## CI = STRATIFICATION & RERANDOMIZATION

In this chapter, we discuss stratification and rerandomization, two strategies that improve <u>covariate balance</u> in RCTs.

### I - STRATIFIED DESIGNS

- For a set of units $i = 1, ..., n$, observe $(X_i, Y_i, W_i)$ where
  - $X_i \in \{1, ..., K\}$ = discrete covariate. We say that $i$ belongs to the $k$-th stratum $\Leftrightarrow X_i = k$
  - $W_i \in \{0, 1\}$ = treatment allocation. In this section, we consider two random allocations: <u>completely randomized</u>, and <u>stratified</u>.
  - $Y_i = W_i Y_i(1) + (1 - W_i) Y_i(0)$
  - $(Y_i(0), Y_i(1))$ = Potential Outcomes

We are interested in estimating the ATE $= \mathbb{E}(Y_i(1) - Y_i(0))$ (under an infinite population model for the P.O.)

- Put $\mu_{jk} := \mathbb{E}(Y_i(j) \mid X_i = k)$
  $$\sigma_{jk}^2 := var(Y_i(j) \mid X_i = k)$$

The conditional means & variances may differ across strata (unless the covariate $X_i$ is independent of the PO)
  - ↳ A stratified design exploits this to reduce the variance of estimators of the ATE.

- Suppose you have collected $n_k$ units with $X_i = k$, and put

---

$\hat{\pi}_k = \dfrac{n_k}{n}$. We consider two randomization procedures:

[CRE] In a CRE, draw a fixed number $n_t$ of units at random from the total (sampled) population of $n$ units. These units are assigned to the treatment group, and the remaining $n_c := n - n_t$ to the control group.

$$\mathbb{P}(\underline{w} \mid \underline{Y}(0), \underline{Y}(1), \underline{X}) = \frac{1}{\binom{n}{n_t}} \quad \forall \underline{w} \text{ s.t. } \sum_{i=1}^{n} w_i = n_t$$

vector notation =
$\underline{w} = (w_1, ..., w_n)^t$       We write $\underline{w} \sim [CRE]$.

A natural choice for the estimation of the ATE is the difference in means estimator

$$\hat{\Delta} = \frac{1}{n_t} \sum_{i=1}^{n} w_i Y_i - \frac{1}{n_c} \sum_{i=1}^{n} (1 - w_i) Y_i.$$

[ST] In a [CRE], unfortunate draws may lead to severe unbalance of the covariates between the treatment and control groups. Taking for example $X_i \in \{man, woman\}$ ($K = 2$), a [CRE] draw may place all men in the treatment group, and all women in the control group. Stratified designs do not allow this to happen by randomizing units stratum per stratum. Specifically, let
  - $n_{1k}$ = # of treated units in the $k$-th stratum
  - $n_{0k} := n_k - n_{1k}$

Put $p_k = \dfrac{n_{1k}}{n_k}$ = proportion of treated units in the $k$-th stratum.

The assignment probability is

$$\mathbb{P}(\underline{w} \mid \underline{Y}(0), \underline{Y}(1), \underline{X}) = \prod_{k=1}^{K} \frac{1}{\binom{n_k}{n_{1k}}}$$

$$\forall \underline{w} \text{ s.t. } \sum_{X_i = k} w_i = n_{1k}$$

We write $[w] \sim [ST]$

The total number of treated units is $n_t = \sum_{k=1}^{K} n_{1k}$

$$= \sum_k p_k n_k \ .$$

Under a [ST] randomisation scheme, we consider the aggregated difference-in-mean estimator:

$$\widehat{\Delta}_{agg} = \sum_{k=1}^{K} \widehat{\pi}_k \, \widehat{\Delta}(k)$$

$$\widehat{\Delta}(k) = \frac{1}{n_{1k}} \sum_{\substack{w_i = 1 \\ X_i = k}} Y_i - \frac{1}{n_{0k}} \sum_{\substack{w_i = 0 \\ X_i = k}} Y_i$$

In the remainder of this section, we compare the bias and variance of $\widehat{\Delta}$ and $\widehat{\Delta}_{agg}$ under the two randomization schemes.

x Remark: Conditionally on the strata sizes, the ATE can be equivalently expressed $ATE = \sum_{k=1}^{K} \widehat{\pi}_k \, \mathbb{E}(Y_i(1) - Y_i(0) \mid X_i = k)$.

Stratum 1

$n_{1k}$ trt units
(proportion $p_k = \dfrac{n_{1k}}{n_k}$)

$n_{0k}$ control units

Stratum K

$k$-th stratum with $n_k$ units
(proportion $\widehat{\pi}_k = \dfrac{n_k}{n}$)

& $\mu_{jk} = \mathbb{E}(Y_i(j) \mid X_i = k)$
$\sigma_{jk}^2 = Var(Y_i(j) \mid X_i = k)$

• In addition,

$$\mu_j := \mathbb{E}\, Y_i(j) = \mathbb{E}_X \, \mathbb{E}(Y_i(j) \mid X) = \sum_k \widehat{\pi}_k \, \mu_{jk}$$

We assume implicitly here that the $n$ individuals represent the whole population of interest, so that $\mathbb{P}(X_i = k) = \dfrac{n_k}{n} = \widehat{\pi}_k$

[Interlude] The assumption above is a consequence of considering non-random $n_1, \dots, n_K$. We may be interested instead in assuming a superpopulation model $\mathbb{P}(X_i = k) = \pi_k$; and selecting $n$ individuals at random from that superpopulation. In this case,

$(n_1, \ldots, n_K) \sim$ Multinomial $(\pi_1, \ldots, \pi_K)$ are random variables, and $\frac{n_k}{n} = \hat{\pi}_k \xrightarrow{a.s} \pi_k$ as $n \to \infty$.

Choosing $\{n_k\}$ as fixed quantities removes some noise from the estimators & we begin our investigations under this assumption. The ATE becomes

$$\text{ATE} = \sum_{k=1}^{K} \hat{\pi}_k (\mu_{1k} - \mu_{0k}).$$

Later on, we consider random $(n_1, \ldots, n_K)$ and quantify the impact on the variance of our estimators.

• $\sigma_j^2 = \text{Var } Y_i(j)$

$= \underbrace{\mathbb{E}_X \text{Var}(Y_i(j)|X_i)}_{\text{var}(Y_i(j)|X_i=k)=\sigma_{jk}^2} + \underbrace{\text{Var}_X \mathbb{E}(Y_i(j)|X_i)}_{\mathbb{E}(Y_i(j)|X_i=k)=\mu_{jk}}$

$\underbrace{= \sum_{k=1}^{K} \hat{\pi}_k \sigma_{jk}^2}_{} \quad \underbrace{= \sum_{k=1}^{K} \hat{\pi}_k (\mu_{jk} - \mu_j)^2}_{}$

thus $\boxed{\sigma_j^2 = \sum_{k=1}^{K} \hat{\pi}_k \sigma_{jk}^2 + \sum_{k=1}^{K} \hat{\pi}_k (\mu_{jk} - \mu_j)^2}$

↖ Unconditional variance of the P.O. $Y_i(j)$

• The table on the next page summarizes the properties of the two estimators of the ATE in terms of their means (& bias) and variances.

