

DIMACS/CCICADA Workshop
Rutgers University

February 4, 2011

Understanding and Improving Propensity Score Methods

Zhiqiang Tan

Department of Statistics
Rutgers University

<http://www.stat.rutgers.edu/~ztan>

Outline

Understanding

- Introduction (to a medical study)
- Causal inference & missing-data problems
- Outcome regression vs propensity score weighting

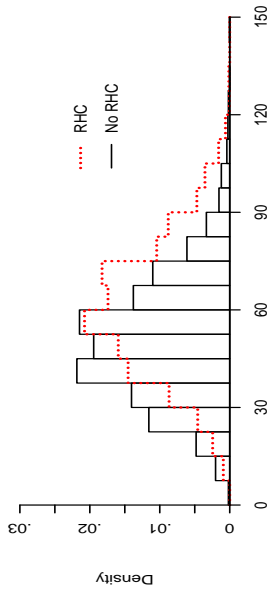
Improvements

- Doubly robust estimation
- Likelihood approach
- Doubly robust likelihood estimator

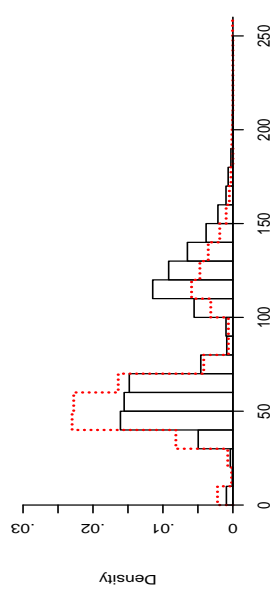
Effect of right heart catheterization

- Right heart catheterization (RHC) is performed daily in hospitals since 1970s.
- The benefit of RHC had **not** been demonstrated in a successful **randomized clinical trial**.
- Connors et al.'s (1996) **observational study** raised the concern that RHC might **not** benefit critically ill patients and might in fact cause harm.
- Data were collected on 5735 critically ill patients admitted to the ICUs of five medical centers:
 - Treatment: No-RHC or RHC
 - Outcome: 30-day survival
 - Covariates: 75 covariates (specified by doctors)
- What is **the effect** of RHC on survival?

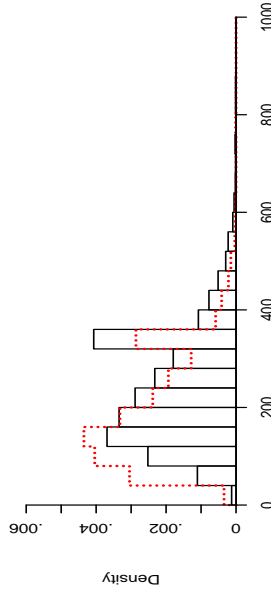
Raw histogram of aps



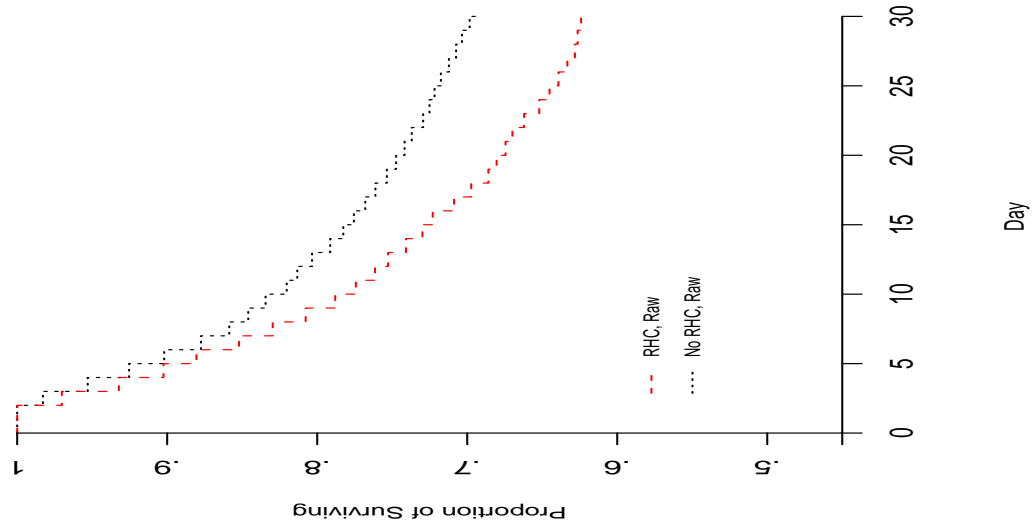
Raw histogram of meanbp



Raw histogram of pafi



Thirty-day survival curves



Causal inference (Rubin 1974, 1978, etc)

- X : pre-treatment **covariates** (measured)
- (Y_0, Y_1) : **potential outcome** that would be observed had a patient received No-RHC or RHC
- T : **treatment** variable taking value 0 or 1 if a patient actually received No-RHC or RHC
- $Y = (1 - T) Y_0 + T Y_1$: **observed outcome**
- Observed data & missing data

id	X_1	...	X_{75}	T	Y	Y_0	Y_1
1	*	*	*	0	10	10	??
2	*	*	*	1	7	??	7
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
n	*	*	*	0	12	12	??

- Average causal effects:

$E(Y_1)$ versus $E(Y_0)$

or other comparisons ...

Missing-data problems

- X : a vector of explanatory variables
- Y : an outcome variable
- R : non-missing indicator
= 1 if Y is observed or 0 if Y is missing
- Data structure

id	X_1	...	X_{75}	R	Y
1	*	*	*	0	??
2	*	*	*	1	7
⋮	⋮	⋮	⋮	⋮	⋮
n	*	*	*	0	??

- Assumption of **ignorability**:

$$R \perp Y | X$$

- Objective: to estimate $\mu = E(Y)$

Survey sampling with a superpopulation

- X : a vector of auxiliary variables
(measured on the entire population)
- Y : a study variable
- R : sampling indicator
= 1 if Y is sampled or 0 otherwise
- Data structure

id	X_1	...	X_{75}	R	Y
1	*	*	*	0	??
2	*	*	*	1	7
⋮	⋮	⋮	⋮	⋮	⋮
N	*	*	*	0	??

- By design:

$$R \perp Y | X$$

- Objective: to estimate $\mu = E(Y)$ [or $\sum_{i=1}^N Y_i$]

No-confounding estimation

- Assumption of **no (unmeasured) confounding**:

$$T \perp (Y_0, Y_1) \mid X,$$

as is the case for randomized clinical trials.

