

CI = UNCONFOUNDEDNESS

- Consider an RCT with $\{Y_i(0), Y_i(1)\} \perp W_i$

Potential outcomes with a population model $(Y_i(0), Y_i(1)) \sim \mathbb{P}$
 $\sim B(\pi)$
 $\in \{0, 1\}$

for $i=1, \dots, n$ units, $Y_i = W_i Y_i(1) + (1-W_i) Y_i(0)$,
 $n_1 = \sum_{i=1}^n W_i$ and $n_0 = n - n_1$.

- In CI: RANDOMIZED CONTROL TRIALS (p. 4), we introduced

$$\hat{\Delta} = \frac{1}{n_1} \sum_{i=1}^n W_i Y_i - \frac{1}{n_0} \sum_{i=1}^n (1-W_i) Y_i$$

(DIFFERENCE ESTIMATOR)

as a natural estimator of the ATE $\Delta^\infty = \mathbb{E}(Y_i(1) - Y_i(0))$.

- Alternatively, noting that $\mathbb{E}\left(\sum_{i=1}^n W_i\right) = n\pi$ and $\mathbb{E}\left(\sum_{i=1}^n (1-W_i)\right) = n(1-\pi)$, we may replace n_0 (resp. n_1) in $\hat{\Delta}$ by $n(1-\pi)$ (resp. $n\pi$) and define:

$$\tilde{\Delta} = \frac{1}{n} \sum_{i=1}^n \frac{W_i Y_i}{\pi} - \frac{1}{n} \sum_{i=1}^n \frac{(1-W_i) Y_i}{1-\pi}$$

A natural estimator of Δ^∞ to consider as well.
 = the Inverse Probability Weighting (IPW) estimator of Δ^∞ in a Bernoulli RCT.

As we will see in the next sections, IPW-like estimators play a central role when estimating causal effects in observational studies (under unconfoundedness) ↴

I - UNCONFOUNDEDNESS

[From Wager (2020)]

The simplest generalization of one RCT = multiple RCT
 Suppose that we have $x=1, \dots, K$ RCT i.e.

$$\{Y_i(0), Y_i(1)\} \perp W_i \mid X_i \quad (*)$$

For example, letting $x = \text{'Paris'}, \text{'Berlin'}, \text{'Melbourne'}$, condition (*) means that three RCTs were conducted in these three different cities.

We write $e(x) = \mathbb{P}(W=1 \mid X=x)$

If the same Bern(π) randomization scheme was used in each of the three cities, then $e(x) = \pi$ for $x = \text{'Paris'}, \text{'Berlin'}$ and 'Melbourne' .

We still want $\Delta^\infty = \mathbb{E}\{Y_i(1) - Y_i(0)\}$.

- $\forall x \in X = \{1, \dots, K\}$, let

$\forall x$, consider the diff in means estimator $\hat{\Delta}(x)$

$$\hat{\Delta}(x) = \frac{1}{n_{x1}} \sum_{\substack{W_i=1 \\ X_i=x}} Y_i - \frac{1}{n_{x0}} \sum_{\substack{W_i=0 \\ X_i=x}} Y_i$$

Agg. Diff. in Means (ADM)

$$\hat{\Delta}_{ADM} = \sum_{x \in X} \frac{n_x}{n} \hat{\Delta}(x) \quad \begin{matrix} n_x = n_{x1} + n_{x0} \\ n = \sum_{x \in X} n_x \end{matrix}$$

Put $\pi_x = P(X_i = x)$

$$\hat{\pi}_x = n_x / n$$

(3)

Q: How accurate is $\hat{\Delta}_{ADM}$?

