

SUPPLEMENTARY APPENDIX FOR
“MOSTLY HARMLESS SIMULATIONS? USING MONTE CARLO
STUDIES FOR ESTIMATOR SELECTION”

BY ARUN ADVANI, TORU KITAGAWA, AND TYMON SŁOCZYŃSKI

MARCH 21, 2019

D Stylised Simulations: Design

Here we provide further details on the parameters and procedures for the stylised simulations described in Subsection 3.2 of Advani *et al.* (2019).

D.1 Details of Simulations for Scenario 2

For each sample we generate 1,000 observations, and for each observation draw a covariate x from a truncated standard normal distribution with the left truncation point at -4 and the right truncation point at 6 .

Propensity score $e(x)$ is then constructed as

$$e(x) = .4 + .1x. \tag{1}$$

For each observation we draw a random number from a standard uniform distribution, and assign treated status, $D = 1$, if $e(x)$ exceeds that random number.

We next generate an unobservable ϵ drawn from a normal distribution with mean zero. Since Scenario 2 is the heteroskedastic case, the standard deviation for those not treated is $\sigma_0 = .5$, while for those who are treated it is $\sigma_1 = 1.5$.¹

Finally, the outcome Y is generated as

$$Y = 3 + .5D + .5X + \epsilon, \tag{2}$$

and hence ATT is equal to $.5$.

This completes the generation of a Scenario 2 sample, which can then be used to implement the two EMCS procedures described in Section 2 of Advani *et al.* (2019). For each EMCS design, we consider 1,000 samples and 1,000 replications per sample.

In the placebo design, we additionally require some choice of π and λ , where λ determines the degree of covariate overlap between the ‘placebo treated’ and ‘placebo control’

¹In the benchmark case (Scenario 1), mentioned at the end of Subsubsection 3.2.1 in Advani *et al.* (2019), $\sigma_0 = \sigma_1 = .5$.

observations and π determines the proportion of the ‘placebo treated’. We choose π to ensure that the proportion of the ‘placebo treated’ observations in each placebo EMCS replication is equal to the proportion of treated units in the sample. We follow Huber *et al.* (2013) in choosing $\lambda = 1$. We also use a linear model to estimate the propensity score, as this corresponds to the true model in equation (1).

In the structured design, we first estimate the mean and variance of X in a given sample, conditional on treatment status. We also regress Y on D and X , excluding the interaction of D and X . Next, in the simulated dataset, X is drawn from a normal distribution with mean and variance conditional on treatment status and equal to the estimates above. Whenever the draw of X lies outside the support observed in the data, conditional on treatment status, the observation is replaced with the limit point of the support. Finally, the simulated outcome, Y , is generated in two steps. In the first step, we calculate its conditional mean based on the estimated coefficients from the regression above. In the second step, the simulated outcome is determined as a draw from a normal distribution with the conditional mean determined above and the variance that is equal to the variance of the residuals in the regression model estimated on the original data.² Again, we replace extreme values with the limit of the support, conditional on treatment status.

We use two estimators in our stylised simulations: linear regression (OLS) and inverse probability weighting (IPW). In the latter case, we first estimate the propensity score using a linear model, as this corresponds to the true model in equation (1), and then use inverse weighting with normalised weights to estimate the ATT.

D.2 Details of Simulations for Scenario 3

A similar procedure to that detailed in the previous subsection is followed. Two changes are made. First, we now have homoskedasticity so $\sigma_0 = \sigma_1 = .5$. Second, in each sample, we now generate the outcome Y as

$$Y = 3 + .5D + .5X + .5XD + \epsilon, \tag{3}$$

and hence ATT is equal to $.5 + .5 \cdot \mathbb{E}(X|D = 1)$.³

The source of misspecification of the structured design in Scenario 3 is in its failure to account for the interaction of D and X when generating the simulated outcomes.

²Thus, by using a single value of variance for both treated and control units, we fail to account for heteroskedasticity of the potential outcome equations. This is the source of misspecification of the structured design in Scenario 2.

³In practice, we estimate $\mathbb{E}(X|D = 1)$ using the mean of X for the treated observations in 1,000 samples from the true data generating process. As a result, ATT is equal to (approximately) .625.