- Identification:

$$E(Y|T = 1, X) = E(Y_1|T = 1, X) \stackrel{Y_1 \perp T|X}{=} E(Y_1|X)$$
$$E[E(Y|T = 1, X)] = E(Y_1)$$

- **Outcome regression (OR)** model [\hookrightarrow Imputation]

$$E(Y|T = t, X) = m(t, X; \alpha_t), \quad t = 0, 1.$$

Let $\hat{m}(t, X) = m(t, X; \hat{\alpha}_t)$. An estimator of $E(Y_1)$ is

$$\hat{\mu}_{\text{OR}} = \frac{1}{n} \sum_{i=1}^n \hat{m}(1, X_i).$$

- **Propensity score (PS)** model (Rosenbaum and Rubin 1983)

$$P(T = 1|X) = \pi(X; \gamma).$$

Let $\hat{\pi}(X) = \pi(X; \hat{\gamma})$. An estimator of $E(Y_1)$ is

$$\hat{\mu}_{\text{IPW}} = \frac{1}{n} \sum_{i=1}^n \frac{T_i}{\hat{\pi}(X_i)} Y_i.$$

Illustration Left: $\#\{Y = 1\}$, Right: $\#\{Y = 0\}$

	$R = 0$ (missing)	$R = 1$
$X = 0$	(?, ?) 80	(11, 9) 20
$X = 1$	(?, ?) 40	(37, 23) 60
	? , ? 120	48, 32 80

- Outcome regression: $E(Y|R, X)$

	$R = 0$	$R = 1$
$X = 0$?	55.0%
$X = 1$?	61.7%

$$E(Y) : [(100) 55.0\% + (100) 61.7\%]/200 = 58.3\%$$

- Propensity score weighting: $P(R|X)$

	$R = 0$	$R = 1$
$X = 0$	80/100	20/100
$X = 1$	40/100	60/100

$$E(Y) : [11/(20/100) + 37/(60/100)]/200 = 58.3\%$$

Illustration Left: $\#\{Y = 1\}$, Right: $\#\{Y = 0\}$

	$T = 0$	$T = 1$
$X = 0$	(52, 28) 80	(11, 9) 20
$X = 1$	(30, 10) 40	(37, 23) 60
	82, 38 120	48, 32 80

- Outcome regression: $E(Y|T, X)$

	$T = 0$	$T = 1$
$X = 0$	65.0%	55.0%
$X = 1$	75.0%	61.7%

$$E(Y_0) : [(100) 65.0\% + (100) 75.0\%]/200 = 70.0\%$$

$$E(Y_1) : [(100) 55.0\% + (100) 61.7\%]/200 = 58.3\%$$

- Propensity score weighting: $P(T|X)$

	$T = 0$	$T = 1$
$X = 0$	80/100	20/100
$X = 1$	40/100	60/100

$$E(Y_0) : [52/(80/100) + 30/(40/100)]/200 = 70.0\%$$

$$E(Y_1) : [11/(20/100) + 37/(60/100)]/200 = 58.3\%$$

Propensity score paradox (Robins and Ritov 1997)

- The likelihood is

$$\prod_{i=1}^n P(\{X_i\}) \times \prod_{i=1}^n \left[(1 - \pi(X_i; \gamma))^{1-T_i} \pi(X_i; \gamma)^{T_i} \right] \\ \times \prod_{i=1}^n \left[P(Y_i|T_i = 0, X_i; \alpha_0)^{1-T_i} P(Y_i|T_i = 1, X_i; \alpha_1)^{T_i} \right].$$

- But, $E(Y_1) = \int E(Y|T = 1, x) dP(x)$.
- By the (strong) likelihood principle, inference should not depend on γ .
- Yet, if $\pi(X)$ is known, then

$$\frac{1}{n} \sum_{i=1}^n \frac{T_i}{\pi(X_i)} Y_i$$

is unbiased & regular, regardless of $\dim(X)$, whereas $\hat{\mu}_{OR}$ suffers the curse of dimensionality.

- If X is categorical, then

$$\frac{1}{n} \sum_{i=1}^n \frac{T_i}{\hat{\pi}(X_i)} Y_i \stackrel{def}{=} \hat{\mu}_{IPW} = \hat{\mu}_{OR}$$

agree with each other.

Propensity score peculiarity

- Outcome regression (OR) model

$$E(Y|T = 1, X) = m(1, X; \alpha_1).$$

Compute $\hat{\alpha}_1$ as a solution to

$$0 = \frac{1}{n} \sum_{i=1}^n \frac{\partial m(1, X_i; \alpha_1) / \partial \alpha_1}{\text{var}(Y|T = 1, X_i)} T_i (Y_i - m(1, X_i; \alpha_1)).$$

Had I used the true value $m(1, X)$:

$$\text{asy.var} \left[\frac{1}{n} \sum_{i=1}^n m(1, X_i) \right] \leq \text{asy.var} \left[\frac{1}{n} \sum_{i=1}^n \hat{m}(1, X_i) \right].$$

- Propensity score (PS) model

$$P(T = 1|X) = \pi(X; \gamma).$$

Compute $\hat{\gamma}$ as a solution to

$$0 = \frac{1}{n} \sum_{i=1}^n \frac{\partial \pi(X_i; \gamma) / \partial \gamma}{\pi(X_i; \gamma)(1 - \pi(X_i; \gamma))} (T_i - \pi(X_i; \gamma)).$$

Had I used the true value $\pi(X)$:

$$\text{asy.var} \left[\frac{1}{n} \sum_{i=1}^n \frac{T_i}{\pi(X_i)} Y_i \right] \geq \text{asy.var} \left[\frac{1}{n} \sum_{i=1}^n \frac{T_i}{\hat{\pi}(X_i)} Y_i \right].$$

Propensity score peculiarity (cont'd)

- Had I used **a bigger OR model**: $\alpha_1 \subset \alpha_1^\dagger$

$$\text{asy.var}(\hat{\mu}_{\text{OR}}) \leq \text{asy.var}(\hat{\mu}_{\text{OR}}^\dagger).$$

Had I used **bigger and bigger OR models**:

$$\alpha_1 \subset \alpha_1^\dagger \subset \alpha_1^{\ddagger} \subset \dots$$

$$\begin{aligned} \text{asy.var}(\hat{\mu}_{\text{OR}}) &\leq \text{asy.var}(\hat{\mu}_{\text{OR}}^\dagger) \leq \text{asy.var}(\hat{\mu}_{\text{OR}}^{\ddagger}) \leq \dots \\ &\leq \text{asy.var. bound} \end{aligned}$$

- Had I used **a bigger PS model**: $\gamma \subset \gamma^\dagger$

$$\text{asy.var}(\hat{\mu}_{\text{OR}}) \geq \text{asy.var}(\hat{\mu}_{\text{OR}}^\dagger).$$