Put $\Delta^\infty = \sum_{x \in X} \pi_x \Delta(x)$; $\Delta(x) = \mathbb{E}\{Y_i(1) - Y_i(0) \mid X_i = x\}$
(the Conditional ATE) (CATE)

Then

$$\hat{\Delta}_{ADM} - \Delta^\infty = \left[\begin{aligned} & \sum_{x \in X} \pi_x (\hat{\Delta}(x) - \Delta(x)) \\ & + \sum_{x \in X} (\hat{\pi}_x - \pi_x) \Delta(x) \\ & + \sum_{x \in X} (\hat{\pi}_x - \pi_x) (\hat{\Delta}(x) - \Delta(x)) \end{aligned} \right] = O_p(1/n)$$

Black term. Assuming $\text{var}(Y_i(w) \mid X_i = x) = \sigma^2(x)$
($w = 0, 1$),

$$\text{then } n_x^{1/2} (\hat{\Delta}(x) - \Delta(x)) \xrightarrow{d} \mathcal{N}\left(0, \frac{\sigma^2(x)}{e(x)(1-e(x))}\right)$$

$$\Rightarrow \sum_x \pi_x (\hat{\Delta}(x) - \Delta(x)) \approx \mathcal{N}\left(0, \frac{1}{n} \sum_x \pi_x^2 \frac{\sigma^2(x)}{\pi_x e(x)(1-e(x))}\right)$$

The asymptotic variance is

$$\sum_x \pi_x \frac{\sigma^2(x)}{e(x)(1-e(x))} = \mathbb{E}\left\{ \frac{\sigma^2(X)}{e(X)(1-e(X))} \right\}$$

Blue term = $\sum_x \frac{n_x}{n} \Delta(x) - \sum_x \Delta(x) \mathbb{P}(X_i = x)$

$$= \frac{1}{n} \sum_x n_x \mathbb{E}\{Y_i(1) - Y_i(0) \mid X_i = x\} - \mathbb{E}\{\Delta(X)\}$$

$$= \frac{1}{n} \sum_{i=1}^n \Delta(X_i) - \mathbb{E}\{\Delta(X)\}$$

$\hookrightarrow \text{CLT} \approx \mathcal{N}\left(0, \frac{\text{var } \Delta(X)}{n}\right)$

Putting pieces together:

$$n^{1/2} (\hat{\Delta}_{ADM} - \Delta^\infty) \xrightarrow{d} \mathcal{N}(0, V_{ADM})$$

where

$$V_{ADM} = \text{Var}\{\Delta(X)\} + \mathbb{E}\left\{ \frac{\sigma^2(X)}{e(X)(1-e(X))} \right\}$$

$\uparrow X \in \{1, \dots, K\}$ has discrete support

What about X_i with continuous support?

$\rightarrow \{Y_i(0), Y_i(1)\} \perp W_i \mid X_i$ is called UNCONFOUNDEDNESS
introduced by Rosenbaum & Rubin (1984)

$\rightarrow e(x) = P(W_i = 1 \mid X_i = x)$ is called the PROPENSITY SCORE (PS)

x Remark: Unconfoundedness $\Rightarrow \{Y_i(0), Y_i(1)\} \perp W_i \mid e(X_i)$
(see p. 64/65 in CI: ELEMENTS OF CAUSAL INFERENCE)

Algorithmic implication: cut into buckets along $e(X_i)$
and use $\hat{\Delta}_{ADM}$ (even though there is still a bit of confounding in each of the buckets).

II. INVERSE PROBABILITY WEIGHTING (IPW)

A conceptual implication of the previous remark is that the propensity score plays a central role. There are many ways to use PS for ATE estimation. One way is IPW.

(4)

$$\hat{\Delta}_{IPW}^* = \frac{1}{n} \sum_{i=1}^n \left(\frac{W_i Y_i}{e(X_i)} - \frac{(1-W_i) Y_i}{1-e(X_i)} \right) \quad \text{= oracle (Horvitz-Thompson)} \quad (5)$$

Compare with $\hat{\Delta}$ on page 1 defined in the context of a single RCT: π is replaced by the PS $e(X_i) = P(W_i=1 | X_i=1)$, a consequence of $\{Y_i(0), Y_i(1)\} \perp W_i | X_i=x$

$\hat{\Delta}_{IPW}^*$ is an oracle estimator since the PS is assumed to be known. In practice, we consider

$$\hat{\Delta}_{IPW} = \frac{1}{n} \sum_{i=1}^n \left(\frac{W_i Y_i}{\hat{e}(X_i)} - \frac{(1-W_i) Y_i}{1-\hat{e}(X_i)} \right)$$

where $\hat{e}(\cdot)$ is an estimate of $e(\cdot)$.