E Stylised Simulations: Detailed Results

Table E1: Simulation results for Scenario 1 in Section 3 of Advani *et al.* (2019)

	Absolute bias	RMSE	SD
Original samples			
IPW	.000	.034	.034
OLS	.000	.032	.032
Placebo			
IPW	.002 (.001)	.044 (.002)	.044 (.002)
OLS	.001 (.001)	.042 (.002)	.042 (.002)
Structured			
IPW	.007 (.005)	.035 (.002)	.034 (.001)
OLS	.001 (.001)	.033 (.001)	.033 (.001)

Notes: Results for ‘Original samples’ correspond to the true values of all features of interest (absolute bias, RMSE, and SD) in the original data generating process. Measures of absolute bias and RMSE are centred around the true value of ATT, reported in Appendix D. All calculations are based on 1,000 samples. For each of these 1,000 samples, ‘Placebo’ and ‘Structured’ generate 1,000 new replications using the placebo and structured approaches described in Section 2 of Advani *et al.* (2019). In each case, we report both the mean and the standard deviation (in brackets) of EMCS estimates of all features of interest across all replications. Estimates of absolute bias and RMSE are centred around 0 for placebo and around the model-implied value for structured. For ease of interpretation, RMSE and SD are reported instead of MSE and variance (as elsewhere in the paper).

Table E2: **Simulation results for Scenario 2 in Section 3 of Advani *et al.* (2019)**

	Absolute bias	RMSE	SD
Original samples			
IPW	.003	.079	.079
OLS	.002	.080	.080
Placebo			
IPW	.002	.044	.044
	(.001)	(.002)	(.002)
OLS	.001	.042	.042
	(.001)	(.002)	(.002)
Structured			
IPW	.016	.070	.067
	(.012)	(.005)	(.003)
OLS	.010	.067	.066
	(.008)	(.003)	(.003)

Notes: Results for ‘Original samples’ correspond to the true values of all features of interest (absolute bias, RMSE, and SD) in the original data generating process. Measures of absolute bias and RMSE are centred around the true value of ATT, reported in Appendix D. All calculations are based on 1,000 samples. For each of these 1,000 samples, ‘Placebo’ and ‘Structured’ generate 1,000 new replications using the placebo and structured approaches described in Section 2 of Advani *et al.* (2019). In each case, we report both the mean and the standard deviation (in brackets) of EMCS estimates of all features of interest across all replications. Estimates of absolute bias and RMSE are centred around 0 for placebo and around the model-implied value for structured. For ease of interpretation, RMSE and SD are reported instead of MSE and variance (as elsewhere in the paper).

Table E3: **Simulation results for Scenario 3 in Section 3 of Advani *et al.* (2019)**

	Absolute bias	RMSE	SD
Original samples			
IPW	.001	.044	.044
OLS	.081	.089	.037
Placebo			
IPW	.002 (.001)	.043 (.002)	.043 (.002)
OLS	.001 (.001)	.042 (.002)	.042 (.002)
Structured			
IPW	.011 (.009)	.040 (.004)	.038 (.001)
OLS	.003 (.003)	.037 (.001)	.036 (.001)

Notes: Results for ‘Original samples’ correspond to the true values of all features of interest (absolute bias, RMSE, and SD) in the original data generating process. Measures of absolute bias and RMSE are centred around the true value of ATT, reported in Appendix D. All calculations are based on 1,000 samples. For each of these 1,000 samples, ‘Placebo’ and ‘Structured’ generate 1,000 new replications using the placebo and structured approaches described in Section 2 of Advani *et al.* (2019). In each case, we report both the mean and the standard deviation (in brackets) of EMCS estimates of all features of interest across all replications. Estimates of absolute bias and RMSE are centred around 0 for placebo and around the model-implied value for structured. For ease of interpretation, RMSE and SD are reported instead of MSE and variance (as elsewhere in the paper).

F Empirical Application: Detailed Results

Table F1: Simulation results for Subsection 5.1 of Advani *et al.* (2019)

	Absolute bias	RMSE	SD
Original samples			
Doubly-robust regression	142	1,019	1,010
IPW	51	1,102	1,101
Kernel matching	818	1,382	1,115
OLS	306	740	674
Oaxaca–Blinder	35	716	715
NN matching	16	1,209	1,209
Bias-adjusted NN matching	102	1,411	1,408
Placebo			
Doubly-robust regression	280 (206)	1,817 (236)	1,785 (234)
IPW	323 (161)	1,960 (225)	1,928 (221)
Kernel matching	68 (50)	1,244 (141)	1,242 (139)
OLS	385 (287)	901 (152)	776 (37)
Oaxaca–Blinder	420 (311)	927 (171)	783 (40)
NN matching	241 (172)	2,423 (326)	2,407 (318)
Bias-adjusted NN matching	319 (392)	3,531 (11,085)	3,509 (11,086)

