

591 **A Zero-shot Rademacher complexity and Proof of Theorem 1**

592 **A.1 Problem setup and assumptions**

593 Let  $w \in \mathcal{W} \subseteq \mathbb{R}^e$  denote an intervention and  $x \in \mathcal{X} \subseteq \mathbb{R}^d$  denote an individual that received it. Assume  
 594 the outcome to predict is a scalar  $y \in [0,1]$ . The hypothesis class is  $\mathcal{F} = \{f : (w,x) \rightarrow y\}$ . The dataset has  
 595  $n$  interventions with  $m$  independent units which received each intervention, *i.e.* first  $n$  *i.i.d.* draws from  
 596  $P_W$  and then  $m$  *i.i.d.* draws from  $P_X$  for each  $w^{(j)}$ . During training we have access to noisy estimate  $\tilde{y} =$   
 597  $y + \xi$  where  $\xi$  is an independent noise with  $\mathbb{E}\xi = 0$  and  $|\xi| \leq \epsilon$  almost surely. We are tested directly on  $y$ .

598 The ERM is

$$\hat{f} = \min_f \hat{L}(f) = \min_f \frac{1}{nm} \sum_{j=1}^n \sum_{i=1}^m (f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2.$$

599 The test error is

$$L(f) = \mathbb{E}_{w,x,y} (f(w,x) - y)^2$$

600 and let  $f^* = \min_f L(f)$ .

601 We are interested in bounding the excess error  $L(\hat{f}) - L(f^*)$ .

602 Our key assumption is that interventions with similar attributes ( $w$ ) have similar effects in expectation.  
 603 More concretely, we assume that all hypotheses in our family are smooth with respect to  $w$ :

**Assumption 2.**

$$\forall f \in \mathcal{F}, \mathbb{E}_{w,x} \left[ \left\| \frac{\partial f}{\partial w} \right\|_2^2 \right] \leq \beta^2.$$

604 Furthermore, we assume that  $P_W$  satisfies a Poincaré-type inequality:

605 **Assumption 3.** For some constant  $C$  that only depends on  $P_W$ , for any smooth function  $F$ ,

$$\text{Var}_w [F(w)] \leq C \mathbb{E} [\|\nabla_w F(w)\|_2^2].$$

606 For example,  $P_W$  can be any of the following distributions:

- 607 • Multivariate Gaussian:  $w \in \mathbb{R}^e \sim \mathcal{N}(\mu, \Sigma)$  for some vector  $\mu \in \mathbb{R}^e$  and positive semi-definite  
 608 matrix  $\Sigma \in \mathbb{R}^{e \times e}$ ;
- 609 •  $w \in \mathbb{R}^e$  has independent coordinates; each coordinate has the symmetric exponential  
 610 distribution  $1/2e^{-|t|}$  for  $t \in \mathbb{R}$ .
- 611 •  $P_W$  is a mixture over base distributions satisfying Poincaré inequalities, and their pair-wise  
 612 chi-squared distances are bounded.
- 613 •  $P_W$  is a mixture of isotropic Gaussians in  $\mathbb{R}^e$ .
- 614 •  $P_W$  is the uniform distribution over  $\mathcal{W} \subset \mathbb{R}^e$ , which is open, connected, and bounded with  
 615 Lipschitz boundary.

616 We note that isotropic Gaussian can approximate any smooth densities in  $\mathbb{R}^e$  [39] (since RBF kernels  
 617 are universal), showing that Assumption 3 is fairly general.

618 We define a novel notion of function complexity specialized to the zero-shot setting. Intuitively, it  
 619 measure how well we can fit random labels, which is first drawing  $n$  interventions and  $m$  recipients  
 620 for each intervention. For examples of concrete upper bound on zero-shot Rademacher complexity  
 621 see section [A.4](#)

$$R_{nm}(F) = \frac{1}{nm} \mathbb{E}_{w,x,\sigma} \sup_f \sum_{j=1}^n \sum_{i=1}^m \sigma_i^j f(w^{(j)}, x_i^{(j)}) \quad (8)$$

622 where  $\sigma_i^j$  are independently randomly drawn from  $\{-1,1\}$ .

623 **A.2 Formal theorem statement**

624 **Theorem 4.** Under Assumptions [2](#)[3](#) with probability  $1 - \delta$ ,

$$L(\hat{f}) \leq L(f^*) + 8(1 + \epsilon)R_{nm}(\mathcal{F}) + 8\sqrt{\frac{(1 + \epsilon)R_{nm}(\mathcal{F})\log(1/\delta)}{n}} \\ + (1 + \epsilon)\sqrt{\frac{(32C\beta^2 + \frac{2(1 + \epsilon)^2}{m})\log(1/\delta)}{n} + \frac{2\log(1/\delta)}{3n}}.$$

625 **A.3 Proof of the main theorem**

626 We define the population loss on the noisy label  $\tilde{L}(f) = \mathbb{E}_{w,x,\tilde{y}}(f(w,x) - \tilde{y})^2$ . Due to independence  
627 of  $\xi$ ,  $\mathbb{E}_{w,x,y,\xi}(f(w,x) - y - \xi)^2 = \mathbb{E}_{w,x,y}(f(w,x) - y)^2 + \mathbb{E}[\xi^2] = L(f) + \mathbb{E}[\xi^2]$  for any  $f$ , so  
628  $L(\hat{f}) - L(f^*) = \tilde{L}(\hat{f}) - \tilde{L}(f^*)$ . We shall focus on bounding the latter.

629 We first need a lemma that bounds the supremum of an empirical process indexed by a bounded  
630 function class.

631 **Lemma 5** (Theorem 2.3 of [6](#)). Assume that  $X_j$  are identically distributed according to  $P$ ,  $\mathcal{G}$  is  
632 a countable set of functions from  $\mathcal{X}$  to  $\mathbb{R}$  and, and all  $g \in \mathcal{G}$  are  $P$ -measurable, square-integrable,  
633 and satisfy  $\mathbb{E}[g] = 0$ . Suppose  $\sup_{g \in \mathcal{G}} \|g\|_\infty \leq 1$ , and we denote  $Z = \sup_g \left| \sum_{j=1}^n g(X_j) \right|$ . Suppose  
634  $\sigma^2 \geq \sup_{g \in \mathcal{G}} \text{Var}(g(X_j))$  almost surely, then for all  $t \geq 0$ , we have

$$\Pr \left[ Z \geq \mathbb{E}Z + \sqrt{2t(n\sigma^2 + 2\mathbb{E}Z)} + \frac{t}{3} \right] \leq e^{-t}.$$

635 We apply Lemma [5](#) with  $X_j = (w^{(j)}, x_1^j, \dots, x_m^j, \tilde{y}_1^j, \dots, \tilde{y}_m^j)$ ,  $g(X_j) =$   
636  $\left( \frac{1}{m} \sum_i (f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2 - \tilde{L}(f) \right)$ ,  $\sigma^2 = \sup_{f \in \mathcal{F}} (\text{Var}(\frac{1}{m} \sum_i (f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2))$ ,  
637  $t = \log(1/\delta)$ . Since  $f - \tilde{y} \in [-1, 1]$ ,  $g \in [-1, 1]$ . With probability  $1 - \delta$ ,

$$n \sup_f |\hat{L}(f) - \tilde{L}(f)| \leq n \mathbb{E} \sup_f |\hat{L}(f) - \tilde{L}(f)| + \sqrt{2 \log \frac{1}{\delta} \left( n \sigma^2 + 2n \mathbb{E} \sup_f |\hat{L}(f) - \tilde{L}(f)| \right)} + \frac{1}{3} \log \frac{1}{\delta}.$$

638 Multiplying both sides by  $1/n$ , and using  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ ,

$$\sup_f |\hat{L}(f) - \tilde{L}(f)| \leq \mathbb{E} \sup_f |\hat{L}(f) - \tilde{L}(f)| + 2 \sqrt{\frac{\mathbb{E} \sup_f |\hat{L}(f) - \tilde{L}(f)| \log(1/\delta)}{n}} + \sqrt{\frac{2\sigma^2 \log(1/\delta)}{n}} + \frac{\log(1/\delta)}{3n}. \quad (9)$$

639 The next lemma bounds the variance  $\sigma^2$  in equation [9](#).

**Lemma 6.**

$$\forall f \in \mathcal{F}, \text{Var}_{w^{(j)}, x_{1 \dots m}^j, \tilde{y}_{1 \dots m}^j} \left[ \frac{1}{m} \sum_{i=1}^m (f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2 \right] \leq 4(1 + \epsilon)^2 C \beta^2 + \frac{(1 + \epsilon)^4}{4m}.$$

640 *Proof of Lemma [6](#)* Using the law of total variance, if we write

$$g(w^{(j)}, x_{1 \dots m}^j, \tilde{y}_{1 \dots m}^j) = \frac{1}{m} \sum_{i=1}^m (f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2,$$

641 then

$$\text{Var}[g] = \text{Var}_w [\mathbb{E}_{x,\tilde{y}} [g(w,x,\tilde{y}) | w]] + \mathbb{E}_h [\text{Var}_{x,\tilde{y}} [g(w,x,\tilde{y}) | w]] \quad (10)$$

642 To bound the first term of equation [10](#), we use Poincaré-type inequalities in Assumption [3](#). For each  
643 of the example distributions, we show that they indeed satisfy Assumption [3](#).

644 **Lemma 7.** *Each of the example distributions in Assumption 3 satisfies a Poincaré-type inequality.*

645 *Proof.* • When  $P_W$  is the uniform distribution over  $\mathcal{W} \in \mathbb{R}^e$ , which is open, connected, and  
 646 bounded with Lipschitz boundary, we use Poincaré–Wirtinger inequality [57] on the smooth  
 647 function  $\mathbb{E}[g | w]$ : For some constant  $C$  that only depends on  $P_W$ ,

$$\text{Var}_w[\mathbb{E}[g | w]] \leq C \mathbb{E}[\|\nabla_w \mathbb{E}[g | w]\|_2^2]. \quad (11)$$

648  $C$  is the Poincaré constant for the domain  $\mathcal{W}$  in  $L_2$  norm. It can be bounded by  $1/\lambda_1$  where  
 649  $\lambda_1$  is the first eigenvalue of the negative Laplacian of the manifold  $\mathcal{W}$  [83]. Many previous  
 650 works study the optimal Poincaré constants for various domains [43]. For example, when  $w$  is  
 651 uniform over  $\mathcal{W}$  which is a bounded, convex, Lipschitz domain with diameter  $d$ ,  $C \leq d/\pi$  [56].

652 We can apply probabilistic Poincaré inequalities over non-Lebesgue measure  $P_W$ :

653 • When  $w \sim \mathcal{N}(\mu, \Sigma)$ , we use the Gaussian Poincaré inequality (see e.g. Theorem 3.20 of [5]  
 654 and using change of variables),

$$\text{Var}[F(w)] \leq \mathbb{E}[\langle \Sigma \nabla_w F(w), \nabla_w F(w) \rangle].$$

655 We apply this with  $F(w) = \mathbb{E}[g | w]$ . Since  $\mathbb{E}[v^\top Av] = \mathbb{E}[\text{Tr}(v^\top Av)] = \mathbb{E}[\text{Tr}(Avv^\top)] =$   
 656  $\text{Tr}(A\mathbb{E}[vv^\top]) \leq \|A\|_2 \mathbb{E}[\|v\|_2^2]$ ,

$$\text{Var}_w[\mathbb{E}[g | w]] \leq \|\Sigma\|_2 \mathbb{E}[\|\nabla_w \mathbb{E}[g | w]\|_2^2],$$

657 which satisfies equation (11) with  $C = \|\Sigma\|_2$ .