* Result = $E(\hat{\Delta}_{IPW}^*) = \Delta^\infty$

proof: $E\left(\frac{W_i Y_i}{e(X_i)} - \frac{(1-W_i) Y_i}{1-e(X_i)}\right)$

$$= E\left(\frac{W_i Y_i(1)}{e(X_i)} - \frac{(1-W_i) Y_i(0)}{1-e(X_i)}\right) \quad \text{(SUTVA)}$$

$$= E\left(E\left(\frac{\quad}{\quad} \mid X_i\right)\right) \quad \text{unconfoundedness}$$

$$= E\left\{\frac{E(W_i | X_i) E(Y_i(1) | X_i)}{e(X_i)} - \dots\right\}$$

$$= E\{Y_i(1) - Y_i(0)\}$$

$$= \Delta^\infty$$

In addition, $n^{1/2}(\hat{\Delta}_{IPW}^* - \Delta^\infty) \xrightarrow{d} \mathcal{N}(0, V_{IPW})$ (6)

where $V_{IPW} = \text{var}\left[\left(\frac{W_i}{e(X_i)} - \frac{1-W_i}{1-e(X_i)}\right) Y_i\right]$ since

the oracle estimator $\hat{\Delta}_{IPW}^*$ is a sum of iid random variables. We derive next an alternative expression for V_{IPW} :

Put \bullet $E(Y_i(w) | X_i) = \mu_{(w)}(X_i)$

\bullet $\text{var}(Y_i(w) | X_i) = \text{var}(\varepsilon_i(w) | X_i) = \sigma^2(X_i)$

\bullet $c(x) := e(x) \mu_{(0)}(x) + (1-e(x)) \mu_{(1)}(x)$

just for convenience; it will simplify calculations of V_{IPW} .

It allows us to write $\mu_{(0)}(X_i)$

$$Y_i(0) = c(X_i) - (1-e(X_i)) \Delta(X_i) + \varepsilon_i(0)$$

$$Y_i(1) = c(X_i) + e(X_i) \Delta(X_i) + \varepsilon_i(1)$$

$\mu_{(0)}(X_i)$

$\Delta(X_i) = \mu_{(1)}(X_i) - \mu_{(0)}(X_i)$

$(= E(Y_i(1) - Y_i(0) | X_i))$

\Downarrow

$$V_{IPW} = \text{var}\left\{\left(\frac{W_i}{e(X_i)} - \frac{1-W_i}{1-e(X_i)}\right) c(X_i)\right.$$

$$\left. + \Delta(X_i) + \left(\frac{W_i \varepsilon_i(1)}{e(X_i)} + \frac{(1-W_i) \varepsilon_i(0)}{1-e(X_i)}\right)\right\}$$

these three terms are uncorrelated

x Remarks: (a) We established consistency and a CLT

(9)

for the oracle estimator $\hat{\Delta}_{IPW}^*$. Assuming that

(i) we have OVERLAP $\eta \leq e(x) \leq 1 - \eta \quad \forall x \in \mathcal{X}$

(ii) $\exists M$ s.t. $|Y_i| \leq M$

(iii) $\sup_{x \in \mathcal{X}} |\hat{e}(x) - e(x)| = O_p(a_n) \rightarrow 0$,

then one can show that $|\hat{\Delta}_{IPW} - \hat{\Delta}_{IPW}^*| = O_p\left(\frac{a_n M}{\eta}\right)$.

In other words, under assumptions (i)-(ii)-(iii),

$\hat{\Delta}_{IPW}^*$ is consistent $\Rightarrow \hat{\Delta}_{IPW}$ is consistent.

To prove this, we use that $\hat{e}(x)$ becomes uniformly bounded away from 0 and 1 as $n \rightarrow \infty$ under (iii), which in turn implies that $1/\hat{e}(x)$ and $1/1-\hat{e}(x)$ decay at the same $O_p(a_n)$ rate as $\hat{e}(x)$.