Notes: Results for ‘Original samples’ correspond to the true values of all features of interest (absolute bias, RMSE, and SD) in the original data generating process. Measures of absolute bias and RMSE are centred around the true value of ATT, equal to \$1,794. All calculations are based on 1,000 samples. For each of these 1,000 samples, ‘Placebo’ and ‘Structured’ generate 1,000 new replications using the placebo and structured approaches described in Section 2 of Advani *et al.* (2019). Similarly, ‘Bootstrap’ generates 1,000 nonparametric bootstrap replications by sampling with replacement the same number of observations as the original data. In each of the three cases, we report both the mean and the standard deviation (in brackets) of EMCS estimates of all features of interest across all replications. Estimates of absolute bias and RMSE are centred around 0 for placebo, around the model-implied value for structured, and around the point estimate in the original sample for bootstrap. For ease of interpretation, RMSE and SD are reported instead of MSE and variance (as in Table 1 in Advani *et al.* (2019)). The ‘minimum’ value for each feature, as reported in Table 1 in Advani *et al.* (2019), is its lowest value among our estimators in the original data generating process (*i.e.* the lowest value in the ‘Original samples’ panel). The minimum value of absolute bias is 16; for MSE, it is 512,322 (or $\simeq 716^2$); for variance, it is 454,278 (or $\simeq 674^2$).

Table F1: **Simulation results for Subsection 5.1 of Advani *et al.* (2019) (cont.)**

	Absolute bias	RMSE	SD
Structured			
Doubly-robust regression	620 (492)	1,402 (340)	1,261 (113)
IPW	591 (488)	1,436 (331)	1,310 (124)
Kernel matching	408 (371)	1,458 (254)	1,426 (145)
OLS	558 (476)	1,125 (359)	1,006 (105)
Oaxaca–Blinder	690 (495)	1,192 (389)	997 (103)
NN matching	626 (492)	1,660 (311)	1,533 (122)
Bias-adjusted NN matching	620 (491)	1,634 (312)	1,509 (119)
Bootstrap			
Doubly-robust regression	128 (122)	1,197 (203)	1,186 (193)
IPW	86 (84)	1,305 (235)	1,301 (231)
Kernel matching	652 (495)	1,789 (307)	1,610 (189)
OLS	24 (18)	906 (68)	906 (68)
Oaxaca–Blinder	25 (19)	961 (93)	961 (93)
NN matching	552 (414)	1,653 (325)	1,518 (250)
Bias-adjusted NN matching	703 (637)	3,126 (3,560)	2,980 (3,562)

Notes: Results for ‘Original samples’ correspond to the true values of all features of interest (absolute bias, RMSE, and SD) in the original data generating process. Measures of absolute bias and RMSE are centred around the true value of ATT, equal to \$1,794. All calculations are based on 1,000 samples. For each of these 1,000 samples, ‘Placebo’ and ‘Structured’ generate 1,000 new replications using the placebo and structured approaches described in Section 2 of Advani *et al.* (2019). Similarly, ‘Bootstrap’ generates 1,000 nonparametric bootstrap replications by sampling with replacement the same number of observations as the original data. In each of the three cases, we report both the mean and the standard deviation (in brackets) of EMCS estimates of all features of interest across all replications. Estimates of absolute bias and RMSE are centred around 0 for placebo, around the model-implied value for structured, and around the point estimate in the original sample for bootstrap. For ease of interpretation, RMSE and SD are reported instead of MSE and variance (as in Table 1 in Advani *et al.* (2019)). The ‘minimum’ value for each feature, as reported in Table 1 in Advani *et al.* (2019), is its lowest value among our estimators in the original data generating process (*i.e.* the lowest value in the ‘Original samples’ panel). The minimum value of absolute bias is 16; for MSE, it is 512,322 (or $\simeq 716^2$); for variance, it is 454,278 (or $\simeq 674^2$).

Table F2: Simulation results for Subsection 5.2 of Advani *et al.* (2019)

	Absolute bias	RMSE	SD
Original samples			
Doubly-robust regression	1,222	1,566	980
IPW	1,081	1,514	1,060
Kernel matching	1,356	1,652	944
OLS	1,111	1,237	545
Oaxaca–Blinder	954	1,106	559
NN matching	1,122	1,732	1,320
Bias-adjusted NN matching	1,101	1,847	1,484
Placebo			
Doubly-robust regression	263 (196)	1,823 (198)	1,794 (197)
IPW	203 (132)	2,026 (197)	2,013 (198)
Kernel matching	78 (88)	1,487 (267)	1,483 (263)
OLS	379 (284)	909 (153)	791 (36)
Oaxaca–Blinder	408 (304)	930 (171)	797 (37)
NN matching	219 (156)	2,480 (253)	2,467 (250)
Bias-adjusted NN matching	290 (226)	3,233 (1,852)	3,214 (1,853)