Had I used **bigger and bigger PS models**:

$$\gamma \subset \gamma^\dagger \subset \gamma^{\ddagger} \subset \dots$$

$$\begin{aligned} \text{asy.var}(\hat{\mu}_{\text{IPW}}) &\geq \text{asy.var}(\hat{\mu}_{\text{IPW}}^\dagger) \geq \text{asy.var}(\hat{\mu}_{\text{IPW}}^{\ddagger}) \geq \dots \\ &\geq \text{asy.var. bound} \end{aligned}$$

- The **asymptotic variance bound** for regular estimators of $\mu_1 = E(Y_1)$ is $n^{-1}E(\tau_1^2)$, where

$$\begin{aligned} \tau_1 &= m(1, X) - \mu_1 + \frac{T}{\pi(X)}(Y - m(1, X)) \\ &= \frac{T}{\pi(X)}Y - \mu_1 - \left(\frac{T}{\pi(X)} - 1\right)m(1, X). \end{aligned}$$

What does this mean? (Tan 2007)

- Caveats

- Assumption: each model is **correctly** specified.

- **Asymptotic** variance: fix each model, let $n \rightarrow \infty$

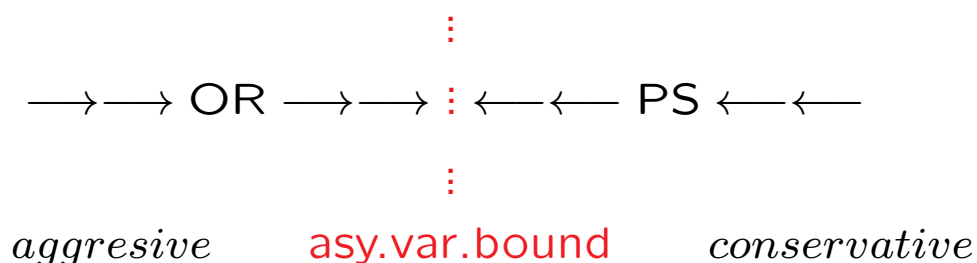
- This “shows” $\hat{\mu}_{IPW}$ is more variable than $\hat{\mu}_{OR}$:

$$\text{asy.var}(\hat{\mu}_{OR}) \leq \text{asy.var.bound} \leq \text{asy.var}(\hat{\mu}_{IPW}).$$

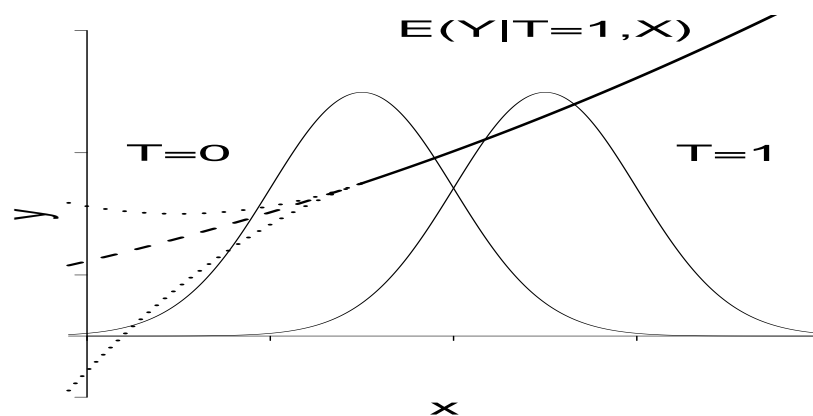
- Does this mean $\hat{\mu}_{IPW}$ is not as good as $\hat{\mu}_{OR}$?

Were a statistician asked to choose a **correct** OR model and a **correct** PS model ...

- Which side of the asy.var.bound?



- The asy.var.bound becomes large or even ∞ whenever $\pi(X) \approx 0$ in some region of X .
 - It is **difficult** to infer $E(Y_1)$, because few subjects with $\pi(X) \approx 0$ receive $T = 1$ (i.e., Y_1 observed).
 - The difficulty holds for OR and PS, although the **symptoms** can be different.



- OR model is built on the “truncated” data

$$\{(X_i, Y_i) : T_i = 1, 1 \leq i \leq n\},$$

and hence makes **extrapolation** to predict $m(1, X)$ in the region of X such that $\pi(X) \approx 0$.

- PS model is built on the “full” data

$$\{(X_i, T_i) : 1 \leq i \leq n\}.$$

Inverse PS weighting of Y for treated subjects with $\pi(X) \approx 0$ leads to **large variance**.

Summary of what's been discussed

OR and PS are two approaches with different characteristics, and one does not dominate the other.

- If an OR model is correctly specified, then $\hat{\mu}_{\text{OR}}$ is consistent and has $\text{asy.var} \leq \text{asy.var.bound}$.
- If a PS model is correctly specified, then $\hat{\mu}_{\text{PS}}$ is consistent and has $\text{asy.var} \geq \text{asy.var.bound}$.
- OR approach suffers the problem of **implicitly making extrapolation** in the region of X with $\pi(X) \approx 0$, due to data truncation.
- PS approach tends to yield large weights for treated subjects with $\pi(X) \approx 0$, **explicitly indicating uncertainty** in the estimate.

Doubly robust estimation

- Given a PS model, a class of estimators of $E(Y_1)$ is

$$\frac{1}{n} \sum_{i=1}^n \frac{T_i}{\widehat{\pi}(X_i)} Y_i - \frac{1}{n} \sum_{i=1}^n \left(\frac{T_i}{\widehat{\pi}(X_i)} - 1 \right) h(X_i),$$

where $h(X)$ is an arbitrary function.

- Consider the estimator (Robins et al. 1994)

$$\widehat{\mu}_{\text{AIPW}} = \frac{1}{n} \sum_{i=1}^n \frac{T_i}{\widehat{\pi}(X_i)} Y_i - \frac{1}{n} \sum_{i=1}^n \left(\frac{T_i}{\widehat{\pi}(X_i)} - 1 \right) \widehat{m}(1, X_i),$$

depending on both OR and PS models.

