

317 6 Appendix

318 6.1 Proofs of theorems

319 **Lemma 3.1.** *Suppose that $\mathbb{E}[Y | \text{Pa}(Y) = a] \neq \mathbb{E}[Y | \text{Pa}(Y) = a']$ whenever $a \neq a'$. Then a*
 320 *representation Φ is invariant across all valid environments if and only if $\mathbb{E}[Y^e | \Phi(T^e, X^e)] =$*
 321 *$\mathbb{E}[Y | \text{Pa}(Y)]$ for all valid environments.*

322 *Proof.* The if direction is immediate.

323 To establish the only if direction, we first show that Φ must contain at least $\text{Pa}(Y)$, in the sense
 324 $\mathbb{E}[Y | \Phi(X)] = \mathbb{E}[Y | \text{Pa}(Y) \cup Z]$ for some set Z . We proceed by contradiction. Suppose that
 325 conditioning on Φ is equivalent to conditioning on only $\text{Pa}(Y) \setminus \{P\} \cup Z$, where P is a parent of Y .
 326 We now create two environments by setting $P = p$ and $P = p'$. Since P is a parent of Y this follows
 327 from the second rule of do calculus [Pea00],

$$\mathbb{E}[Y | \text{Pa}(Y) \setminus \{P\} \cup Z; \text{do}(P = p)] = \mathbb{E}[Y | \text{Pa}(Y) \setminus \{P\} \cup Z, P = p]$$

328 and

$$\mathbb{E}[Y | \text{Pa}(Y) \setminus \{P\} \cup Z; \text{do}(P = p')] = \mathbb{E}[Y | \text{Pa}(Y) \setminus \{P\} \cup Z, P = p'].$$

329 The equality $\mathbb{E}[Y | \text{Pa}(Y) \setminus \{P\} \cup Z, P = p] = \mathbb{E}[Y | \text{Pa}(Y) \setminus \{P\} \cup Z, P = p']$ holds only if P is
 330 conditionally independent of Y given $\text{Pa}(Y) \setminus \{P\} \cup Z$. Since P is a parent of Y , by the first assumption
 331 of the lemma, the equality does not hold. It follows that $\mathbb{E}[Y | \text{Pa}(Y) \setminus \{P\} \cup Z; \text{do}(P = p)] \neq$
 332 $\mathbb{E}[Y | \text{Pa}(Y) \setminus \{P\} \cup Z; \text{do}(P = p')]$. That is, if conditioning on Φ was equivalent to conditioning
 333 on less information than $\text{Pa}(Y) \cup Z$, then Φ would not be invariant across all valid environments.

334 It remains to show that Φ does not contain any more information than $\text{Pa}(Y)$.

335 Φ cannot contain any descendants of the outcome. Suppose that Φ depends on some descendant D of
 336 Y in the sense that there is at least one environment and $d \neq d'$ where $\mathbb{E}[Y | \Phi(X \setminus D, D = d)] \neq$
 337 $\mathbb{E}[Y | \Phi(X \setminus D, D = d')]$. Then, construct a new environment e by randomly intervening and setting
 338 $\text{do}(D = d)$ or $\text{do}(D = d')$, each with probability 0.5. In this new environment, there is no relationship
 339 between Y and D . Accordingly, $\mathbb{E}[Y^e | \Phi(X^e \setminus D^e, D^e = d)] = \mathbb{E}[Y^e | \Phi(X^e \setminus D^e, D^e = d')]$.
 340 Thus, the conditional expectations are not equal (as functions of d) in the two environments—a
 341 contradiction.

342 Next, we show that, Φ need not to contain the non-parent ancestors A of the outcome, because
 343 $\mathbb{E}[Y | \{A\} \cup \text{Pa}(Y)] = \mathbb{E}[Y | \text{Pa}(Y)]$ by the Markov property of the causal graph, where A is any
 344 non-ancestor variables. Since Φ contains $\text{Pa}(Y)$, it follows that Φ does not depend on any non-parent
 345 ancestor A .