(b) Rewriting $q(w, x) = 1 / \mathbb{P}(W=w | X=x)$,

we see that $\mathbb{E}(W q(W, X)) = \mathbb{E}((1-W) q(W, X)) = 1$.

In addition,

$$\begin{aligned} \mathbb{E} Y_i(1) &= \mathbb{E} \left(\frac{W_i Y_i}{e(X_i)} \right) = \mathbb{E} (W_i Y_i q(W_i, X_i)) \\ &= \frac{\mathbb{E}(W_i Y_i q(W_i, X_i))}{\mathbb{E}(W_i q(W_i, X_i))} \leftarrow = 1 \end{aligned}$$

and similarly for $\mathbb{E} Y_i(0)$. This suggests defining a version

of the IPW estimator using normalized weights:

(10)

$$\begin{aligned} \hat{\Delta}_{IPW}^* &= \sum_{i|W_i=1} \left\{ \frac{q(W_i, X_i)}{\sum_{j|W_j=1} q(W_j, X_j)} \right\} Y_i - \sum_{i|W_i=0} \left\{ \frac{q(W_i, X_i)}{\sum_{j|W_j=0} q(W_j, X_j)} \right\} Y_i \\ &= \sum_{i|W_i=1} q'(W_i, X_i) Y_i - \sum_{i|W_i=0} q'(W_i, X_i) Y_i \end{aligned}$$

(Robins '98)

Normalized weights are often preferred over original weights, especially when $e(x)$ is close to 0 or 1, which increases the variance of the IPW estimator.

III. BALANCING WEIGHTS

IPW is a special case of a general class of balancing weights. We introduced p.3 the CATE

$$\Delta(x) = \mathbb{E}\{Y_i(1) - Y_i(0) | X=x\}$$

Usually, the CATE is not computed for a single x , but averaged over a target distribution of the covariates. Let

$\rightarrow f(x)$ = marginal density of X over the whole pop.

$\rightarrow g(x)$ = density of X of a target population, possibly different from X .

Ex: Medical treatment applied to the whole population ($X \sim f$) vs Medical treatment applied to those individuals who need it the most ($X \sim g$)

$\hookrightarrow h(x) = \frac{g(x)}{f(x)}$ reweights observations to represent 11
the target population.

We introduce next a new class of estimators: the ATE over the target population g (Li, Morgan, Zaslavsky 2018)

$$\Delta_h^\infty := \mathbb{E}_g [Y_i(1) - Y_i(0)] = \frac{\mathbb{E}\{h(X)\Delta(X)\}}{\mathbb{E}\{h(X)\}} \quad \text{X w/ f}$$

$$= \frac{\int \Delta(x) f(x) h(x) dx}{\int f(x) h(x) dx}$$

Let $f_w(x) = \mathbb{P}(X=x | W=x)$
(working w.l.o.g. with discrete distributions here)
= density of X in the $W=w$ group.

Then

$$f_1(x) = \frac{\mathbb{P}(X=x, W=1)}{\mathbb{P}(W=1)} = \frac{\mathbb{P}(W=1 | X=x) \mathbb{P}(X=x)}{\mathbb{P}(W=1)} \propto f(x) e(x)$$

and similarly $f_0(x) \propto f(x)(1-e(x))$

For a specific $h(x)$, we can estimate Δ_h^∞ using weights $w_0(x), w_1(x)$ defined such that

$$w_0(x) f_0(x) = w_1(x) f_1(x) = \underbrace{h(x) f(x)}_{\text{target } g(x)}$$

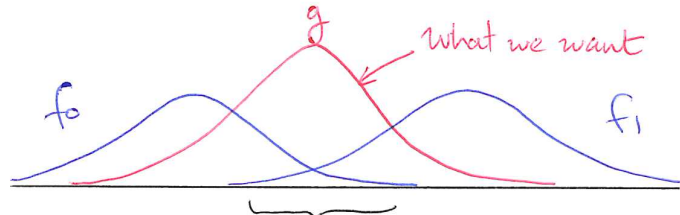
\uparrow
 \uparrow

reweight each $W=w$ group

$$\Rightarrow \begin{cases} w_0(x) \propto \frac{f(x)h(x)}{f(x)(1-e(x))} = \frac{h(x)}{1-e(x)} \\ w_1(x) \propto \frac{f(x)h(x)}{f(x)e(x)} = \frac{h(x)}{e(x)} \end{cases} \quad \text{12}$$