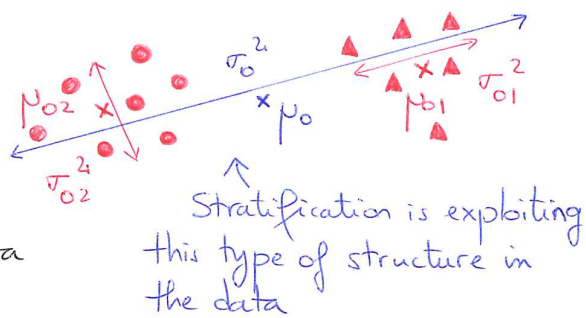| | |
|---|---|
| [ST] Mean $\hat{\Delta}_{agg}$ | $\mathbb{E}\,\hat{\Delta}_{agg} = \sum_k \hat{\pi}_k (\mu_{1k} - \mu_{0k})$ <br> (zero bias) |
| [ST] Var $\hat{\Delta}_{agg}$ | $n \text{ var } \hat{\Delta}_{agg} = \sum_k \hat{\pi}_k \left( \frac{\sigma_{1k}^2}{p_k} + \frac{\sigma_{0k}^2}{1-p_k} \right)$ |
| [CRE] Mean $\hat{\Delta}$ | $\mathbb{E}\,\hat{\Delta} = \mu_1 - \mu_0 = \sum_k \hat{\pi}_k (\mu_{1k} - \mu_{0k})$ <br> (zero bias) |
| [CRE] Var $\hat{\Delta}$ | $n \text{ var } \hat{\Delta} = \frac{n}{n_t} \left( \sum_k \hat{\pi}_k \sigma_{1k}^2 + \sum_k \hat{\pi}_k (\mu_{1k} - \mu_1)^2 \right)$ <br> $+ \frac{n}{n_c} \left( \sum_k \hat{\pi}_k \sigma_{0k}^2 + \sum_k \hat{\pi}_k (\mu_{0k} - \mu_0)^2 \right)$ |

[ See Appendix A ]

• Remarks

(i) Both $\hat{\Delta}$ and $\hat{\Delta}_{agg}$ are unbiased for the ATE

(ii) When $p_k = p \; \forall k$ ( all strata have the same assignment probability ), $n_t = pn$ and

$\begin{cases} n \text{ var } \hat{\Delta}_{agg} = \sum_k \hat{\pi}_k \left( \frac{\sigma_{1k}^2}{p} + \frac{\sigma_{0k}^2}{1-p} \right) \\[4mm] n \text{ var } \hat{\Delta} = \sum_k \hat{\pi}_k \left( \frac{\sigma_{1k}^2}{p} + \frac{\sigma_{0k}^2}{1-p} \right) \\[2mm] \qquad\qquad + \sum_k \frac{\hat{\pi}_k}{p} (\mu_{1k} - \mu_1)^2 + \sum_k \frac{\hat{\pi}_k}{1-p} (\mu_{0k} - \mu_0)^2 \end{cases}$

$\Rightarrow \operatorname{var} \widehat{\Delta} > \operatorname{var} \widehat{\Delta}_{agg}$ unless all means $\mu_{1k} = \mu_1$, $\mu_{0k} = \mu_0$ are equal across all strata. In this case, the P.O. are uncorrelated with $X$ and stratification does not reduce the variance compared to a simple [CRE] procedure. On the other hand, the more heterogeneity there is across strata and the more variance reduction we get.



Stratification is exploiting this type of structure in the data

In practice, the strata are constructed from treatment-insensitive vector of covariates. For example, $X_i \in \{man, woman\}$ and we randomize individuals within the m and f populations. With a constant $p_k = p$ across the two strata, this ensures that men and women are well balanced across the treatment and control groups. The [CRE] may suffer from some unlikely imbalance for a specific random draw.

With continuous covariate, we may start by discretizing them first, or consider rerandomization (see later)

x <u>Remark</u> = We may consider a [ST] randomization scheme, but analyze the experiment omitting stratum information, pooling all the data together and using a difference in means estimator. Doing so yields the following expressions for the mean and variance:

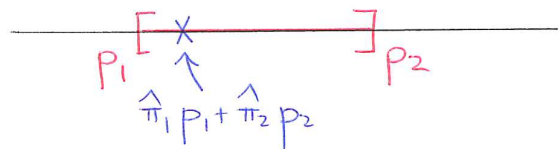| [ST] Mean $\widehat{\Delta}$ | $\mathbb{E}\,\widehat{\Delta} = \sum_k \widehat{\pi}_k \left( \dfrac{n\,p_k}{n_t}\,\mu_{1k} - \dfrac{n(1-p_k)}{n_c}\,\mu_{0k} \right)$ |
|---|---|
| [ST] Var $\widehat{\Delta}$ | $n\operatorname{var}\widehat{\Delta} = \sum_k \widehat{\pi}_k \left( \left(\dfrac{n}{n_t}\right)^2 p_k\,\sigma_{1k}^2 + \left(\dfrac{n}{n_c}\right)(1-p_k)\,\sigma_{0k}^2 \right)$ |

[Appendix B]

↘ the difference in means estimator computed under a [ST] randomization scheme is biased in general, unless $p_k = p \;\forall k$, or $\mu_{1k} = \mu_1$ and $\mu_{0k} = \mu_0 \;\forall k$

↘ however, the variance of $\widehat{\Delta}$ may be smaller than that of $\widehat{\Delta}_{agg}$ → bias / variance tradeoff. Recall that $n\operatorname{var}\widehat{\Delta}_{agg} = \sum_k \widehat{\pi}_k \left( \dfrac{\sigma_{1k}^2}{p_k} + \dfrac{\sigma_{0k}^2}{1-p_k} \right)$

Comparing this expression with $n\operatorname{var}\widehat{\Delta}$ indicates that one must evaluate the respective contributions of $\left(\dfrac{n}{n_t}\right)^2 p_k$ and $\dfrac{1}{p_k}$ in front of $\sigma_{1k}^2$ (and similarly for $\sigma_{0k}^2$).

Note that $\left(\frac{n}{n_t}\right)^2 p_k > \frac{1}{p_k} \iff \frac{p_k}{\sum_\ell p_\ell \hat{\pi}_\ell} > 1$

Take $k = 2$, $p_1 < p_2$. Since $\hat{\pi}_1 + \hat{\pi}_2 = 1$,
$p_1 \hat{\pi}_1 + p_2 \hat{\pi}_2$ lies in the interval $[p_1, p_2]$



and necessarily one $\frac{p_k}{\sum_\ell p_\ell \hat{\pi}_\ell}$ is larger than $1$, and the

other one smaller than $1$.

↘ Select $\hat{\pi}_1$ small so that $\hat{\pi}_1 p_1 + \hat{\pi}_2 p_2$
lies close to $p_2$.

Then

$$\begin{cases} \dfrac{p_2}{\sum p_\ell \hat{\pi}_\ell} \text{ close to } 1, \text{ slightly larger than } 1 \\[3mm] \dfrac{p_1}{\sum p_\ell \hat{\pi}_\ell} \text{ close to } 0 \end{cases}$$

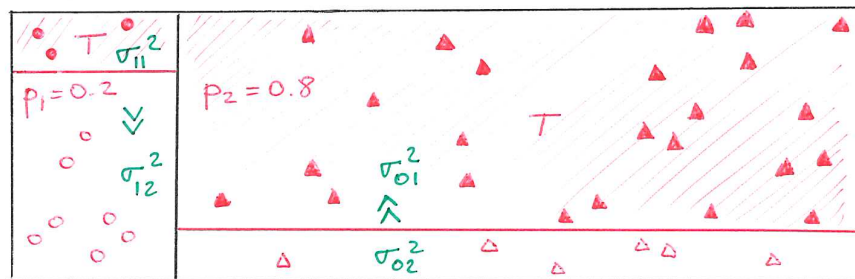To find cases where $\hat{\Delta}$ has smaller variance than $\hat{\Delta}_{agg}$,
we need to select

$\sigma_{11}^2$ large when $\left(\frac{n}{n_t}\right)^2 p_1 < \frac{1}{p_1} \iff \frac{p_1}{\sum_\ell p_\ell \hat{\pi}_\ell} < 1$

$\sigma_{12}^2$ small when $\left(\frac{n}{n_t}\right)^2 p_2 > \frac{1}{p_2} \iff \frac{p_2}{\sum_\ell p_\ell \hat{\pi}_\ell} > 1$

& the other way around for $\sigma_{01}^2$ and $\sigma_{02}^2$.