658 • When  $w \in \mathbb{R}^e$  has independent coordinates  $w_1, \dots, w_e$  and each coordinate has the symmetric  
 659 exponential distribution  $1/2e^{-|t|}$  for  $t \in \mathbb{R}$ , we first bound a single dimension using Lemma  
 660 4.1 of [45], which says for any function  $k \in L^1$ ,

$$\text{Var}(k(w_i)) \leq 4\mathbb{E}[k'(w_i)^2]$$

661 which, combined with the Efro-Stein inequality (Theorem 3.1 of [5]),

$$\text{Var}(F(w)) = \mathbb{E} \sum_{i=1}^e \text{Var}(F(w) | w_1, \dots, w_{i-1}, w_{i+1}, \dots, w_n),$$

662 yields:

$$\text{Var}(F(w)) \leq 4\mathbb{E}[\|F'(w)\|_2^2]$$

663 which satisfies equation (11) with  $C = 4$ .

664 Lastly, we consider the case where  $P_W$  is a mixture over base distributions satisfying Poincaré  
 665 inequalities. We first consider the case where the pair-wise chi-squared distances are bounded. Next,  
 666 we show that mixture of isotropic Gaussians satisfies Poincaré inequality without further condition  
 667 on pair-wise chi-squared distances.

668 • When  $\{P_W^q\}_{q \in \mathcal{Q}}$  is a family of distributions, each satisfying Poincaré inequality with constant  
 669  $C^q$ , and  $P_W$  is any mixture over  $\{P_W^q\}_{q \in \mathcal{Q}}$  with density  $\mu$ , let  $K_P(\mu) = \text{ess sup}_q C^q$ ,  
 670 which is an upper bound on the base Poincaré constants almost surely, and  $K_{\chi^2}(\mu) =$   
 671  $\mathbb{E}_{q, q' \sim \mu} [(1 + \chi^2(P_W^q || P_W^{q'}))^p]^{1/p}$ , which is an upper bound on the pairwise  $\chi^2$ -divergence.  
 672 Using Theorem 1 of [8] we get that  $P_W$  satisfies Poincaré inequality with constant  $C$  such  
 673 that  $C \leq K_P(\mu)(p^* + K_{\chi^2}^p(\mu))$  where  $p^*$  is the dual exponent of  $p$  satisfying  $1/p + 1/p^* = 1$ .

674 As an example, when base distributions are from the same exponential family and  
 675 the natural parameter space is affine, such as mixture of Poisson or Multinomial dis-  
 676 tributions, the pair-wise chi-squared distances are bounded (under some additional  
 677 conditions) and hence the mixture satisfies Poincaré inequality. More formally, let  
 678  $p_\theta(x) = \exp(T(x)^\top \theta - F(\theta) + k(x))$  where  $\theta \in \Theta$  is the natural parameter space and  $A(\theta)$   
 679 is the log partition function. Lemma 1 in [54] shows that

$$\chi^2(p_{\theta_1} || p_{\theta_2}) = e^{(A(2\theta_2 - \theta_1) - (2A(\theta_2) - A(\theta_1)))} - 1,$$

680 which is bounded as long as  $2\theta_2 - \theta_1 \in \Theta$ . This is satisfied for mixture of 1-D  
681 Poisson distributions which can be written as  $p(w|\lambda) = \frac{1}{w!} \exp(w \log \lambda - \lambda)$  with  
682 natural parameter space  $\mathbb{R}$ , and mixture of  $e$ -dimensional Multinomial distributions  
683  $p(w|\pi) = \exp\left(\langle w, \log\left(\frac{\pi}{1 - \sum_{i=1}^{e-1} \pi_i}\right)\rangle + \log\left(1 - \sum_{i=1}^{e-1} \pi_i\right)\right)$  with natural parameter  
684 space  $R^{e-1}$ . When applied to Gaussian family the natural parameters are

$$\theta_q = \begin{pmatrix} \Sigma_q^{-1} \mu_q \\ \text{vec}\left(-\frac{1}{2} \Sigma_q^{-1}\right) \end{pmatrix}.$$

685 Since the covariance has to be positive definite matrices,  $2\theta_q - \theta_{q'}$  may not be a set of valid  
686 natural parameter. We deal with this in the next case.

687 • When  $\{P_W^q\}_{q \in \mathcal{Q}}$  is a mixture of isotropic Gaussians, each with mean  $\mu_q \in \mathbb{R}^e$  and covariance  
688  $\Sigma_q = \sigma_q^2 I_e$ , each satisfying Poincaré inequality with constant  $C^q$  (in the single-Gaussian  
689 case above we know that  $C^q \leq \sigma_q^2$ ),  $P_W$  also satisfies Poincaré inequality. We prove this via  
690 induction. The key lemma is below:

691 **Lemma 8** (Corollary 1 of [64]). *Suppose measure  $p_0$  is absolutely continuous with respect  
692 to measure  $p_1$ , and  $p_0, p_1$  satisfy Poincaré inequality with constants  $C_0, C_1$  respectively, then  
693 for all  $\alpha \in [0, 1]$  and  $\beta = 1 - \alpha$ , mixture measure  $p = \alpha p_0 + \beta p_1$  satisfies Poincaré inequality  
694 with  $C \leq \max\{C_0, C_1(1 + \alpha \chi_1)\}$  where  $\chi_1 = \int \frac{dp_0}{dp_1} dp_0 - 1$ .*

695 We sort the components in the order of non-decreasing  $\sigma_q^2$ , and add in each component one  
696 by one. For each new component  $i = 2, \dots, |\mathcal{Q}|$ , we apply the above lemma with  $p_0$  being  
697 mixture of  $P_W^1, \dots, P_W^{i-1}$  and  $p_1$  being the new component  $P_W^i$ . We only need to prove that  
698  $\chi_1$  is bounded at every step. Suppose  $p_0 = \sum_{j=1}^{i-1} \alpha_j P_W^j$  with  $\sum_{j=1}^{i-1} \alpha_j = 1$ ,  $p_1 = P_W^i$ , and  
699  $P_W^j = \frac{1}{(2\pi)^{e/2} \sigma_j^e} \exp\left\{-\frac{1}{2}(w - \mu_j)^\top \Sigma_j^{-1} (w - \mu_j)\right\}$ . Therefore

$$\begin{aligned} \chi_1 + 1 &= \int \frac{dp_0}{dp_1} dp_0 = \int_w \frac{p_0(w)^2}{p_1(w)} dw \\ &= \int_w \frac{\sum_{j=1}^{i-1} \frac{\alpha_j^2}{\sigma_j^{2e}} \exp\left\{-\frac{\|w - \mu_j\|^2}{\sigma_j^2}\right\} + \sum_{j=1}^{i-1} \sum_{j' \neq j} \frac{2\alpha_j \alpha_{j'}}{\sigma_j^e \sigma_{j'}^e} \exp\left\{-\frac{\|w - \mu_j\|^2}{2\sigma_j^2} - \frac{\|w - \mu_{j'}\|^2}{2\sigma_{j'}^2}\right\}}{\frac{(2\pi)^{e/2}}{\sigma_i^e} \exp\left\{-\frac{\|w - \mu_i\|^2}{2\sigma_i^2}\right\}} dw \end{aligned}$$

700 The convergence condition of the above integral is  $2\sigma_i^2 \geq 2\sigma_j^2$  for all  $j < i$  which is satisfied  
701 when  $\sigma_i^2 \geq \sigma_j^2$ .

702

□

703 Next we observe that

$$\nabla_w \mathbb{E}[g|w] = \nabla_w \int_{x, \tilde{y}} (f(w, x) - \tilde{y})^2 p(x, \tilde{y}) dx d\tilde{y} = 2 \int_{x, \tilde{y}} (f(w, x) - \tilde{y}) \frac{\partial f}{\partial w} p(x, \tilde{y}) dx d\tilde{y} = 2 \mathbb{E}\left[(f(w, x) - \tilde{y}) \frac{\partial f}{\partial w}\right].$$

704 Since  $|f(w, x) - \tilde{y}| \leq 1 + \epsilon$  almost surely,  $\mathbb{E}\left[\left\|\frac{\partial f}{\partial w}\right\|_2^2\right] \leq \beta^2$ ,

$$\mathbb{E}_h\left[\|\nabla_w \mathbb{E}[g|w]\|_2^2\right] = 4 \mathbb{E}\left[\left\|(f(w, x) - y) \frac{\partial f}{\partial w}\right\|_2^2\right] \leq 4(1 + \epsilon)^2 \beta^2.$$

705 Therefore

$$\text{Var}_w[\mathbb{E}[g|w]] \leq C \mathbb{E}\left[\|\nabla_w \mathbb{E}[g|w]\|_2^2\right] \leq 4(1 + \epsilon)^2 C \beta^2.$$

706 To bound the second term of equation (10), we use concentration of mean of  $m$  i.i.d. random variables.

707 Conditioned on  $w^{(j)}$ , each of the loss  $(f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2$  are *i.i.d.* and bounded in  $[0, (1 + \epsilon)^2]$ .  
 708 Hence each variable has variance upper bound  $((1 + \epsilon)^2 - 0)^2/4 = (1 + \epsilon)^4/4$  and the mean has  
 709 variance upper bound  $(1 + \epsilon)^4/4m$ .

710 Therefore  $\text{Var}[g] \leq 4(1 + \epsilon)^2 C \beta^2 + (1 + \epsilon)^4/4m$ .  $\square$

*Proof of Theorem 4*

$$L(\hat{f}) - L(f^*) \leq 2 \sup_{f \in \mathcal{F}} |\tilde{L}(f) - \hat{L}(f)|$$

711

$$\leq 2 \mathbb{E} \sup_f |\tilde{L}(f) - \hat{L}(f)| + 4 \sqrt{\frac{\mathbb{E} \sup_f |\hat{L}(f) - \tilde{L}(f)| \log(1/\delta)}{n}} + \sqrt{\frac{(32(1 + \epsilon)^2 C \beta^2 + \frac{2(1 + \epsilon)^4}{m}) \log(1/\delta)}{n}} + \frac{2 \log(1/\delta)}{3n} \quad (12)$$

712 by equation 9 and Lemma 6

713 We now show that  $\mathbb{E} \sup_f |\tilde{L}(f) - \hat{L}(f)| \leq 2(1 + \epsilon) R_{nm}(F)$ . This is similar to the argument for  
 714 classical Rademacher complexity

$$\begin{aligned} & \mathbb{E}_{w,x,\tilde{y}} \sup_f \left( \frac{1}{nm} \sum_{i,j} (f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2 - \mathbb{E}_{w,x,\tilde{y}} (f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2 \right) \\ & \leq \frac{1}{nm} \mathbb{E}_{S,S'} \sup_f \left( \sum_{i,j} [(f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2 - (f(w'^{(j)}, x_i'^{(j)}) - \tilde{y}_i'^{(j)})^2] \right) \\ & = \frac{1}{nm} \mathbb{E}_{S,S',\sigma} \sup_f \left( \sum_{i,j} [\sigma_i^j (f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2 - \sigma_i^j (f(w'^{(j)}, x_i'^{(j)}) - \tilde{y}_i'^{(j)})^2] \right) \\ & \leq \frac{1}{nm} \mathbb{E}_{S,\sigma} \sup_f \left( \sum_{i,j} \sigma_i^j (f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2 \right) + \frac{1}{nm} \mathbb{E}_{S',\sigma} \sup_f \left( \sum_{i,j} \sigma_i^j (f(w'^{(j)}, x_i'^{(j)}) - \tilde{y}_i'^{(j)})^2 \right) \\ & = 2R_{nm}(\tilde{\mathcal{L}}). \end{aligned}$$

715 where the first inequality uses Jensen's inequality and convexity of *sup*.

716 Now we prove the equivalent of Talagrand's contraction lemma to show that  $R_{nm}(\tilde{\mathcal{L}}) \leq 2R_{nm}(F)$ .