346

347 **Theorem 3.2.** *Let L be a loss function such that the minimizer of the associated risk is a conditional*
 348 *expectation, and let Φ be a representation that elicits a predictor Q^{inv} that is invariant for all*
 349 *valid distributions. Assuming there is no mediators between the treatment and the outcome, then*
 350 *$\psi = \mathbb{E}[Q^{\text{inv}}(1, X) - Q^{\text{inv}}(0, X) | T = 1]$.*

351 *Proof.* We assume the technical condition of lemma 3.1, that $\mathbb{E}[Y | \text{Pa}(Y) = a] \neq$
 352 $\mathbb{E}[Y | \text{Pa}(Y) = a']$ whenever $a \neq a'$. This is without loss of generality because violations of
 353 this condition will not lead to different causal effects.

354 By the assumption on the loss function, the elicited invariant predictor is $\mathbb{E}[Y | \Phi(T, X)]$. Lemma 3.1
 355 shows that $\mathbb{E}[Y | \Phi(T, X)] = \mathbb{E}[Y | \text{Pa}(Y)]$. We further observe that the non-treatment parents of
 356 Y are sufficient to block backdoor paths. It follows the ATT can be expressed as the following.

$$\begin{aligned} \Psi &= \mathbb{E}[\mathbb{E}[Y | T = 1, \text{Pa}(Y) \setminus \{T\}] - \mathbb{E}[Y | T = 0, \text{Pa}(Y) \setminus \{T\}]] | T = 1 \\ &= \mathbb{E}[\mathbb{E}[Y | \Phi(1, X)] - \mathbb{E}[Y | \Phi(0, X)] | T = 1] \end{aligned}$$

357

358 **Theorem 3.3.** *Suppose $0 < P(T^e = 1 | X^e) < 1$ with probability 1, then $0 < P(T^e = 1 | \Phi(X^e)) < 1$*
 359 *with probability 1.*

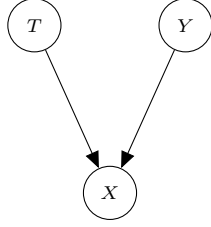


Figure 6: V-structure graph. We denote the bias induced by conditioning on X as V-bias.

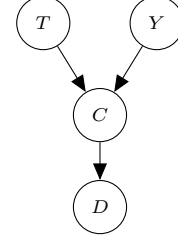


Figure 7: Y-structure graph. We denote the bias induced by conditioning on D as Y-bias.

360 *Proof.* The proof follows directly from Theorem 1 in [D'A+20]. The intuition is that the richer the
 361 covariate set is, the more likely it is to predict the treatment assignment accurately [D'A+20]. The
 362 covariate representation $\Phi(X^e)$ by definition contains less information than X^e , therefore $\Phi(X^e)$
 363 satisfies overlap if X^e satisfies overlap.

364 Consider the DGP with binary variables $\{X, Y, T\}$ illustrated in figure 6, where X is causally
 365 influence by Y and T .

366 **Theorem 6.1.** Let cov denote the covariance between two variables, we define collider bias at
 367 $X = c$ as $\Delta(X = c) = cov(T, Y | X = c) - cov(T, Y)$, and collider bias of X as $\Delta(X) =$
 368 $|P(X = 1)\Delta(X = 1) + P(X = 0)\Delta(X = 0)|$. Let $\Phi(T, X)$ be a random variable, where
 369 $P(\Phi(T, X) = X) \geq 0.5$. Suppose $P(X = 1) = 0.5$, and $\Delta(X = 1)$ has the same sign as
 370 $\Delta(X = 0)$, conditioning on X induce more collider bias than conditioning its coarsening $\Phi(T, X)$:

$$\Delta(\Phi(T, X)) \leq \Delta(X)$$

371
 372 *Proof.* The proof follows corollary 2.1 in [NDO19].