---

Taking $\begin{cases} p_1 = 0.2 = 1 - p_2 \\ \hat{\pi}_1 = 0.1 = 1 - \hat{\pi}_2 \\ 10 = \sigma_{11}^2 > \sigma_{12}^2 = 1 \\ 2 = \sigma_{01}^2 < \sigma_{02}^2 = 8 \end{cases}$ does the job :

$\begin{cases} n \text{ var } \hat{\Delta} = 25.9 \\ n \text{ var } \hat{\Delta}_{Agg} = 42.3 \end{cases}$



$\hat{\pi}_1 = 0.1$          $\hat{\pi}_2 = 0.9$

The var differ
on page 6 since
the $\omega$ are $\sim [ST]$
or $[CRE]$

x Remark = If all $p_k = p$ $\forall k$,

$\hat{\Delta}_{Agg} = \sum_{k=1}^{K} \hat{\pi}_k \left( \frac{1}{n_{1k}} \sum_{\substack{\omega_i = 1 \\ X_i = k}} Y_i - \frac{1}{n_{0k}} \sum_{\substack{\omega_i = 0 \\ X_i = k}} Y_i \right)$

$= \sum_{k=1}^{K} \frac{1}{n} \left( \underbrace{\frac{n_k}{n_{1k}}}_{\frac{1}{p}} \sum_{\substack{\omega_i = 1 \\ X_i = k}} Y_i - \underbrace{\left(\frac{n_k}{n_{0k}}\right)}_{\frac{1}{1-p}} \sum_{\substack{\omega_i = 0 \\ X_i = k}} Y_i \right)$

$= \underbrace{\left( \frac{1}{np} \right.}_{\frac{1}{n_t}} \sum_{k=1}^{K} \sum_{\substack{\omega_i = 1 \\ X_i = k}} Y_i - \underbrace{\left( \frac{1}{n(1-p)} \right)}_{\frac{1}{n_c}} \sum_{k=1}^{K} \sum_{\substack{\omega_i = 0 \\ X_i = k}} Y_i = \hat{\Delta}$

x <u>Remark</u> = Mean and variances can be computed <u>unconditionally</u> on $(n_1, \ldots, n_K)$, assuming a multinomial distribution $\mathrm{Mult}(\pi_1, \ldots, \pi_K)$ for $(n_1, \ldots, n_K)$ [ Since all $X_i$ are assumed independent with $\mathbb{P}(X_i = k) = \pi_k$ ]. The ATE is $\sum_{k=1}^{K} \pi_k (\mu_{1k} - \mu_{0k})$.

| | |
|---|---|
| [ST] Mean $\widehat{\Delta}_{agg}$ | $\mathbb{E}\,\widehat{\Delta}_{agg} = \sum_{k=1}^{K} \pi_k (\mu_{1k} - \mu_{0k})$ (unbiased) |
| [ST] Var $\widehat{\Delta}_{agg}$ | $n\,\mathrm{var}\,\widehat{\Delta}_{agg} = \sum_k \pi_k \left( \dfrac{\sigma_{1k}^2}{p_k} + \dfrac{\sigma_{0k}^2}{1 - p_k} \right)$ $+ \sum_k \pi_k (\mu_{1k} - \mu_{0k})^2$ $- \left[ \sum_k \pi_k (\mu_{1k} - \mu_{0k}) \right]^2$ |
| [CRE] Mean $\widehat{\Delta}$ | $\mathbb{E}\,\widehat{\Delta} = \sum_{k=1}^{K} \pi_k (\mu_{1k} - \mu_{0k})$ (unbiased) |
| [CRE] Var $\widehat{\Delta}$ [$(n_t, n_c)$ remain fixed] | $n\,\mathrm{var}\,\widehat{\Delta} = \left( \dfrac{n}{n_t} \right) \sum_{k=1}^{K} \pi_k \left[ \sigma_{1k}^2 + (\mu_{1k} - \mu_1)^2 \right]$ $+ \left( \dfrac{n}{n_c} \right) \sum_{k=1}^{K} \pi_k \left[ \sigma_{0k}^2 + (\mu_{0k} - \mu_0)^2 \right]$ |

[ See Appendix C ]

↓ Both $\widehat{\Delta}_{agg}$ and $\widehat{\Delta}$ remain unbiased

---

↘ Comparing the variance terms with those on page 6, we see that taking into account the variability in the $n_k$'s increases the variance of the estimators (as expected). The green terms p.11 represent the additional terms.

↘ The expression of $\mathrm{var}\,\widehat{\Delta}_{agg}$ under [ST] was also derived on p.3 in <u>CI : UNCONFOUNDEDNESS</u>.

Similar calculations can be derived for $\widehat{\Delta}$ under [ST]. For example, one can show that

$$\mathbb{E}\,\widehat{\Delta} = \sum_k \pi_k \left( \frac{p_k}{\sum_\ell p_\ell \pi_\ell} \mu_{1k} - \frac{1 - p_k}{\sum_\ell (1 - p_\ell) \pi_\ell} \mu_{0k} \right) + O(n^{-1})$$

[ See Appendix D ]
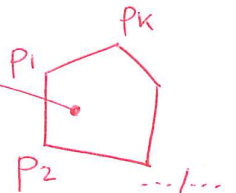
Under $\underline{W} \sim$ [ST]
[ $n_t, n_c$ are random here, unlike under $\underline{W} \sim$ [CRE] ]

Note that the asymptotic term differs from the ATE, since all coefficients $\dfrac{p_k}{\sum p_\ell \pi_\ell}$ differ from the value 1; unless all $p_k$ are equal.

↘ To see this geometrically, note that $\sum p_\ell \pi_\ell$ belongs to the interior of the convex hull of $\{p_1, \ldots, p_K\}$ for $0 < \pi_\ell < 1$ $\forall \ell$

x Remark : Finite population approach.

The mean and variance calculations can be derived similarly assuming fixed potential outcomes $\{Y_i(0), Y_i(1)\}$. One can show that

$$n \operatorname{var} \hat{\Delta}_{agg} = \sum_{k=1}^{K} \hat{\pi}_k \left( \frac{s_{ik}^2}{p_k} + \frac{s_{0k}^2}{1-p_k} - s_{\Delta k}^2 \right)$$

$W \sim [ST]$

where

• $\hat{\pi}_k = \frac{n_k}{n}$ ; $n_k$ fixed (non-random) ; $p_k = \frac{n_{1k}}{n_k}$

• $s_{jk}^2 := \frac{1}{n_k-1} \sum_{X_i=k} \left( Y_i(j) - \bar{Y}_k(j) \right)^2$

$$\bar{Y}_k(j) := \frac{1}{n_k} \sum_{X_i=k} Y_i(j)$$

• $s_{\Delta k}^2 := \frac{1}{n_k-1} \sum_{X_i=k} \left( Y_i(1) - Y_i(0) - \bar{\Delta}_k \right)^2$

$$\bar{\Delta}_k := \frac{1}{n_k} \sum_{X_i=k} \left( Y_i(1) - Y_i(0) \right)$$

↑ See also p.7 in CI : RANDOMIZED CONTROL TRIALS
for the expression of the variance of the simple difference in means in the finite population approach.

x Remark :
$(n_1, .., n_K)$ fixed = sample $n_k$ units from category $k$, where the $K$ categories are known in advance.
$(n_1, .., n_K)$ random = sample $n$ units & then create strata based on covariates. The # of units in each stratum is not known in advance.

---

## II - RERANDOMIZATION

Li, Ding & Rubin (2018) : " Although complete randomization ensures covariate balance on average, the chance of observing significant differences between treatment and control covariate distributions increases with many covariates. Rerandomization discards randomizations that do not satisfy a predetermined covariate balance criterion, generally resulting in better covariate balance and more precise estimates of causal effects. "

In this section, we present asymptotic results for the difference in means estimator of the ATE under a rerandomization scheme using a Mahalanobis criterion : For some covariates $X_i$, compute $\hat{\Delta}_X := \frac{1}{n_t} \sum_{W_i=1} X_i - \frac{1}{n_c} \sum_{W_i=0} X_i$ and accept the assignment vector $(W_1, .., W_n)$ if the normalized covariate mean difference between the treatment and control groups $\frac{\hat{\Delta}_X^2}{\operatorname{var} \hat{\Delta}_X}$ is less than some predefined threshold $\varepsilon$.

Specifically, we provide LT for $\boxed{\sqrt{n}(\hat{\Delta} - ATE) \,\Big|\, \frac{\hat{\Delta}_X^2}{\operatorname{var} \hat{\Delta}_X} \leq \varepsilon}$
(Limit Theorems)

↘ Assuming a finite & infinite population model ...
↘ Based on the work of Morgan & Rubin (2012, 2015),
  Li, Ding & Rubin (2018)
↘ Extend to stratified designs Wang, Wang, Liu (2021).

We proceed as follows:

(i) Provide CLT in the finite & infinite population models for $\hat{\Delta}$

(ii) Provide Limit Theorems for the difference in means estimator under rerandomization with Mahalanobis distance

(iii) Provide Limit Theorems under a stratified design with rerandomization.