717 Note that the squared loss is  $2(1 + \epsilon)$ -Lipschitz since  $\left| \frac{\partial (f - \tilde{y})^2}{\partial f} \right| = 2|f - \tilde{y}| \leq 2(1 + \epsilon)$ . We use the  
 718 following lemma to prove this:

719 **Lemma 9** (Lemma 5 of [49]). *Suppose  $\{\phi_i\}, \{\psi_i\}, i = 1, \dots, N$  are two sets of functions on  $\Theta$  such  
 720 that for each  $i$  an  $\theta, \theta' \in \Theta$ ,  $|\phi_i(\theta) - \phi_i(\theta')| \leq |\psi_i(\theta) - \psi_i(\theta')|$ . Then for all functions  $c: \Theta \rightarrow \mathbb{R}$ ,*

$$\mathbb{E}_\sigma \left[ \sup_\theta \left\{ c(\theta) + \sum_{i=1}^N \sigma_i \phi_i(\theta) \right\} \right] \leq \mathbb{E}_\sigma \left[ \sup_\theta \left\{ c(\theta) + \sum_{i=1}^N \sigma_i \psi_i(\theta) \right\} \right]$$

721 For any set of  $w, x$ , we apply Lemma 9 with  $\Theta = \mathcal{F}$ ,  $\theta = f$ ,  $N = nm$ ,  $\phi_{ij}(f) = (f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2$ ,  
 722  $\psi_{ij}(f) = 2(1 + \epsilon)f(w^{(j)}, x_i^{(j)})$ , and  $c(\theta) = 0$ . Since  $|(f - \tilde{y})^2 - (f' - \tilde{y})^2| \leq 2(1 + \epsilon)|f - f'|$ , so the  
 723 condition for Lemma 9 hold. We take expectation over  $w, x$  and divide both sides by  $nm$  to get

$$\frac{1}{nm} \mathbb{E}_{w,x,\sigma} \sup_f \sum_{j=1}^n \sum_{i=1}^m \sigma_i^j (f(w^{(j)}, x_i^{(j)}) - \tilde{y}_i^{(j)})^2 \leq \frac{2(1 + \epsilon)}{nm} \mathbb{E}_{w,x,\sigma} \sup_f \sum_{j=1}^n \sum_{i=1}^m \sigma_i^j f(w^{(j)}, x_i^{(j)})$$

724 which means  $R_{nm}(\mathcal{L}) \leq 2(1 + \epsilon) R_{nm}(F)$ . Substituting this into inequality 12 finishes the proof.

725  $\square$

726 **A.4 Zero-shot Rademacher complexity bound for the linear hypothesis class**

727 Consider the linear classifier  $F = \{(w_1^\top w + w_2^\top x : \|w_1\|_2 \leq B, \|w_2\|_2 \leq C)\}$ . Suppose  $\|w\|_2 \leq 1$  and  
 728  $\|x\|_2 \leq 1$ .

$$\begin{aligned}
 R_{nm}(F) &= \frac{1}{nm} \mathbb{E}_{\sigma, w, x} \sup_w \left\{ \langle w_1, \sum_{ij} \sigma_i^j w^{(j)} \rangle + \langle w_2, \sum_{ij} \sigma_i^j x_i^{(j)} \rangle \right\} \\
 &= \frac{1}{nm} \left( B_1 \mathbb{E}_{\sigma, w} \left\| \sum_{ij} \sigma_i^j w^{(j)} \right\|_2 + B_2 \mathbb{E}_{\sigma, x} \left\| \sum_{ij} \sigma_i^j x_i^{(j)} \right\|_2 \right) \\
 &\leq \frac{1}{nm} \left( B_1 \sqrt{m \sum_j \|w^{(j)}\|_2^2} + B_2 \sqrt{\sum_{ij} \|x_i^{(j)}\|_2^2} \right) \\
 &= (B_1 + B_2) / \sqrt{nm}.
 \end{aligned}$$

729 We observe that the bound is the same as the standard Rademacher complexity for  $nm$  independent  
 730 samples, which is interesting. The relationship between standard and zero-shot Rademacher  
 731 complexity for other function classes is an important future direction.

732 **B Extended Related Work**

733 Our approach to zero-shot prediction of intervention effects is related to recent advances in  
 734 heterogenous treatment effect (HTE) estimation, zero-shot learning, and meta-learning.

735 **B.1 Heterogenous treatment effect (HTE) estimation**

736 **Conditional average treatment effect (CATE) estimation.** A number of approaches have been  
 737 developed to predict the effect of an existing intervention on an individual or subgroup, based  
 738 on historical data from individuals who received it. This problem is often referred to in the  
 739 literature as heterogeneous treatment effect (HTE) estimation [26, 11], to denote that the goal is  
 740 to detect heterogeneities in how individuals respond to an intervention. A more specific instance  
 741 of HTE estimation, which we focus on here, is conditional average treatment effect (CATE)  
 742 estimation [76, 42], in which the goal is to predict the effect of a treatment *conditioned* on an  
 743 individual’s features. A variety of methods and specific models have been developed to achieve this  
 744 goal [26, 32, 21, 28, 76, 66, 1, 84, 24, 86, 25, 14, 12, 42, 34, 11, 2], and we refer to Bica et al. and Curth  
 745 et al. for a detailed review of these methods [4, 14]. These methods estimate CATE for an existing  
 746 intervention, based on historical data from individuals who received it and those that did not.

747 While these approaches have a number of useful applications, they do not address CATE for  
 748 novel interventions which did not exist during training (zero-shot). Our primary contribution is a  
 749 meta-learning framework to leverage these existing CATE estimators for zero-shot predictions. In  
 750 the CaML framework (Figure 2), each task corresponds to predicting CATE for a single intervention.  
 751 We synthesize a task by sampling a natural experiment for each intervention, and then use any existing  
 752 CATE estimator to generate a noisy target label for our the task (Step 2: estimate pseudo-outcomes).  
 753 We rely on pseudo-outcome estimates as training labels because prior work has shown that training  
 754 on observed outcomes directly leads to biased CATE estimates [9, 42, 34], a result which we find  
 755 holds true in our experiments as well (see T-learner and S-learner w/ meta-learning in Tables 2 and 3).

756 **Pseudo-outcome estimators.** Prior work has developed a variety of methods to estimate CATE pseudo-  
 757 outcomes, which are noisy but unbiased estimates of CATE, such as the X-learner [42], R-learner [53],  
 758 DR-learner [34], and RA-learner [14]. Moreover, the outputs of any other CATE estimation method,  
 759 such as methods which directly estimate CATE via an end-to-end neural network [32, 66, 68] are an  
 760 equally valid choice of pseudo-outcome. The literature on pseudo-outcome estimation is growing con-  
 761 tinuously as new estimators are being developed [19, 38]. Typically, these estimators are specific to a *sin-*  
 762 *gle binary intervention*, for which a set of nuisance models are trained and used to compute the pseudo-  
 763 outcomes. As such, applying meta-learning algorithms to these pseudo-outcomes requires synthesizing  
 764 a natural experiment for each intervention, which corresponds to a single task in the CaML framework.

765 **Multi-cause estimators.** Our methods to address zero-shot CATE estimation for combinations  
766 of interventions are distinct from multi-cause estimators for combinations of binary or categorical  
767 interventions [78, 58, 62]. Recent work has shown that these methods can predict the effects of new  
768 combinations of interventions [48], when every intervention in the combination has been observed  
769 at some point during. However, these methods do not estimate CATE for novel interventions which did  
770 not exist during training. By contrast, CaML estimates CATE for zero-shot intervention combinations  
771 in which none of the interventions in the combo was ever observed during training (Appendix Table C).

## 772 B.2 Zero-shot learning

773 Zero-shot learning (ZSL) has traditionally aimed to reason over new concepts and classes [80, 60]  
774 which did not exist during training time. While ZSL has primarily focused on natural language  
775 processing and computer vision [77], recent interest has been sparked in generalizing over novel  
776 interventions (zero-shot) in the biomedical domain [61, 27] in which data can be cheaply collected  
777 for hundreds or thousands of possible interventions [87, 71, 17]. However, general-purpose zero-shot  
778 causal methods have been largely unexplored. Notable exceptions include GranITE [23] and SIN [23],  
779 which each extend a specific CATE estimation [53, 42] method to incorporate intervention features  
780 ( $W$ ). However, these approaches have significant drawbacks, which we discuss in Section 2.

## 781 B.3 Meta-learning

782 Meta-learning, or *learning to learn*, aims to train models which can quickly adapt to new settings  
783 and tasks. The key idea is to enable a model to gain experience over multiple learning episodes - in  
784 which episodes typically correspond to distinct tasks - to accelerate learning in subsequent learning  
785 episodes [30]. The meta-learning literature is rich and spans multiple decades [72, 65, 63, 3], with  
786 recent interest focused on model-agnostic methods to train deep learning models to quickly adapt to  
787 new tasks [18, 59, 52]. A common focus in the meta-learning literature is few-shot learning, in which  
788 a model must adapt to a new task given a small support set of labeled examples. By contrast, we focus  
789 on the zero-shot setting, in which no such support set exists. However, we hypothesize that the typical  
790 meta-learning problem formulation and training algorithms may also improve zero-shot performance.  
791 Thus, CaML’s problem formulation and algorithm inspiration from the meta-learning literature,  
792 particularly the Reptile algorithm [52] and its application to other tasks in causal inference [67].  
793 Our experimental results show that this meta-learning formulation improves CaML’s performance,  
794 compared to a standard multi-task learning strategy.

## 795 C Experimental details

### 796 C.1 Experimental setup

797 Here, we provide more details about the experimental setup for each investigated setting. This serves  
798 to complement the high-level overview given in Table I. Experiments were run using Google Cloud  
799 Services. Deep learning-based methods (*i.e.*, CaML and its ablations, S-learner w/ meta-learning,  
800 T-learner w/ meta-learning, SIN, GraphITE, FlexTENET, TARNET, and DragonNet) were run on  
801 n1-highmem-64 machines with 4x NVIDIA T4 GPU devices. The remaining baselines (RA-learner,  
802 R-learner, X-learner, and T-learner) were run on n1-highmem-64 machines featuring 64 CPUs.

803 **Fair comparison.** We perform hyper-parameter optimization with random search for all models, with  
804 the meta-testing dataset predetermined and held out. To avoid “hyperparameter hacking”, hyperparam-  
805 eters ranges are consistent between methods wherever possible, and were chosen using defaults similar  
806 to prior work [33, 23]. Choice of final model hyper-parameters was determined using performance met-  
807 rics (specific to each dataset) computed on the meta-validation dataset, using the best hyper-parameters  
808 over 48 runs (6 servers x 4 NVIDIA T4 GPUs per server x 2 runs per GPU ) (Appendix C.4). All table  
809 results are computed as the mean across 8 runs of the final model with distinct random seeds.

#### 810 C.1.1 Claims dataset

811 **Interventions ( $W$ ):** We consider drug prescriptions consisting of either one drug, or two drugs  
812 prescribed in combination. We observed 745 unique single drugs, and 22,883 unique drug pairs,  
813 excluding interventions which occurred less than 500 times. Time of intervention corresponds to the  
814 *first* day of exposure. To obtain intervention information, we generated pre-trained drug embeddings

815 from a large-scale biomedical knowledge graph [7] (see Appendix C.5). Drugs correspond to nodes  
 816 in the knowledge graph, which are linked to other nodes (e.g. genes, based on the protein target of  
 817 the drug). Drug combination embeddings are the sum of the embeddings for their constituent drugs.

818 **Control group.** A challenge in such causal analyses of clinical settings is defining a control group. We  
 819 randomly sample 5% (1.52M patients) to use as controls, with a 40/20/40 split between meta-train/meta-  
 820 val/meta-test. When sampling a natural experiment for a given intervention, we select all patients from  
 821 this control group that did not receive such an intervention. An additional challenge is defining time  
 822 of intervention for the control group. It is not possible to naively sample a random date, because there  
 823 are large quiet periods in the claims dataset in which no data is logged. We thus sample a date in which  
 824 the control patient received a *random drug*, and thus our measure of CATE estimates the *increase* in  
 825 side effect likelihood from the drug(s)  $W$ , compared to another drug intervention chosen at random.