- $\widehat{\mu}_{\text{AIPW}}$ is **locally efficient**: it attains asy.var.bound if both OR **and** PS models are correctly specified.
- $\widehat{\mu}_{\text{AIPW}}$ is **doubly robust**: it remains consistent if either OR **or** PS model is correctly specified.
- Comparison of $\widehat{\mu}_{\text{AIPW}}$ vs $\widehat{\mu}_{\text{OR}}$:
$$\text{asy.var}(\widehat{\mu}_{\text{AIPW}}) \geq \text{asy.var}(\widehat{\mu}_{\text{OR}}).$$
- Comparison of $\widehat{\mu}_{\text{AIPW}}$ vs $\widehat{\mu}_{\text{IPW}}$:
$$\text{asy.var}(\widehat{\mu}_{\text{AIPW}}) \stackrel{??}{\leq} \text{asy.var}(\widehat{\mu}_{\text{IPW}}).$$

A more desirable DR estimator

- Consider the estimator (Tan 2006)

$$\tilde{\mu}_{\text{REG}} = \frac{1}{n} \sum_{i=1}^n \frac{T_i}{\hat{\pi}(X_i)} Y_i - \tilde{\beta}^{(1)} \frac{1}{n} \sum_{i=1}^n \left(\frac{T_i}{\hat{\pi}(X_i)} - 1 \right) \hat{m}(1, X_i),$$

where $\tilde{\beta}^{(1)}$ is the first element of

$$\tilde{\beta} = \left[\frac{1}{n} \sum_{i=1}^n \hat{\xi}_i \hat{\zeta}_i^\top \right]^{-1} \left[\frac{1}{n} \sum_{i=1}^n \hat{\xi}_i \hat{\eta}_i \right]$$

and

$$\hat{\eta}_i = \frac{T_i}{\hat{\pi}(X_i)} Y_i,$$

$$\hat{\xi}_i = \left(\frac{T_i}{\hat{\pi}(X_i)} - 1 \right) \left(\hat{m}(1, X_i), \frac{\frac{\partial \hat{\pi}}{\partial \gamma^\top}(X_i)}{1 - \hat{\pi}(X_i)} \right)^\top,$$

$$\hat{\zeta}_i = \frac{T_i}{\hat{\pi}(X_i)} \left(\hat{m}(1, X_i), \frac{\frac{\partial \hat{\pi}}{\partial \gamma^\top}(X_i)}{1 - \hat{\pi}(X_i)} \right)^\top.$$

- $\tilde{\mu}_{\text{REG}}$ is **locally efficient** and **doubly robust**.
- $\tilde{\mu}_{\text{REG}}$ is **intrinsically efficient**: if PS model is correct, it attains the smallest asy.var among

$$\frac{1}{n} \sum_{i=1}^n \frac{T_i}{\hat{\pi}(X_i)} Y_i - b^{(1)} \frac{1}{n} \sum_{i=1}^n \left(\frac{T_i}{\hat{\pi}(X_i)} - 1 \right) \hat{m}(1, X_i),$$

where $b^{(1)}$ is an arbitrary coefficient.

- Comparison of $\tilde{\mu}_{\text{REG}}$ vs $\hat{\mu}_{\text{OR}}$.
- Comparison of $\tilde{\mu}_{\text{REG}}$ vs $\hat{\mu}_{\text{IPW}}$.

Likelihood approach (Tan 2006)

- Rewrite the likelihood as

$$L_1 \times L_2 = \prod_{i=1}^n \left[(1 - \pi(X_i; \gamma))^{1-T_i} \pi(X_i; \gamma)^{T_i} \right] \\ \times \prod_{i=1}^n \left[G_0(\{X_i, Y_i\})^{1-T_i} G_1(\{X_i, Y_i\})^{T_i} \right]$$

where G_0 is the joint dist'n of (X, Y_0) and G_1 is the joint dist'n of (X, Y_1) .

- G_0 and G_1 induce the same marginal of X :

$$\int h(x) dG_0(x, y_0) = \int h(x) dG_1(x, y_1)$$

for each bounded function $h(x)$.

- We retain a finite subset of constraints and ignore other constraints on (G_0, G_1) .

- Let $\hat{h}(x) = \left\{ 1 - \hat{\pi}(x), (1 - \hat{\pi}(x))\hat{m}(1, x), \frac{\partial \hat{\pi}}{\partial \gamma^\top}(x) \right\}^\top$.
Maximize $L_2(G_0, G_1)$ subject to

$$\int dG_1 = 1, \\ \int (1 - \hat{\pi}(x)) dG_1 = \int (1 - \hat{\pi}(x)) dG_0, \\ \int (1 - \hat{\pi}(x))\hat{m}(1, x) dG_1 = \int (1 - \hat{\pi}(x))\hat{m}(1, x) dG_0, \\ \int \frac{\partial \hat{\pi}}{\partial \gamma}(x) dG_1 = \int \frac{\partial \hat{\pi}}{\partial \gamma}(x) dG_0.$$

Likelihood approach (cont'd)

- Postulate the following model for (T, X, Y) :

- T is generated according to

$$P(T = 0) = \int (1 - \hat{\pi}(x)) dG_0,$$

$$P(T = 1) = \int \hat{\pi}(x) dG_1.$$

- (X, Y) given $T = 0$ or 1 is generated from the biased-sampling distribution (Vardi 1985)

$$\frac{1 - \hat{\pi}(x)}{\int (1 - \hat{\pi}(x')) dG_0} dG_0,$$

$$\frac{\hat{\pi}(x)}{\int \hat{\pi}(x') dG_1} dG_1.$$

- Let $\varpi(X; \lambda) = \hat{\pi}(X) + \lambda^\top \hat{h}(X)$ and

$$\hat{\lambda} = \operatorname{argmax}_\lambda$$

$$\frac{1}{n} \sum_{i=1}^n \left[T_i \log(\varpi(X_i; \lambda)) + (1 - T_i) \log(1 - \varpi(X_i; \lambda)) \right].$$

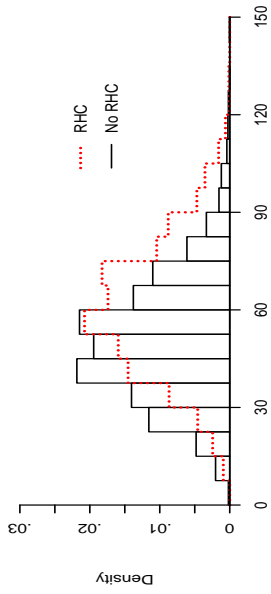
Then

$$\hat{G}_1(\{X_i, Y_i\}) = \frac{n^{-1}}{\varpi(X_i; \hat{\lambda})}, \quad \text{if } T_i = 1,$$

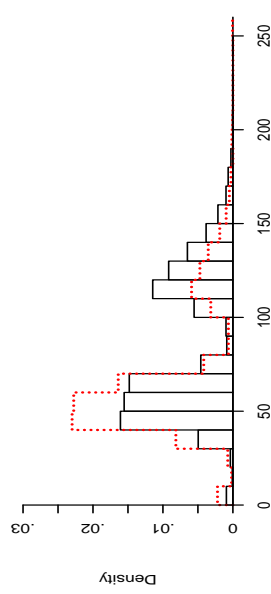
$$\hat{G}_0(\{X_i, Y_i\}) = \frac{n^{-1}}{1 - \varpi(X_i; \hat{\lambda})}, \quad \text{if } T_i = 0.$$

The dist'n (\hat{G}_0, \hat{G}_1) are **weighted histograms**.