373 **Corollary 2.1.** We refer to collider bias in the V substructure embedded in the Y structure as
 374 'embedded V-bias' and denote it as $\Delta(C = c)$. For the covariance effect scale, Y-bias $\Delta(D = d)$
 375 relates to embedded V-bias through the following formula:

$$\Delta(D = d) = \frac{p(D = d | C = 1) - p(D = d | C = 0)}{\{P(D = d)\}^2} \cdot \left[\frac{p(D = d | C = 1)\{P(C = 1)\}^2 \cdot \Delta(C = 1) - p(D = d | C = 0)\{P(C = 0)\}^2 \cdot \Delta(C = 0)}{p(D = d | C = 0)\{P(C = 0)\}^2 \cdot \Delta(C = 0)} \right].$$

376 With the corollary above, let D denote $\Phi(T, X)$, let C denote the collider X in figure 6. The bias
 377 induced by conditioning on D is less than the bias induced by conditioning on C .

$$\begin{aligned} \Delta(D = 1) &= \frac{2\alpha - 1}{0.25} (0.25\alpha \cdot \Delta(C = 1) - 0.25(1 - \alpha) \cdot \Delta(C = 0)) \\ &= (2\alpha - 1)(\alpha \cdot \Delta(C = 1) - (1 - \alpha) \cdot \Delta(C = 0)) \\ \Delta(D = 0) &= \frac{1 - 2\alpha}{0.25} (0.25(1 - \alpha) \cdot \Delta(C = 1) - 0.25\alpha \cdot \Delta(C = 0)) \\ &= (1 - 2\alpha)((1 - \alpha) \cdot \Delta(C = 1) - \alpha \cdot \Delta(C = 0)) \\ \Delta(C) &= |0.5 \cdot \Delta(C = 0) + 0.5 \cdot \Delta(C = 1)| \\ \Delta(D) &= |0.5 \cdot \Delta(D = 0) + 0.5 \cdot \Delta(D = 1)| \\ \Delta(D) &= |0.5((1 - 2\alpha)((1 - \alpha) \cdot \Delta(C = 1) - \alpha \cdot \Delta(C = 0)) \\ &\quad + 0.5(2\alpha - 1)((\alpha \cdot \Delta(C = 1) - (1 - \alpha) \cdot \Delta(C = 0)))| \\ &= |0.5(2\alpha - 1)^2 \cdot \Delta(C = 1) + 0.5(2\alpha - 1)^2 \cdot \Delta(C = 0)| \\ &\leq \Delta(C) \end{aligned}$$

378 6.2 The Case of Mediators

379 In the most part of the paper, we assumed no mediators between treatment and outcome. What
 380 happens to the interpretation of the learned parameter if the adjustment set contains mediators?
 381 Intuitively, NICE retains the direct link between the treatment and the outcome. Specifically, if
 382 there are no mediators, the parameter reduces to ATT. If there are mediators but no confounders, the

383 parameter reduces to the Natural Direct Effect [Pea00]. If there are mediators and confounders, the
384 NICE estimand is a non-standard causal target that we call the natural direct effect on the treated
385 (NDET).

386 Conceptually, NDET describes the expected change in outcome Y for the treated population, induced
387 by changing the value of T , while keeping all mediating factors M , constant at whatever value they
388 would have obtained under $\text{do}(t)$. The main point is that NDET provides answers to questions such
389 as, “does this treatment have a substantial direct effect on this outcome?”. Substantively, NDET is the
390 natural direct effect, adjusted for confounders.