Notes: Results for ‘Original samples’ correspond to the true values of all features of interest (absolute bias, RMSE, and SD) in the original data generating process. Measures of absolute bias and RMSE are centred around the true value of ATT, equal to \$0. All calculations are based on 1,000 samples. For each of these 1,000 samples, ‘Placebo’ and ‘Structured’ generate 1,000 new replications using the placebo and structured approaches described in Section 2 of Advani *et al.* (2019). Similarly, ‘Bootstrap’ generates 1,000 nonparametric bootstrap replications by sampling with replacement the same number of observations as the original data. In each of the three cases, we report both the mean and the standard deviation (in brackets) of EMCS estimates of all features of interest across all replications. Estimates of absolute bias and RMSE are centred around 0 for placebo, around the model-implied value for structured, and around the point estimate in the original sample for bootstrap. For ease of interpretation, RMSE and SD are reported instead of MSE and variance (as in Table 2 in Advani *et al.* (2019)). The ‘minimum’ value for each feature, as reported in Table 2 in Advani *et al.* (2019), is its lowest value among our estimators in the original data generating process (*i.e.* the lowest value in the ‘Original samples’ panel). The minimum value of absolute bias is 954; for MSE, it is 1,222,627 (or $\simeq 1,106^2$); for variance, it is 296,671 (or $\simeq 545^2$).

Table F2: Simulation results for Subsection 5.2 of Advani *et al.* (2019) (cont.)

	Absolute bias	RMSE	SD
Structured			
Doubly-robust regression	1,009 (398)	1,440 (327)	1,065 (70)
IPW	1,027 (401)	1,500 (326)	1,120 (84)
Kernel matching	827 (405)	1,411 (301)	1,156 (86)
OLS	1,023 (395)	1,295 (351)	858 (60)
Oaxaca–Blinder	1,072 (389)	1,327 (351)	851 (57)
NN matching	1,041 (403)	1,704 (303)	1,364 (84)
Bias-adjusted NN matching	1,010 (398)	1,642 (301)	1,318 (76)
Bootstrap			
Doubly-robust regression	155 (127)	1,152 (189)	1,136 (180)
IPW	84 (77)	1,262 (233)	1,258 (231)
Kernel matching	430 (369)	1,452 (284)	1,352 (208)
OLS	23 (17)	849 (45)	849 (44)
Oaxaca–Blinder	21 (16)	853 (43)	853 (43)
NN matching	643 (530)	1,789 (450)	1,615 (321)
Bias-adjusted NN matching	839 (804)	3,356 (1,025)	3,180 (932)

Notes: Results for ‘Original samples’ correspond to the true values of all features of interest (absolute bias, RMSE, and SD) in the original data generating process. Measures of absolute bias and RMSE are centred around the true value of ATT, equal to \$0. All calculations are based on 1,000 samples. For each of these 1,000 samples, ‘Placebo’ and ‘Structured’ generate 1,000 new replications using the placebo and structured approaches described in Section 2 of Advani *et al.* (2019). Similarly, ‘Bootstrap’ generates 1,000 nonparametric bootstrap replications by sampling with replacement the same number of observations as the original data. In each of the three cases, we report both the mean and the standard deviation (in brackets) of EMCS estimates of all features of interest across all replications. Estimates of absolute bias and RMSE are centred around 0 for placebo, around the model-implied value for structured, and around the point estimate in the original sample for bootstrap. For ease of interpretation, RMSE and SD are reported instead of MSE and variance (as in Table 2 in Advani *et al.* (2019)). The ‘minimum’ value for each feature, as reported in Table 2 in Advani *et al.* (2019), is its lowest value among our estimators in the original data generating process (*i.e.* the lowest value in the ‘Original samples’ panel). The minimum value of absolute bias is 954; for MSE, it is 1,222,627 (or $\simeq 1,106^2$); for variance, it is 296,671 (or $\simeq 545^2$).

