more important covariates can be used instead of cluster size. We have done some work on how stratification may affect the HT-PPS estimator, but we plan to expand on it. We also plan on extending our results to include blocking too.

References

- P. M. Aronow and J. Middleton. A class of unbiased estimators of the average treatment effect in randomized experiments. *Journal of Causal Inference*, 1(1), 2013. doi: 10.1515/jci-2012-0009.
- A. Beath, F. Christina, and R. Enikolopov. Empowering women through development aid: Evidence from a field experiment in Afghanistan. *American Political Science Review*, 107(03): 540 – 557, 2013. doi: 10.1017/S0003055413000270. URL <http://www.jstor.org/stable/43654923>.
- Y. G. Berger and Y. Till. *Sampling with Unequal Probabilities*, volume 29A, chapter 2, pages 39 – 54. Elsevier B. V., 2009.
- K. R. W. Brewer and M. Hanif. *Sampling with Unequal Probabilities (Lecture Notes in Statistics)*. Springer, 1982.
- M. L. Bruce, T. R. T. Have, C. F. Reynolds, H. C. Schulberg, B. H. Mulsant, G. K. Brown, G. J. McAvay, J. L. Pearson, and G. S. Alexopoulos. Reducing suicidal ideation and depression symptoms in depressed older primary care patients: A randomized controlled trial. *JAMA*, 291(09):1081 – 1091, 2004. doi: 10.1001/jama.291.9.1081. URL <https://jamanetwork.com/journals/jama/fullarticle/198310>.
- W. G. Cochran. *Sampling Techniques*. Wiley, 1977.
- J. Cornfield. Randomized by group: A formal analysis. *American Journal of Epidemiology*, 108(02):100 – 102, 1978. doi: 10.1093/oxfordjournals.aje.a112592.
- A. Donner. Some aspects of the design and analysis of cluster randomization trials. *Journal of the Royal Statistical Society*, 47(01):95 – 113, 1998. doi: 10.1111/1467-9876.00100.
- A. Donner and N. Klar. Pitfalls of and controversies in cluster randomization trials. *American Journal of Public Health*, 94(03):416 – 422, 2004. doi: 10.2105/AJPH.94.3.416.

- Lorenzo Fattorini. An adaptive algorithm for estimating inclusion probabilities and performing the Horvitz–Thompson criterion in complex designs. *Computational Statistics*, 24(4):623, Mar 2009. ISSN 1613-9658. doi: 10.1007/s00180-009-0149-9. URL <https://doi.org/10.1007/s00180-009-0149-9>.
- M. H. Gail, S. D. Mark, R. J. Carroll, S. B. Green, and D. Pee. On design considerations and randomization-based inference for community intervention trials. *Statistics in Medicine*, 15(11): 1069 – 1092, 1996. doi: 10.1002/(SICI)1097-0258(19960615)15:11<1069::AID-SIM220>3.0.CO;2-Q.
- A. S. Gerber and D. P. Green. The effects of canvassing, telephone calls, and direct mail on voter turnout: A field experiment. *American Political Science Review*, 94(3):653 – 663, September 2000. doi: 10.2307/2585837.
- J. Hajek. Discussion of 'An essay on the logical foundations of survey sampling, part one,' by D. Basu. In V. P. Godambe and D. A. Sprout, editors, *Foundations of Statistical Inference*, page 236. Holt, Rhinehart, and Winston, 1971.
- B. B. Hansen, P. R. Rosenbaum, and D. S. Small. Clustered treatment assignments and sensitivity to unmeasured biases in observational studies. *Journal of the American Statistical Society*, 109(505), 2014. doi: 10.1080/01621459.2013.863157.
- M. H. Hansen and W. N. Hurwitz. On the theory of sampling from finite populations. *The Annals of Mathematical Statistics*, 14(04):333 – 362, 1943. doi: 10.1214/aoms/1177731356.
- T. V. Hanurav. Optimum utilization of auxiliary information: π ps sampling of two units from a stratum. *Journal of the Royal Statistical Society*, pages 374 – 391, 1967.
- R. Hayes and L. Moulton. *Cluster Randomised Trials*. Chapman & Hall/CRC, 2009.
- P. W. Holland. Statistics and causal inference. *American Statistical Association*, 81(396):945 – 960, 1986.
- D. G. Horvitz and D. J. Thompson. A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, 47(260):663 – 685, 1952. doi: 10.1080/01621459.1952.10483446.
- K. Imai, G. King, and C. Nall. The essential role of pair matching in cluster-randomized experiments, with application to the Mexican universal health insurance evaluation. *Statistical Science*, 24(01):29 – 53, 2009. doi: 10.1214/08-STS274.
- G. W. Imbens. Experimental design for unit and cluster randomized trials. 2011. URL http://cyrussamii.com/wp-content/uploads/2011/06/Imbens_June_8_paper.pdf.
- S. M. Kerry, F. P. Cappuccio, L. Emmett, J. Plange-Rhule, and J. B. Eastwood. Reducing selection bias in a cluster randomized trial in West African villages. *Clinical Trials*, 02(02):125 – 129, 2005. doi: 10.1191/1740774505cn074oa.

- G. King, E. Gakidou, N. Ravishankar, R. T. Moore, J. Lakin, M. Vargas, M. M. Tllez-Rojo, vila J. E. H., M. H. vila, and H. H. Llamas. A “politically robust” experiment design for public policy evaluation, with application to the Mexican universal health insurance program. *Journal of Policy Analysis and Management*, 26(03):479 – 506, 2007. doi: 10.1002/pam.20279.
- J. D. Lewsey. Comparing completely and stratified randomized design in cluster randomized trials when the stratifying factor is the cluster size: A simulation study. *Statistics in Medicine*, 23(06): 897 – 905, 2004. doi: 10.1002/sim.1665.
- S. L. Lohr. *Sampling: Design and Analysis*. Duxbury Press, 2nd edition, 2010.
- J. A. Middleton and P. M. Aronow. Unbiased estimation of the average treatment effect in cluster-randomized experiments. *Statistics, Politics and Policy*, 6:39 – 75, 2015. doi: 10.1515/spp-2013-0002.
- L. W. Miratrix, J. S. Sekhon, and B. Yu. Adjusting treatment effect estimates by post-stratification in randomized experiments. *Journal of Royal Statistical Society*, 75(02):369 – 396, 2013. doi: 10.1111/j.1467-9868.2012.01048.x.
- E. L. Paluck. Reducing intergroup prejudice and conflict using the media: A field experiment in Rwanda. *Journal of Personality and Social Psychology*, 96(03):574, 2009. doi: 10.1037/a0011989.
- C. C. Raver, S. M. Jones, C. Li-Grining, F. Zhai, M. W. Metzger, and B. Soloman. Targeting children’s behavior problems in preschool classrooms: A cluster-randomized controlled trial. *Journal of Consulting and Clinical Psychology*, 77(02):302, 2009. doi: 10.1037/a0015302.
- D. B. Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(05):688, 1974. doi: 10.1037/h0037350.
- B. K. Sinha. On sampling schemes to realize preassigned sets of inclusion probabilities of first two orders. *Calcutta Statistical Association Bulletin*, 22:89 – 110, 1973. doi: 10.1177/0008068319740103.
- D. S. Small, T. R. T. Have, and P. R. Rosenbaum. Randomization inference in a group-randomized trial of treatments for depression: Covariate adjustment, noncompliance, and quantile effects. *Journal of the American Statistical Association*, 103(481):271 – 279, 2008. doi: 10.1198/016214507000000897.
- Jerzy Splawa-Neyman, D. M. Dabrowska, and T. P. Speed. On the application of probability theory to agricultural experiments. essay on principles. section 9. *Statistical Science*, 5(4):465–472, 1923. ISSN 08834237. URL <http://www.jstor.org/stable/2245382>. (Translated in 1990).
- A. Sunter. Solutions to the problem of unequal probability sampling without replacement. *International Statistical Review*, 54(01):33 – 50, 1986. doi: 10.2307/1403257.

- A. B. Sunter. List sequential sampling with equal or unequal probabilities without replacement. *Applied Statistics*, 26(3):261–268, 1977.
- K. Vijayan. An exact π ps sampling scheme-generalization of a method of Hanurav. *Journal of the Royal Statistical Society*, pages 556 – 566, 1968.
- L. Wantchekon. Clientelism and voting behavior: Evidence from a field experiment in Benin. *World Politics*, 55(03):399 – 422, 2003. doi: 10.1353/wp.2003.0018.