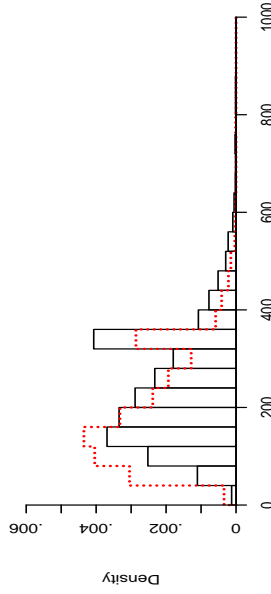
Raw histogram of aps



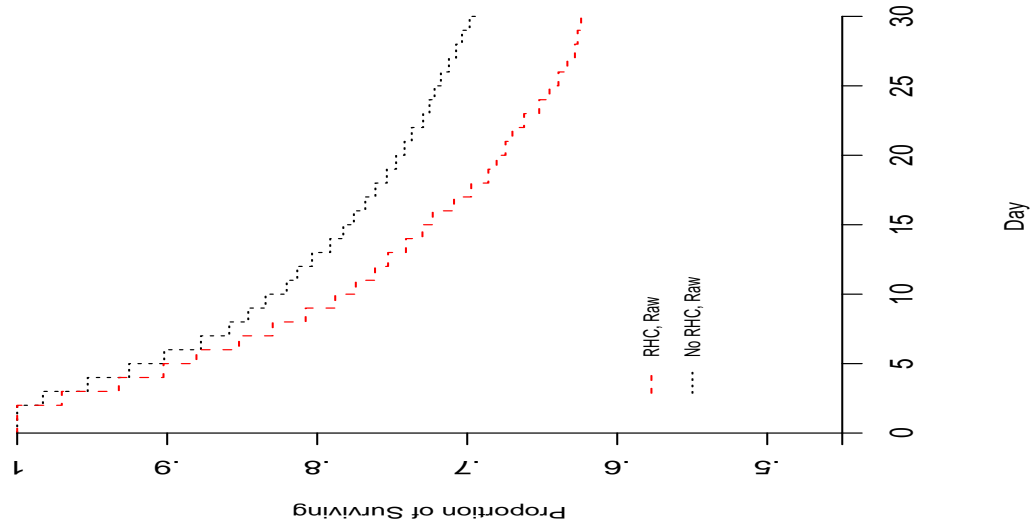
Raw histogram of meanbp



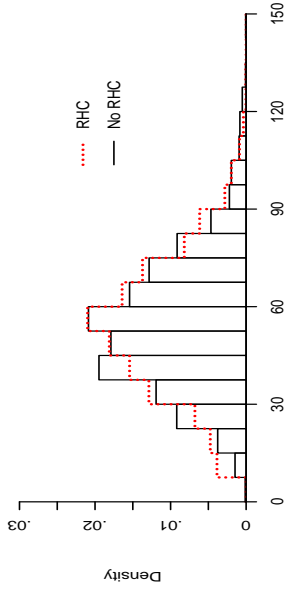
Raw histogram of pafi



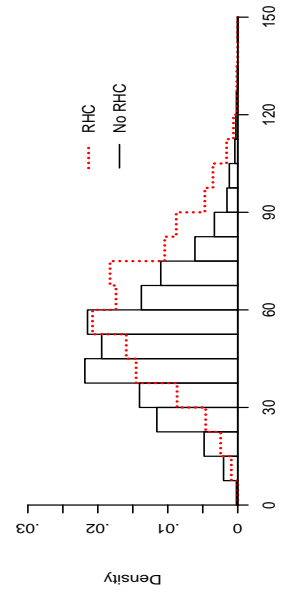
Thirty-day survival curves



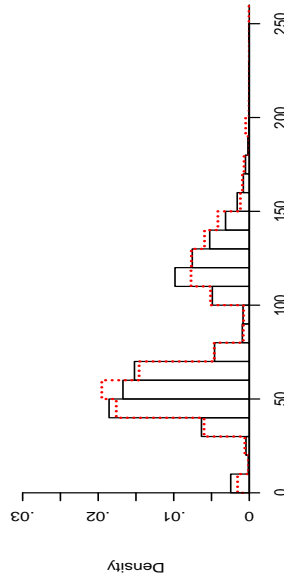
Weighted histogram of aps



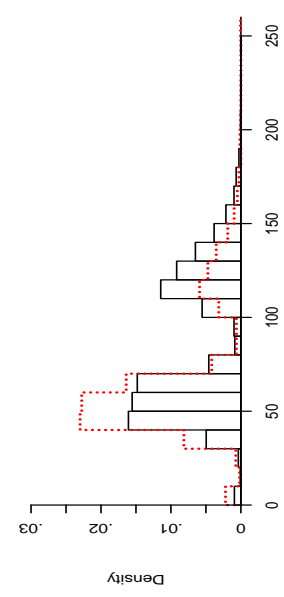
Raw histogram of aps



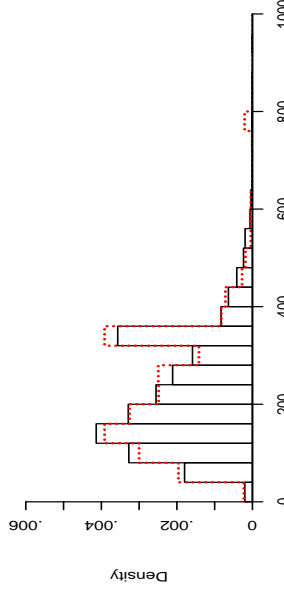
Weighted histogram of meanbp



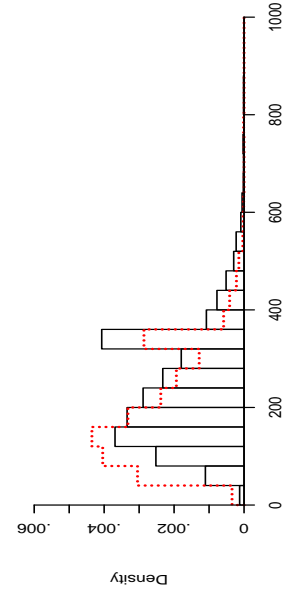
Raw histogram of meanbp



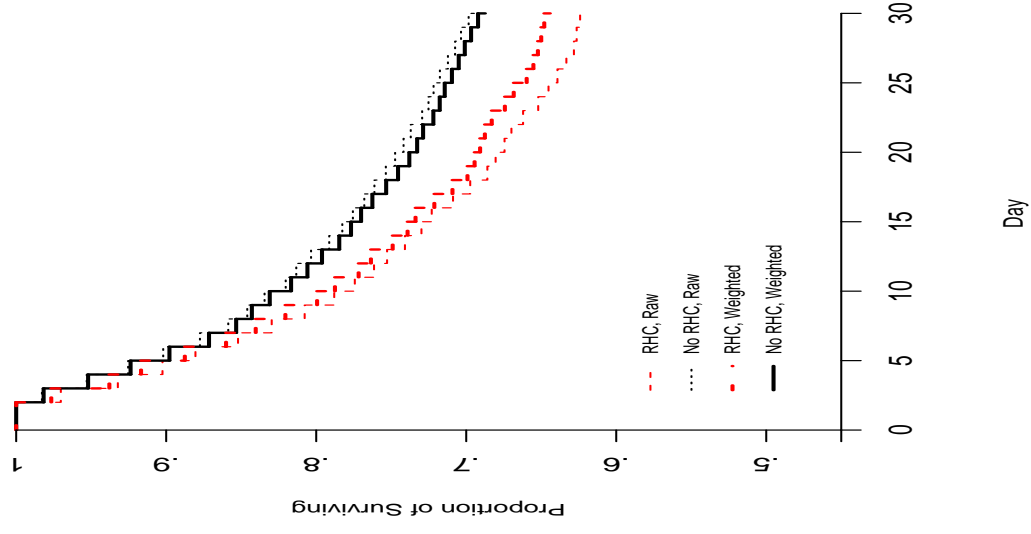
Weighted histogram of pafi



Raw histogram of pafi



Thirty-day survival curves



Likelihood approach (cont'd)

- The resulting estimator of μ

$$\hat{\mu}_{\text{LIK}} = \int y d\hat{G}_1 = \frac{1}{n} \sum_{i=1}^n \left[\frac{T_i Y_i}{\varpi(X_i; \hat{\lambda})} \right].$$

- $\hat{\mu}_{\text{LIK}}$ is **locally efficient**: it attains asy.var.bound if both OR **and** PS models are correctly specified.
- $\hat{\mu}_{\text{LIK}}$ is **intrinsically efficient**: if PS model is correct, it is asymptotically at least efficient as $\hat{\mu}_{\text{IPW}}$ and $\hat{\mu}_{\text{AIPW}}$.

- $\hat{\lambda}$ is a solution to

$$0 = \frac{1}{n} \sum_{i=1}^n \left[\frac{T_i - \varpi(X_i; \lambda)}{\varpi(X_i; \lambda)(1 - \varpi(X_i; \lambda))} \hat{h}(X_i) \right].$$

- $\hat{\lambda}$ also satisfies

$$1 = \int d\hat{G}_1 = \frac{1}{n} \sum_{i=1}^n \left[\frac{T_i}{\varpi(X_i; \hat{\lambda})} \right].$$

- $\varpi(X_i; \hat{\lambda})$ with $T_i = 1$, $1 \leq i \leq n$, are bounded from below by n^{-1} .
- $\hat{\mu}_{\text{LIK}}$ satisfies **sample boundedness**: it lies within the range of $\{Y_i : T_i = 1, 1 \leq i \leq n\}$, even if OR and/or PS model are misspecified.

DR Likelihood estimator (Tan 2010)

- Let

$$\hat{v}(x) = (1, \hat{m}(1, x))^{\top},$$