- Ex: $h(x) = 1$; the target population is $f(x)$; weights are $(w_0(x), w_1(x)) = \left(\frac{1}{1-e(x)}, \frac{1}{e(x)}\right)$ and $\Delta_h^\infty = \Delta^\infty = \mathbb{E}\{Y_i(1) - Y_i(0)\}$ IPW
- Ex: $h(x) = e(x)$; the target population is $f(x)e(x) \propto f_1(x)$ i.e. the treated subpopulation. The weights are $(w_0(x), w_1(x)) = \left(\frac{e(x)}{1-e(x)}, 1\right)$ and the estimand is $\Delta_h^\infty = \text{ATT} = \mathbb{E}\{Y_i(1) - Y_i(0) | W=1\}$
- Ex: $h(x) = 1-e(x)$, the target population is the control subpopulation. The weights are $(w_0(x), w_1(x)) = \left(1, \frac{1-e(x)}{e(x)}\right)$ and $\Delta_h^\infty = \text{ATNT} = \mathbb{E}\{Y_i(1) - Y_i(0) | W=0\}$
- Ex: $h(x) = \text{indicator function} = \mathbb{1}(\alpha < e(x) < 1-\alpha)$ leads to an ATE for a subpopulation $\alpha \in (0, 1/2)$ with overlap of the covariates between the two groups.
- Ex: $h(x) = e(x)(1-e(x))$ yields $(w_0, w_1) = (e(x), 1-e(x))$. The estimand $\Delta_h^\infty = \mathbb{E}\{e(x)(1-e(x))\Delta(x)\} / \mathbb{E}\{e(x)(1-e(x))\}$ is called the ATE for the overlap population.

Indeed, $h(x)$ is maximal when $e(x) = 1/2$; i.e. for those individuals which have an equal chance to be allocated to a treatment or control group.



grey area : in medical applications, unsure if this patient should be given this new treatment or not \Rightarrow we want to focus on these patients the most, and this is exactly what $h(x)$ is doing.

IV - AUGMENTED IPW

We refer to Chapter 3 in Wager (2020) for a formal account on AIPW.

BALANCING WEIGHTS

(From a keynote talk of Betsy Ogburn at EURO CIM '24)

[REF] Augmented balancing weights as linear regression
D. Brun-Smith, O. Duker, A. Feller & B. Ogburn.

• Consider iid tuples $(X_i, T_i, Y_i(0), Y_i(1))$
 $\begin{matrix} \mathbb{R}^d & \{0,1\} & \mathbb{R} & \mathbb{R} \\ \text{(cov)} & \text{(trt alloc)} & & \end{matrix}$

↳ We are interested in estimating $\Delta = \mathbb{E}[Y(1) - Y(0)]$
& focus on the counterfactual mean $\mu_1 := \mathbb{E}[Y(1)]$.

• Usual assumptions → Ignorability $T \perp \{Y(0), Y(1)\} \mid X$
→ Overlap $e(x) = \mathbb{P}(T=1 \mid X=x) > 0$

• Two common strategies:

→ IPW
 $\mu_1 = \mathbb{E}\left(\frac{T}{e(X)} Y\right)$
 under ignorability

OUTCOME MODELLING
 $\mu(x, t) = \mathbb{E}(Y \mid X=x, T=t)$

$\hat{\Delta}_{REG} = \frac{1}{n} \sum_{i=1}^n \hat{\mu}(X_i, 1) - \hat{\mu}(X_i, 0)$

$$\hat{\Delta}_{IPW} = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{T_i}{\hat{e}(X_i)} Y_i - \frac{(1-T_i)}{1-\hat{e}(X_i)} Y_i \right\}$$

↳ Doubly Robust estimators combine both approaches.