A Properties of the SRS estimators

A.1 Horvitz-Thompson estimator

The *Horvitz-Thompson (HT-SRS) estimator for the population mean under treatment t* is defined as:

$$\hat{\mu}_{t,\text{HT,SRS}} = \ell \sum_{c=1}^{\ell} \frac{S_c T_{ct} n_c}{\#T_t} \frac{1}{n} \sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c}. \quad (36)$$

The *HT-SRS estimator for the PATE* is then the difference of the HT-SRS estimator for the population mean under treatment and under control:

$$\hat{\delta}_{\text{HT,SRS}} = \hat{\mu}_{1,\text{HT,SRS}} - \hat{\mu}_{0,\text{HT,SRS}}. \quad (37)$$

The HT-SRS estimator for PATE is not location-invariant since

$$\begin{aligned} \hat{\delta}_{\text{HT,SRS}}(a + \mathbf{y}) &= \hat{\mu}_{1,\text{HT,SRS}}(a + \mathbf{y}) - \hat{\mu}_{0,\text{HT,SRS}}(a + \mathbf{y}) \\ &= \ell \sum_{c=1}^{\ell} \frac{S_c T_{c1} n_c}{\#T_1} \frac{1}{n} \sum_{k=1}^{n_c} \frac{(a + y_{kc1}) S_{kc}}{s_c} - \ell \sum_{c=1}^{\ell} \frac{S_c T_{c0} n_c}{\#T_0} \frac{1}{n} \sum_{k=1}^{n_c} \frac{(a + y_{kc0}) S_{kc}}{s_c} \\ &= a \left(\ell \sum_{c=1}^{\ell} \frac{S_c T_{c1} n_c}{\#T_1} \frac{1}{n} \sum_{k=1}^{n_c} \frac{S_{kc}}{s_c} - \ell \sum_{c=1}^{\ell} \frac{S_c T_{c0} n_c}{\#T_0} \frac{1}{n} \sum_{k=1}^{n_c} \frac{S_{kc}}{s_c} \right) \\ &\quad + \left(\ell \sum_{c=1}^{\ell} \frac{S_c T_{c1} n_c}{\#T_1} \frac{1}{n} \sum_{k=1}^{n_c} \frac{y_{kc1} S_{kc}}{s_c} - \ell \sum_{c=1}^{\ell} \frac{S_c T_{c0} n_c}{\#T_0} \frac{1}{n} \sum_{k=1}^{n_c} \frac{y_{kc0} S_{kc}}{s_c} \right) \\ &= a \left(\frac{\ell \#N_1}{n \#T_1} - \frac{\ell \#N_0}{n \#T_0} \right) + \hat{\delta}_{\text{HT,SRS}}(\mathbf{y}), \end{aligned} \quad (38)$$

where

$$\#N_t = \sum_{c=1}^{\ell} S_c T_{ct} n_c \quad (39)$$

represents all units given treatment t in the sampled clusters. Lack of location-invariance will not affect the unbiasedness of the estimator, even for linearly transformed outcomes, since

$$\mathbb{E}(\#N_t) = \sum_{c=1}^{\ell} \mathbb{E}(S_c T_{ct} n_c) = \frac{n \#T_t}{\ell}. \quad (40)$$

However, this problem presents itself in the variance calculation:

$$\begin{aligned} \text{Var} \left(\hat{\delta}_{\text{HT,SRS}}(a + \mathbf{y}) \right) &= a^2 \left(\frac{\ell}{n} \right)^2 \left[\text{Var} \left(\frac{\#N_1}{\#T_1} \right) + \text{Var} \left(\frac{\#N_0}{\#T_0} \right) - 2 \text{Cov} \left(\frac{\#N_1}{\#T_1}, \frac{\#N_0}{\#T_0} \right) \right] \\ &\quad + 2a \frac{\ell}{n} \left[\text{Cov} \left(\frac{\#N_1}{\#T_1}, \hat{\delta}_{\text{HT,SRS}} \right) - \text{Cov} \left(\frac{\#N_0}{\#T_0}, \hat{\delta}_{\text{HT,SRS}} \right) \right] \\ &\quad + \text{Var}(\hat{\delta}_{\text{HT,SRS}}). \end{aligned} \quad (41)$$

A.2 Difference-in-means estimator

The *sample mean under treatment t* is:

$$\hat{\mu}_{t,\text{DIM,SRS}} \equiv \frac{\sum_{c=1}^{\ell} S_c T_{ct} \sum_{k=1}^{n_c} y_{kct} S_{kc}}{\sum_{c=1}^{\ell} S_c T_{ct} S_c}. \quad (42)$$

The *difference-in-means (DIM) estimator for the PATE* is the difference of the sample means under treatment and under control:

$$\hat{\delta}_{\text{DIM,SRS}} = \hat{\mu}_{1,\text{DIM,SRS}} - \hat{\mu}_{0,\text{DIM,SRS}}. \quad (43)$$

Using the relationship

$$\mathbb{E} \left(\frac{u}{v} \right) = \frac{1}{\mathbb{E}(v)} \left[\mathbb{E}(u) - \text{Cov} \left(\frac{u}{v}, v \right) \right], \quad (44)$$

it can be shown that the DIM estimator is actually estimating the quantity

$$\mathbb{E}(\hat{\mu}_{t,\text{DIM},\text{SRS}}) = \frac{1}{\sum_{c=1}^{\ell} s_c} \left[\sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \frac{s_c}{n_c} y_{kct} - \ell \text{Cov} \left(\hat{\mu}_{t,\text{DIM},\text{SRS}}, \sum_{c=1}^{\ell} \frac{S_c T_{ct} s_c}{\#T_t} \right) \right]. \quad (45)$$

A.3 Des Raj estimate of PATE under SRS

Middleton and Aronow [2015] advocate the Des Raj (DR) estimator, and they define the *DR estimator for the population mean under treatment t* as:

$$\hat{\mu}_{t,\text{DR},\text{SRS}} = \ell \sum_{c=1}^{\ell} \frac{S_c T_{ct} n_c}{\#T_t n} \left[\sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c} - \frac{\theta}{n_c} \left(n_c - \frac{n}{\ell} \right) \right] \quad (46)$$

where θ is a regression coefficient. *The DR estimator for the PATE* is then the difference of the DR estimator for the population mean under treatment and under control:

$$\hat{\delta}_{\text{DR},\text{SRS}} = \hat{\mu}_{1,\text{DR},\text{SRS}} - \hat{\mu}_{0,\text{DR},\text{SRS}}. \quad (47)$$

We show here that the DR estimator is biased when the regression coefficient θ needs to be estimated:

$$\begin{aligned} \mathbb{E}(\hat{\delta}_{\text{DR},\text{SRS}}) &= \delta - \ell \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Cov} \left(\frac{S_c T_{c1}}{\#T_1}, \hat{\theta} \right) \\ &\quad + \ell \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Cov} \left(\frac{S_c T_{c0}}{\#T_0}, \hat{\theta} \right). \end{aligned} \quad (48)$$

B Properties of HT-PPS estimator

We prove the results for Lemma 1 and Theorem 2. For more detailed derivations, we provide an expanded supplemental appendix.

B.1 Linearity in potential outcomes for HT-PPS mean estimator

$$\begin{aligned}
\hat{\mu}_{t,\text{HT,PPS}}(a + \mathbf{y}) &= \sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{(a + y_{kct}) S_{kc}}{s_c} \\
&= a \sum_{c=1}^{\ell} \frac{S_c T_{c1}}{\#T_1} \sum_{k=1}^{n_c} \frac{S_{kc}}{s_c} + \sum_{c=1}^{\ell} \frac{S_c T_{c1}}{\#T_1} \sum_{k=1}^{n_c} \frac{y_{kc1} S_{kc}}{s_c} \\
&= a + \hat{\mu}_{t,\text{HT,PPS}}(\mathbf{y})
\end{aligned} \tag{49}$$

B.2 Expectation of HT-PPS estimator for PATE

$$\begin{aligned}
\mathbb{E}(\hat{\delta}_{\text{HT-PPS}}) &= \mathbb{E} \left(\sum_{c=1}^{\ell} \frac{S_c T_{c1}}{\#T_1} \sum_{k=1}^{n_c} \frac{y_{kc1} S_{kc}}{s_c} - \sum_{c=1}^{\ell} \frac{S_{c'} T_{c'0}}{\#T_0} \sum_{k^*=1}^{n_{c'}} \frac{y_{k^*c'0} S_{k^*c'}}{s_{c'}} \right) \\
&= \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \frac{y_{kc1}}{s_c} \mathbb{E} \left(S_c \mathbb{E} \left(\frac{T_{c1}}{\#T_1} \middle| \mathbf{S} \right) \mathbb{E}(S_{kc} | \mathbf{S}) \right) \\
&\quad - \sum_{c'=1}^{\ell} \sum_{k^*=1}^{n_{c'}} \frac{y_{k^*c'0}}{s_{c'}} \mathbb{E} \left(S_{c'} \mathbb{E} \left(\frac{T_{c'0}}{\#T_0} \middle| \mathbf{S} \right) \mathbb{E}(S_{k^*c'} | \mathbf{S}) \right) \\
&= \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \frac{y_{kc1}}{n_c s} \mathbb{E}(S_c) - \sum_{c'=1}^{\ell} \sum_{k^*=1}^{n_{c'}} \frac{y_{k^*c'0}}{n_{c'} s} \mathbb{E}(S_{c'}) \\
&= \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \frac{y_{kc1}}{n} - \sum_{c'=1}^{\ell} \sum_{k^*=1}^{n_{c'}} \frac{y_{k^*c'0}}{n} \\
&= \mu_1 - \mu_0 = \delta.
\end{aligned} \tag{50}$$

B.3 Variance of HT-PPS estimator for PATE

From the property

$$\text{Var}(\hat{\delta}) = \text{Var}(\hat{\mu}_1 - \hat{\mu}_0) = \text{Var}(\hat{\mu}_1) + \text{Var}(\hat{\mu}_0) - 2\text{Cov}(\hat{\mu}_1, \hat{\mu}_0). \quad (51)$$

each term is expanded upon to derive the variance of the HT-PPS estimator for PATE and obtain a variance estimator.