## II.1. Limit Theorems for $\hat{\Delta}$

○ **Set-up** = For a set of units $i = 1, \dots, n$ we observe

$(X_i, Y_i, W_i)$ where

• $X_i \in \mathbb{R}$ = covariate

• $W_i \in \{0, 1\}$ ; $\sum_{i=1}^{n} W_i = n_t$

$(W_1, \dots, W_n) \sim$ Completely Randomized

• $Y_i = W_i Y_i(1) + (1 - W_i) Y_i(0)$

We consider both a finite & infinite population approach:

**[finite]** $\quad ATE = \Delta^n = \frac{1}{n} \sum_{i=1}^{n} (Y_i(1) - Y_i(0))$

**[ ∞ ]** $\quad ATE = \Delta^\infty = \mathbb{E}(Y_i(1) - Y_i(0))$

$\hat{\Delta} = \frac{1}{n_t} \sum_{i=1}^{n} W_i Y_i - \frac{1}{n_c} \sum_{i=1}^{n} (1 - W_i) Y_i$

We discard covariate information $X_i$ in this section

○ **Goal:**

**[finite]** Establish $\sqrt{n}(\hat{\Delta} - \Delta^n) \xrightarrow{d} ?$

**[ ∞ ]** Establish $\sqrt{n}(\hat{\Delta} - \Delta^\infty) \xrightarrow{d} ?$

---

○ **[finite]** Put

• $s_j^2 = \frac{1}{n-1} \sum_{i=1}^{n} (Y_i(j) - \bar{Y}(j))^2$

• $s_{01} = \frac{1}{n-1} \sum_{i=1}^{n} (Y_i(1) - \bar{Y}(1))(Y_i(0) - \bar{Y}(0))$

• $s_\Delta^2 = \frac{1}{n-1} \sum_{i=1}^{n} (Y_i(1) - Y_i(0) - \Delta^n)^2$

One can show that $\hat{\Delta}$ has mean $\Delta^n$ and variance

$var\, \hat{\Delta} = \frac{s_0^2}{n_c} + \frac{s_1^2}{n_t} - \frac{s_\Delta^2}{n}$ (see also p.13)

$= \frac{s_0^2}{n_c} + \frac{s_1^2}{n_t} - \frac{1}{n}(s_0^2 + s_1^2 - 2s_{01})$

since $s_\Delta^2 = \frac{1}{n-1} \sum (Y_i(1) - Y_i(0) - (\bar{Y}(1) - \bar{Y}(0)))^2$

$= \frac{1}{n-1} \sum (Y_i(1) - \bar{Y}(1) - (Y_i(0) - \bar{Y}(0)))^2$

$= s_0^2 + s_1^2 - 2s_{01}$

$= \left(\frac{1}{n_c} - \frac{1}{n}\right)s_0^2 + \left(\frac{1}{n_t} - \frac{1}{n}\right)s_1^2 + 2\frac{s_{01}}{n}$

put $\hat{\rho}_{01} = \frac{s_{01}}{s_0 s_1}$ = sample correlation coefficient

$= \frac{n_t}{n\,n_c} s_0^2 + \frac{n_c}{n\,n_t} s_1^2 + \frac{2}{n}\hat{\rho}_{01} s_0 s_1$

Assuming $\frac{n_t}{n} \to p \qquad \frac{n_c}{n} \to 1-p$ as $n \to \infty$,

& that $s_0^2, s_1^2, \hat{\rho}_{01} \to \sigma_0^2, \sigma_1^2, \rho_{01}$,

$$n\, var\, \hat{\Delta} \longrightarrow V_\Delta(p) = \frac{p}{1-p}\sigma_0^2 + \frac{1-p}{p}\sigma_1^2 + 2\rho_{01}\sigma_0\sigma_1 \quad (*)$$

• [∞] All derivations can be made conditionally on a design $\underline{\omega} = (\omega_1, \ldots, \omega_n)$.

$$\text{var}(\hat{\Delta} \mid \underline{\omega}) = \text{var}\left(\sum_{i=1}^{n} \underbrace{\left\{ \frac{\omega_i Y_i(1)}{n_E} - \frac{(1-\omega_i) Y_i(0)}{n_C} \right\}}_{=: T_i} \mid \underline{\omega}\right)$$

$$= \sum_{i=1}^{n} \text{var}(T_i \mid \underline{\omega}) + \sum_{j \neq i} \text{cov}(T_i, T_j \mid \underline{\omega})$$

independent

$$= \sum_{i=1}^{n} \left\{ \frac{\omega_i^2}{n_E^2} \sigma_1^2 + \frac{(1-\omega_i)^2}{n_C^2} \sigma_0^2 \right.$$

$$\left. - \frac{2}{n_C n_E} \omega_i(1-\omega_i) \text{cov}(Y(0), Y(1)) \right.$$

$$= \frac{\sigma_1^2}{n_E} + \frac{\sigma_0^2}{n_C}.$$

Note that unconditionally on $\underline{\omega}$,

$$\text{var}\,\hat{\Delta} = \text{var}\,\underbrace{\mathbb{E}(\hat{\Delta} \mid \underline{\omega})}_{= \mu_1 - \mu_0 \atop = 0} + \mathbb{E}\,\text{var}(\hat{\Delta} \mid \underline{\omega}) = \text{var}(\hat{\Delta} \mid \underline{\omega})$$

(**)

$$n\,\text{var}\,\hat{\Delta} \longrightarrow V_\Delta^\infty(p) := \frac{\sigma_1^2}{P} + \frac{\sigma_0^2}{1-P}.$$

**Theorem:**

$$[\text{finite}] \quad \sqrt{n}(\hat{\Delta} - \Delta^n) \xrightarrow{d} \mathcal{N}(0, V_\Delta(p))$$

$$[\infty] \quad \sqrt{n}(\hat{\Delta} - \Delta^\infty) \xrightarrow{d} \mathcal{N}(0, V_\Delta^\infty(p))$$

where $V_\Delta(p)$, $V_\Delta^\infty(p)$ are defined in (*) and (**) p. 16/17

---

× Remarks : (i) The results p. 17 complement the discussion page 12 in CI : UNCONFOUNDEDNESS & justify the normal approximation made by Neyman.

(ii) When the P.O. are uncorrelated, $\rho_{01} = 0$.
Putting $\sigma_1^2 := \lambda \sigma_0^2$ ; $\lambda > 0$,

$$\begin{pmatrix} V_\Delta^\circ(p) = \left( \frac{P}{1-P} + \frac{1-P}{P} \lambda \right) \sigma_0^2 \\[2mm] V_\Delta^\infty(p) = \left( \frac{1}{1-P} + \frac{1}{P} \lambda \right) \sigma_0^2 \end{pmatrix}$$

We immediately see that $V_\Delta^\circ(p) < V_\Delta^\infty(p)$ $\forall p \in (0,1)$.

(iii) When the P.O. are perfectly correlated, $\rho_{01} = 1$
[sharp null or constant effect across units]

$$\begin{pmatrix} \overline{V}_\Delta(p) = \left( \frac{P}{1-P} + \frac{1-P}{P} \lambda + 2\sqrt{\lambda} \right) \sigma_0^2 \\[2mm] V_\Delta^\infty(p) = \left( \frac{1}{1-P} + \frac{1}{P} \lambda \right) \sigma_0^2 \end{pmatrix}$$

One can easily see that $\overline{V}_\Delta(p) \leq V_\Delta^\infty(p)$ with equality iff $\lambda = 1 \Leftrightarrow \sigma_0^2 = \sigma_1^2$.

**Summary =**
• If the P.O. are uncorrelated $\quad V_\Delta(p) < V_\Delta^\infty(p) \quad \forall p \in (0,1)$

• If the P.O. are perfectly linearly correlated $\quad V_\Delta(p) \leq V_\Delta^\infty(p) \quad \forall p \in (0,1)$ with equality iff $\sigma_0^2 = \sigma_1^2$

For $X_i \in \mathbb{R}$ (univariate), let

$$\hat{\Delta}_x = \frac{1}{n_t} \sum_{i=1}^n W_i X_i - \frac{1}{n_c} \sum_{i=1}^n (1 - W_i) X_i$$

denote the covariate mean difference between the treatment and control group, given the allocation $\underline{W} = (W_1, \ldots, W_n)$.