$$\hat{h}(x) = \left\{ (1 - \hat{\pi}(x)) \hat{v}^{\top}(x), \frac{\partial \hat{\pi}}{\partial \gamma^{\top}}(x) \right\}^{\top}.$$

Recall $\varpi(X; \lambda) = \hat{\pi}(X) + \lambda^{\top} \hat{h}(X).$

- Let $\tilde{\lambda}$ be a solution to

$$0 = \frac{1}{n} \sum_{i=1}^n \left[\left(\frac{T_i}{\varpi(X_i; \lambda)} - 1 \right) \hat{v}(X_i) \right],$$

$$0 = \frac{1}{n} \sum_{i=1}^n \left[\frac{T_i - \varpi(X_i; \lambda)}{\varpi(X_i; \lambda)(1 - \varpi(X_i; \lambda))} \frac{\partial \hat{\pi}}{\partial \gamma}(X_i) \right],$$

subject to $\varpi(X_i; \lambda) > 0$ if $T_i = 1$.

- The resulting estimator of μ

$$\tilde{\mu}_{\text{LIK}} = \frac{1}{n} \sum_{i=1}^n \left[\frac{T_i Y_i}{\varpi(X_i; \tilde{\lambda})} \right].$$

- $\tilde{\lambda}$ satisfies

$$1 = \frac{1}{n} \sum_{i=1}^n \left[\frac{T_i}{\varpi(X_i; \tilde{\lambda})} \right],$$

$$\frac{1}{n} \sum_{i=1}^n \hat{m}(1, X_i) = \frac{1}{n} \sum_{i=1}^n \left[\frac{T_i}{\varpi(X_i; \tilde{\lambda})} \hat{m}(1, X_i) \right].$$

- $\tilde{\mu}_{\text{LIK}}$ satisfies **sample boundedness**.
 - $\tilde{\mu}_{\text{LIK}}$ is **doubly robust**.
 - $\tilde{\mu}_{\text{LIK}}$ is **locally efficient** and **intrinsically efficient**.
- Issues: existence and computation of $\tilde{\lambda}$.

Asymptotic expansions under misspecified models

- $\hat{\lambda}$ converges to a constant λ^* . Then

$$\hat{\mu}_{\text{LIK}} = \tilde{E} \left[\frac{TY}{\varpi(X; \lambda^*)} \right] - \hat{C}^\top \hat{B}^{-1} \tilde{E} \left[\frac{T - \varpi(X; \lambda^*)}{\varpi(X; \lambda^*)(1 - \varpi(X; \lambda^*))} \hat{h}(X) \right] + o_p(n^{-1/2}),$$

where

$$\hat{B} = \tilde{E} \left[\frac{(T - \varpi(X; \lambda^*))^2}{\varpi^2(X; \lambda^*)(1 - \varpi(X; \lambda^*))^2} \hat{h}(X) \hat{h}^\top(X) \right],$$

$$\hat{C} = \tilde{E} \left[\frac{T}{\varpi^2(X; \lambda^*)} Y \hat{h}(X) \right].$$

If $\lambda^* = 0$, then the **first-order term** reduces to $\hat{\mu}_{\text{REG}}$.