B.3.1 Variance of HT-PPS estimator for population mean

Using the law of total variance,

$$\begin{aligned} \text{Var}(\hat{\mu}_{t,\text{HT,PPS}}) &= \text{Var} \left(\sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c} \right) \\ &= \text{Var} \left[\mathbb{E} \left(\sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c} \middle| \mathbf{S}, \mathbf{T} \right) \right] \\ &\quad + \mathbb{E} \left[\text{Var} \left(\sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c} \middle| \mathbf{S}, \mathbf{T} \right) \right]. \end{aligned} \quad (52)$$

The first terms can be further simplified:

$$\begin{aligned} &\text{Var} \left[\mathbb{E} \left(\sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c} \middle| \mathbf{S}, \mathbf{T} \right) \right] \\ &= \sum_{c=1}^{\ell} \mu_{ct}^2 \text{Var} \left(\frac{S_c T_{ct}}{\#T_t} \right) + \sum_{c=1}^{\ell} \sum_{c' \neq c} \mu_{ct} \mu_{c't} \text{Cov} \left(\frac{S_c T_{ct}}{\#T_t}, \frac{S_{c'} T_{c't}}{\#T_t} \right) \\ &= \sum_{c=1}^{\ell} \mu_{ct}^2 \left[\mathbb{E} \left(\frac{S_c T_{ct}}{\#T_t} \right) - \mathbb{E} \left(\frac{S_c T_{ct}}{\#T_t} \right)^2 \right] \\ &\quad + \sum_{c=1}^{\ell} \sum_{c' \neq c} \mu_{ct} \mu_{c't} \left[\mathbb{E} \left(\frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t^2} \right) - \mathbb{E} \left(\frac{S_c T_{ct}}{\#T_t} \right) \mathbb{E} \left(\frac{S_{c'} T_{c't}}{\#T_t} \right) \right] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \mu_{ct}^2 - \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{n_c^2}{n^2} \mu_{ct}^2 \\
&\quad + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{n_c^2}{n^2} \mu_{ct} \mu_{c't} \\
&= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sigma_{t,bet}^2 + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \left[\sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2 \right] \tag{53}
\end{aligned}$$

where $\sigma_{t,bet}^2$ is the weighted variance of cluster means. Simplifying the second term:

$$\begin{aligned}
&\mathbb{E} \left[\text{Var} \left(\sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c} \middle| \mathbf{S}, \mathbf{T} \right) \right] = \sum_{c=1}^{\ell} \text{Var}(\hat{\mu}_{ct} | \mathbf{S}, \mathbf{T}) \mathbb{E} \left(\frac{S_c T_{ct}}{\#T_t^2} \right) \\
&= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Var}(\hat{\mu}_{ct}) \\
&= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \left(1 - \frac{s_c}{n_c} \right) \frac{\sigma_{ct}^2}{s_c}. \tag{54}
\end{aligned}$$

The variance for the HT-PPS mean estimator is

$$\begin{aligned}
\text{Var}(\hat{\mu}_{t,\text{HT,SRS}}) &= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \mu_{ct}^2 + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2 \\
&\quad + \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Var}(\hat{\mu}_{ct}) \tag{55}
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sigma_{t,bet}^2 + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \left[\sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2 \right] \\
&\quad + \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \left(1 - \frac{s_c}{n_c} \right) \frac{\sigma_{ct}^2}{s_c}. \tag{56}
\end{aligned}$$

B.3.2 SYG estimator for variance

The SYG variance estimator is

$$\begin{aligned}\widehat{\text{Var}}(\hat{\mu}_t) &= \frac{1}{2} \sum_{c=1}^{\ell} \sum_{c \neq c'} \left[\frac{s(s-1)}{\pi_{cc'} \#T_t (\#T_t - 1)} \frac{n_c n_{c'}}{n^2} - \frac{1}{\#T_t^2} \right] S_c T_{ct} S_{c'} T_{c't} (\hat{\mu}_{ct} - \hat{\mu}_{c't})^2 \\ &\quad + \sum_{c=1}^{\ell} \frac{S_c T_{ct} n_c}{\#T_t n} \widehat{\text{Var}}(\hat{\mu}_{ct})\end{aligned}\tag{57}$$

where

$$\widehat{\text{Var}}(\hat{\mu}_{ct}) = \left(1 - \frac{s_c}{n_c}\right) \frac{\hat{\sigma}_{ct}^2}{s_c}.\tag{58}$$

The $\hat{\sigma}_{ct}^2$ is the sample variance of outcomes, which is unbiased for the population variance σ_{ct}^2 . We will now show that the SYG variance is unbiased for $\text{Var}(\hat{\mu}_t)$.

$$\begin{aligned}&\mathbb{E} \left(\widehat{\text{Var}}(\hat{\mu}_t) \right) \\ &= \mathbb{E} \left(\mathbb{E} \left(\frac{1}{2} \sum_{c=1}^{\ell} \sum_{c \neq c'} \left[\frac{s(s-1)}{\pi_{cc'}} \frac{n_c n_{c'}}{n^2} \frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t (\#T_t - 1)} - \frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t^2} \right] [\hat{\mu}_{ct} - \hat{\mu}_{c't}]^2 \middle| \mathbf{S}, \mathbf{T} \right) \right) \\ &\quad + \mathbb{E} \left(\mathbb{E} \left(\sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{ct}}{\#T_t} \widehat{\text{Var}}(\hat{\mu}_{ct}) \middle| \mathbf{S}, \mathbf{T} \right) \right) \\ &= \sum_{c=1}^{\ell} \sum_{c \neq c'} \left[\frac{s(s-1)}{\pi_{cc'}} \frac{n_c n_{c'}}{n^2} \mathbb{E} \left(\frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t (\#T_t - 1)} \right) - \mathbb{E} \left(\frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t^2} \right) \right] [\mu_{ct}^2 + \text{Var}(\hat{\mu}_{ct}) - \mu_{ct} \mu_{c't}] \\ &\quad + \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Var}(\hat{\mu}_{ct}) \mathbb{E} \left(\frac{S_c T_{ct}}{\#T_t} \right) \\ &= \sum_{c=1}^{\ell} \sum_{c \neq c'} \left[\frac{s(s-1)}{\pi_{cc'}} \frac{n_c n_{c'}}{n^2} \mathbb{E} \left(\frac{\mathbb{E}(S_c S_{c'} T_{ct} T_{c't} | \#T_t)}{\#T_t (\#T_t - 1)} \right) - \mathbb{E} \left(\frac{\mathbb{E}(S_c S_{c'} T_{ct} T_{c't} | \#T_t)}{\#T_t^2} \right) \right] \\ &\quad \cdot [\mu_{ct}^2 + \text{Var}(\hat{\mu}_{ct}) - \mu_{ct} \mu_{c't}] + \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Var}(\hat{\mu}_{ct}) \mathbb{E} \left(\frac{\mathbb{E}(S_c T_{ct} | \#T_t)}{\#T_t} \right) \\ &= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \mu_{ct}^2 + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \sum_{c \neq c'} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2\end{aligned}$$

$$+ \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Var}(\hat{\mu}_{ct}). \quad (59)$$

This is equal to eq. (55).