We consider a rerandomization criterion based on Mahalanobis distance $\dfrac{\hat{\Delta}_x^2}{\text{var } \hat{\Delta}_x}$

↗ the variance of $\hat{\Delta}_x$ depends on the chosen framework = finite or ∞ pop.

↖ Extends to multivariate $X_i \in \mathbb{R}^p$ using $\hat{\Delta}_x^t \Sigma_x^{-1} \hat{\Delta}_{xj}$
$\Sigma_x =$ covariance matrix of $X_i$

For some $\varepsilon > 0$, put $M_\varepsilon := \left\{ \dfrac{\hat{\Delta}_x^2}{\text{var } \hat{\Delta}_x} \leq \varepsilon \right\}$.

To establish the asymptotic distributions of

$$\begin{pmatrix} \sqrt{n}(\hat{\Delta} - \Delta^n) \\ \sqrt{n}(\hat{\Delta} - \Delta^\infty) \end{pmatrix} \mid M_\varepsilon,$$

we first need the joint asymptotic distribution of

$$\sqrt{n}\begin{pmatrix} \hat{\Delta} - \Delta^n \\ \hat{\Delta}_x \end{pmatrix} \text{ (finite)} \quad \text{and} \quad \sqrt{n}\begin{pmatrix} \hat{\Delta} - \Delta^\infty \\ \hat{\Delta}_x \end{pmatrix} (\infty),$$

generalizing the results of section II.1.

Note that in both cases $\mathbb{E}\,\hat{\Delta}_x = 0$ holds.

---

* **Notation =**

**[finite]** non-random

$$S_{xj} := \frac{1}{n-1} \sum_{i=1}^n (Y_i(j) - \overline{Y}(j))(X_i - \overline{X})$$

$\xrightarrow{\text{some value}} c_{xj}$

$$S_x^2 := \frac{1}{n-1} \sum_{i=1}^n (X_i - \overline{X})^2 \longrightarrow \sigma_x^2$$

$$S_j^2 \to \sigma_j^2 \\ \hat{\rho}_{01} \to \rho_0 \quad\} \text{ see p. 16}$$

**[∞]**

$\begin{pmatrix} X \\ Y(0) \\ Y(1) \end{pmatrix}$ has covariance matrix $\begin{pmatrix} \sigma_x^2 & c_{x0} & c_{x1} \\ * & \sigma_0^2 & c_{01} \\ * & * & \sigma_1^2 \end{pmatrix}$

**Theorem =** Suppose that $\dfrac{n_t}{n} \to p$ as $n \to \infty$. Then

**[finite]** $\sqrt{n}\begin{pmatrix} \hat{\Delta} - \Delta^n \\ \hat{\Delta}_x \end{pmatrix} \xrightarrow{d} \mathcal{N}\left(0, \begin{pmatrix} V_\Delta & V_{\Delta x} \\ * & V_x \end{pmatrix}\right)$

**[∞]** $\sqrt{n}\begin{pmatrix} \hat{\Delta} - \Delta^\infty \\ \hat{\Delta}_x \end{pmatrix} \xrightarrow{d} \mathcal{N}\left(0, \begin{pmatrix} V_\Delta^\infty & V_{\Delta x} \\ * & V_x \end{pmatrix}\right)$

where → $V_\Delta$ and $V_\Delta^\infty$ are defined in (*) and (**) p.16/17

→ $V_{\Delta x} = \dfrac{c_{x0}}{1-p} + \dfrac{c_{x1}}{p}$

$V_x = \dfrac{\sigma_x^2}{p(1-p)}$ ⎫ are common to the two frameworks.

proof = The finite population CLT is used in Li, Ding & Rubin (2018), and proved in Li & Ding (2017) under their regularity condition (18).

In the $[\infty]$ framework,

$\bullet \; \mathrm{cov}(\hat{\Delta}, \hat{\Delta}_x) = \mathrm{cov}\left(\sum_i \left\{\frac{\omega_i Y_i(1)}{n_t} - \frac{(1-\omega_i) Y_i(0)}{n_c}\right\}, \right.$

$$\left. \sum_i \left\{\frac{\omega_i X_i}{n_t} - \frac{(1-\omega_i) X_i}{n_c}\right\}\right)$$

$$= \sum_{i,j} \mathrm{Cov}\left(\frac{\omega_i Y_i(1)}{n_t} - \frac{(1-\omega_i) Y_i(0)}{n_c}, \right.$$

$$\left. \frac{\omega_j X_j}{n_t} - \frac{(1-\omega_j) X_j}{n_c}\right)$$

$$= \sum_i \left\{\frac{\omega_i}{n_t^2} \mathrm{cov}(Y_i(1), X_i) + \frac{1-\omega_i}{n_c^2} \mathrm{Cov}(Y_i(0), X_i)\right\}$$

$$= \frac{C_{x1}}{n_t} + \frac{C_{x0}}{n_c}.$$

$\bullet \; \mathrm{var}\,\hat{\Delta}_x = \left(\frac{1}{n_c} + \frac{1}{n_t}\right)\sigma_x^2 = \frac{n}{n_c n_t}\sigma_x^2$

---

__Theorem__   Li, Ding & Rubin (2018) for [finite]:

$[\text{finite}] \quad \mathrm{var}\,\sqrt{n}(\hat{\Delta} - \Delta^n)\,|\,\mathcal{M}_\varepsilon \longrightarrow V_\Delta(p)\left(1 - (1-\upsilon_{1,\varepsilon})R^2\right)$

$[\,\infty\,] \quad \mathrm{var}\,\sqrt{n}(\hat{\Delta} - \Delta^\infty)\,|\,\mathcal{M}_\varepsilon \longrightarrow V_\Delta^\infty(p)\left(1 - (1-\upsilon_{1,\varepsilon})R_\infty^2\right)$

where

$$R^2 = \frac{V_{\Delta x} V_x^{-1} V_{x\Delta}}{V_\Delta} \quad \& \quad R_\infty^2 = \frac{V_{\Delta x} V_x^{-1} V_{x\Delta}}{V_\Delta^\infty}$$

$$\upsilon_{p,\varepsilon} = \frac{\mathbb{P}(\chi^2_{p+2} \le \varepsilon)}{\mathbb{P}(\chi^2_p \le \varepsilon)} \quad ; \quad X \in \mathbb{R}^P, \; \varepsilon = \text{threshold}$$

$\leftarrow$ The results of Li, Ding & Rubin (2018) extend naturally to the infinite population framework with $V_\Delta$ replaced by $V_\Delta^\infty$.

---

Note that when $X \in \mathbb{R}$ (univariate)

$$R^2 = \frac{\left(\frac{C_{x0}}{1-p} + \frac{C_{x1}}{p}\right)^2}{\frac{\sigma_x^2}{p(1-p)}\left(\frac{p}{1-p}\sigma_0^2 + \frac{1-p}{p}\sigma_1^2 + 2\rho_{01}\sigma_1\sigma_0\right)}$$

$$R_\infty^2 = \frac{\left(\frac{C_{x0}}{1-p} + \frac{C_{x1}}{p}\right)^2}{\frac{\sigma_x^2}{p(1-p)}\left(\frac{\sigma_0^2}{1-p} + \frac{\sigma_1^2}{p}\right)}$$

$f(p) = $ common component of $R^2$ and $R_\infty^2$.

---

__Corollary =__

$[\text{finite}]$   Asymptotic variance is $V_\Delta(p) - (1-\upsilon_{1,\varepsilon})f(p)$

$[\,\infty\,]$   $\quad$ —"— $\quad$ is $V_\Delta^\infty(p) - (1-\upsilon_{1,\varepsilon})f(p)$

Same variance reduction in absolute value in the two frameworks.

---

x __Remark:__ The higher $R^2$ and the more the variance reduction, where $R^2$ corresponds to the squared correlation coefficient between $\hat{\Delta}$ and $\hat{\Delta}_x$ i.e. the proportion of variance of $\hat{\Delta}$ explained by $\hat{\Delta}_x$.