Ex: AIPW

$$\hat{\Delta}_{AIPW} = \frac{1}{n} \sum_{i=1}^n \left\{ \hat{\mu}(X_i, 1) + \frac{T_i}{\hat{e}(X_i)} (Y_i - \hat{\mu}(X_i, 1)) \right\} - \left\{ \hat{\mu}(X_i, 0) + \frac{1-T_i}{1-\hat{e}(X_i)} (Y_i - \hat{\mu}(X_i, 0)) \right\}$$

Issues with inverse PS weights: small errors in estimating the PS are magnified when taking its inverse. & obs with small proba of treatment dominate (expected, but again small errors are demultiplied).

Alternative approach = estimate the inverse PS directly.

Put $w(X) = \frac{1}{e(X)}$. We have =

$$\begin{aligned} \mu_1 &= \mathbb{E}[Y(1)] = \mathbb{E} \mathbb{E}[Y \mid X, T=1] = \mathbb{E} \mu(X, 1) \\ &= \mathbb{E}[T w(X) Y] \\ &= \mathbb{E}[T w(X) \mu(X, 1)] \end{aligned}$$

$$\mathbb{E} \mu(X, 1) = \mathbb{E}[T w(X) \mu(X, 1)]$$

In fact, we can show that $w(x) = \frac{1}{e(x)}$ is the only functional satisfying this equality for all measurable functions $m(\cdot, 1)$ [Riesz Representer]

We can use this property to characterize $w(x)$ as the unique solution to the following optimization problem (3)

$$\min_w \sup_{f \text{ measurable}} \{ \mathbb{E}[T w(X) f(X)] - \mathbb{E}[f(X)] \}$$

In practice, restrict to a class \mathcal{F} of functionals & derive the optimal weights balancing all functions in that class. Let:

$$\text{Imbalance}_{\mathcal{F}}(w) := \sup_{f \in \mathcal{F}} \{ \mathbb{E}[T w(X) f(X)] - \mathbb{E}[f(X)] \}$$

Introduce a regularization hyperparameter $\delta > 0$ to ensure a unique min & for bias-variance tradeoff.

$$w^*(\cdot) = \operatorname{argmin}_w \{ \text{Imbalance}_{\mathcal{F}}(w) + \delta \|w\|^2 \}$$

δ typically selected via cross-validation.

* Guarantees = Hirshberg & Wager (2021) show that if $m(\cdot, 1) \in \mathcal{F}$, then imbalance is small & that $\mathbb{E}[T w^*(X) Y]$ can be used as an approx. unbiased estimate of p_1 .
 [aka imbalance _{\mathcal{F}} used to bound the estimator's bias]

Ogburn et al restrict analysis to the space of linear functionals $\mathcal{F} = \{ f(x) = \theta^t x, \theta, x \in \mathbb{R}^d \}$. (4)

any norm on \mathbb{R}^d $\rightarrow \|\theta\| \leq 1$
 Authors are more general & consider rich sets of covariates [& potentially co-dim space \mathcal{F} using RKHS]

$$\begin{aligned} \text{Imbalance}_{\mathcal{F}}(w) &= \sup_{\|\theta\| \leq 1} \theta^t [\mathbb{E} T w(X) X - \mathbb{E} X] \\ &= \| \mathbb{E} T w(X) X - \mathbb{E} X \|_*$$

where $\|\cdot\|_*$ is the dual norm of $\|\cdot\|$.

$\underline{\mathbb{E}} X = l_2$ is the dual norm of l_2
 l_∞ is the dual norm of l_1 .

↳ Consider l_2 balancing.

The dual representation shows that we are minimising the l_2 norm of the difference $\mathbb{E} T w(X) X - \mathbb{E} X$.