B.3.3 Covariance of HT-PPS estimator for population means

Note that:

$$\begin{aligned} \hat{\mu}_{1,\text{HT,PPS}} \hat{\mu}_{0,\text{HT,PPS}} &= \left(\sum_{c=1}^{\ell} \frac{S_c T_{c1}}{\#T_1} \sum_{k=1}^{n_c} \frac{y_{kc1} S_{kc}}{s_c} \right) \left(\sum_{c'=1}^{\ell} \frac{S_{c'} T_{c'0}}{\#T_0} \sum_{k^*=1}^{n_{c'}} \frac{y_{k^*c'0} S_{k^*c'}}{s_{c'}} \right) \\ &= \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \sum_{c' \neq c}^{n_{c'}} \sum_{k^*=1}^{n_{c'}} \frac{y_{kc1} y_{k^*c'0}}{s_c s_{c'}} \frac{S_c T_{c1} S_{c'} T_{c'0} S_{k^*c'}}{\#T_1 \#T_0}. \end{aligned} \quad (60)$$

Therefore,

$$\begin{aligned} \text{Cov}(\hat{\mu}_{1,\text{HT,PPS}}, \hat{\mu}_{0,\text{HT,PPS}}) &= \mathbb{E}(\hat{\mu}_{1,\text{HT,PPS}} \hat{\mu}_{0,\text{HT,PPS}}) - \mathbb{E}(\hat{\mu}_{1,\text{HT,PPS}}) \mathbb{E}(\hat{\mu}_{0,\text{HT,PPS}}) \\ &= \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \sum_{c' \neq c}^{n_{c'}} \sum_{k^*=1}^{n_{c'}} \frac{y_{kc1} y_{k^*c'0}}{s_c s_{c'}} \mathbb{E} \left[S_c S_{c'} \mathbb{E} \left(\frac{T_{c1} T_{c'0}}{\#T_1 \#T_0} \middle| \mathbf{S} \right) \mathbb{E}(S_{kc} S_{k^*c'} | \mathbf{S}) \right] - \mu_{c1} \mu_{c'0} \\ &= \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{\pi_{cc'}}{s(s-1)} - \frac{n_c n_{c'}}{n^2} \right] \mu_{c1} \mu_{c'0} - \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mu_{c1} \mu_{c0}. \end{aligned} \quad (61)$$

B.3.4 Covariance bound

The covariance is bounded by

$$\begin{aligned} \widehat{\text{Cov}}_C(\hat{\mu}_1, \hat{\mu}_0) &= \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[1 - \frac{n_c n_{c'}}{n^2} \frac{s(s-1)}{\pi_{cc'}} \right] \frac{S_c T_{ct} S_{c'} T_{c't}}{\#T_1 \#T_0} \hat{\mu}_{c1} \hat{\mu}_{c'0} \\ &\quad - \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c1}}{\#T_1} \hat{\mu}_{c1}^2 - \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c0}}{\#T_0} \hat{\mu}_{c0}^2 \\ &\quad + \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c1}}{\#T_1} \widehat{\text{Var}}(\hat{\mu}_{c1}) + \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c0}}{\#T_0} \widehat{\text{Var}}(\hat{\mu}_{c0}). \end{aligned} \quad (62)$$

Taking expectation:

$$\begin{aligned}
\mathbb{E} \left(\widehat{\text{Cov}}_C(\hat{\mu}_1, \hat{\mu}_0) \right) &= \mathbb{E} \left[\mathbb{E} \left(\sum_{c=1}^{\ell} \sum_{c' \neq c} \left[1 - \frac{n_c n_{c'}}{n^2} \frac{s(s-1)}{\pi_{cc'}} \right] \frac{S_c T_{ct} S_{c'} T_{c't}}{\#T_1 \#T_0} \hat{\mu}_{c1} \hat{\mu}_{c'0} \middle| \mathbf{S}, \mathbf{T} \right) \right] \\
&\quad - \mathbb{E} \left[\mathbb{E} \left(\frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c1}}{\#T_1} \hat{\mu}_{c1}^2 \middle| \mathbf{S}, \mathbf{T} \right) \right] - \mathbb{E} \left[\mathbb{E} \left(\frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c0}}{\#T_0} \hat{\mu}_{c0}^2 \middle| \mathbf{S}, \mathbf{T} \right) \right] \\
&\quad + \mathbb{E} \left[\mathbb{E} \left(\frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c1}}{\#T_1} \widehat{\text{Var}}(\hat{\mu}_{c1}) \middle| \mathbf{S}, \mathbf{T} \right) \right] + \mathbb{E} \left[\mathbb{E} \left(\frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c0}}{\#T_0} \widehat{\text{Var}}(\hat{\mu}_{c0}) \middle| \mathbf{S}, \mathbf{T} \right) \right] \\
&= \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[1 - \frac{n_c n_{c'}}{n^2} \frac{s(s-1)}{\pi_{cc'}} \right] \mu_{c1} \mu_{c'0} \mathbb{E} \left(\frac{S_c T_{ct} S_{c'} T_{c't}}{\#T_1 \#T_0} \right) \\
&\quad - \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} [\mu_{c1}^2 + \text{Var}(\hat{\mu}_{c1})] \mathbb{E} \left(\frac{S_c T_{c1}}{\#T_1} \right) - \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} [\mu_{c0}^2 + \text{Var}(\hat{\mu}_{c0})] \mathbb{E} \left(\frac{S_c T_{c0}}{\#T_0} \right) \\
&\quad + \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Var}(\hat{\mu}_{c1}) \mathbb{E} \left(\frac{S_c T_{c1}}{\#T_1} \right) + \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Var}(\hat{\mu}_{c0}) \mathbb{E} \left(\frac{S_c T_{c0}}{\#T_0} \right) \\
&= \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{\pi_{cc'}}{s(s-1)} - \frac{n_c n_{c'}}{n^2} \right] \mu_{c1} \mu_{c'0} - \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mu_{c1}^2 - \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mu_{c0}^2. \tag{63}
\end{aligned}$$

We next show that eq. (63) is no larger than eq. (61), using Young's inequality.

Lemma 4 (Young's Inequality) *If a, b are nonnegative real numbers and p, q are positive real numbers such that $\frac{1}{p} + \frac{1}{q} = 1$, then*

$$ab \leq \frac{a^p}{p} + \frac{b^q}{q}. \tag{64}$$

Take $p = q = 2$, then

$$\begin{aligned}
\text{COV}(\hat{\mu}_{1, \text{HT-PPS}}, \hat{\mu}_{0, \text{HT-PPS}}) &= \sum_{c=1}^{\ell} \sum_{c' \neq c} \left(\frac{\pi_{cc'}}{s(s-1)} - \frac{n_c n_{c'}}{n^2} \right) \mu_{c1} \mu_{c'0} - \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mu_{c1} \mu_{c0} \\
&\geq \sum_{c=1}^{\ell} \sum_{c' \neq c} \left(\frac{\pi_{cc'}}{s(s-1)} - \frac{n_c n_{c'}}{n^2} \right) \mu_{c1} \mu_{c'0}
\end{aligned}$$

$$\begin{aligned}
& -\frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mu_{c1}^2 - \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mu_{c0}^2 \\
& = \text{COV}_C(\hat{\mu}_{1,\text{HT-PPS}}, \hat{\mu}_{0,\text{HT-PPS}}).
\end{aligned} \tag{65}$$

C Supplementary appendix on HT-PPS properties

C.1 Useful indicator properties under PPS

We begin this section by computing expectations, variances, and covariances of indicators under PPSWOR sampling of clusters. Define $\pi_{cc'} \equiv E(S_c S_{c'}) = P(S_c = 1, S_{c'} = 1)$ as the probability of sampling both cluster c and c' .

$$\mathbb{E}(S_c | \#T_t) = \frac{n_c s}{n} \quad (66)$$

$$\begin{aligned} \mathbb{E}(S_c^2 T_{ct}^2 | \#T_t) &= \mathbb{E}(S_c T_{ct} | \#T_t) = \mathbb{E}(S_c \mathbb{E}(T_{ct} | \mathbf{S}) | \#T_t) \\ &= \mathbb{E}\left(S_c \frac{\#T_t}{s} \middle| \#T_t\right) = \frac{\#T_t}{s} \mathbb{E}(S_c | \#T_t) \\ &= \frac{n_c \#T_t}{n} \end{aligned} \quad (67)$$

$$\begin{aligned} \mathbb{E}(S_c S_{c'} T_{ct} T_{c't} | \#T_t) &= \mathbb{E}(S_c S_{c'} \mathbb{E}(T_{ct} T_{c't} | \mathbf{S}) | \#T_t) \\ &= \mathbb{E}\left(S_c S_{c'} \frac{\#T_t}{s} \frac{\#T_t - 1}{s - 1} \middle| \#T_t\right) = \frac{\#T_t (\#T_t - 1)}{s(s - 1)} \mathbb{E}(S_c S_{c'} | \#T_t) \\ &= \frac{\#T_t (\#T_t - 1)}{s(s - 1)} \pi_{cc'} \end{aligned} \quad (68)$$

$$\begin{aligned} \mathbb{E}(S_c S_{c'} T_{c1} T_{c'0} | \#T_1, \#T_0) &= \mathbb{E}(S_c S_{c'} \mathbb{E}(T_{c1} T_{c'0} | \mathbf{S}) | \#T_1, \#T_0) \\ &= \frac{\#T_1}{s} \frac{\#T_0}{s - 1} \mathbb{E}(S_c S_{c'} | \#T_1, \#T_0) \\ &= \frac{\#T_1 \#T_0}{s(s - 1)} \pi_{cc'} \end{aligned} \quad (69)$$

$$\mathbb{E}\left(\frac{S_c^2 T_{ct}^2}{\#T_t^2}\right) = \mathbb{E}\left(\frac{1}{\#T_t^2} \mathbb{E}(S_c T_{ct} | \#T_t)\right) = \frac{n_c}{n} \mathbb{E}\left(\frac{1}{\#T_t}\right) \quad (70)$$

$$\begin{aligned} \mathbb{E}\left(\frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t^2}\right) &= \mathbb{E}\left(\frac{1}{\#T_t^2} \mathbb{E}(S_c S_{c'} T_{ct} T_{c't} | \#T_t)\right) \\ &= \frac{\pi_{cc'}}{s(s - 1)} \mathbb{E}\left(1 - \frac{1}{\#T_t}\right) \end{aligned} \quad (71)$$

$$\text{Var}(S_c T_{ct} | \#T_t) = \frac{n_c \#T_t}{n} \left(1 - \frac{n_c \#T_t}{n}\right) \quad (72)$$