In general, Li, Ding & Rubin (2018) showed that the limit distribution of $\sqrt{n}(\hat{\Delta} - \Delta^n)\,|\,\mathcal{M}_\varepsilon$ of two independent random variables: standard normal + truncated RV:

__Theorem__ (Theorem 1 in Li, Ding & Rubin (2018) )

[finite] $\sqrt{n} (\hat{\Delta} - \Delta^n) | \mathcal{M}_\varepsilon \xrightarrow{d} \sqrt{V_\Delta} \left( \sqrt{1-R^2} \, \mathcal{N}(0,1) + \sqrt{R^2} \, L_{p,\varepsilon} \right)$

where $\circ L_{p,\varepsilon} \sim \chi_{p,\varepsilon} S \sqrt{\beta_p}$

$\searrow \chi^2_{p,\varepsilon} \sim \chi^2_p | \chi^2_p \le \varepsilon$

$\searrow S = \pm 1$ w.p. $1/2$

$\searrow \beta_p \sim \text{Beta} \left( \frac{1}{2}, \frac{p-1}{2} \right)$

— Truncated RV
$\delta$ has variance $v_{p,\varepsilon}$
appearing p. 21 —

[∞] $\sqrt{n} (\hat{\Delta} - \Delta^\infty) | \mathcal{M}_\varepsilon \xrightarrow{d} \sqrt{V_\emptyset^\infty} \left( \sqrt{1-R^2_\infty} \, \mathcal{N}(0,1) + \sqrt{R^2_\infty} \, L_{p,\varepsilon} \right)$

↖ The authors treat the finite population case explicitly only.

When the covariates are uncorrelated with the P.O., the $R^2$ coefficient vanish and we recover the usual CLT (with no reduction in variance).

When $R^2 \ne 0$, the asymptotic distribution of the difference in mean estimator is not normal. When computing confidence intervals based on the asymptotic distribution, this must be taken into account.

Let $\circ \ z_\alpha = \alpha$-th quantile of $\mathcal{N}(0,1)$

$\circ \ u_\alpha(R^2, p, \varepsilon) = \alpha$-th quantile of
$$\sqrt{1-R^2} \, \mathcal{N}(0,1) + \sqrt{R^2} \, L_{p,\varepsilon}$$
$u_\alpha(R^2)$ for simplicity

Under $\mathcal{M}_\varepsilon$, the $(1-\alpha)$ confidence interval for $\sqrt{n}(\hat{\Delta} - \Delta^n)$ is $\left[ u_{\frac{\alpha}{2}}(R^2) \sqrt{V_\Delta} \ , \ u_{1-\frac{\alpha}{2}}(R^2) \sqrt{V_\Delta} \right]$

Under [CRE], the corresponding interval is
$$\left[ z_{\frac{\alpha}{2}} \sqrt{V_\Delta} \ , \ z_{1-\frac{\alpha}{2}} \sqrt{V_\Delta} \right]$$

Li, Ding & Rubin (2018) showed that the [CRE] interval is never less in length than under $\mathcal{M}_\varepsilon$.
⇒ Not taking re-randomization at the time of the analysis leads to conservative confidence bounds for $\hat{\Delta}$.

x __Remark__: When the effect is additive, $R^2$ is the squared multiple correlation between $X$ and $Y(0)$. As noted in Li, Ding & Rubin (2018), " when $\varepsilon$ is close to zero, [...] the asymptotic sampling variance $V_\Delta(p) \left( 1 - (1 - v_{p,\varepsilon}) R^2 \right)$ reduces to $V_\Delta(p) (1 - R^2)$, which is identical to the asymptotic sampling variance of the regression-adjusted estimator under the CRE (see Lin (2013)). Therefore, rerandomization accomplishes covariate adjustments in the design stage, whereas regression accomplishes covariate adjustment in the analysis stage "

Stratified sampling was introduced in Section I. A natural way to extend rerandomization to stratified designs is to rerandomize stratum per stratum until a criterion is reached in all strata. Define

**[finite]** ( stratum k ) **[∞]**

$$\Delta_k^n = \frac{1}{n_k} \sum_{X_i=k} \{Y_i(1) - Y_i(0)\} \; ; \; \Delta_k^\infty = \mathbb{E}(Y_i(1) - Y_i(0) \mid X_i = k)$$
$$= \mu_{1k} - \mu_{0k}$$

$$\Delta^n = \sum_{k=1}^{K} \widehat{\pi}_k \Delta_k^n \; ; \; \Delta^\infty = \sum_{k=1}^{K} \widehat{\pi}_k \Delta_k^\infty$$

where $\widehat{\pi}_k = \frac{n_k}{n}$ ; $n_k = \#$ observations in stratum $k$

$$\widehat{\Delta} = \sum_{k=1}^{K} \widehat{\pi}_k \widehat{\Delta}(k) \text{ with } \widehat{\Delta}(k) = \frac{1}{n_{1k}} \sum_{S_i=k} w_i Y_i - \frac{1}{n_{0k}} \sum_{S_i=k} (1-w_i) Y_i$$
$$= \text{DM estimator in } k\text{-th stratum}$$

$$\widehat{\Delta}_x = \sum_{k=1}^{K} \widehat{\pi}_k \widehat{\Delta}_x(k) \text{ with } \widehat{\Delta}_x(k) = \frac{1}{n_{1k}} \sum_{S_i=k} w_i X_i - \frac{1}{n_{0k}} \sum_{S_i=k} (1-w_i) X_i$$

$S_i = k \Leftrightarrow$ unit $i$ belong to stratum $k$
$X_i = $ covariate on which balance is required.
(slight change of notation from section I)

$p_k = \frac{n_{1k}}{n_k}$

$V_{x,k} = \frac{\sigma_{x,k}^2}{p_k(1-p_k)} \; ; \; V_{\Delta x, k} = \frac{C_{x0,k}}{1-p_k} + \frac{C_{x1,k}}{p_k}$

---

$V_{\Delta,k}^\infty = \frac{\sigma_{1,k}^2}{p_k} + \frac{\sigma_{0,k}^2}{1-p_k}$

$V_{\Delta,k} = \frac{p_k}{1-p_k} \sigma_{0,k}^2 + \frac{1-p_k}{p_k} \sigma_{1,k}^2 + 2 \rho_{01,k} \sigma_{0,k} \sigma_{1,k}$ **[finite]**

$R_k^2 = \frac{V_{\Delta x,k} V_{x,k}^{-1} V_{x\Delta,k}}{V_{\Delta,k}} \; ; \; R_{\infty,k}^2 = \frac{V_{\Delta x,k} V_{x,k}^{-1} V_{x\Delta,k}}{V_{\Delta,k}^\infty}$

We consider a rerandomization procedure where each stratum is rerandomized independently, based on the stratum-specific Mahalanobis distances. For a stratum $k$, specify a threshold $\varepsilon_k$, and compute $\frac{\widehat{\Delta}_x^2(k)}{\text{var } \widehat{\Delta}_x(k)}$.

Accept the stratum $k$ assignment vector if this quantity is less than $\varepsilon_k$. In the multivariate case, compute $\widehat{\Delta}_x^t(k) \; \text{cov}^{-1} \widehat{\Delta}_x(k) \; \widehat{\Delta}_x(k)$. For $\varepsilon = (\varepsilon_1, \dots, \varepsilon_K)$,

put $\boxed{\mathcal{M}_\varepsilon = \left\{ (w_1, \dots, w_n) \mid \frac{\widehat{\Delta}_x^2(k)}{\text{var } \widehat{\Delta}_x(k)} \leq \varepsilon_k \; \forall k \right\}}$

↑ Stratum-specific Mahalanobis distance

Wang, Wang, Liu (2021)

Then

**[finite]** $\sqrt{n} \begin{pmatrix} \widehat{\Delta} - \Delta^n \\ \widehat{\Delta}_x \end{pmatrix} \xrightarrow{d} \mathcal{N}(0, V)$ ⟵ only

where $V = \sum_{k=1}^{K} \widehat{\pi}_k \left( \begin{array}{c|c} V_{\Delta,k} & V_{\Delta x,k} \\ \hline V_{x\Delta,k} & V_{x,k} \end{array} \right)$

**[∞]** $\sqrt{n} \begin{pmatrix} \widehat{\Delta} - \Delta^\infty \\ \widehat{\Delta}_x \end{pmatrix} \xrightarrow{d} \mathcal{N}(0, V^\infty)$

where $V^\infty = \sum_{k=1}^{K} \widehat{\pi}_k \left( \begin{array}{c|c} V_{\Delta,k}^\infty & V_{\Delta x,k} \\ \hline V_{x\Delta,k} & V_{x,k} \end{array} \right)$

Under some general conditions, Wang, Wang, Liu (2021) (27)
established for the finite population framework that

$$\sqrt{n}\left(\hat{\Delta} - \Delta^n\right) \mid \mathcal{M}_\varepsilon \sim \sqrt{\sum_{k=1}^{K} \hat{\pi}_k V_{\Delta, k}(1 - R_k^2)} \; \mathcal{N}(0,1)$$

$K$ fixed, $n_k \to \infty \; \forall k$ $\quad + \sqrt{\sum_{k=1}^{K} \hat{\pi}_k V_{\Delta, k} R_k^2} \; L_{P, \varepsilon_k}^k$,

where $L_{P, \varepsilon_k}^k$ are independent RVs $\sim L_{P, \varepsilon_k}$.