- $\tilde{\lambda}$ converges to a constant λ^\dagger . Then

$$\tilde{\mu}_{\text{LIK}} = \tilde{E} \left[\frac{TY}{\varpi(X; \lambda^\dagger)} \right] - \hat{C}^\top \tilde{B}^{-1} \tilde{E} \left[\left(\begin{array}{c} \left(\frac{T}{\varpi(X; \lambda^\dagger)} - 1 \right) \hat{v}(X) \\ \frac{T - \varpi(X; \lambda^\dagger)}{\varpi(X; \lambda^\dagger)(1 - \varpi(X; \lambda^\dagger))} \hat{h}_2(X) \end{array} \right) \right] + o_p(n^{-1/2}),$$

where

$\tilde{B} =$

$$\tilde{E} \left[\left(\begin{array}{cc} \frac{T}{\varpi^2(X; \lambda^\dagger)} \hat{h}_1(X) \hat{v}^\top(X) & \frac{(T - \varpi(X; \lambda^\dagger))^2}{\varpi^2(X; \lambda^\dagger)(1 - \varpi(X; \lambda^\dagger))^2} \hat{h}_1(X) \hat{h}_2^\top(X) \\ \frac{T}{\varpi^2(X; \lambda^\dagger)} \hat{h}_2(X) \hat{v}^\top(X) & \frac{(T - \varpi(X; \lambda^\dagger))^2}{\varpi^2(X; \lambda^\dagger)(1 - \varpi(X; \lambda^\dagger))^2} \hat{h}_2(X) \hat{h}_2^\top(X) \end{array} \right) \right]$$

with

$$\hat{h}_1(x) = (1 - \hat{\pi}(x)) \hat{v}^\top(x) = (1 - \hat{\pi}(x)) (1, \hat{m}(1, x))^\top,$$

$$\hat{h}_2(x) = \frac{\partial \hat{\pi}}{\partial \gamma}(x).$$

If $\lambda^\dagger = 0$, then the **first-order term** reduces (almost) to $\tilde{\mu}_{\text{REG}}$.

DR Likelihood estimator (cont'd)

- Two-step estimator

- Compute $\hat{\lambda} = (\hat{\lambda}_1^\top, \hat{\lambda}_2^\top)^\top$, partitioned according to $\hat{h}(x) = \{(1 - \hat{\pi}(x))\hat{v}^\top(x), \frac{\partial \hat{\pi}}{\partial \gamma^\top}(x)\}^\top$.

- Compute $\tilde{\lambda}_{\text{step2}} = (\tilde{\lambda}_{1,\text{step2}}^\top, \hat{\lambda}_2^\top)^\top$, where

$$\tilde{\lambda}_{1,\text{step2}} = \operatorname{argmax}_{\lambda_1} \frac{1}{n} \sum_{i=1}^n \left[T_i \frac{\log\{\varpi(X_i; \lambda_1, \hat{\lambda}_2) / \varpi(X_i; \hat{\lambda})\}}{1 - \hat{\pi}(X_i)} - \lambda_1^\top \hat{v}(X_i) \right].$$

The objective function is **strictly concave** and **bounded from above** under mild conditions.

- For step 2, $\tilde{\lambda}_{1,\text{step2}}$ is a solution to

$$0 = \frac{1}{n} \sum_{i=1}^n \left[\left(\frac{T_i}{\varpi(X_i; \lambda_1, \hat{\lambda}_2)} - 1 \right) \hat{v}(X_i) \right].$$

- The resulting estimator of μ

$$\tilde{\mu}_{\text{LIK}} = \frac{1}{n} \sum_{i=1}^n \left[\frac{T_i Y_i}{\varpi(X_i; \tilde{\lambda})} \right].$$

- $\tilde{\mu}_{\text{LIK2}}$ satisfies **sample boundedness**.

- $\tilde{\mu}_{\text{LIK2}}$ is **doubly robust**.

- $\tilde{\mu}_{\text{LIK2}}$ is **locally efficient** and **intrinsically efficient**.

Simulation study

- The same design as in Kang and Schafer (2007)

$$X = (X_1, X_2, X_3, X_4)^\top,$$

$$Y = 210 + 27.4X_1 + 13.7X_2 + 13.7X_3 + 13.7X_4 + \epsilon,$$

$$T = 1\{U \leq \text{expit}(-X_1 + .5X_2 - .25X_3 - .1X_4)\}.$$

- Missing-data setup: observed data (X, T, TY)
Causal-inference setup: observed data (X, T, Y)

- 22 estimators are labelled as follows.

(1–3) $\hat{\mu}_{\text{LIK,OLS}}$, $\hat{\mu}_{\text{REG,OLS}}$, $\tilde{\mu}_{\text{REG,OLS}}$;

(4–7) $\hat{\mu}_{\text{AIPW (ratio)}}$, $\hat{\mu}_{\text{OLS,ext}}$, $\hat{\mu}_{\text{WLS}}$, $\tilde{\mu}_{\text{RV}}$;

(8–12) $\hat{\mu}_{\text{IPW,ext (ratio)}}$, $\tilde{\mu}_{\text{IPW,ext2}}$,
 $\hat{\mu}_{\text{AIPW,ext (ratio)}}$, $\hat{\mu}_{\text{WLS,ext (ratio)}}$, $\hat{\mu}_{\text{WLS,ext2}}$;

(13–15) $\hat{\mu}_{\text{TIPW (ratio)}}$, $\hat{\mu}_{\text{TML}}$, $\hat{\mu}_{\text{TAIPW (ratio)}}$;

(16–22) $\hat{\mu}_{\text{AIPW,lik}}$, $\tilde{\mu}_{\text{LIK2,OLS}}$,
 $\hat{\mu}_{\text{WLS,lik}}$, $\hat{\mu}_{\text{WLS,lik2}}$, $\tilde{\mu}_{\text{LIK2,WLS}}$,
 $\tilde{\mu}_{\text{RV,lik}}$, $\tilde{\mu}_{\text{LIK2,RV}}$.