Conditional on the sampled clusters, units are sampled within a cluster using simple random sampling, and this secondary sampling stage is independent of cluster treatment assignment. Thus, the expectation of within-cluster sampling indicators are independent of the cluster treatment indicators. Moreover, within-cluster samples are drawn independently across clusters, and so for distinct units k and k' in the same cluster c or distinct units k and k^* in different clusters c and c' :

$$\mathbb{E}(S_{kc}|\mathbf{S}) = \frac{s_c}{n_c} \quad (73)$$

$$\mathbb{E}(S_{kc}S_{k'c}|\mathbf{S}) = \frac{s_c(s_c - 1)}{n_c(n_c - 1)} \quad (74)$$

$$\mathbb{E}(S_{kc}S_{k^*c'}|\mathbf{S}) = \frac{s_c s_{c'}}{n_c n_{c'}} \quad (75)$$

$$\text{Var}(S_{kc}|\mathbf{S}) = \frac{s_c}{n_c} \left(1 - \frac{s_c}{n_c}\right) \quad (76)$$

$$\begin{aligned} \text{Cov}(S_{kc}, S_{k'c}|\mathbf{S}) &= \mathbb{E}(S_{kc}S_{k'c}|\mathbf{S}) - \mathbb{E}(S_{kc}|\mathbf{S})\mathbb{E}(S_{k'c}|\mathbf{S}) \\ &= -\frac{s_c}{n_c} \frac{1}{n_c - 1} \left(1 - \frac{s_c}{n_c}\right) \end{aligned} \quad (77)$$

C.2 Location invariance of HT-PPS estimator for PATE

Since,

$$\begin{aligned} \hat{\mu}_{t,\text{HT,PPS}}(a + \mathbf{y}) &= \sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{(a + y_{kct})S_{kc}}{s_c} \\ &= a \sum_{c=1}^{\ell} \frac{S_c T_{c1}}{\#T_1} \sum_{k=1}^{n_c} \frac{S_{kc}}{s_c} + \sum_{c=1}^{\ell} \frac{S_c T_{c1}}{\#T_1} \sum_{k=1}^{n_c} \frac{y_{kc1}S_{kc}}{s_c} \\ &= a + \hat{\mu}_{t,\text{HT,PPS}}(\mathbf{y}) \end{aligned} \quad (78)$$

the HT-PPS estimator for PATE is location-invariant:

$$\hat{\delta}_{\text{HT,PPS}}(a + \mathbf{y}) = \hat{\mu}_{1,\text{HT,PPS}}(a + \mathbf{y}) - \hat{\mu}_{0,\text{HT,PPS}}(a + \mathbf{y}) = \hat{\delta}_{\text{HT,PPS}}(\mathbf{y}). \quad (79)$$

C.3 Expectation of HT-PPS estimator for PATE

$$\begin{aligned}
\mathbb{E}(\hat{\delta}_{\text{HT-PPS}}) &= \mathbb{E} \left(\sum_{c=1}^{\ell} \frac{S_c T_{c1}}{\#T_1} \sum_{k=1}^{n_c} \frac{y_{kc1} S_{kc}}{s_c} - \sum_{c=1}^{\ell} \frac{S_{c'} T_{c'0}}{\#T_0} \sum_{k^*=1}^{n_{c'}} \frac{y_{k^*c'0} S_{k^*c'}}{s_{c'}} \right) \\
&= \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \frac{y_{kc1}}{s_c} \mathbb{E} \left(\frac{S_c T_{c1} S_{kc}}{\#T_1} \right) - \sum_{c'=1}^{\ell} \sum_{k^*=1}^{n_{c'}} \frac{y_{k^*c'0}}{s_{c'}} \mathbb{E} \left(\frac{S_{c'} T_{c'0} S_{k^*c'}}{\#T_0} \right) \\
&= \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \frac{y_{kc1}}{s_c} \mathbb{E} \left(S_c \mathbb{E} \left(\frac{T_{c1}}{\#T_1} \middle| \mathbf{S} \right) \mathbb{E}(S_{kc} | \mathbf{S}) \right) \\
&\quad - \sum_{c'=1}^{\ell} \sum_{k^*=1}^{n_{c'}} \frac{y_{k^*c'0}}{s_{c'}} \mathbb{E} \left(S_{c'} \mathbb{E} \left(\frac{T_{c'0}}{\#T_0} \middle| \mathbf{S} \right) \mathbb{E}(S_{k^*c'} | \mathbf{S}) \right) \\
&= \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \frac{y_{kc1}}{n_c S} \mathbb{E}(S_c) - \sum_{c'=1}^{\ell} \sum_{k^*=1}^{n_{c'}} \frac{y_{k^*c'0}}{n_{c'} S} \mathbb{E}(S_{c'}) \\
&= \sum_{c=1}^{\ell} \frac{n_c}{n} \mu_{c1} - \sum_{c'=1}^{\ell} \frac{n_{c'}}{n} \mu_{c'0} \\
&= \mu_1 - \mu_0 = \delta.
\end{aligned} \tag{80}$$

C.4 Variance of HT-PPS estimator for PATE

From the property

$$\text{Var}(\hat{\delta}) = \text{Var}(\hat{\mu}_1 - \hat{\mu}_0) = \text{Var}(\hat{\mu}_1) + \text{Var}(\hat{\mu}_0) - 2\text{Cov}(\hat{\mu}_1, \hat{\mu}_0). \tag{81}$$

each term is expanded upon to derive the variance of the HT-PPS estimator for PATE and obtain a variance estimator.

C.4.1 Variance of HT-PPS estimator for population mean

Using the law of total variance,

$$\begin{aligned}
\text{Var}(\hat{\mu}_{t,\text{HT,PPS}}) &= \text{Var} \left(\sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c} \right) \\
&= \text{Var} \left[\mathbb{E} \left(\sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c} \middle| \mathbf{S}, \mathbf{T} \right) \right] \\
&\quad + \mathbb{E} \left[\text{Var} \left(\sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c} \middle| \mathbf{S}, \mathbf{T} \right) \right]. \tag{82}
\end{aligned}$$

The first terms can be further simplified:

$$\begin{aligned}
\text{Var} \left[\mathbb{E} \left(\sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c} \middle| \mathbf{S}, \mathbf{T} \right) \right] &= \text{Var} \left[\sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{y_{kct}}{s_c} \mathbb{E}(S_{kc} | \mathbf{S}, \mathbf{T}) \right] \\
&= \sum_{c=1}^{\ell} \text{Var} \left(\mu_{ct} \frac{S_c T_{ct}}{\#T_t} \right) + \sum_{c=1}^{\ell} \sum_{c' \neq c} \text{Cov} \left(\mu_{ct} \frac{S_c T_{ct}}{\#T_t}, \mu_{c't} \frac{S_{c'} T_{c't}}{\#T_t} \right) \\
&= \sum_{c=1}^{\ell} \mu_{ct}^2 \text{Var} \left(\frac{S_c T_{ct}}{\#T_t} \right) + \sum_{c=1}^{\ell} \sum_{c' \neq c} \mu_{ct} \mu_{c't} \text{Cov} \left(\frac{S_c T_{ct}}{\#T_t}, \frac{S_{c'} T_{c't}}{\#T_t} \right) \\
&= \sum_{c=1}^{\ell} \mu_{ct}^2 \left[\mathbb{E} \left(\frac{S_c T_{ct}}{\#T_t} \right)^2 - \mathbb{E} \left(\frac{S_c T_{ct}}{\#T_t} \right)^2 \right] \\
&\quad + \sum_{c=1}^{\ell} \sum_{c' \neq c} \mu_{ct} \mu_{c't} \left[\mathbb{E} \left(\frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t^2} \right) - \mathbb{E} \left(\frac{S_c T_{ct}}{\#T_t} \right) \mathbb{E} \left(\frac{S_{c'} T_{c't}}{\#T_t} \right) \right] \\
&= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \mu_{ct}^2 - \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{n_c^2}{n^2} \mu_{ct}^2 \\
&\quad + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{n_c^2}{n^2} \mu_{ct} \mu_{c't} \\
&= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \mu_{ct}^2 + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2 \\
&= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \mu_{ct}^2 - \mathbb{E} \left(\frac{1}{\#T_t} \right) \mu_t^2
\end{aligned}$$