The sample balance property is

$$\sum_{i=1}^n w(X_i) T_i X_{ij} \approx \sum_{i=1}^n X_{ij} \quad j=1, \dots, d$$

Want: balancing for each component

covariate distribution of the treatment group

covariate distribution of the target group (whole sample)

Let $\hat{w}_\delta^{l_2}(\cdot)$ denote the solution to the (sample) optimization problem. The DR estimator of ρ_1 is (5)

$$\frac{1}{n} \sum_{i=1}^n \left\{ \hat{p}_\lambda^{RR}(X_i, 1) + \hat{w}_\delta^{l_2}(X_i) T_i (Y_i - \hat{p}_\lambda^{RR}(X_i, 1)) \right\} \quad (*)$$

Outcome model fit using RR with penalty λ

$$\hat{p}_\lambda^{RR}(x, 1) = x^t \hat{\beta}_\lambda^{RR}$$

↳ regularization when d is large

↳ bias/variance tradeoff issue = we bring confounding back \Rightarrow RR does not perform well in high dim settings for estimating a causal effect in the presence of confounding

Augmented Term = Bias Correction

Required when the outcome model is regularized.

The main result of Ogburn et al is to show that in this setting [RR(λ) for outcome model + l_2 norm with hyp δ for w]

this expression collapses to a form that can be represented as a single output regression

$$(*) = \frac{1}{n} \sum_{i=1}^n X_i^t \hat{\beta}_{\lambda, \delta}^{AUG}$$

In fact, $\hat{\beta}_{\lambda, \delta}^{AUG}$ is the solution of a (generalized) ridge regression. (see later)

To further characterize the solution of the generalized RR, we consider an important special case: we compute $\hat{w}_\delta^{l_2}(\cdot)$ with exact balancing. (6)

↳ The most common optimization problem in that case is

$$\min_{w \in \mathbb{R}^n} \|w\|_2^2$$

more unknowns than eq when $d < n$

such that $\sum_{i=1}^n w_i T_i X_{ij} = \sum_{i=1}^n X_{ij} \quad j=1, \dots, d$

The Lagrangian of this convex opt. problem is

$$\mathcal{L}(w, v) = \sum_{i=1}^n w_i^2 + \sum_{j=1}^d v_j \left(\sum_{i=1}^n X_{ij} - \sum_{i=1}^n w_i T_i X_{ij} \right)$$

KKT conditions: optimal (w^*, v^*) parameters must satisfy:

- primary constraints $\sum_{i=1}^n X_{ij} = \sum_{i=1}^n w_i^* T_i X_{ij} \equiv \tilde{X}_{ij}^*$
- $\nabla_w \mathcal{L}(w, v) = 0 \Rightarrow 2w_i^* = \sum_{j=1}^d v_j T_i X_{ij} \quad i=1, \dots, n$

Matrix notation

$$X = \begin{pmatrix} -X_1^t & - \\ & -X_n^t & - \end{pmatrix}_{(n \times d)} \quad \tilde{X} = \begin{pmatrix} -\tilde{X}_1^t & - \\ & -\tilde{X}_n^t & - \end{pmatrix}_{(n \times d)} \quad v^* = \begin{pmatrix} v_1^* \\ \vdots \\ v_d^* \end{pmatrix} \quad w^* = \begin{pmatrix} w_1^* \\ \vdots \\ w_n^* \end{pmatrix}$$

$$\underline{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^n \quad \text{KKT} \quad \begin{cases} 2w^* = \tilde{X} v^* \\ X^t \underline{1} = \tilde{X}^t w^* \end{cases}$$

$$\Rightarrow v^* = 2(\tilde{X}^t \tilde{X})^{-1} X^t \underline{1}$$

$$w^* = \tilde{X} (\tilde{X}^t \tilde{X})^{-1} X^t \underline{1} \quad \leftarrow \text{linear in } X$$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n \hat{p}_0^{RR}(X_i, 1) = \frac{1}{n} (w^*)^t Y = \frac{1}{n} \underline{1}^t X (\tilde{X}^t \tilde{X})^{-1} \tilde{X}^t Y$$

x Remarks: (i) The weighting estimator of μ_j that exactly balances features is numerically equal to the OLS linear reg. estimator. Indeed, assuming $\mu(x, 1) = x^t \beta$,

$$\hat{\beta}^{OLS} \leftarrow OLS(Y_i \sim X_i | T_i = 1)$$

Fit one model for the trt group & one for the control group, separately.

$$\hat{\beta}^{OLS} = (\tilde{X}^t \tilde{X})^{-1} \tilde{X}^t Y$$

restricting the sample to $T_i = 1$ or truncating entries in X and using \tilde{X} yield the same numerical result.