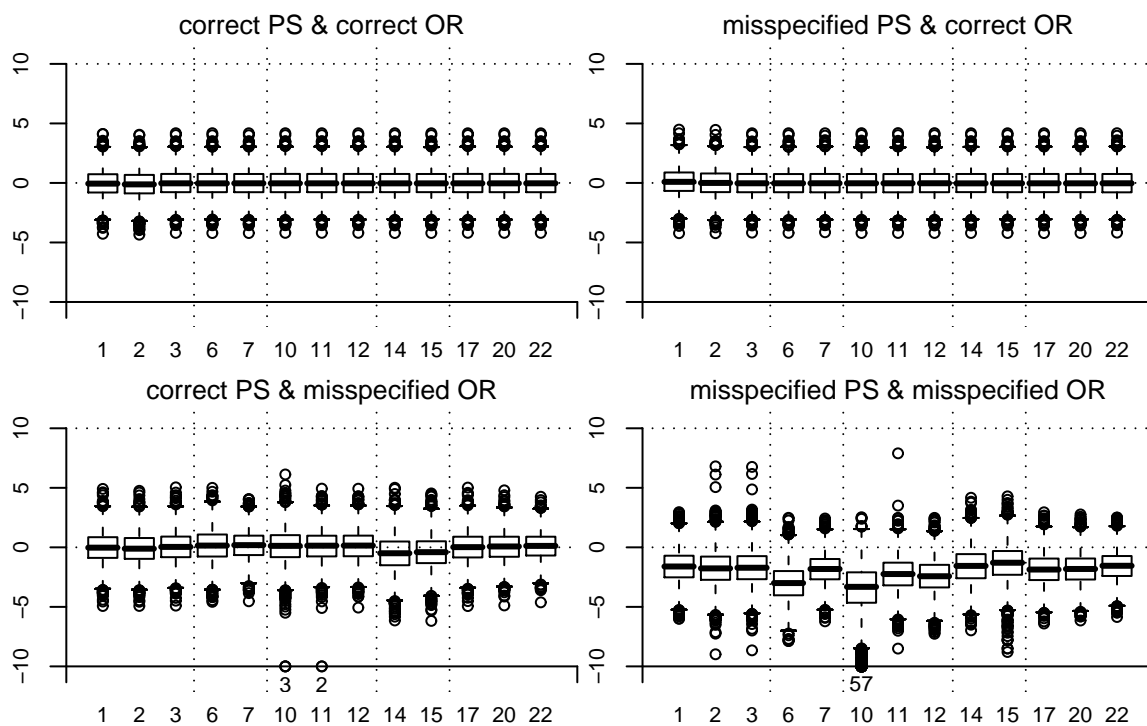


Figure 1: Boxplots of estimators of $\mu_1 - 210$ ($n = 1000$)

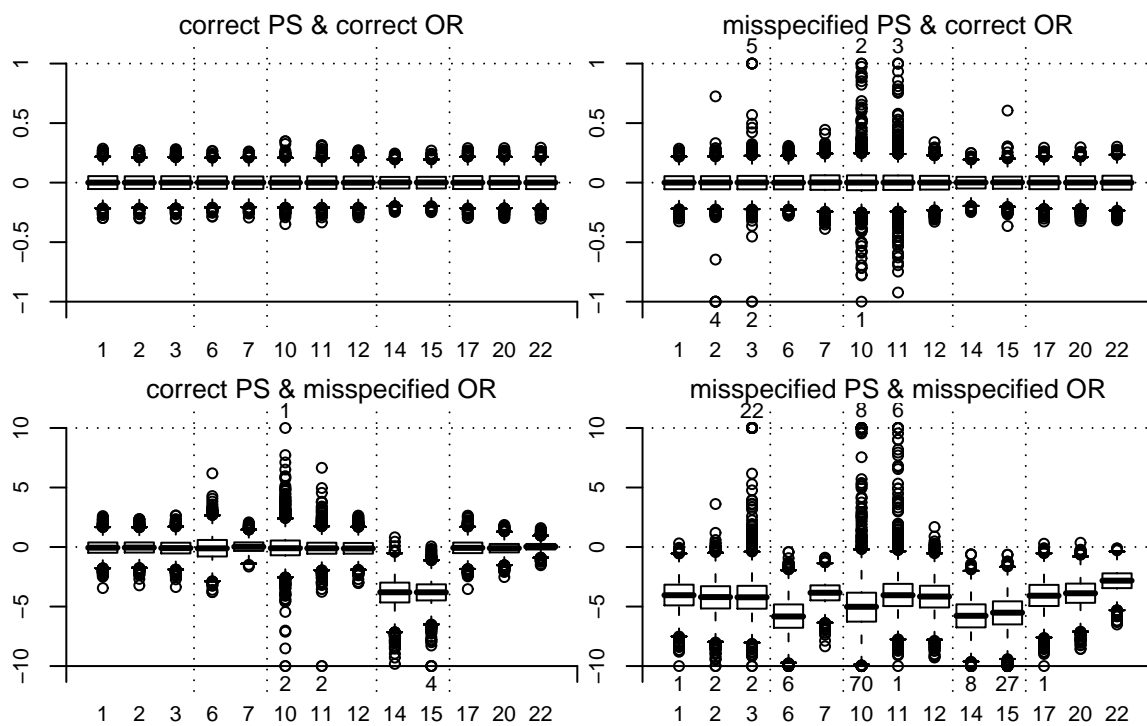


Figure 2: Boxplots of estimators of $\mu_1 - \mu_0$ ($n = 1000$)

Table 1: Numerical comparison of estimators

Estimator	1	2	3	6	7	10	11	12	14	15	17	20	22
Estimators of μ_1 in missing data setup													
C-PS&M-OR	1.23	1.28	1.20	1.32	1.07	1.39	1.25	1.24	1.51	1.35	1.24	1.17	1.00
	1.21	1.21	1.20	1.36	1.07	1.75	1.33	1.20	1.80	1.49	1.21	1.13	1.00
M-PS&M-OR	1.32	1.50	1.34	1.62	1.12	1.87	1.31	1.30	1.32	1.12	1.27	1.20	1.00
	1.10	1.31	1.27	2.82	1.24	4.86	1.78	1.99	1.21	1.04	1.31	1.27	1.00
Estimators of $\mu_1 - \mu_0$ in causal inference setup													
C-PS&M-OR	1.97	1.79	2.76	4.50	1.85	5.00	3.27	3.16	19.8	18.1	2.65	1.87	1.00
	3.80	3.58	4.05	9.08	2.33	11.7	5.42	4.17	134	130	4.07	2.58	1.00
M-PS&M-OR	1.77	1.77	2.51	2.75	1.46	2.39	1.87	1.82	2.91	2.42	1.88	1.74	1.00
	2.02	2.22	73.3	3.99	1.77	3.42	2.08	2.16	3.97	3.72	2.07	1.87	1.00

C-PS (or M-PS): correct (or misspecified) propensity score model; M-OR: misspecified outcome regression model; Each cell gives the ratios of mean squared errors for $n = 200$ (upper) and $n = 1000$ (lower).