391 Formally, NDET is

$$\psi = \mathbb{E}_{M|T=1} [\mathbb{E}[Y | M; \text{do}(T = 1)] - \mathbb{E}[Y | M; \text{do}(T = 0)] | T = 1]. \quad (6.1)$$

392 With adjustment set W , the causal effect can be expressed through a parameter of the observational
393 distribution:

$$\psi = \mathbb{E}_{M,W} [\mathbb{E}[Y | T = 1, M, W] - \mathbb{E}[Y | T = 0, M, W] | T = 1]. \quad (6.2)$$

394 Importantly, the mediators M and the confounders W show up in the same way in (6.2). Accordingly,
395 we don’t need to know which observed variables are mediators and which are confounders to compute
396 the parameter. Under the NICE procedure, we condition on all parents of Y , including possible
397 mediators. Thus, the NICE estimand is the NDET.

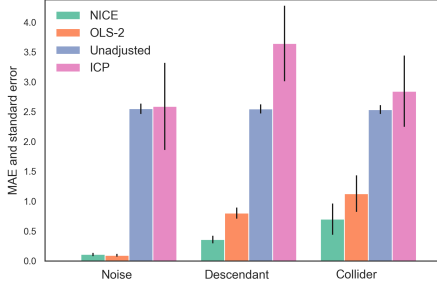


Figure 8: Models performance under the scrambled and heteroskedastic setting

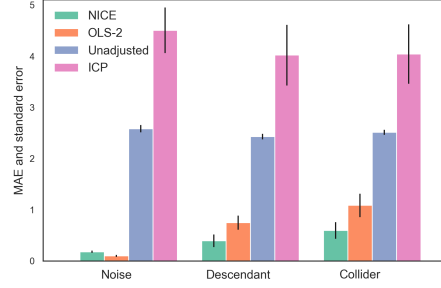


Figure 9: Models performance under the scrambled and homoscedastic setting

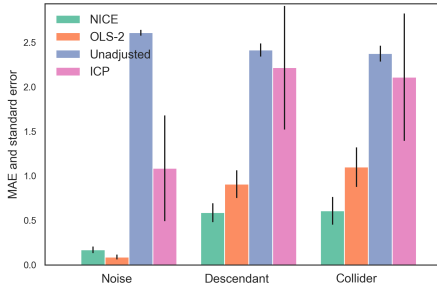


Figure 10: Models performance under the unscrambled and heteroskedastic setting

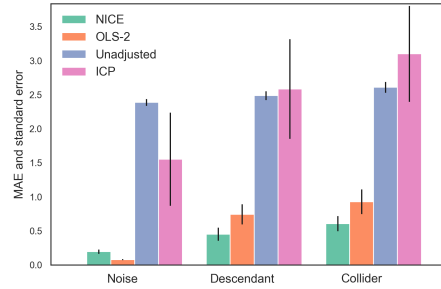


Figure 11: Models performance under the unscrambled and homoscedastic setting

398 **6.3 Details of the experiments**

399 **Experiment 1**

400 We simulate data with three causal graphs in Figure 3. With a slight abuse of notation, each
 401 intervention e generates a new environment e with interventional distribution $P(X^e, T^e, Y^e)$. T^e
 402 is the binary treatment and Y^e is the outcome. X^e is a 10-dimensional covariate set that differs
 403 across DGPs. $X^e = (X_1^e, X_2^e)$, where X_1^e is a five-dimensional confounder. X_2^e is either noise, a
 404 descendant, or a collider in each DGP. We examine the models' performance under two types of
 405 variations: 1) whether the observed covariates are scrambled versions of the true covariates. 2)
 406 whether the treatment effects are heteroskedastic across environments. The data generating process is
 407 illustrated below.

$$\begin{aligned}
 X_1^e &\leftarrow \mathcal{N}(0, e^2) \\
 P^e &\leftarrow \text{sigmoid}(X_1^e \cdot w_{xt^e} + \mathcal{N}(0, 1)) \\
 T^e &\leftarrow \text{Bern}(P^e) \\
 \tau &\leftarrow 5 + \mathcal{N}(0, \sigma^2) \\
 Y^e &\leftarrow X_1^e \cdot w_{xy^e} + T^e \cdot \tau + \mathcal{N}(0, e^2)
 \end{aligned}$$