826 **Outcome ( $Y$ ):** We focus on the side effect pancytopenia: a deficiency across all three blood cell  
 827 lines (red blood cells, white blood cells, and platelets). Pancytopenia is life-threatening, with a 10-20%  
 828 mortality rate [36, 41], and is a rare side effect of many common medications [40] (e.g. arthritis and  
 829 cancer drugs), which in turn require intensive monitoring of the blood work. Our outcome is defined  
 830 as the (binary) occurrence of pancytopenia within 90 days of intervention exposure.

831 **Features ( $X$ ):** Following prior work [22], patient medical history features were constructed by  
 832 time-binned counts of each unique medical code (diagnosis, procedure, lab result, drug prescription)  
 833 before the drug was prescribed. In total, 443,940 features were generated from the following time bins:  
 834 0-24 hours, 24-48 hours, 2-7 days, 8-30 days, and 31-90 days, 91-365 days, and 365+ days prior. All  
 835 individuals in the dataset provided by the insurance company had at least 50 unique days of claims data.

836 **Metrics:** We rely on best practices for evaluating CATE estimators as established established by recent  
 837 work [81, 10], which recommend to assess treatment rules by comparing subgroups across different  
 838 quantiles of estimated CATE. We follow the high vs. others RATE (rank-weighted average treatment  
 839 effect) approach from Yadlowsky et. al [81], which computes the difference in average treatment effect  
 840 (ATE) of the top  $u$  percent of individuals (ranked by predicted CATE), versus all individuals:

$$RATE@u = \mathbb{E} \left[ Y(1) - Y(0) \mid F_S(S(X)) \geq 1 - u \right] - \mathbb{E} \left[ Y(1) - Y(0) \right], \quad (13)$$

841 where  $S(\cdot)$  is a priority score which ranks samples lowest to highest predicted CATE, and  $F_S(\cdot)$  is the  
 842 cumulative distribution function (CDF) of  $S(X_i)$ . For instance, RATE @ 0.99 would be the difference  
 843 between the top 1% of the samples (by estimated CATE) vs. the average treatment effect (ATE) across  
 844 all samples, which we would expect to be high if the CATE estimator is accurate. The real-world use  
 845 case of our model would be preventing drug prescription a small subset of high-risk individuals. Thus,  
 846 more specifically, for each task  $j$ , intervention  $w_j$  in the meta-dataset, and meta-model  $\Psi_\theta$  (our priority  
 847 score  $S(\cdot)$ ), we compute  $RATE@u$  for each  $u$  in  $[0.999, 0.998, 0.995, 0.99]$  across individuals who  
 848 received the intervention.

849 We now summarize how to estimate RATE performance metrics for a single intervention (task).  
 850 As RATE performance is calculated separately per-intervention we are concerned with a single  
 851 intervention, we use the simplified notation (i.e.  $Y_i(1)$  instead of  $Y_i(w)$ ) from Section 3. Due to the  
 852 fundamental problem of causal inference (we can only observe  $Y_i(0)$  or  $Y_i(1)$  for a given sample),  
 853 the true RATE, as defined above, cannot be directly observed.

854 We follow the method outlined in Section 2.2 and 2.4 of Yadlowsky et. al. [81] in which we compute  
 855  $\hat{\Gamma}_i$ , a (noisy but unbiased) estimate for CATE which is in turn used to estimate RATE:

$$\mathbb{E} \left[ \hat{\Gamma}_i \mid X_i \right] \approx \tau(X_i) = \mathbb{E} \left[ Y_i(1) - Y_i(0) \mid X_i \right]. \quad (14)$$

856 Our data is observational, and as such we can estimate  $\hat{\Gamma}_i$  using a direct non-parametric estimator [75]:

$$\hat{\Gamma}_i = W_i(Y_i - \hat{m}(X_i, 0)) + (1 - W_i)(\hat{m}(X_i, 1) - Y_i) \quad (15)$$

$$m(x, w) = \mathbb{E}[Y_i(w) \mid X_i = x] \quad (16)$$

857 where  $m(x, w)$  is a model that predicts the outcome. Here  $\hat{m}(x, w)$  represent nonparametric estimates  
 858 of  $m(x, w)$ , respectively, which we obtain by fitting a cross-fitting a model to the intervention natural  
 859 experiment over 5-folds. We use random forest models for  $\hat{m}(x, w)$ , as they perform well (achieving  
 860  $\geq 0.90$  ROC AUC across all meta-testing tasks for predicting outcomes) and are robust to choice of  
 861 hyperparameters.

862 RATE can then be estimated via sample-averaging estimator. Specifically, we compute the difference  
863 between the average value of  $\hat{\Gamma}_i$  for those in the top  $u$  percent of individuals (based on our meta-model’s  
864 predictions), compared to the average  $\hat{\Gamma}_i$  across all individuals. For further discussion on estimating  
865 RATE, we refer readers to [81]. Note that estimates of RATE are *unbounded*: RATE can be less than  
866 0 (due to predictions inversely relating to CATE).

867 Finally, because our meta-testing dataset consists of individuals treated with drugs *known* in the  
868 medical literature to cause pancytopenia (identified by filtering drugs using the side effect database  
869 SIDER [40]), observational metrics of recall and precision are also a rough *proxy* for successful  
870 CATE estimation. Thus, as secondary metrics, we also compute *Recall@u* and *Precision@u* for  
871 the same set of thresholds as RATE, where a positive label is defined as occurrence of pancytopenia  
872 after intervention. We find that these metrics are highly correlated to RATE in our performance results.

873 **Training & Evaluation:** For each method, we ran a hyperparameter search with 48 random  
874 configurations (48 due to running 8 jobs in parallel on 6 servers each) that were drawn uniformly  
875 from a pre-defined hyperparameter search space (see Appendix C.4). Methods that can be trained on  
876 multiple tasks to then be applied to tasks unseen during training (*i.e.*, CaML and its ablations, S-learner  
877 w/ meta-learning, T-learner w/ meta-learning, SIN, GraphITE) were trained for 24 hours (per run) on  
878 the meta-training tasks. Model selection was performed on the meta-validation tasks by maximizing  
879 the mean RATE@0.998 across meta-validation tasks. Then, the best hyperparameter configuration  
880 was used to fit 8 repetition runs across 8 different random seeds. Each repetition model was then tested  
881 on the meta-testing tasks, where for all metrics averages across the testing tasks are reported. To make  
882 the setting of multi-task models comparable with single-task models that were trained on meta-testing  
883 tasks (requiring a train and test split of each meta-testing task), the evaluation of all models was  
884 computed on the test split of the meta-testing tasks, respectively. Single-task baselines (FlexTENET,  
885 TARNet, and DragonNet, RA-learner, R-learner, X-learner, and T-learner) were given access to the  
886 meta-testing tasks during training. Specifically, model selection was performed on the meta-validation  
887 tasks, while the best hyperparameter configuration was used to train 8 repetition models (using 8  
888 random seeds) on the train split of each meta-testing task. For the final evaluation, each single-task  
889 model that was fit on meta-testing task  $i$  was tested on the test split of the same meta-testing task  $i$ ,  
890 and the average metrics were reported across meta-testing tasks.

## 891 C.1.2 LINCS

892 **Interventions ( $W$ ):** Interventions in the LINCS dataset consist of a single perturbagen (small  
893 molecule). For intervention information, we used the molecular embeddings for each perturbagen  
894 using the RDKit featurizer. The same cell line-perturbagen combinations are tested with different  
895 perturbagen dosages and times of exposure. [44]. To maintain consistency in experimental conditions  
896 while also ensuring that the dataset is sufficiently large for training a model, we filter for most  
897 frequently occurring dosage and time of exposure in the dataset, which are  $10\mu M$  and 24 hours,  
898 respectively. We use data from 10,322 different perturbagens.

899 **Control group.** For each perturbagen (at a given timepoint and dose), we use cell lines which did  
900 not receive that intervention as the control group.

901 **Outcomes ( $Y$ ):** We measure gene expression across the top-50 and top-20 landmark differentially  
902 expressed genes (DEGs) in the LINCS dataset. Accurately predicting in gene expression in these  
903 DEGs is most crucial to the drug discovery process.

904 **Features ( $X$ ):** We use 19,221 features from the Cancer Cell Line Encyclopedia (CCLE) [20] to  
905 describe each cell-line, based on historical gene expression values in a different lab environment. Our  
906 dataset consisted of 99 unique cell lines (after filtering for cell-lines with CCLE features).

907 **Metrics:** A key advantage of experiments on cells is that at evaluation time we can observe both  $Y(0)$   
908 and  $Y(1)$  for the same cell line  $X$ , through multiple experiments on clones of the same cell-line in  
909 controlled lab conditions. In the LINCS dataset,  $Y(0)$  is also measured for all cells which received  
910 an intervention. Thus, we can directly compute the Precision Estimation of Heterogenous Effects  
911 (PEHE) on all treated cells in our meta-testing dataset. PEHE is a standard metric for CATE estimation  
912 performance [28], analogous to mean squared error (MSE).

$$PEHE = \frac{1}{N} \sum_{i=1}^N (\tau_i - \hat{\tau}_i)^2 \quad (17)$$

913 **Training & Evaluation:** For each method, we ran a hyperparameter search with 48 random  
914 configurations (48 due to running 8 jobs in parallel on 6 servers each) that were drawn uniformly  
915 from a pre-defined hyperparameter search space (see Appendix C.4). Methods that can be trained  
916 on multiple tasks to then be applied to tasks unseen during training (*i.e.*, CaML and its ablations,  
917 S-learner w/ meta-learning, T-learner w/ meta-learning, SIN) were trained for 12 hours (per run) on  
918 the meta-training tasks. Model selection was performed on the meta-validation tasks by minimizing  
919 the overall PEHE for the Top-20 most differentially expressed genes (DEGs) across meta-validation  
920 tasks. Then, the best hyperparameter configuration was used to fit 8 repetition runs across 8 different  
921 random seeds. Each repetition model was then tested on the meta-testing tasks, where for all metrics  
922 averages across the testing tasks are reported.

## 923 C.2 Selecting holdout interventions for meta-validation and meta-testing

### 924 C.2.1 Claims.

925 In the 30.4 million patient insurance claims dataset, each intervention task in meta-train/meta-  
926 val/meta-testing corresponds to a natural experiment of multiple patients, with some interventions (*e.g.*  
927 commonly prescribed drugs) having millions of associated patients who were prescribed the drug. One  
928 challenge is that in this setting, there is overlap in subjects between the natural experiments sampled by  
929 CaML, which can lead to data leakage between training and testing. For instance, if a patient received  
930 Drug 1 (in meta-test) and Drug 2 (meta-train), they would appear in both natural experiments, resulting  
931 in data leakage.

932 We take a conservative approach and exclude all patients who have ever received a meta-testing drug in  
933 their lifespan from the natural experiments for meta-val/meta-train. Similarly, we exclude all patients  
934 who received a meta-validation drug from meta-training.

935 This approach means we must take great care in selecting meta-testing drugs. Specifically, we must  
936 trade off between selecting drugs that are important (covering enough patients) while not diminishing  
937 the training dataset size. For instance selecting a commonly prescribed (*e.g.* aspirin) for meta-testing  
938 would deplete our meta-training dataset by over 50% of patients. Thus we only selected meta-test/meta-  
939 validation drugs which were prescribed to between 1,000,000 and 100K patients in our dataset, after  
940 filtering for only drugs which known to cause Pancytopenia [40] (using the SIDER database). From  
941 this subset of drugs, we randomly selected 10 meta-testing drugs and 2 meta-validation drugs, resulting  
942 in a total meta-testing/meta-validation pool of 4.1 million patients and 685K patients respectively.