Table F3: **Simulation results for Subsection 5.3 of Advani *et al.* (2019)**

	Absolute bias	RMSE	SD
Original samples			
Doubly-robust regression	68	1,573	1,572
IPW	565	1,682	1,585
Kernel matching	540	1,649	1,559
OLS	1,069	1,131	371
Oaxaca–Blinder	171	583	558
NN matching	442	2,374	2,333
Bias-adjusted NN matching	102	1,837	1,835
Placebo			
Doubly-robust regression	174 (134)	1,721 (156)	1,709 (153)
IPW	187 (117)	2,162 (172)	2,153 (171)
Kernel matching	144 (140)	1,925 (281)	1,917 (277)
OLS	208 (160)	651 (72)	600 (26)
Oaxaca–Blinder	298 (224)	753 (116)	664 (30)
NN matching	175 (136)	2,705 (242)	2,699 (238)
Bias-adjusted NN matching	175 (137)	1,942 (174)	1,931 (171)

Notes: Results for ‘Original samples’ correspond to the true values of all features of interest (absolute bias, RMSE, and SD) in the original data generating process. Measures of absolute bias and RMSE are centred around the true value of ATT, equal to $-\$405$. All calculations are based on 1,000 samples. For each of these 1,000 samples, ‘Placebo’ and ‘Structured’ generate 500 new replications using the placebo and structured approaches described in Section 2 of Advani *et al.* (2019). Similarly, ‘Bootstrap’ generates 500 nonparametric bootstrap replications by sampling with replacement the same number of observations as the original data. In each of the three cases, we report both the mean and the standard deviation (in brackets) of EMCS estimates of all features of interest across all replications. Estimates of absolute bias and RMSE are centred around 0 for placebo, around the model-implied value for structured, and around the point estimate in the original sample for bootstrap. For ease of interpretation, RMSE and SD are reported instead of MSE and variance (as in Table 3 in Advani *et al.* (2019)). The ‘minimum’ value for each feature, as reported in Table 3 in Advani *et al.* (2019), is its lowest value among our estimators in the original data generating process (*i.e.* the lowest value in the ‘Original samples’ panel). The minimum value of absolute bias is 68; for MSE, it is 340,300 (or $\simeq 583^2$); for variance, it is 137,574 (or $\simeq 371^2$).

Table F3: **Simulation results for Subsection 5.3 of Advani *et al.* (2019) (cont.)**

	Absolute bias	RMSE	SD
Structured			
Doubly-robust regression	198 (116)	819 (78)	816 (65)
IPW	194 (121)	1,228 (119)	1,224 (116)
Kernel matching	149 (116)	910 (107)	913 (106)
OLS	405 (253)	631 (172)	470 (20)
Oaxaca–Blinder	202 (114)	509 (54)	503 (21)
NN matching	140 (107)	1,205 (113)	1,210 (111)
Bias-adjusted NN matching	200 (117)	942 (81)	939 (71)
Bootstrap			
Doubly-robust regression	283 (252)	1,347 (396)	1,304 (363)
IPW	159 (303)	1,177 (625)	1,157 (570)
Kernel matching	381 (365)	1,276 (512)	1,187 (456)
OLS	15 (12)	415 (23)	415 (23)
Oaxaca–Blinder	22 (17)	620 (33)	620 (33)
NN matching	953 (808)	2,245 (967)	1,974 (726)
Bias-adjusted NN matching	689 (574)	1,946 (555)	1,763 (435)

Notes: Results for ‘Original samples’ correspond to the true values of all features of interest (absolute bias, RMSE, and SD) in the original data generating process. Measures of absolute bias and RMSE are centred around the true value of ATT, equal to $-\$405$. All calculations are based on 1,000 samples. For each of these 1,000 samples, ‘Placebo’ and ‘Structured’ generate 500 new replications using the placebo and structured approaches described in Section 2 of Advani *et al.* (2019). Similarly, ‘Bootstrap’ generates 500 nonparametric bootstrap replications by sampling with replacement the same number of observations as the original data. In each of the three cases, we report both the mean and the standard deviation (in brackets) of EMCS estimates of all features of interest across all replications. Estimates of absolute bias and RMSE are centred around 0 for placebo, around the model-implied value for structured, and around the point estimate in the original sample for bootstrap. For ease of interpretation, RMSE and SD are reported instead of MSE and variance (as in Table 3 in Advani *et al.* (2019)). The ‘minimum’ value for each feature, as reported in Table 3 in Advani *et al.* (2019), is its lowest value among our estimators in the original data generating process (*i.e.* the lowest value in the ‘Original samples’ panel). The minimum value of absolute bias is 68; for MSE, it is 340,300 (or $\approx 583^2$); for variance, it is 137,574 (or $\approx 371^2$).

References

- ADVANI, A., T. KITAGAWA, AND T. SŁOCZYŃSKI (2019): “Mostly Harmless Simulations? Using Monte Carlo Studies for Estimator Selection,” unpublished manuscript.
- HUBER, M., M. LECHNER, AND C. WUNSCH (2013): “The Performance of Estimators Based on the Propensity Score,” *Journal of Econometrics*, 175, 1–21.