408 X_2^e equals $\mathcal{N}(0, e^2)$ in setting a), X_2^e equals $e * Y^e + \mathcal{N}(0, q)$ in setting b), and X_2^e equals
 409 $e * Y^e + T^e + \mathcal{N}(0, q)$ in setting c). For the four variants, in the scrambled setting $\mathcal{N}(0, \sigma^2) = \mathcal{N}(0, e^2)$,
 410 in the un-scrambled setting $\mathcal{N}(0, \sigma^2) = \mathcal{N}(0, 1)$. In the environment-level heteroskedastic setting
 411 $\tau \leftarrow 5 + \mathcal{N}(0, e^2)$. In the environment-level homoscedastic setting $\tau \leftarrow 5 + \mathcal{N}(0, 1)$. The
 412 performance under the four variants are illustrated in Figure 8, Figure 9, Figure 10, and Figure 11.

413 **Experiment 2**

414 We validate NICE for the non-linear case on the benchmark dataset SpeedDating. SpeedDating was
 415 collected to study the gender difference in mate selection [Fis+06]. The study recruited university
 416 students to participate in speed dating, and collected objective and subjective information such as
 417 'undergraduate institution' and 'perceived attractiveness'. It has 8378 entries and 185 covariates.
 418 ACIC 2019's simulation samples subsets of the covariates to simulate the treatment T and outcome

419 Y . Specifically, it provides four modified DGPs: Mod1: parametric models; Mod2: complex models;
 420 Mod3: parametric models with poor overlap; Mod4: complex models with treatment heterogeneity.
 421 Each modification includes three versions: low, med, high, indicating an increasing number of
 422 covariates included in the models for T and Y .

423 We compare the estimation quality of the within-sample SATT and out-of-sample CATE over 10
 424 bootstraps in Table 3 and Table 4. The main paper reports the models performance under the low
 setting. We now report results for the med and high setting.

Within-sample				
	ϵ_{att}			
	MOD1	MOD2	MOD3	MOD4
med				
TARNet	.04 ± .02	.09 ± .1	.23 ± .16	.12 ± .09
+NICE	.02 ± .02	.06 ± .03	.07 ± .04	.02 ± .02
Dragon	.12 ± .12	.09 ± .1	.23 ± .16	.13 ± .07
+NICE	.04 ± .02	.06 ± .03	.07 ± .04	.07 ± .08
high				
TARNet	.06 ± .05	.18 ± .12	.10 ± .08	.03 ± .03
+NICE	.01 ± .01	.06 ± .03	.05 ± .02	.05 ± .06
Dragon	.18 ± .12	.18 ± .12	.10 ± .08	.05 ± .03
+NICE	.02 ± .01	.06 ± .03	.05 ± .02	.09 ± .09

Table 3: NICE perform well relative to the baselines if the adjustment set does not contain bad controls. The table reports MAE and bootstrap standard deviation of the SATT estimation. The model is trained and evaluated on all three environments.

425

Out-of-sample				
	$\sqrt{\epsilon_{PEHE}}$			
	MOD1	MOD2	MOD3	MOD4
med				
TARNet	.14 ± .04	.13 ± .03	.11 ± .05	.09 ± .06
+NICE	.06 ± .02	.06 ± .01	.08 ± .05	.07 ± .02
Dragon	.15 ± .03	.13 ± .03	.22 ± .15	.07 ± .02
+NICE	.04 ± .01	.04 ± .01	.08 ± .02	.07 ± .04
high				
TARNet	.14 ± .06	.14 ± .08	.13 ± .03	.08 ± .10
+NICE	.05 ± .01	.07 ± .01	.06 ± .01	.07 ± .08
Dragon	.14 ± .04	.14 ± .05	.15 ± .04	.09 ± .08
+NICE	.06 ± .01	.06 ± .02	.05 ± .01	.05 ± .06

Table 4: NICE perform well relative to the baselines if the adjustment set does not contain bad controls. The table reports PEHE and bootstrap standard deviation of the out-of-distribution CATE estimation. The model is trained on two environments and evaluated on the third.