The asymptotic variance of $\sqrt{n}\left(\hat{\Delta} - \Delta^n\right)$ under $\mathcal{M}_\varepsilon$ is

$$\sum_{k=1}^{K} \hat{\pi}_k V_{\Delta, k} \left\{ 1 - (1 - \nu_{P, \varepsilon_k}) R_k^2 \right\}.$$

x Remark: These expressions should be compared with the
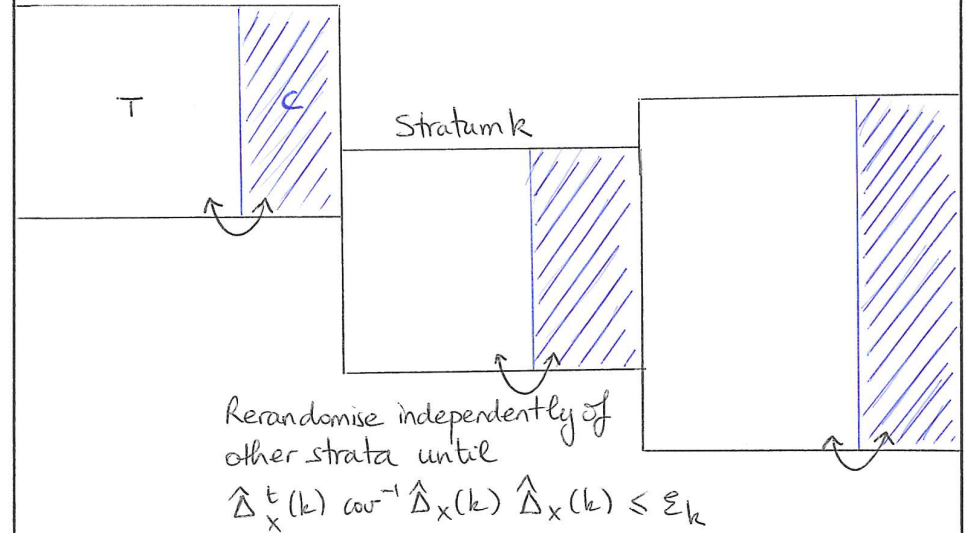asymptotic distribution of the aggregated difference in
means estimator under [ST]:

$$\sqrt{n}\left(\hat{\Delta} - \Delta^n\right) \xrightarrow{d} \mathcal{N}\left(0, \sum_{k=1}^{K} \hat{\pi}_k V_{\Delta, k}\right).$$

Alternatively, we may consider an overall Mahalanobis
distance $\hat{\Delta}_x^t \; cov^{-1} \hat{\Delta}_x \; \hat{\Delta}_x$ and accept the full
allocation $(w_1, \ldots, w_n)$ if this quantity is less than
some value $\varepsilon$.

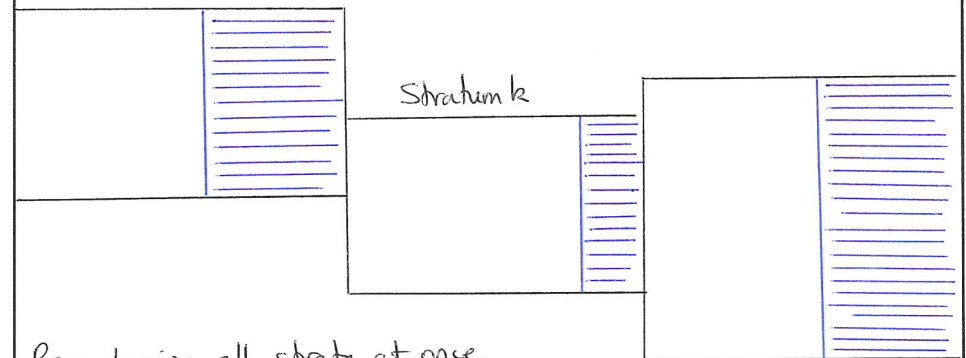↖ Here, strata are all rerandomized at one,
and all accepted at the same time, which
differs from the previous rerandomization scheme.

[Scheme 1]  Stratum-Specific Distance (28)



Rerandomise independently of
other strata until
$$\hat{\Delta}_x^t(k) \; cov^{-1} \hat{\Delta}_x(k) \; \hat{\Delta}_x(k) \leq \varepsilon_k$$

[Scheme 2]  Overall Distance



Rerandomize all strata at once
until $\hat{\Delta}_x^t \; cov^{-1} \hat{\Delta}(k) \; \hat{\Delta}_x \leq \varepsilon.$

Wang, Wang, Liu (2021) showed that for identical thresholds
$\varepsilon_1 = \cdots = \varepsilon_K = \varepsilon$ in [Scheme 1], the asymptotic
variance of $\sqrt{n}\left(\hat{\Delta} - \Delta^n\right)$ is smaller than for [Scheme 2].

Specifically,

$$\sum_{k=1}^{K} \widehat{\pi}_k V_{\Delta,k} \left\{ 1 - (1 - \sigma_{P,\varepsilon}) R_k^2 \right\} \quad \text{Scheme 1}$$

$$\text{Scheme 2} \quad \leqslant V_\Delta \left\{ 1 - (1 - \sigma_{P,\varepsilon}) \right\} R^2$$

where

$$\left( \begin{array}{c|c} V_\Delta & V_{\Delta X} \\ \hline V_{X\Delta} & V_X \end{array} \right) = \sum_{k=1}^{K} \widehat{\pi}_k \left( \begin{array}{c|c} V_{\Delta,k} & V_{\Delta X,k} \\ \hline V_{X\Delta,k} & V_{X,k} \end{array} \right)$$

and $R^2 = \dfrac{V_{\Delta X} V_X^{-1} V_{X\Delta}}{V_\Delta}$

---

[Appendix A]

× Consider first $\widehat{\Delta}_{agg}$ with $W \sim [ST]$.

Mean $\quad \mathbb{E} \widehat{\Delta}_{agg} = \sum_k \widehat{\pi}_k \, \mathbb{E} \, \widehat{\Delta}(k)$ ; $\mathbb{E} \widehat{\Delta}(k) = \mu_{1k} - \mu_{0k}$

Variance $\quad \text{var} \, \widehat{\Delta}_{agg} = \sum_k \widehat{\pi}_k^2 \, \text{var} \, \widehat{\Delta}(k)$

with $\text{var} \, \widehat{\Delta}(k) = \dfrac{\sigma_{1k}^2}{n_{1k}} + \dfrac{\sigma_{0k}^2}{n_{0k}}$

$$\text{var} \, \widehat{\Delta}_{agg} = \sum_k \frac{n_k}{n} \widehat{\pi}_k \left( \frac{\sigma_{1k}^2}{n_{1k}} + \frac{\sigma_{0k}^2}{n_{0k}} \right)$$

$$= \frac{1}{n} \sum_k \widehat{\pi}_k \left( \frac{\sigma_{1k}^2}{p_k} + \frac{\sigma_{0k}^2}{1 - p_k} \right)$$

× We turn our attention to $\widehat{\Delta} = \dfrac{1}{n_t} \sum W_i Y_i - \dfrac{1}{n_c} \sum (1 - W_i) Y_i$

Mean $\qquad\qquad\qquad\qquad W \sim [CRE]$

$$\mathbb{E} \widehat{\Delta} = \frac{n}{n_t} \underbrace{\mathbb{E}(WY)}_{\substack{\| \\ \mathbb{E}W \, \mathbb{E}Y(1) \text{ since } W \perp Y(1) \\ = \frac{n_t}{n} \mu_1}} - \frac{n}{n_c} \mathbb{E}\left[ (1-W) Y \right] = \mu_1 - \mu_0$$

Variance: $\text{var} \, \widehat{\Delta} = \dfrac{\text{var} \, Y(1)}{n_t} + \dfrac{\text{var} \, Y(0)}{n_c}$

with $\sigma_j^2 = \text{var} \, Y(j) = \sum_k \widehat{\pi}_k \sigma_{jk}^2 + \sum_k \widehat{\pi}_k (\mu_{jk} - \mu_j)^2$

## [Appendix B]

Consider $\hat{\Delta} = \frac{1}{n_E} \sum_{i=1}^{n} w_i Y_i - \frac{1}{n_C} \sum_{i=1}^{\hat{n}} (1-w_i) Y_i$,

where the $w_i \sim [ST]$.