943 To evaluate on unseen pairs of drugs on the same hold-out test dataset, we additionally created a second  
944 pairs testing dataset from the 5 most frequently occurring combinations from the meta-testing dataset.  
945 This allowed us to train a single model on the same meta-train split and evaluate on both single drug  
946 and drug pair interventions without occurrence of data leakage. Designing a larger evaluation of  
947 pairs was not possible because while pairs of drugs are commonly prescribed as intervention, each  
948 particular pair of drugs is a rare event, and accurately evaluating CATE estimation performance (for  
949 a rare outcome such as Pancytopenia) requires amassing a natural experiment with at least several  
950 thousand patients who received the same intervention.

### 951 C.2.2 LINCS.

952 The goal in selecting holdout interventions for the meta-validation and meta-testing sets was to ensure  
953 that they consisted of both cell lines and tasks (small molecules) that had not been seen previously  
954 at the time of training (*i.e.* zero-shot on cell lines and tasks).

955 Using a random data splitting approach would result in large portions (up to 50%) of the data being  
956 unused to comply with the zero-shot requirements on cell lines and tasks. One approach to tackle  
957 this was to reserve only those tasks in the held-out sets which had been tested on the fewest cell lines.  
958 This preserved the maximum amount of data but resulted in an average of just 1 cell line per task in  
959 the meta-testing and meta-validation sets, which would not be fair to the non-zero shot baselines.

960 To address these issues, we designed a new data split procedure that exploits the structure of how  
961 tasks and cell lines are paired. To do so, We first clustered tasks by the cell lines they are tested on.  
962 We then identified a set of 600 drugs that had all been tested on a shared set of roughly 20 cell lines.  
963 We divided the cell lines and tasks within this set into the meta-validation and meta-testing set, while  
964 enforcing zero-shot constraints on both. This resulted in roughly 10 cell lines per intervention in both

965 the meta-validation and meta-testing sets, while still maintaining a reasonably large size of 11 distinct  
966 cell lines and 300 distinct tasks in both sets. All remaining tasks and cell lines were reserved for the  
967 training set. (See Table 8)

### 968 C.3 Understanding CaML’s performance

969 Our comparison to CATE estimators which are restricted to single interventions (Grey, Table 2.B.3)  
970 shows that a key reason for CaML’s strong performance is the ability to jointly learn across from many  
971 intervention datasets, in order to generalize to unseen intervention.

972 Additionally, in both the Claims and LINCS settings, we conduct two key ablation studies to further  
973 understand the underlying reason for CaML’s strong performance results.

974 In our first ablation experiment (w/o meta-learning), we trained the CaML model without employing  
975 meta-learning, instead using the standard empirical risk minimization (ERM) technique [73]. This can  
976 be seen as a specific implementation of the CaML algorithm (refer to Algorithm 1) when  $k = 1$  [52].  
977 The results of this experiment showed a varying degree of performance deterioration across our  
978 primary tests. In the Claims settings, we observed a decrease in the RATE performance metric by  
979 15%-22% (refer to Table 2), while in the LINCS settings, the PEHE performance metric decreased  
980 by approximately 0.01 (see Table 3). These results indicate that the absence of meta-learning affects  
981 the model’s performance, although the impact varies depending on the specific setting. An important  
982 detail to consider is that the Claims data experiments dealt with substantially larger datasets, each  
983 comprising hundreds of thousands of patients per intervention. This extensive scale of data potentially  
984 amplifies the benefits of using meta-learning in the CaML model for the Claims dataset. The larger  
985 dataset enables the model to adapt to a given task over a larger set of iterations without reusing the  
986 same data, thereby enhancing the efficacy of meta-learning.

987 Our second ablation (w/o RA-learner) assesses the sensitivity of CaML’s performance to different  
988 pseudo-outcome estimation strategies. A key aspect of CaML is *flexibility* in choice of any pseudo-  
989 outcome estimator to infer CATE, in contrast to prior work which uses specific CATE estimation strate-  
990 gies [23, 33]. We find that CaML performance benefits strongly from flexibility of pseudo-outcome esti-  
991 mator choice. We assess this by using an alternative pseudo-outcome estimator. Firstly, we find that this  
992 ablation results in much noisier model training. For instance, the standard deviation in RATE across the  
993 8 random seeds increases by  $20\times$  when using the alternative pseudo-outcome estimator in the claims set-  
994 ting. Moreover, the alternative pseudo-outcome estimator typically worsens performance, decreasing  
995 RATE by up to 6% in the Claims setting, and increasing PEHE by 20%-21% in the LINCS setting (Table  
996 3). We note that this ablation performs slightly better at the 0.99 threshold, which may be a result of the  
997 high variance in this ablation. Specific choice of alternative pseudo-outcome estimator for this ablation  
998 varies by setting. We use the R-learner [53] for Claims as it also achieves strong single task performance  
999 (Table 2, grey) on Claims data. However, R-learner is restricted to single-dimensional outcomes, and  
1000 thus for LINCS (in which outcomes are 50 and 20 dimensional), we use the PW-learner instead [14].

### 1001 C.4 Hyperparameter space

#### 1002 C.4.1 Claims dataset hyperparameter space

1003 We list the hyperparameter search spaces for the medical claims dataset in the following tables. Table 9  
1004 represents the search space for CaML. The SIN baseline consists of two stages, Stage 1 and Stage 2. For  
1005 the Stage 1 model, we searched the identical hyperparameter search space as for CaML (Table 9). For  
1006 Stage 2, we used the hyperparameters displayed in Table 10. The search space for the GraphITE baseline  
1007 is displayed in Table 11. For the S-learner and T-learner w/ meta-learning baselines, we use the same  
1008 hyperparameter space as for CaML (Table 9) with the only major difference that these baselines  
1009 predict the outcome  $Y$  instead of  $\hat{\tau}$ . For all deep learning-based methods, we employed a batch size of  
1010 8,192, except for GraphITE, where we were restricted to using a batch size of 512 due to larger memory  
1011 requirements. Single-task neural network baselines (FlexTENet, TARNet, and DragonNet) are shown  
1012 in Tables 12, 13, and 14, respectively. For the remaining baselines, i.e., the model-agnostic CATE estima-  
1013 tors, the (shared) hyperparameter search space is shown in Table 15. Finally, applied L1 regularization to  
1014 the encoder layer of the customizable neural network models (that were not reused as external packages),  
1015 i.e., SIN learner, GraphITE, T-learner w/ meta-learning, and S-learner w/ meta-learning, and CaML.

## 1016 C.4.2 LINCS hyperparameter space

1017 We list the hyperparameter search spaces for LINCS in the following tables. CaML is shown in Table 16.  
1018 SIN Stage 1 used the same search space as CaML (Table 16). The search space of SIN Stage 2 is shown  
1019 in Table 17. S learner and T-learner w/ meta-learning used the same search space as CaML. The search  
1020 space of GraphITE is shown in Table 18. All methods that were applied to LINCS used a batch size of 20.

## 1021 C.5 More details on intervention information

1022 Here we give more details about the intervention information used for the medical claims dataset.  
1023 In order to perform zero-shot generalization, we acquired information about a specific intervention  
1024 through the use of pretrained embeddings. We generated these embeddings on the Precision Medicine  
1025 Knowledge Graph [7] that contains drug nodes as well as 9 other node types. We extracted embeddings  
1026 for 7957 drugs from the knowledge graph.

1027 To extract rich neighborhood information from the knowledge graph we used Stargraph [47], which  
1028 is a coarse-to-fine representation learning algorithm. StarGraph generates a subgraph for each node  
1029 by sampling from its neighbor nodes (all nodes in the one-hop neighborhood) and anchor nodes (a  
1030 preselected subset of nodes appearing in the multihop neighborhood). In our case the anchor nodes  
1031 were the 2% of graph nodes with the highest degree. For the scoring function we used the augmented  
1032 version of TripleRE [85] presented in the StarGraph article [47].

1033 We performed a hyperparameter optimization to compare different models and determine the one  
1034 we used to calculate our final embeddings (see Table C.5). The hyperparameter search was random  
1035 with the objective of minimizing the loss function used in training on held out data. The search range  
1036 for each of the parameters is displayed in C.5. Since certain parameters did not seem to influence the  
1037 final score as much we decided to use them as constants and focus on optimizing the hyperparameters  
1038 in the table. Therefore the number of sampled anchors was set to 20 and  $u = 0.1$  in the augmented  
1039 TripleRE function, the values matching those seen in Stargraph [46].

1040 Our final embeddings were 256-dimensional, the learning rate was  $2e-4$ , the drop-ratio was  $5e-3$ . We  
1041 used the self-adversarial negative sampling loss with  $\gamma = 8$  and we sampled 4 neighbor nodes for each  
1042 subgraph.

1043 To additionally evaluate the quality of the embeddings we assigned classes to drug combinations  
1044 and then scored them using multiple clustering metrics. We were interested to see if embeddings  
1045 of drug combinations used for similar purposes would be embedded closer together than other drug  
1046 combinations. For the class label of single drugs we used the first level of the Anatomical Therapeutic  
1047 Chemical (ATC) code, which represents one of the 14 anatomical or pharmacological groups. Since  
1048 certain medications have more than one ATC code, we took the mode of all labels for a specific drug.  
1049 For multiple drugs we combined all distinct first level values and took the mode of them as the label. We  
1050 used the Silhouette metric, Calinski Harabasz index and Davies Bouldin index as well as the average  
1051 classification accuracy over 10 runs of training a random forest classifier on a random sample of 80%  
1052 of the dataset and evaluating on the remaining 20%. Out of all tested embeddings the hyperparameter  
1053 optimized StarGraph embeddings performed best (exceeding 93% in the classification accuracy metric).

## 1054 C.6 Pseudo-outcome estimation

1055 In our experiments, we estimate pseudo-outcomes  $\tilde{\tau}$  for a given intervention  $w$  using the  
1056 RA-learner [14]:

$$\tilde{\tau} = W(Y - \hat{\mu}_0(X)) + (1 - W)(\hat{\mu}_1(X) - Y) \quad (18)$$

1057 where  $\hat{\mu}_w$  is an estimate of  $\mu_w(X) = \mathbb{E}_{\mathcal{P}}[Y | X = x, W = w]$ .

1058 Furthermore, in both settings we only estimate CATE for treated individuals. We focus on treated  
1059 individuals in the Claims setting because we care about the risk of an adverse event for prescribing  
1060 a sick patients drugs that may cure their sickness, not the adverse event risk of prescribing healthy  
1061 patients drugs (which is of less clinical interest). In the LINCS setting, we focus on treated cells as  
1062 for these cell-lines  $Y(0)$  is also measured from a cloned cell-line under similar laboratory conditions,  
1063 which allows us to directly estimate CATE prediction performance using the PEHE metric. As we  
1064 focus on treated samples, the RA-learner can be simplified to  $\tilde{\tau} = Y - \hat{\mu}_0(X)$ . We estimate  $\hat{\mu}_0(X)$

1065 using a random forest model in the Claims setting, whereas in the LINCS setting we use the point  
1066 estimate from the untreated control cell line’s gene expression.

## 1067 C.7 Baselines

1068 Here we provide more details on the baselines used in our experiments.

1069 *Trained on test task:* These baselines leverage CATE estimators which can only be trained on a single  
1070 task (typically these are the strongest baselines, when there is a large enough dataset for a single task).  
1071 Thus, we train a single model for each meta-testing task on its train split, and evaluate performance  
1072 on its test split. We use a number of strong baselines for CATE estimation developed by prior work  
1073 including both model-agnostic and end-to-end deep learning approaches: T-learner. Specifically,  
1074 we use the model-agnostic CATE estimators: [42], X-learner [42], RA-learner [14], R-learner [53].  
1075 We additionally use the end-to-end deep learning estimators DragonNet [68], TARNet [66], and  
1076 FlexTENet [15], using implementations from [15]. For model-agnostic CATE estimators, we use  
1077 random forest models following prior work [12, 76].