426 To examine whether NICE helps reduce collider bias, we simulated 20 copies of a collider: $X_{co}^e =$
 427 $T^e + Y^e + \mathcal{N}(0, e^2)$, where $e \in \{0.01, 0.2, 1\}$ and include it in the covariate set. table 5 compares
 428 Dragonnet trained under standard empirical risk minimization framework and trained under NICE.
 429 NICE reduces collider bias across simulation setups.

430 Experiment 3

431 In the third experiment, we consider data generated according figure 2. Notably, in the setup, we
 432 observe $P(A, X, T, Y, Z)$, where $A = \{A_1, A_t, A_y\}$ and $X = \{X_t, X_y\}$. Here A is a 50 dimensional
 433 covariate, X a 30 dimensional covariate, and Z a 50 dimensional covariate. Z is causally affected by
 434 A and Y .

435 We compare NICE against a neural network model similar to the structure of TARNet. The model
 436 architecture is the same as the models in the SpeedDating experiment, except the hidden layer

Within-sample

		ϵ_{att}			
		Mod1	Mod2	Mod3	Mod4
low					
Dragon		.32 ± .16	.39 ± .14	.32 ± .02	.50 ± .04
+NICE		.11 ± .05	.08 ± .07	.15 ± .04	.08 ± .05
med					
Dragon		.39 ± .15	.29 ± .13	.37 ± .13	.27 ± .10
+NICE		.08 ± .04	.17 ± .10	.09 ± .03	.06 ± .04
high					
Dragon		.36 ± .11	.35 ± .09	.49 ± .16	.28 ± .06
+NICE		.09 ± .06	.15 ± .08	.10 ± .08	.14 ± .09

Table 5: NICE reduces estimation bias in the presence of colliders. The table reports the MAE and bootstrap standard deviation of SATT. The model is trained and evaluated on three environments.

Out-of-sample

		ϵ_{pehe}			
		Mod1	Mod2	Mod3	Mod4
low					
TARNet		.18 ± .05	.42 ± .03	.25 ± .04	.36 ± .12
+NICE		.08 ± .02	.07 ± .01	.08 ± .02	.08 ± .03
Dragon		.25 ± .06	.49 ± .05	.29 ± .06	.45 ± .06
+NICE		.09 ± .01	.09 ± .03	.09 ± .03	.09 ± .04
med					
TARNet		.41 ± .08	.28 ± .08	.35 ± .06	.21 ± .03
+NICE		.09 ± .02	.08 ± .02	.07 ± .02	.09 ± .01
Dragon		.40 ± .08	.32 ± .04	.47 ± .11	.25 ± .06
+NICE		.07 ± .01	.10 ± .05	.08 ± .03	.07 ± .01
high					
TARNet		.24 ± .06	.26 ± .10	.28 ± .08	.25 ± .10
+NICE		.06 ± .02	.09 ± .02	.07 ± .02	.10 ± .05
Dragon		.34 ± .14	.36 ± .03	.37 ± .08	.28 ± .07
+NICE		.09 ± .03	.13 ± .04	.07 ± .01	.12 ± .05

Table 6: NICE reduces estimation bias in the presence of colliders. The table reports the CATE and standard deviation of CATE. The model is trained on two environments and evaluated on a third environment.

437 size is 200 for the shared representation. For the exact data generating process and the detailed
 438 implementation of the models, see the associated codebase.