⇒ Apply $\hat{\beta}^{OLS}$ to the whole (target) pop:

$$X \hat{\beta}^{OLS} = X (\tilde{X}^t \tilde{X})^{-1} \tilde{X}^t Y, \text{ and}$$

$$\frac{1}{n} \sum_{i=1}^n X_i \hat{\beta}^{OLS} = \frac{1}{n} 1^t X (\tilde{X}^t \tilde{X})^{-1} \tilde{X}^t Y \quad \square$$

(ii) This justifies the notation $\hat{\beta}^{OLS} = \hat{\beta}_{\lambda=0}^{RR}$.

(iii) This result is long known, see [Robins et al \(2007\)](#)

⇒ When both the reg function and the weights are linear, the OLS estimator is DR.

Back to our problem, [Ogburn et al](#) show that $\hat{\beta}_{\lambda, \delta}^{AUG}$ [under a ridge penalty (λ) for the outcome model and a ridge penalty (δ) for the weights] can be expressed as a shrank version of $\hat{\beta}^{OLS}$. Specifically,

$$\hat{\beta}_{\lambda, \delta}^{RR} = \left(\frac{\sigma_j^2}{\sigma_j^2 + \lambda} \right) \hat{\beta}_j^{OLS} + \hat{w}_s^{L2} \Rightarrow \hat{\beta}_j^{AUG} = \left(\frac{\sigma_j^2}{\sigma_j^2 + \delta_j} \right) \hat{\beta}_j^{OLS}$$

with

$$\delta_j = \frac{\delta \lambda}{\sigma_j^2 + \lambda + \delta}$$

Under the assumption that $\tilde{X}^t \tilde{X}$ is diagonal $\text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ Close form expression

$$\text{Ridge}(\lambda) + \text{Ridge}(\delta) = \text{Ridge}$$

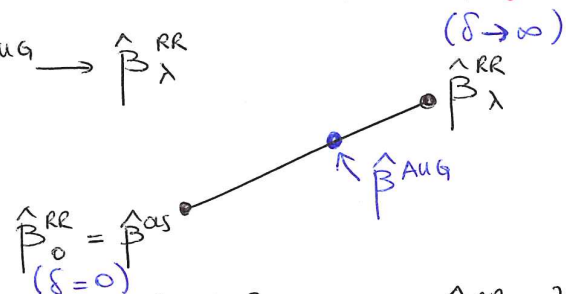
↳ $\lambda = 0$ [OLS] ⇒ $\hat{\beta}^{AUG} = \hat{\beta}^{OLS}$ [Robins et al \(2007\)](#)

↳ $\delta = 0$ ⇒ $\hat{\beta}^{AUG} = \hat{\beta}^{OLS}$

⚠ In practice, the Ridge + Ridge procedure with cross validation spits at the OLS estimator, since $\delta = 0$ is selected.

↳ Most researchers are unaware that they are simply doing OLS

↳ $\delta \rightarrow \infty$ ⇒ $\hat{\beta}^{AUG} \rightarrow \hat{\beta}_{\lambda}^{RR}$ ($\delta \rightarrow \infty$)



Also, can show that $\|Y - \tilde{X} \hat{\beta}^{AUG}\|_2^2 \leq \|Y - \tilde{X} \hat{\beta}_{\lambda}^{RR}\|_2^2$
 ⇒ $\hat{\beta}^{AUG}$ has smaller in-sample training error than $\hat{\beta}_{\lambda}^{RR}$.

⇒ The bias correction term in the DR procedure ⑨
brings the overall solution back closer to the OLS
"undersmoothing" / "underfitting".

↳ the present framework automatically selects the
amount of undersmoothing required for the
estimator to achieve \sqrt{n} convergence (not true for the
Ridge estimator due to regularisation & optimization
of the MSE).

In practice, the authors suggest to cross validate $\hat{\beta}_\lambda^{RR}$
to select λ , and then take $\delta = \lambda$ and $\gamma = \lambda^2$
(data driven & seems to work well in practice).