1078 *Zero-shot.* These baselines use CATE estimators which incorporate intervention information ( $W$ ) and  
1079 are capable of multi-task learning. We train these baselines on all meta-training tasks. These baselines  
1080 have no access to the meta-testing tasks during training. We found in preliminary experiments that  
1081 in some cases, baseline models trained with vanilla ERM would not even converge. To allow for  
1082 fair comparison to baselines, we allow for all zero-shot baselines to be trained using Reptile (by  
1083 training using the same optimization strategy as Algorithm 1, while allowing for training with ERM  
1084 by including  $k = 1$  in the hyperparameter search space).

1085 Firstly, we use GraphITE [23] and Structured Intervention Networks [33]. These are, to the best  
1086 of our knowledge, the only methods from prior work which are (in principle) capable of zero-shot  
1087 generalization. We use existing implementations provided by the authors [33].

1088 Additionally, we implement two strong baselines which estimate CATE by modeling potential  
1089 outcomes, rather than via pseudo-outcomes. These are variants of the S-learner and T-learner [42] with  
1090 meta-learning, which use the intervention information as input, rather than one-hot encoded vectors of  
1091 the different interventions—such that they also have zero-shot capability. Specifically, we train MLPs  
1092 using the same architecture as CaML to estimate the response function from observed outcomes:

$$\mu(x, w) = \mathbb{E}_{\mathcal{P}} \left[ Y \mid X = x, W = w \right] \quad (19)$$

1093 and estimate CATE by

$$\hat{\tau}_w(x) = \hat{\mu}(x, w) - \hat{\mu}(x, \mathbf{0}) \quad (20)$$

1094 Where  $w$  denotes the corresponding intervention information  $w$  for an intervention, and  $\mathbf{0}$  denotes  
1095 a null intervention vector. In the LINCS setting, we represent  $\mathbf{0}$  as a vector of zeros, whereas in the  
1096 Claims setting we represent  $\mathbf{0}$  as the mean embedding of all drugs (as the estimand is the increase in  
1097 adverse event likelihood compared to a randomly chosen drug). The difference between the T-learner  
1098 and the S-learner is that the T-learner estimates two models, one for control units and one for treated  
1099 units. By contrast, the S-learner estimates a shared model across all units.

	RATE @ $u$ ( $\uparrow$ )			Recall @ $u$ ( $\uparrow$ )			Precision @ $u$ ( $\uparrow$ )		
	0.999	.998	0.995	0.999	0.998	0.995	0.999	0.998	0.995
Random	0.00±<0.001	0.00±<0.001	0.00±<0.001	0.00±<0.001	0.00±<0.001	0.01±<0.001	0.00±<0.001	0.00±<0.001	0.00±<0.001
T-learner	0.32±<0.001	0.26±<0.001	0.16±<0.001	0.10±<0.001	0.18±<0.001	0.26±<0.001	0.12±<0.001	0.29±<0.001	0.18±<0.001
X-learner	0.06±<0.001	0.05±<0.001	0.04±<0.001	0.03±<0.001	0.04±<0.001	0.08±<0.001	0.12±<0.001	0.07±<0.001	0.06±<0.001
R-learner	0.19±<0.001	0.17±<0.001	0.12±<0.001	0.08±<0.001	0.10±<0.001	0.19±<0.001	0.26±<0.001	0.21±<0.001	0.15±<0.001
RA-learner	0.47±0.001	0.37±<0.001	0.23±<0.001	0.14±<0.001	0.17±<0.001	0.38±<0.001	0.45±<0.001	0.42±<0.001	0.26±<0.001
DragonNet	0.09±0.037	0.07±0.030	0.05±0.019	0.04±0.013	0.02±0.008	0.04±0.012	0.10±0.027	0.10±0.045	0.07±0.023
TARNet	0.15±0.011	0.12±0.011	0.07±0.006	0.05±0.004	0.05±0.003	0.08±0.006	0.12±0.008	0.18±0.013	0.15±0.012
FlexTENet	0.10±0.015	0.09±0.016	0.06±0.008	0.04±0.006	0.04±0.006	0.07±0.009	0.17±0.017	0.12±0.018	0.08±0.010
GraphITE	0.19±0.024	0.12±0.013	0.05±0.004	0.03±0.002	0.07±0.009	0.08±0.010	0.10±0.008	0.23±0.027	0.14±0.015
SIN	0.00±0.002	0.00±0.001	0.00±0.001	0.00±0.001	0.00±0.001	0.00±0.001	0.02±0.002	0.01±0.002	0.01±0.001
S-learner w/ meta-learning	0.21±0.032	0.16±0.028	0.09±0.020	0.05±0.012	0.08±0.013	0.11±0.022	0.15±0.035	0.25±0.034	0.18±0.031
T-learner w/ meta-learning	0.40±0.012	0.31±0.010	0.18±0.007	0.11±0.004	0.15±0.006	0.22±0.008	0.32±0.013	0.45±0.013	0.35±0.011
CaML - w/o meta-learning	0.39±0.012	0.31±0.006	0.18±0.008	0.11±0.006	0.15±0.005	0.22±0.007	0.32±0.014	0.45±0.010	0.35±0.006
CaML - w/o RA-learner	0.45±0.058	0.36±0.066	0.22±0.067	0.14±0.041	0.16±0.020	0.24±0.019	0.35±0.016	0.51±0.076	0.41±0.082
CaML (ours)	<b>0.48±0.010</b>	<b>0.38±0.007</b>	<b>0.23±0.003</b>	0.13±0.002	<b>0.18±0.004</b>	<b>0.27±0.005</b>	<b>0.38±0.006</b>	<b>0.54±0.012</b>	<b>0.43±0.008</b>
									<b>0.26±0.078</b>
									<b>0.16±0.003</b>

Table 4: Performance results for the Claims dataset (predicting pancytopenia onset from drug exposure using patient medical history). This table extends Table 2 with standard deviations.

	RATE @ $u$ ( $\uparrow$ )			Recall @ $u$ ( $\uparrow$ )			Precision @ $u$ ( $\uparrow$ )			
	0.999	.998	0.995	0.99	0.999	0.998	0.995	0.999	0.998	0.995
Random	0.00±<0.001	0.00±<0.001	0.00±<0.001	0.00±<0.001	0.0±<0.001	0.0±<0.001	0.01±<0.001	0.01±<0.0014	0.01±<0.001	0.01±<0.001
T-learner	0.10±<0.001	0.07±<0.001	0.05±<0.001	0.04±<0.001	0.05±<0.001	0.07±<0.001	0.11±<0.001	0.10±<0.001	0.08±<0.001	0.06±<0.001
X-learner	0.00±<0.001	-0.01±<0.001	0.00±<0.001	0.00±<0.001	0.00±<0.001	0.00±<0.001	0.01±<0.001	0.00±<0.001	0.00±<0.001	0.00±<0.001
R-learner	-0.01±<0.001	-0.01±<0.001	-0.01±<0.001	0.00±<0.001	0.00±<0.001	0.00±<0.001	0.04±<0.001	0.00±<0.001	0.00±<0.001	0.00±<0.001
RA-learner	0.28±<0.001	0.26±<0.001	0.17±<0.001	0.10±<0.001	0.10±<0.001	0.19±<0.001	0.30±<0.001	0.30±<0.001	0.28±<0.001	0.18±<0.001
DragonNet	-0.01±0.002	0.00±0.009	0.00±0.004	0.00±0.003	0.00±<0.001	0.00±0.003	0.00±0.005	0.00±<0.001	0.00±0.010	0.00±0.004
TARNet	0.04±0.046	0.03±0.030	0.02±0.013	0.02±0.012	0.01±0.011	0.02±0.015	0.04±0.013	0.05±0.046	0.04±0.032	0.03±0.013
FlexTENet	0.02±0.024	0.02±0.019	0.04±0.012	0.03±0.013	0.01±0.009	0.03±0.018	0.08±0.012	0.02±0.027	0.03±0.020	0.04±0.012
S-learner w/ meta-learning	0.27±0.173	0.16±0.118	0.08±0.052	0.04±0.030	0.09±0.055	0.10±0.070	0.13±0.084	0.29±0.180	0.18±0.123	0.09±0.055
T-learner w/ meta-learning	0.27±0.173	0.16±0.118	0.08±0.052	0.04±0.030	0.09±0.055	0.10±0.070	0.13±0.084	0.29±0.180	0.18±0.123	0.09±0.055
GraphITE	0.25±0.088	0.15±0.054	0.06±0.025	0.03±0.011	0.08±0.024	0.10±0.034	0.11±0.045	0.27±0.091	0.16±0.057	0.07±0.027
SIN	0.00±0.008	0.00±0.014	0.00±0.008	0.00±0.005	0.00±0.005	0.00±0.008	0.02±0.015	0.00±0.007	0.01±0.014	0.01±0.009
CaML - w/o meta-learning	0.45±0.070	<b>0.38</b> ±0.057	0.21±0.017	0.13±0.008	0.19±0.019	0.28±0.026	0.38±0.025	0.49±0.070	<b>0.41</b> ±0.057	0.23±0.017
CaML - w/o RA-learner	0.40±0.101	0.33±0.034	<b>0.24</b> ±0.014	0.15±0.010	0.18±0.025	0.28±0.010	0.42±0.024	0.44±0.099	0.36±0.033	<b>0.26</b> ±0.014
CaML (ours)	<b>0.47</b> ±0.084	0.37±0.044	0.23±0.022	<b>0.15</b> ±0.013	<b>0.20</b> ±0.015	<b>0.30</b> ±0.016	<b>0.43</b> ±0.024	<b>0.51</b> ±0.079	0.40±0.044	0.25±0.023

Table 5: Performance results for the medical claims dataset, in which the task is to predict the effect of a *pair* of drugs the drug on pancytopenia occurrence. Mean and standard deviation between runs is reported. Single-task methods were trained on the meta-testing tasks (best model underlined). Methods that were capable of training across multiple tasks were trained on meta-training tasks and applied to previously unseen meta-testing tasks (best model in bold). CaML outperforms the strongest baseline that had access to testing tasks on 12 out of 12 metrics, and outperforms all zero-shot baselines. Notably, due to the small sample size for natural experiments with combinations of drugs, *the RATE estimation process is very noisy* which is reflected in high variability of the measured RATE. Here, the secondary metrics (Recall and Precision) that are not affected, additionally assert the dominance of CaML over all baselines.

	Split	# of Patients
Allopurinol	Test	815,921
Pregabalin	Test	636,995
Mirtazapine	Test	623,980
Indomethacin	Test	560,380
Colchicine	Test	370,397
Hydralazine	Test	363,070
Hydroxychloroquine	Test	324,750
Methotrexate	Test	323,387
Memantine	Test	306,832
Fentanyl	Test	261,000
Etodolac	Val	438,854
Azathioprine	Val	100,000

Table 6: Held-out test and validation drugs for our single-drug meta-testing and meta-validation datasets for our Claims evaluation in Table 2. Drugs are unseen (excluded) during training. All drugs are known to cause pancytopenia [40].

	Split	# of Patients
Allopurinol + Hydralazine	Test	7,859
Methotrexate + Hydroxychloroquine	Test	25,716
Pregabalin + Fentanyl	Test	5,424
Indomethacin + Colchicine	Test	42,846
Mirtazapine + Memantine	Test	10,215

Table 7: Held-out test pairs of drugs for our meta-testing and meta-validation datasets in Appendix Table B.3. Both drugs are unseen (excluded) during training. All drugs are known to cause pancytopenia [40].

Split	# Perturbagens	# Cell-Lines	Mean #Cell Lines/Task
Meta-training	9717	77	5.79
Meta-validation	304	11	9.99
Meta-testing	301	11	10.77

Table 8: Composition of the meta-training, meta-validation and meta-testing sets for the LINC6 dataset. No cell lines or drugs (tasks) were shared across any of the splits.