$$\begin{aligned}
& + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \mu_t^2 \\
& = \mathbb{E} \left(\frac{1}{\#T_t} \right) \left[\sum_{c=1}^{\ell} \frac{n_c}{n} \mu_{ct}^2 - \mu_t^2 \right] + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \left[\sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2 \right] \\
& = \mathbb{E} \left(\frac{1}{\#T_t} \right) \left[\sum_{c=1}^{\ell} \frac{n_c}{n} (\mu_{ct}^2 - 2\mu_t \mu_{ct} + \mu_t^2) \right] + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \left[\sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2 \right] \\
& = \mathbb{E} \left(\frac{1}{\#T_t} \right) \left[\sum_{c=1}^{\ell} \frac{n_c}{n} (\mu_{ct} - \mu_t)^2 \right] + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \left[\sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2 \right] \\
& = \mathbb{E} \left(\frac{1}{\#T_t} \right) \sigma_{t,bet}^2 + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \left[\sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2 \right] \tag{83}
\end{aligned}$$

where $\sigma_{t,bet}^2$ is the weighted variance of cluster means. Simplifying the second term:

$$\begin{aligned}
& \mathbb{E} \left[\text{Var} \left(\sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c} \middle| \mathbf{S}, \mathbf{T} \right) \right] = \sum_{c=1}^{\ell} \text{Var}(\hat{\mu}_{ct} | \mathbf{S}, \mathbf{T}) \mathbb{E} \left(\frac{S_c T_{ct}}{\#T_t^2} \right) \\
& = \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Var}(\hat{\mu}_{ct}) \\
& = \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \left[\text{Var} \left(\sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c} \right) + \text{Cov} \left(\sum_{k=1}^{n_c} \frac{y_{kct} S_{kc}}{s_c}, \sum_{k' \neq k} \frac{y_{k'ct} S_{k'c}}{s_c} \right) \right] \\
& = \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \left[\sum_{k=1}^{n_c} \frac{y_{kct}^2}{s_c n_c} \left(1 - \frac{s_c}{n_c} \right) - \sum_{k=1}^{n_c} \sum_{k' \neq k} \frac{y_{kct} y_{k'ct}}{s_c n_c (n_c - 1)} \left(1 - \frac{s_c}{n_c} \right) \right] \\
& = \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{1}{n(n_c - 1) s_c} \left(1 - \frac{s_c}{n_c} \right) \left[(n_c - 1) \sum_{k=1}^{n_c} y_{kct}^2 - \sum_{c=1}^{n_c} \sum_{c' \neq c} y_{kct} y_{k'ct} \right] \\
& = \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{1}{n(n_c - 1) s_c} \left(1 - \frac{s_c}{n_c} \right) \left[(n_c - 1) \sum_{k=1}^{n_c} y_{kct}^2 - \sum_{k=1}^{n_c} \sum_{k'=1}^{n_c} y_{kct} y_{k'ct} + \sum_{k=1}^{n_c} y_{kct}^2 \right] \\
& = \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{1}{n(n_c - 1) s_c} \left(1 - \frac{s_c}{n_c} \right) \left[n_c \sum_{k=1}^{n_c} y_{kct}^2 - \left(\sum_{k=1}^{n_c} y_{kct} \right)^2 \right] \\
& = \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \left(1 - \frac{s_c}{n_c} \right) \frac{\sigma_{ct}^2}{s_c}. \tag{84}
\end{aligned}$$

The variance for the HT-PPS mean estimator is

$$\begin{aligned} \text{Var}(\hat{\mu}_{t,\text{HT,SRS}}) &= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \mu_{ct}^2 + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2 \\ &\quad + \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Var}(\hat{\mu}_{ct}) \end{aligned} \quad (85)$$

$$\begin{aligned} &= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sigma_{t,\text{bet}}^2 + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \left[\sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2 \right] \\ &\quad + \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \left(1 - \frac{s_c}{n_c} \right) \frac{\sigma_{ct}^2}{s_c}. \end{aligned} \quad (86)$$

C.4.2 Covariance of HT-PPS estimator for population means

Note that:

$$\begin{aligned} \hat{\mu}_{1,\text{HT,PPS}} \hat{\mu}_{0,\text{HT,PPS}} &= \left(\sum_{c=1}^{\ell} \frac{S_c T_{c1}}{\#T_1} \sum_{k=1}^{n_c} \frac{y_{kc1} S_{kc}}{s_c} \right) \left(\sum_{c'=1}^{\ell} \frac{S_{c'} T_{c'0}}{\#T_0} \sum_{k^*=1}^{n_{c'}} \frac{y_{k^*c'0} S_{k^*c'}}{s_{c'}} \right) \\ &= \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \sum_{c'=1}^{n_{c'}} \sum_{k^*=1}^{n_{c'}} \frac{y_{kc1} y_{k^*c'0}}{s_c s_{c'}} \frac{S_c T_{c1} S_{kc} S_{c'} T_{c'0} S_{k^*c'}}{\#T_1 \#T_0} \\ &= \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \sum_{c' \neq c} \sum_{k^*=1}^{n_{c'}} \frac{y_{kc1} y_{k^*c'0}}{s_c s_{c'}} \frac{S_c T_{c1} S_{kc} S_{c'} T_{c'0} S_{k^*c'}}{\#T_1 \#T_0}. \end{aligned} \quad (87)$$

The last equality comes from the fact that a cluster can only be given one treatment. Therefore,

$$\begin{aligned} \text{COV}(\hat{\mu}_{1,\text{HT,PPS}}, \hat{\mu}_{0,\text{HT,PPS}}) &= \mathbb{E}(\hat{\mu}_{1,\text{HT,PPS}} \hat{\mu}_{0,\text{HT,PPS}}) - \mathbb{E}(\hat{\mu}_{1,\text{HT,PPS}}) \mathbb{E}(\hat{\mu}_{0,\text{HT,PPS}}) \\ &= \mathbb{E} \left(\sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \sum_{c' \neq c} \sum_{k^*=1}^{n_{c'}} \frac{y_{kc1} y_{k^*c'0}}{s_c s_{c'}} \frac{S_c T_{c1} S_{kc} S_{c'} T_{c'0} S_{k^*c'}}{\#T_1 \#T_0} \right) - \mu_1 \mu_0 \\ &= \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \sum_{c' \neq c} \sum_{k^*=1}^{n_{c'}} \frac{y_{kc1} y_{k^*c'0}}{s_c s_{c'}} \mathbb{E} \left[\mathbb{E} \left(\frac{S_c T_{c1} S_{kc} S_{c'} T_{c'0} S_{k^*c'}}{\#T_1 \#T_0} \mid \mathbf{S} \right) \right] - \mu_1 \mu_0 \\ &= \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \sum_{c' \neq c} \sum_{k^*=1}^{n_{c'}} \frac{y_{kc1} y_{k^*c'0}}{s_c s_{c'}} \mathbb{E} \left[S_c S_{c'} \mathbb{E} \left(\frac{T_{c1} T_{c'0}}{\#T_1 \#T_0} \mid \mathbf{S} \right) \mathbb{E}(S_{kc} S_{k^*c'} \mid \mathbf{S}) \right] - \mu_1 \mu_0 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{s(s-1)} \sum_{c=1}^{\ell} \sum_{k=1}^{n_c} \sum_{c' \neq c} \sum_{k^*=1}^{n_{c'}} \frac{y_{kc1} y_{k^*c'0}}{n_c n_{c'}} \mathbb{E}(S_c S_{c'}) - \mu_1 \mu_0 \\
&= \frac{1}{s(s-1)} \sum_{c=1}^{\ell} \sum_{c' \neq c} \pi_{cc'} \mu_{c1} \mu_{c'0} - \mu_1 \mu_0 \\
&= \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{\pi_{cc'}}{s(s-1)} - \frac{n_c n_{c'}}{n^2} \right] \mu_{c1} \mu_{c'0} - \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mu_{c1} \mu_{c0}. \tag{88}
\end{aligned}$$