Here, $\{Y_i(0), Y_i(1)\} \perp w_i \mid X_i = k$, $k = 1, \dots, K$, but not unconditionally.

• **Mean** $\mathbb{E}\hat{\Delta} = \frac{n}{n_E} \mathbb{E} w Y(1) - \frac{n}{n_C} \mathbb{E}(1-w) Y(0)$

with $\mathbb{E} w Y(1) = \mathbb{E}_X \underbrace{\mathbb{E} w Y(1) \mid X}_{} = \sum_k \hat{\pi}_k p_k \mu_{1k}$.

since $\mathbb{E} w Y(1) \mid X = k$
$$= \mathbb{E}(w \mid X = k)\, \mathbb{E}(Y(1) \mid X = k)$$
$$= \frac{n_{1k}}{n_k} \mu_{1k} = p_k \mu_{1k}$$

& similarly for $\mathbb{E}(1-w) Y(0) = \sum_k \hat{\pi}_k (1-p_k) \mu_{0k}$

$$\Rightarrow \mathbb{E}\hat{\Delta} = \sum_k \hat{\pi}_k \left( \frac{n p_k}{n_E} \mu_{1k} - \frac{n(1-p_k)}{n_C} \mu_{0k} \right)$$

• **Variance** Put $z_k := \sum_{X_i = k} \left\{ \frac{w_i Y_i}{n_E} - \frac{(1-w_i) Y_i}{n_C} \right\}$

$z_1, \dots, z_K$ are independent, and $\hat{\Delta} = \sum_{k=1}^{K} z_k$. $\leftarrow$ $n_k$ terms

$$\text{var } z_k = \frac{n_{1k}}{n_E^2} \sigma_{1k}^2 + \frac{n_{0k}}{n_C^2} \sigma_{0k}^2 \quad \rightarrow = (1-p_k) \hat{\pi}_k \left(\frac{n}{n_C}\right)^2 \frac{1}{n}$$

$$\Rightarrow \text{var } \hat{\Delta} = \sum_k \left( \frac{n_{1k}}{n_E^2} \sigma_{1k}^2 + \frac{n_{0k}}{n_C^2} \sigma_{0k}^2 \right)$$

$$\hookrightarrow = \frac{n_{1k}}{n_k} \frac{n_k}{n} n \frac{1}{n_E^2} = p_k \hat{\pi}_k \left(\frac{n}{n_E}\right)^2 \frac{1}{n}$$

## [Appendix C]

• $\text{var } \hat{\Delta}_{agg} = \mathbb{E} \text{ var}\left(\hat{\Delta}_{agg} \mid \{n_{1k}, n_{0k}\}\right) + \text{var } \mathbb{E}\left(\hat{\Delta}_{agg} \mid \{n_{1k}, n_{0k}\}\right)$

$\swarrow$ p. 6 $\qquad\qquad$ $\swarrow$ p. 6

$$= \mathbb{E}\left( \frac{1}{n} \sum_{k=1}^{K} \hat{\pi}_k \left( \frac{\sigma_{1k}^2}{p_k} + \frac{\sigma_{0k}^2}{1-p_k} \right) \right) \qquad = \text{var}\left( \sum_{k=1}^{K} \hat{\pi}_k (\mu_{1k} - \mu_{0k}) \right)$$

$$= \frac{1}{n} \sum_{k=1}^{K} \pi_k \left( \frac{\sigma_{1k}^2}{p_k} + \frac{\sigma_{0k}^2}{1-p_k} \right) \qquad = \frac{1}{n^2} \sum_{k, \ell} (\mu_{1k} - \mu_{0k})(\mu_{1\ell} - \mu_{0\ell})$$
$$\times \text{Cov}(n_k, n_\ell)$$

where $\text{Cov}(n_k, n_\ell) = -n \pi_k \pi_\ell$

$\left( (n_1, \dots, n_K) \sim \text{Mult}(\pi_1, \dots, \pi_K) \right)$

& $\text{var } n_k = n \pi_k (1-\pi_k)$.

$$= -\frac{1}{n} \sum_{k \neq \ell} (\mu_{1k} - \mu_{0k})(\mu_{1\ell} - \mu_{0\ell}) \pi_k \pi_\ell + \frac{1}{n} \sum_k (\mu_{1k} - \mu_{0k})^2 \pi_k (1-\pi_k)$$

$$= -\frac{1}{n} \sum_{k, \ell} (\mu_{1k} - \mu_{0k})(\mu_{1\ell} - \mu_{0\ell}) \pi_k \pi_\ell + \frac{1}{n} \sum_k \pi_k (\mu_{1k} - \mu_{0k})^2$$

$$= -\frac{1}{n} \left[ \sum_k (\mu_{1k} - \mu_{0k}) \pi_k \right]^2 + \frac{1}{n} \sum_k \pi_k (\mu_{1k} - \mu_{0k})^2$$

$$= \frac{1}{n} \text{var } \Delta(X), \quad \text{where} \quad \Delta(k) := \mathbb{E}(Y(1) - Y(0) \mid X = k).$$

• $\text{var } \hat{\Delta}$ is derived like in [Appendix A] with $\hat{\pi}_k$ replaced by the population proportions $\pi_k$

$$= \frac{1}{n_E} \left( \sum_k \pi_k (\sigma_{1k}^2 + (\mu_{1k} - \mu_1)^2) \right) + \frac{1}{n_C} \left( \sum_k \pi_k (\sigma_{0k}^2 + (\mu_{0k} - \mu_0)^2) \right)$$

[Appendix D]

• $\mathbb{E}\hat{\Delta} = \mathbb{E}\,\mathbb{E}\,\hat{\Delta} \mid \{n_k\}$  ↘ p.8

$$= \sum_{k=1}^{K} \left[ \mathbb{E}\left(\hat{\pi}_k \frac{n\,p_k}{n_t}\right) \mu_{1k} - \mathbb{E}\left(\hat{\pi}_k \frac{n(1-p_k)}{n_c}\right) \mu_{0k} \right]$$

Here $n_t$ is random since $n_t = \sum_{k=1}^{K} p_k n_k$

$$\mathbb{E}\,n_t = n\left(\sum_{k=1}^{K} p_k \pi_k\right)$$

$\hat{\pi}_k$ = 

$$\mathbb{E}\left(\frac{n_k}{n}\, n\, \frac{p_k}{n_t}\right) = p_k\, \mathbb{E}\left(\frac{n_k}{n_t}\right)$$

$$= p_k\left(\frac{\mathbb{E}\,n_k}{\mathbb{E}\,n_t} + O(n^{-1})\right)$$

$$= \frac{p_k\,\pi_k}{\sum_{\ell=1}^{K} p_\ell \pi_\ell} + O(n^{-1}).$$

$$\Rightarrow \mathbb{E}\hat{\Delta} = \sum_{k=1}^{K} \pi_k\left(\frac{p_k}{\sum_\ell p_\ell \pi_\ell}\,\mu_{1k} - \frac{1-p_k}{\sum_\ell (1-p_\ell)\pi_\ell}\,\mu_{0k}\right)$$

↕                    ↕

Compare these coefficients with the value $1$; the reference point for unbiasedness.



$\sum p_\ell \pi_\ell$ lies in the interior of the convex hull of $\{p_1, \ldots, p_K\} \Rightarrow$ all coefs $\frac{p_k}{\sum p_\ell \pi_\ell} \neq 1$ unless all $p_k$ are equal.