Hyperparameter	Search range
Num. of layers	{2,4,6}
Dim. of hidden layers	{128,256}
Dropout	{0,0.1}
Learning rate	$\{3 \times 10^{-3}, 1 \times 10^{-3}, 3 \times 10^{-4}, 1 \times 10^{-4}\}$
Meta learning rate	{1}
Weight decay	$\{5 \times 10^{-3}\}$
Reptile k	{1,10,50}
L1 regularization coefficient	$\{0, 1 \times 10^{-7}, 5 \times 10^{-7}\}$

Table 9: Hyperparameter search space for CaML (our proposed method) on the medical claims dataset.

Hyperparameter	Search range
Num. of como layers	{2,4,6}
Num. of covariate layers	{2,4,6}
Num. of propensity layers	{2,4,6}
Num. of treatment layers	{2,4,6}
Dim. of hidden como layers	{128,256}
Dim. of hidden covariate layers	{128,256}
Dim. of hidden treatment layers	{128,256}
Dim. of hidden propensity layers	{16,32,64,128}
Dropout	{0,0.1}
Learning rate	$\{3 \times 10^{-3}, 1 \times 10^{-3}, 3 \times 10^{-4}, 1 \times 10^{-4}\}$
Meta learning rate	{1}
Sin Weight decay	$\{0,5 \times 10^{-3}\}$
Pro Weight decay	$\{0,5 \times 10^{-3}\}$
GNN Weight decay	$\{0,5 \times 10^{-3}\}$
Reptile k	{1,10,50}
L1 regularization coefficient	$\{0,1 \times 10^{-7}, 5 \times 10^{-7}\}$

Table 10: Hyperparameter search space for **SIN** on the medical claims dataset. The SIN model consists of two stages, Stage 1 and Stage 2. For the Stage 1 model we searched the identical hyperparameter search space as for CaML (Table 9). For Stage 2, we used the hyperparameters shown in this table.

Hyperparameter	Search range
Num. of covariate layers	{2,4,6}
Num. of treatment layers	{2,4,6}
Dim. of hidden treatment layers	{128,256}
Dim. of hidden covariate layers	{128,256}
Dropout	{0,0.1}
Independence regularization coefficient	{0,0.01,0.1,1.0}
Learning rate	$\{3 \times 10^{-3}, 1 \times 10^{-3}, 3 \times 10^{-4}, 1 \times 10^{-4}\}$
Meta learning rate	{1}
Weight decay	$\{5 \times 10^{-3}\}$
Reptile k	{1,10,50}
L1 regularization coefficient	$\{0,1 \times 10^{-7}, 5 \times 10^{-7}\}$

Table 11: Hyperparameter search space for **GraphITE** on the medical claims dataset.

Hyperparameter	Search range
Num. of out layers	{1,2,4}
Num. of r layers	{2,4,6}
Num. units p out	{32,64,128,256}
Num. units s out	{32,64,128,256}
Num. units s r	{32,64,128,256}
Num. units p r	{32,64,128,256}
Weight decay	$\{5 \times 10^{-3}\}$
Orthogonal penalty	$\{0,1 \times 10^{-5}, 1 \times 10^{-3}, 0.1\}$
Private out	{True, False}
Learning rate	$\{3 \times 10^{-3}, 1 \times 10^{-3}, 3 \times 10^{-4}, 1 \times 10^{-4}\}$

Table 12: Hyperparameter search space for **FlexTENet** on the medical claims dataset.

Hyperparameter	Search range
Num. of out layers	{1,2,4}
Num. of r layers	{2,4,6}
Num. units out	{128,256}
Weight decay	$\{5 \times 10^{-3}\}$
Penalty disc	$\{0, 1 \times 10^{-3}\}$
Learning rate	$\{3 \times 10^{-3}, 1 \times 10^{-3}, 3 \times 10^{-4}, 1 \times 10^{-4}\}$

Table 13: Hyperparameter search space for **TARNet** on the medical claims dataset.

Hyperparameter	Search range
Num. of out layers	{1,2,4}
Num. of r layers	{2,4,6}
Num. units r	{128,256}
Num. units out	{128,256}
Weight decay	$\{5 \times 10^{-3}\}$
Learning rate	$\{3 \times 10^{-3}, 1 \times 10^{-3}, 3 \times 10^{-4}, 1 \times 10^{-4}\}$

Table 14: Hyperparameter search space for **DragonNet** on the medical claims dataset.

Hyperparameter	Search range
Num. of estimators	[50,250]
Max depth	[10,50]
Min sample split	[2,8]
Criterion regress	{squared error, absolute error}
Criterion binary	{gini, entropy}
Max features	{sqrt, log2, auto}

Table 15: Hyperparameter search space for model-agnostic CATE estimators, i.e., **R-learner**, **X-learner**, **RA-learner**, and **T-learner** on the medical claims dataset.

Hyperparameter	Search range
Num. of layers	{2,4,6}
Dim. of hidden layers	{512,1024}
Dropout	{0,0.1}
Learning rate	$\{3 \times 10^{-3}, 1 \times 10^{-3}, 3 \times 10^{-4}, 1 \times 10^{-4}\}$
Meta learning rate	{0.1,0.5,0.9}
Weight decay	{0.1}
Reptile k	{1,2,3}
L1 regularization coefficient	$\{0, 1 \times 10^{-7}, 5 \times 10^{-7}\}$

Table 16: Hyperparameter search space for **CaML** (our proposed method) on the LINC dataset.

Hyperparameter	Search range
Num. of como layers	{2,4,6}
Num. of covariates layers	{2,4,6}
Num. of propensity layers	{2,4,6}
Num. of treatment layers	{2,4,6}
Dim. output	{128,256}
Dim. of hidden treatment layers	{128,256}
Dim. of hidden covariate layers	{128,256}
Dim. of hidden como layers	{128,256}
Dim. of hidden propensity layers	{16,32,64,128}
Model dim.	{512,1024}
Dropout	{0,0.1}
Learning rate	$\{3 \times 10^{-3}, 1 \times 10^{-3}, 3 \times 10^{-4}, 1 \times 10^{-4}\}$
Meta learning rate	{0.1,0.5,0.9}
Sin weight decay	{0.0,0.005}
Pro weight decay	{0.0,0.005}
GNN weight decay	{0.0,0.005}
Weight decay	{0.1}
Reptile k	{1,2,3}
L1 regularization coefficient	$\{0, 1 \times 10^{-7}, 5 \times 10^{-7}\}$

Table 17: Hyperparameter search space for the **SIN** baseline on the LINC5 dataset.

Hyperparameter	Search range
Num. of covariate layers	{2,4,6}
Num. of treatment layers	{2,4,6}
Num. of layers	{2,4,6}
Dim. of hidden covariate layers	{128,256}
Independence regularization coefficient	{0,0.01,0.1,1.0}
Dropout	{0,0.1}
Model dim.	{512,1024}
Learning rate	$\{3 \times 10^{-3}, 1 \times 10^{-3}, 3 \times 10^{-4}, 1 \times 10^{-4}\}$
Meta learning rate	{0.1,0.5,0.9}
Weight decay	{0.1}
Reptile k	{1,2,3}
L1 regularization coefficient	$\{0, 1 \times 10^{-7}, 5 \times 10^{-7}\}$

Table 18: Hyperparameter search space for the **GraphITE** baseline on the LINC5 dataset.

Hyperparameter	Search range
Dropout	[1e-4,1e-1]
Learning rate	[1e-5,1e-3]
Weight decay	[1e-5,1e-2]
Adversarial temperature	[1,10]
Gamma	[0,30]
Num. of sampled neighbors	0-10
Dim. of hidden layers	{ 64, 128, 256, 512}

Table 19: The hyperparameter optimization search ranges used in the selection of the optimal model for the generation of knowledge graph node embeddings that would serve as intervention information for the medical claims dataset.

361 **References**

- 362 [1] Ahmed M Alaa and Mihaela Van Der Schaar. Bayesian inference of individualized treatment  
363 effects using multi-task gaussian processes. *Advances in neural information processing systems*,  
364 30, 2017.
- 365 [2] Susan Athey and Guido Imbens. Recursive partitioning for heterogeneous causal effects.  
366 *Proceedings of the National Academy of Sciences*, 113(27):7353–7360, 2016.
- 367 [3] Yoshua Bengio, Samy Bengio, and Jocelyn Cloutier. *Learning a synaptic learning rule*. Citeseer,  
368 1990.
- 369 [4] Ioana Bica, Ahmed M Alaa, Craig Lambert, and Mihaela Van Der Schaar. From real-world patient  
370 data to individualized treatment effects using machine learning: current and future methods to  
371 address underlying challenges. *Clinical Pharmacology & Therapeutics*, 109(1):87–100, 2021.
- 372 [5] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A*  
373 *nonasymptotic theory of independence*. Oxford university press, 2013.
- 374 [6] Olivier Bousquet. A bennett concentration inequality and its application to suprema of empirical  
375 processes. *Comptes Rendus Mathématique*, 334(6):495–500, 2002.
- 376 [7] Payal Chandak, Kexin Huang, and Marinka Zitnik. Building a knowledge graph to enable  
377 precision medicine. *bioRxiv*, 2022.
- 378 [8] Hong-Bin Chen, Sinho Chewi, and Jonathan Niles-Weed. Dimension-free log-sobolev  
379 inequalities for mixture distributions. *Journal of Functional Analysis*, 281(11):109236, 2021.
- 380 [9] Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen,  
381 Whitney Newey, and James Robins. Double/debiased machine learning for treatment and  
382 structural parameters, 2018.
- 383 [10] Victor Chernozhukov, Mert Demirer, Esther Duflo, and Ivan Fernandez-Val. Generic machine  
384 learning inference on heterogeneous treatment effects in randomized experiments, with an appli-  
385 cation to immunization in india. Technical report, National Bureau of Economic Research, 2018.
- 386 [11] Richard K Crump, V Joseph Hotz, Guido W Imbens, and Oscar A Mitnik. Nonparametric tests  
387 for treatment effect heterogeneity. *The Review of Economics and Statistics*, 90(3):389–405, 2008.
- 388 [12] Alicia Curth, David Svensson, Jim Weatherall, and Mihaela van der Schaar. Really doing great  
389 at estimating cate? a critical look at ml benchmarking practices in treatment effect estimation.  
390 In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks*  
391 *Track (Round 2)*, 2021.
- 392 [13] Alicia Curth and Mihaela van der Schaar. Doing great at estimating cate? on the neglected  
393 assumptions in benchmark comparisons of treatment effect estimators. *arXiv preprint*  
394 *arXiv:2107.13346*, 2021.
- 395 [14] Alicia Curth and Mihaela van der Schaar. Nonparametric estimation of heterogeneous treatment  
396 effects: From theory to learning algorithms. In *International Conference on Artificial Intelligence*  
397 *and Statistics*, pages 1810–1818. PMLR, 2021.
- 398 [15] Alicia Curth and Mihaela van der Schaar. On inductive biases for heterogeneous treatment effect  
399 estimation. *Advances in Neural Information Processing Systems*, 34:15883–15894, 2021.
- 400 [16] Gerald DeJong and Raymond Mooney. Explanation-based learning: An alternative view.  
401 *Machine learning*, 1986.
- 402 [17] Qiaonan Duan, Corey Flynn, Mario Niepel, Marc Hafner, Jeremy L Muhlich, Nicolas F  
403 Fernandez, Andrew D Rouillard, Christopher M Tan, Edward Y Chen, Todd R Golub, et al. Lincs  
404 canvas browser: interactive web app to query, browse and interrogate lincs l1000 gene expression  
405 signatures. *Nucleic acids research*, 42(W1):W449–W460, 2014.