439 **References**

- 440 [Arj+19] M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz. “Invariant risk minimization”.
441 In: *arXiv preprint arXiv:1907.02893* (2019).
- 442 [BP14] E. Bareinboim and J. Pearl. “Transportability from multiple environments with limited
443 experiments: completeness results”. In: *Advances in neural information processing*
444 *systems*. 2014.
- 445 [BV07] J. Bhattacharya and W. B. Vogt. *Do instrumental variables belong in propensity scores?*
446 Tech. rep. National Bureau of Economic Research, 2007.
- 447 [Büh18] P. Bühlmann. “Invariance, causality and robustness”. In: *arXiv preprint*
448 *arXiv:1812.08233* (2018).
- 449 [CEP16] K. Chalupka, F. Eberhardt, and P. Perona. “Multi-level cause-effect systems”. In:
450 *Artificial Intelligence and Statistics*. 2016.
- 451 [CPE14] K. Chalupka, P. Perona, and F. Eberhardt. “Visual causal feature learning”. In: *arXiv*
452 *preprint arXiv:1412.2309* (2014).
- 453 [CH20] C. Cinelli and C. Hazlett. “Making sense of sensitivity: extending omitted variable
454 bias”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 1
455 (2020).
- 456 [D’A+20] A. D’Amour, P. Ding, A. Feller, L. Lei, and J. Sekhon. “Overlap in observational studies
457 with high-dimensional covariates”. In: *Journal of Econometrics* (2020).
- 458 [Fis+06] R. Fisman, S. S. Iyengar, E. Kamenica, and I. Simonson. “Gender differences in mate
459 selection: evidence from a speed dating experiment”. In: *The Quarterly Journal of*
460 *Economics* 2 (2006).
- 461 [GZS19] C. Glymour, K. Zhang, and P. Spirtes. “Review of causal discovery methods based on
462 graphical models”. In: *Frontiers in genetics* (2019).
- 463 [Gru+] S. Gruber, G. Lefebvre, A. PichÃ, and T. Schuster. *Data Challenge*. URL: <https://sites.google.com/view/acic2019datachallenge>.
464
- 465 [Haa43] T. Haavelmo. “The statistical implications of a system of simultaneous equations”. In:
466 *Econometrica, Journal of the Econometric Society* (1943).
- 467 [HDPM18] C. Heinze-Deml, J. Peters, and N. Meinshausen. “Invariant causal prediction for
468 nonlinear models”. In: *Journal of Causal Inference* 2 (2018).
- 469 [Hil11] J. Hill. “Bayesian nonparametric modeling for causal inference”. In: *Journal of Compu-*
470 *tational and Graphical Statistics* (2011).
- 471 [Imb04] G. W. Imbens. “Nonparametric estimation of average treatment effects under exogeneity:
472 a review”. In: *Review of Economics and statistics* 1 (2004).
- 473 [Kao+12] C.-H. Kao, L.-M. Sun, J.-A. Liang, S.-N. Chang, F.-C. Sung, and C.-H. Muo. “Rela-
474 tionship of zolpidem and cancer risk: a taiwanese population-based cohort study”. In:
475 *Mayo clinic proceedings*. 5. Elsevier. 2012.
- 476 [KB14] D. P. Kingma and J. Ba. “Adam: a method for stochastic optimization”. In: *arXiv*
477 *preprint arXiv:1412.6980* (2014).
- 478 [KLK12] D. F. Kripke, R. D. Langer, and L. E. Kline. “Hypnotics’ association with mortality or
479 cancer: a matched cohort study”. In: *BMJ open* 1 (2012).
- 480 [Mag+18] S. Magliacane, T. van Ommen, T. Claassen, S. Bongers, P. Versteeg, and J. M.
481 Mooij. “Domain adaptation by using causal inference to predict invariant conditional
482 distributions”. In: *Advances in Neural Information Processing Systems*. 2018.
- 483 [ML09] K. L. Moore and M. J. van der Laan. “Covariate adjustment in randomized trials with
484 binary outcomes: targeted maximum likelihood estimation”. In: *Statistics in medicine* 1
485 (2009).
- 486 [MM+99] K. Murphy, S. Mian, et al. *Modelling gene expression data using dynamic Bayesian*
487 *networks*. Tech. rep. Citeseer, 1999.
- 488 [NDO19] T. Q. Nguyen, A. Dafoe, and E. L. Ogburn. “The magnitude and direction of collider
489 bias for binary variables”. In: *Epidemiologic Methods* 1 (2019).