C.4.3 SYG estimator for variance

The SYG variance estimator is

$$\begin{aligned}
\widehat{\text{Var}}(\hat{\mu}_t) &= \frac{1}{2} \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{s(s-1)}{\pi_{cc'} \#T_t (\#T_t - 1)} \frac{n_c n_{c'}}{n^2} - \frac{1}{\#T_t^2} \right] S_c T_{ct} S_{c'} T_{c't} (\hat{\mu}_{ct} - \hat{\mu}_{c't})^2 \\
&\quad + \sum_{c=1}^{\ell} \frac{S_c T_{ct}}{\#T_t} \frac{n_c}{n} \widehat{\text{Var}}(\hat{\mu}_{ct}) \tag{89}
\end{aligned}$$

where

$$\widehat{\text{Var}}(\hat{\mu}_{ct}) = \left(1 - \frac{s_c}{n_c} \right) \frac{\hat{\sigma}_{ct}^2}{s_c}. \tag{90}$$

The $\hat{\sigma}_{ct}^2$ is the sample variance of outcomes, which is unbiased for the population variance σ_{ct}^2 . We will now show that the SYG variance is unbiased for $\text{Var}(\hat{\mu}_t)$. This requires the following:

$$\sum_{c' \neq c}^{\ell} n_{c'} = n - n_c \tag{91}$$

and

$$\sum_{c' \neq c}^{\ell} \pi_{cc'} = \sum_{c' \neq c}^{\ell} E(S_c S_{c'}) = E[S_c (s - S_c)] = \frac{n_c s}{n} (s - 1). \tag{92}$$

Therefore, the expectation is

$$\mathbb{E} \left(\widehat{\text{Var}}(\hat{\mu}_t) \right) = \mathbb{E} \left(\frac{1}{2} \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{s(s-1)}{\pi_{cc'}} \frac{n_c n_{c'}}{n^2} \frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t (\#T_t - 1)} - \frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t^2} \right] [\hat{\mu}_{ct} - \hat{\mu}_{c't}]^2 \right)$$

$$\begin{aligned}
& + \sum_{c=1}^{\ell} \frac{n_c S_c T_{ct}}{n \#T_t} \widehat{\mathbf{Var}}(\hat{\mu}_{ct}) \Big) \\
= & \mathbb{E} \left(\mathbb{E} \left(\frac{1}{2} \sum_{c=1}^{\ell} \sum_{c' \neq c}^{\ell} \left[\frac{s(s-1) n_c n_{c'}}{\pi_{cc'}} \frac{S_c S_{c'} T_{ct} T_{c't}}{n^2 \#T_t (\#T_t - 1)} - \frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t^2} \right] [\hat{\mu}_{ct} - \hat{\mu}_{c't}]^2 \Big| \mathbf{S}, \mathbf{T} \right) \right) \\
& + \mathbb{E} \left(\mathbb{E} \left(\sum_{c=1}^{\ell} \frac{n_c S_c T_{ct}}{n \#T_t} \widehat{\mathbf{Var}}(\hat{\mu}_{ct}) \Big| \mathbf{S}, \mathbf{T} \right) \right) \\
= & \mathbb{E} \left(\mathbb{E} \left(\sum_{c=1}^{\ell} \sum_{c' \neq c}^{\ell} \left[\frac{s(s-1) n_c n_{c'}}{\pi_{cc'}} \frac{S_c S_{c'} T_{ct} T_{c't}}{n^2 \#T_t (\#T_t - 1)} - \frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t^2} \right] [\hat{\mu}_{ct}^2 - \hat{\mu}_{ct} \hat{\mu}_{c't}] \Big| \mathbf{S}, \mathbf{T} \right) \right) \\
& + \mathbb{E} \left(\sum_{c=1}^{\ell} \frac{n_c S_c T_{ct}}{n \#T_t} \mathbf{Var}(\hat{\mu}_{ct}) \right) \\
= & \mathbb{E} \left(\sum_{c=1}^{\ell} \sum_{c' \neq c}^{\ell} \left[\frac{s(s-1) n_c n_{c'}}{\pi_{cc'}} \frac{S_c S_{c'} T_{ct} T_{c't}}{n^2 \#T_t (\#T_t - 1)} - \frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t^2} \right] [\mu_{ct}^2 + \mathbf{Var}(\hat{\mu}_{ct}) - \mu_{ct} \mu_{c't}] \right) \\
& + \mathbb{E} \left(\sum_{c=1}^{\ell} \frac{n_c S_c T_{ct}}{n \#T_t} \mathbf{Var}(\hat{\mu}_{ct}) \right) \\
= & \sum_{c=1}^{\ell} \sum_{c' \neq c}^{\ell} \left[\frac{s(s-1) n_c n_{c'}}{\pi_{cc'}} \frac{S_c S_{c'} T_{ct} T_{c't}}{n^2 \#T_t (\#T_t - 1)} \mathbb{E} \left(\frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t (\#T_t - 1)} \right) - \mathbb{E} \left(\frac{S_c S_{c'} T_{ct} T_{c't}}{\#T_t^2} \right) \right] [\mu_{ct}^2 + \mathbf{Var}(\hat{\mu}_{ct}) - \mu_{ct} \mu_{c't}] \\
& + \sum_{c=1}^{\ell} \frac{n_c}{n} \mathbf{Var}(\hat{\mu}_{ct}) \mathbb{E} \left(\frac{S_c T_{ct}}{\#T_t} \right) \\
= & \sum_{c=1}^{\ell} \sum_{c' \neq c}^{\ell} \left[\frac{s(s-1) n_c n_{c'}}{\pi_{cc'}} \frac{1}{n^2 \#T_t (\#T_t - 1)} \mathbb{E} (S_c S_{c'} T_{ct} T_{c't} | \#T_t) \right] - \mathbb{E} \left(\frac{1}{\#T_t^2} \mathbb{E} (S_c S_{c'} T_{ct} T_{c't} | \#T_t) \right) \Big] \\
& \cdot [\mu_{ct}^2 + \mathbf{Var}(\hat{\mu}_{ct}) - \mu_{ct} \mu_{c't}] + \sum_{c=1}^{\ell} \frac{n_c}{n} \mathbf{Var}(\hat{\mu}_{ct}) \mathbb{E} \left(\frac{1}{\#T_t} \mathbb{E} (S_c T_{ct} | \#T_t) \right) \\
= & \sum_{c=1}^{\ell} \sum_{c' \neq c}^{\ell} \left[\frac{n_c n_{c'}}{n^2} - \frac{\pi_{cc'}}{s(s-1)} \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \right] (\mu_{ct}^2 + \mathbf{Var}(\hat{\mu}_{ct}) - \mu_{ct} \mu_{c't}) + \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mathbf{Var}(\hat{\mu}_{ct}) \\
= & \sum_{c=1}^{\ell} \sum_{c' \neq c}^{\ell} \left[\frac{n_c n_{c'}}{n^2} - \frac{\pi_{cc'}}{s(s-1)} \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \right] (\mu_{ct}^2 + \mathbf{Var}(\hat{\mu}_{ct})) \\
& - \sum_{c=1}^{\ell} \sum_{c' \neq c}^{\ell} \left[\frac{n_c n_{c'}}{n^2} - \frac{\pi_{cc'}}{s(s-1)} \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \right] \mu_{ct} \mu_{c't} + \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mathbf{Var}(\hat{\mu}_{ct}) \\
= & \sum_{c=1}^{\ell} \left[\frac{n_c}{n^2} \sum_{c' \neq c}^{\ell} n_{c'} - \frac{1}{s(s-1)} \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \sum_{c' \neq c}^{\ell} \pi_{cc'} \right] (\mu_{ct}^2 + \mathbf{Var}(\hat{\mu}_{ct}))
\end{aligned}$$

$$\begin{aligned}
& - \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{n_c n_{c'}}{n^2} - \frac{\pi_{cc'}}{s(s-1)} \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \right] \mu_{ct} \mu_{c't} + \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \text{Var}(\hat{\mu}_{ct}) \\
&= \sum_{c=1}^{\ell} \left[\frac{n_c}{n} \left(1 - \frac{n_c}{n} \right) - \frac{n_c}{n} \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \right] (\mu_{ct}^2 + \text{Var}(\hat{\mu}_{ct})) \\
& \quad - \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{n_c n_{c'}}{n^2} - \frac{\pi_{cc'}}{s(s-1)} \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \right] \mu_{ct} \mu_{c't} + \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \text{Var}(\hat{\mu}_{ct}) \\
&= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \mu_{ct}^2 - \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mu_{ct}^2 - \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{n_c n_{c'}}{n^2} \mu_{ct} \mu_{c't} \\
& \quad + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} + \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Var}(\hat{\mu}_{ct}) \\
&= \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \mu_{ct}^2 + \mathbb{E} \left(1 - \frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \sum_{c' \neq c} \frac{\pi_{cc'}}{s(s-1)} \mu_{ct} \mu_{c't} - \mu_t^2 \\
& \quad + \mathbb{E} \left(\frac{1}{\#T_t} \right) \sum_{c=1}^{\ell} \frac{n_c}{n} \text{Var}(\hat{\mu}_{ct}). \tag{93}
\end{aligned}$$

This is equal to eq. (85).