- 490 [Pat+17] E. Patorno, R. J. Glynn, R. Levin, M. P. Lee, and K. F. Huybrechts. “Benzodiazepines
491 and risk of all cause mortality in adults: cohort study”. In: *bmj* (2017).
- 492 [Pea00] J. Pearl. *Causality: models, reasoning and inference*. 2000.
- 493 [Pea09] J. Pearl. *Causality*. 2009.
- 494 [PP14] J. Pearl and A. Paz. “Confounding equivalence in causal inference”. In: *Journal of
495 Causal Inference J. Causal Infer.* 1 (2014).
- 496 [PBM16] J. Peters, P. Bühlmann, and N. Meinshausen. “Causal inference by using invariant
497 prediction: identification and confidence intervals”. In: *Journal of the Royal Statistical
498 Society: Series B (Statistical Methodology)* (2016).
- 499 [RJ91] L. D. Robinson and N. P. Jewell. “Some surprising results about covariate adjustment in
500 logistic regression models”. In: *International Statistical Review/Revue Internationale
501 de Statistique* (1991).
- 502 [RC+18] M. Rojas-Carulla, B. Schölkopf, R. Turner, and J. Peters. “Invariant models for causal
503 transfer learning”. In: *The Journal of Machine Learning Research* 1 (2018).
- 504 [Ros02] P. R. Rosenbaum. “Overt bias in observational studies”. In: *Observational studies*. 2002.
- 505 [RR83] P. R. Rosenbaum and D. B. Rubin. “The central role of the propensity score in
506 observational studies for causal effects”. In: *Biometrika* (1983).
- 507 [Rub09] D. B. Rubin. “Should observational studies be designed to allow lack of balance in
508 covariate distributions across treatment groups?” In: *Statistics in Medicine* 9 (2009).
- 509 [SJS16] U. Shalit, F. D. Johansson, and D. Sontag. “Estimating individual treatment effect:
510 generalization bounds and algorithms”. In: *arXiv e-prints arXiv:1606.03976* (2016).
- 511 [SBV19] C. Shi, D. M. Blei, and V. Veitch. “Adapting neural networks for the estimation of
512 treatment effects”. In: *Advances in neural information processing systems* (2019).
- 513 [Shi+06] S. Shimizu, P. O. Hoyer, A. Hyvärinen, and A. Kerminen. “A linear non-gaussian
514 acyclic model for causal discovery”. In: *Journal of Machine Learning Research* Oct
515 (2006).
- 516 [SE17] S. M. Shortreed and A. Ertefaie. “Outcome-adaptive lasso: variable selection for causal
517 inference”. In: *Biometrics* 4 (2017).
- 518 [Spi+00] P. Spirtes, C. Glymour, R. Scheines, S. Kauffman, V. Aimale, and F. Wimberly.
519 “Constructing bayesian network models of gene expression networks from microarray
520 data”. In: (2000).
- 521 [Wei+14] S. Weich, H. L. Pearce, P. Croft, S. Singh, I. Crome, J. Bashford, and M. Frisher. “Effect
522 of anxiolytic and hypnotic drug prescriptions on mortality hazards: retrospective cohort
523 study”. In: *Bmj* (2014).
- 524 [Zha+20] K. Zhang, M. Gong, P. Stojanov, B. Huang, and C. Glymour. “Domain adaptation as
525 a problem of inference on graphical models”. In: *arXiv preprint arXiv:2002.03278*
526 (2020).
- 527 [Zha+16] Q. Zhao, C. Zheng, T. Hastie, and R. Tibshirani. “Comment on causal inference using
528 invariant prediction”. In: *arXiv preprint arXiv: 1501.01332* (2016).