C.4.4 Covariance bound

The covariance is bounded by

$$\begin{aligned}
\widehat{\text{Cov}}_C(\hat{\mu}_1, \hat{\mu}_0) &= \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[1 - \frac{n_c n_{c'}}{n^2} \frac{s(s-1)}{\pi_{cc'}} \right] \frac{S_c T_{ct} S_{c'} T_{c't}}{\#T_1 \#T_0} \hat{\mu}_{c1} \hat{\mu}_{c'0} \\
& \quad - \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c1}}{\#T_1} \hat{\mu}_{c1}^2 - \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c0}}{\#T_0} \hat{\mu}_{c0}^2 \\
& \quad + \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c1}}{\#T_1} \widehat{\text{Var}}(\hat{\mu}_{c1}) + \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c0}}{\#T_0} \widehat{\text{Var}}(\hat{\mu}_{c0}). \tag{94}
\end{aligned}$$

Taking expectation:

$$\begin{aligned}
\mathbb{E}\left(\widehat{\text{Cov}}_C(\hat{\mu}_1, \hat{\mu}_0)\right) &= \mathbb{E}\left[\mathbb{E}\left(\sum_{c=1}^{\ell}\sum_{c'\neq c}\left[1-\frac{n_cn_{c'}}{n^2}\frac{s(s-1)}{\pi_{cc'}}\right]\frac{S_cT_{ct}S_{c'}T_{c't}}{\#T_1\#T_0}\hat{\mu}_{c1}\hat{\mu}_{c'0}\middle|\mathbf{S},\mathbf{T}\right)\right] \\
&\quad - \mathbb{E}\left[\mathbb{E}\left(\frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c}{n}\frac{S_cT_{c1}}{\#T_1}\hat{\mu}_{c1}^2\middle|\mathbf{S},\mathbf{T}\right)\right] - \mathbb{E}\left[\mathbb{E}\left(\frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c}{n}\frac{S_cT_{c0}}{\#T_0}\hat{\mu}_{c0}^2\middle|\mathbf{S},\mathbf{T}\right)\right] \\
&\quad + \mathbb{E}\left[\mathbb{E}\left(\frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c}{n}\frac{S_cT_{c1}}{\#T_1}\widehat{\text{Var}}(\hat{\mu}_{c1})\middle|\mathbf{S},\mathbf{T}\right)\right] + \mathbb{E}\left[\mathbb{E}\left(\frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c}{n}\frac{S_cT_{c0}}{\#T_0}\widehat{\text{Var}}(\hat{\mu}_{c0})\middle|\mathbf{S},\mathbf{T}\right)\right] \\
&= \sum_{c=1}^{\ell}\sum_{c'\neq c}\left[1-\frac{n_cn_{c'}}{n^2}\frac{s(s-1)}{\pi_{cc'}}\right]\mu_{c1}\mu_{c'0}\mathbb{E}\left(\frac{S_cT_{ct}S_{c'}T_{c't}}{\#T_1\#T_0}\right) \\
&\quad - \frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c}{n}\left[\mu_{c1}^2 + \text{Var}(\hat{\mu}_{c1})\right]\mathbb{E}\left(\frac{S_cT_{c1}}{\#T_1}\right) - \frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c}{n}\left[\mu_{c0}^2 + \text{Var}(\hat{\mu}_{c0})\right]\mathbb{E}\left(\frac{S_cT_{c0}}{\#T_0}\right) \\
&\quad + \frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c}{n}\text{Var}(\hat{\mu}_{c1})\mathbb{E}\left(\frac{S_cT_{c1}}{\#T_1}\right) + \frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c}{n}\text{Var}(\hat{\mu}_{c0})\mathbb{E}\left(\frac{S_cT_{c0}}{\#T_0}\right) \\
&= \sum_{c=1}^{\ell}\sum_{c'\neq c}\left[1-\frac{n_cn_{c'}}{n^2}\frac{s(s-1)}{\pi_{cc'}}\right]\mu_{c1}\mu_{c'0}\mathbb{E}\left[\frac{1}{\#T_1\#T_0}\mathbb{E}(S_cS_{c'}T_{ct}T_{c't}|\#T_1,\#T_0)\right] \\
&\quad - \frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c}{n}\left[\mu_{c1}^2 + \text{Var}(\hat{\mu}_{c1})\right]\mathbb{E}\left[\frac{1}{\#T_1}\mathbb{E}(S_cT_{c1}|\#T_1)\right] \\
&\quad - \frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c}{n}\left[\mu_{c0}^2 + \text{Var}(\hat{\mu}_{c0})\right]\mathbb{E}\left[\frac{1}{\#T_0}\mathbb{E}(S_cT_{c0}|\#T_0)\right] \\
&\quad + \frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c}{n}\text{Var}(\hat{\mu}_{c1})\mathbb{E}\left[\frac{1}{\#T_1}\mathbb{E}(S_cT_{c1}|\#T_1)\right] + \frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c}{n}\text{Var}(\hat{\mu}_{c0})\mathbb{E}\left[\frac{1}{\#T_0}\mathbb{E}(S_cT_{c0}|\#T_0)\right] \\
&= \sum_{c=1}^{\ell}\sum_{c'\neq c}\left[\frac{\pi_{cc'}}{s(s-1)} - \frac{n_cn_{c'}}{n^2}\right]\mu_{c1}\mu_{c'0} - \frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c^2}{n^2}\mu_{c1}^2 - \frac{1}{2}\sum_{c=1}^{\ell}\frac{n_c^2}{n^2}\mu_{c0}^2. \tag{95}
\end{aligned}$$

We next show that eq. (95) is no larger than eq. (88), using Young's inequality.

Lemma 5 (Young's Inequality) *If a, b are nonnegative real numbers and p, q are positive real numbers such that $\frac{1}{p} + \frac{1}{q} = 1$, then*

$$ab \leq \frac{a^p}{p} + \frac{b^q}{q}. \tag{96}$$

Take $p = q = 2$, then

$$\begin{aligned}
\text{COV}(\hat{\mu}_{1,\text{HT-PPS}}, \hat{\mu}_{0,\text{HT-PPS}}) &= \sum_{c=1}^{\ell} \sum_{c' \neq c} \left(\frac{\pi_{cc'}}{s(s-1)} - \frac{n_c n_{c'}}{n^2} \right) \mu_{c1} \mu_{c0} - \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mu_{c1} \mu_{c0} \\
&\geq \sum_{c=1}^{\ell} \sum_{c' \neq c} \left(\frac{\pi_{cc'}}{s(s-1)} - \frac{n_c n_{c'}}{n^2} \right) \mu_{c1} \mu_{c0} \\
&\quad - \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mu_{c1}^2 - \frac{1}{2} \sum_{c=1}^{\ell} \frac{n_c^2}{n^2} \mu_{c0}^2 \\
&= \text{COV}_C(\hat{\mu}_{1,\text{HT-PPS}}, \hat{\mu}_{0,\text{HT-PPS}}).
\end{aligned} \tag{97}$$

From eq. (93) and eq. (97), we see that

$$\begin{aligned}
\widehat{\text{Var}}(\hat{\delta}_{\text{HT,PPS}}) &= \frac{1}{2} \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{s(s-1)}{\pi_{cc'} \#T_1 (\#T_1 - 1)} \frac{n_c n_{c'}}{n^2} - \frac{1}{\#T_1^2} \right] S_c T_{c1} S_{c'} T_{c'1} (\hat{\mu}_{c1} - \hat{\mu}_{c'1})^2 \\
&\quad + \frac{1}{2} \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{s(s-1)}{\pi_{cc'} \#T_0 (\#T_0 - 1)} \frac{n_c n_{c'}}{n^2} - \frac{1}{\#T_0^2} \right] S_c T_{c0} S_{c'} T_{c'0} (\hat{\mu}_{c0} - \hat{\mu}_{c'0})^2 \\
&\quad - 2 \sum_{c=1}^{\ell} \sum_{c' \neq c} \left[\frac{\pi_{cc'}}{s(s-1)} - \frac{n_c n_{c'}}{n^2} \right] \frac{s(s-1)}{\pi_{cc'}} \frac{S_c S_{c'} T_{c1} T_{c'0}}{\#T_1 \#T_0} \hat{\mu}_{c1} \hat{\mu}_{c0} \\
&\quad + \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c1}}{\#T_1} \hat{\mu}_{c1}^2 + \sum_{c=1}^{\ell} \frac{n_c}{n} \frac{S_c T_{c0}}{\#T_0} \hat{\mu}_{c0}^2
\end{aligned} \tag{98}$$

is a conservative bound for $\text{Var}(\hat{\delta}_{\text{HT,PPS}})$.