- 406 [18] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast  
407 adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135.  
408 PMLR, 2017.
- 409 [19] Dennis Frauen and Stefan Feuerriegel. Estimating individual treatment effects under unobserved  
410 confounding using binary instruments. *arXiv preprint arXiv:2208.08544*, 2022.
- 411 [20] Mahmoud Ghandi, Franklin W Huang, Judit Jané-Valbuena, Gregory V Kryukov, Christopher C  
412 Lo, E Robert McDonald, 3rd, Jordi Barretina, Ellen T Gelfand, Craig M Bielski, Haoxin Li,  
413 Kevin Hu, Alexander Y Andreev-Drakhlin, Jaegil Kim, Julian M Hess, Brian J Haas, François  
414 Aguet, Barbara A Weir, Michael V Rothberg, Brenton R Paoella, Michael S Lawrence, Rehan  
415 Akbani, Yiling Lu, Hong L Tiv, Prafulla C Gokhale, Antoine de Weck, Ali Amin Mansour, Coyin  
416 Oh, Juliann Shih, Kevin Hadi, Yanay Rosen, Jonathan Bistline, Kavitha Venkatesan, Anupama  
417 Reddy, Dmitriy Sonkin, Manway Liu, Joseph Lehar, Joshua M Korn, Dale A Porter, Michael D  
418 Jones, Javad Golji, Giordano Caponigro, Jordan E Taylor, Caitlin M Dunning, Amanda L  
419 Creech, Allison C Warren, James M McFarland, Mahdi Zamanighomi, Audrey Kauffmann,  
420 Nicolas Stransky, Marcin Imielinski, Yosef E Maruvka, Andrew D Cherniack, Aviad Tsherniak,  
421 Francisca Vazquez, Jacob D Jaffe, Andrew A Lane, David M Weinstock, Cory M Johannessen,  
422 Michael P Morrissey, Frank Stegmeier, Robert Schlegel, William C Hahn, Gad Getz, Gordon B  
423 Mills, Jesse S Boehm, Todd R Golub, Levi A Garraway, and William R Sellers. Next-generation  
424 characterization of the cancer cell line encyclopedia. *Nature*, 569(7757):503–508, May 2019.
- 425 [21] Donald P Green and Holger L Kern. Modeling heterogeneous treatment effects in survey exper-  
426 iments with bayesian additive regression trees. *Public opinion quarterly*, 76(3):491–511, 2012.
- 427 [22] Lin Lawrence Guo, Ethan Steinberg, Scott Lanyon Fleming, Jose Posada, Joshua Lemmon,  
428 Stephen R Pfohl, Nigam Shah, Jason Fries, and Lillian Sung. Ehr foundation models improve  
429 robustness in the presence of temporal distribution shift. *medRxiv*, 2022.
- 430 [23] Shonosuke Harada and Hisashi Kashima. Graphite: Estimating individual effects of graph-  
431 structured treatments. In *Proceedings of the 30th ACM International Conference on Information  
432 & Knowledge Management*, pages 659–668, 2021.
- 433 [24] Negar Hassanpour and Russell Greiner. Counterfactual regression with importance sampling  
434 weights. In *IJCAI*, pages 5880–5887, 2019.
- 435 [25] Negar Hassanpour and Russell Greiner. Learning disentangled representations for counterfactual  
436 regression. In *International Conference on Learning Representations*, 2019.
- 437 [26] Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements  
438 of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- 439 [27] Leon Hetzel, Simon Böhm, Niki Kilbertus, Stephan Günemann, Mohammad Lotfollahi, and  
440 Fabian Theis. Predicting single-cell perturbation responses for unseen drugs. *arXiv preprint  
441 arXiv:2204.13545*, 2022.
- 442 [28] Jennifer L Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational  
443 and Graphical Statistics*, 20(1):217–240, 2011.
- 444 [29] Jennifer L Hill, Jeanne Brooks-Gunn, and Jane Waldfogel. Sustained effects of high participation  
445 in an early intervention for low-birth-weight premature infants. *Developmental psychology*,  
446 39(4):730, 2003.
- 447 [30] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in  
448 neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*,  
449 44(9):5149–5169, 2021.
- 450 [31] Guido W Imbens and Donald B Rubin. *Causal inference in statistics, social, and biomedical  
451 sciences*. Cambridge University Press, 2015.
- 452 [32] Fredrik Johansson, Uri Shalit, and David Sontag. Learning representations for counterfactual  
453 inference. In *International conference on machine learning*, pages 3020–3029. PMLR, 2016.

- 454 [33] Jean Kaddour, Yuchen Zhu, Qi Liu, Matt J Kusner, and Ricardo Silva. Causal effect inference  
455 for structured treatments. *Advances in Neural Information Processing Systems*, 34:24841–24854,  
456 2021.
- 457 [34] Edward H Kennedy. Optimal doubly robust estimation of heterogeneous causal effects. *arXiv*  
458 *preprint arXiv:2004.14497*, 2020.
- 459 [35] Edward H Kennedy. Towards optimal doubly robust estimation of heterogeneous causal effects  
460 (2020). URL <https://arxiv.org/abs>, 2020.
- 461 [36] Jitender Mohan Khunger, S Arulsevi, Uma Sharma, Sunil Ranga, and VH Talib. Pancytopenia—a  
462 clinico haematological study of 200 cases. *Indian journal of pathology & microbiology*,  
463 45(3):375–379, 2002.
- 464 [37] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay  
465 Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds:  
466 A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*  
467 *(ICML)*, 2021.
- 468 [38] Andrei V Konstantinov, Stanislav R Kirpichenko, and Lev V Utkin. Heterogeneous treatment  
469 effect with trained kernels of the nadaraya-watson regression. *arXiv preprint arXiv:2207.09139*,  
470 2022.
- 471 [39] N Kostantinos. Gaussian mixtures and their applications to signal processing. *Advanced signal*  
472 *processing handbook: theory and implementation for radar, sonar, and medical imaging real*  
473 *time systems*, pages 3–1, 2000.
- 474 [40] Michael Kuhn, Ivica Letunic, Lars Juhl Jensen, and Peer Bork. The sider database of drugs and  
475 side effects. *Nucleic acids research*, 44(D1):D1075–D1079, 2016.
- 476 [41] R Kumar, SP Kalra, H Kumar, AC Anand, and H Madan. Pancytopenia—a six year study. *The*  
477 *Journal of the Association of Physicians of India*, 49:1078–1081, 2001.
- 478 [42] Sören R Künnel, Jasjeet S Sekhon, Peter J Bickel, and Bin Yu. Metalearners for estimating  
479 heterogeneous treatment effects using machine learning. *Proceedings of the national academy*  
480 *of sciences*, 116(10):4156–4165, 2019.
- 481 [43] Nikolay Kuznetsov and Alexander Nazarov. Sharp constants in the poincaré, steklov and related  
482 inequalities (a survey). *Mathematika*, 61(2):328–344, 2015.
- 483 [44] Greg Landrum et al. Rdkit: Open-source cheminformatics. 2006.
- 484 [45] Michel Ledoux. Concentration of measure and logarithmic sobolev inequalities. In *Seminaire*  
485 *de probabilites XXXIII*, pages 120–216. Springer, 1999.
- 486 [46] Hongzhu Li, Xiangrui Gao, and Yafeng Deng. Stargraph: A coarse-to-fine representation method  
487 for large-scale knowledge graph, 2022.
- 488 [47] Michelle M Li, Kexin Huang, and Marinka Zitnik. Graph representation learning in biomedicine  
489 and healthcare. *Nature Biomedical Engineering*, pages 1–17, 2022.
- 490 [48] Jing Ma, Ruocheng Guo, Aidong Zhang, and Jundong Li. Multi-cause effect estimation with  
491 disentangled confounder representation. In *Proceedings of the Thirtieth International Joint*  
492 *Conference on Artificial Intelligence*, 2021.
- 493 [49] Ron Meir and Tong Zhang. Generalization error bounds for bayesian mixture algorithms.  
494 *Journal of Machine Learning Research*, 4(Oct):839–860, 2003.
- 495 [50] Stephen L Morgan and Christopher Winship. *Counterfactuals and causal inference*. Cambridge  
496 University Press, 2015.
- 497 [51] Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms.  
498 *arXiv preprint arXiv:1803.02999*, 2018.

- 499 [52] Alex Nichol and John Schulman. Reptile: a scalable metalearning algorithm. *arXiv preprint*  
500 *arXiv:1803.02999*, 2(3):4, 2018.
- 501 [53] Xinkun Nie and Stefan Wager. Quasi-oracle estimation of heterogeneous treatment effects.  
502 *Biometrika*, 108(2):299–319, 2021.
- 503 [54] Frank Nielsen and Richard Nock. On the chi square and higher-order chi distances for  
504 approximating f-divergences. *IEEE Signal Processing Letters*, 21(1):10–13, 2013.
- 505 [55] Hamed Nilforoshan and Eugene Wu. Leveraging quality prediction models for automatic writing  
506 feedback. In *Twelfth International AAAI Conference on Web and Social Media*, 2018.
- 507 [56] Lawrence E Payne and Hans F Weinberger. An optimal poincaré inequality for convex domains.  
508 *Archive for Rational Mechanics and Analysis*, 5(1):286–292, 1960.
- 509 [57] Henri Poincaré. Sur les équations aux dérivées partielles de la physique mathématique. *American*  
510 *Journal of Mathematics*, pages 211–294, 1890.
- 511 [58] Zhaozhi Qian, Alicia Curth, and Mihaela van der Schaar. Estimating multi-cause treatment  
512 effects via single-cause perturbation. *Advances in Neural Information Processing Systems*,  
513 34:23754–23767, 2021.
- 514 [59] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. Rapid learning or feature  
515 reuse? towards understanding the effectiveness of maml. *arXiv preprint arXiv:1909.09157*, 2019.
- 516 [60] Bernardino Romera-Paredes and Philip Torr. An embarrassingly simple approach to zero-shot  
517 learning. In *International conference on machine learning*, pages 2152–2161. PMLR, 2015.
- 518 [61] Yusuf Roohani, Kexin Huang, and Jure Leskovec. Gears: Predicting transcriptional outcomes  
519 of novel multi-gene perturbations. *bioRxiv*, 2022.
- 520 [62] Shiv Kumar Saini, Sunny Dhamnani, Akil Arif Ibrahim, and Prithviraj Chavan. Multiple  
521 treatment effect estimation using deep generative model with task embedding. In *The World*  
522 *Wide Web Conference*, pages 1601–1611, 2019.
- 523 [63] Tim Salimans and Durk P Kingma. Weight normalization: A simple reparameterization to acceler-  
524 ate training of deep neural networks. *Advances in neural information processing systems*, 29, 2016.
- 525 [64] André Schlichting. Poincaré and log–sobolev inequalities for mixtures. *Entropy*, 21(1):89, 2019.
- 526 [65] Jürgen Schmidhuber. *Evolutionary principles in self-referential learning, or on learning how*  
527 *to learn: the meta-meta-... hook*. PhD thesis, Technische Universität München, 1987.
- 528 [66] Uri Shalit, Fredrik D Johansson, and David Sontag. Estimating individual treatment effect:  
529 generalization bounds and algorithms. In *International Conference on Machine Learning*, pages  
530 3076–3085. PMLR, 2017.
- 531 [67] Ankit Sharma, Garima Gupta, Ranjitha Prasad, Arnab Chatterjee, Lovekesh Vig, and Gautam  
532 Shroff. Metaci: Meta-learning for causal inference in a heterogeneous population. *arXiv preprint*  
533 *arXiv:1912.03960*, 2019.
- 534 [68] Claudia Shi, David Blei, and Victor Veitch. Adapting neural networks for the estimation of  
535 treatment effects. *Advances in neural information processing systems*, 32, 2019.
- 536 [69] Yishai Shimoni, Chen Yanover, Ehud Karavani, and Yaara Goldschmid. Benchmarking frame-  
537 work for performance-evaluation of causal inference analysis. *arXiv preprint arXiv:1802.05046*,  
538 2018.
- 539 [70] Aravind Subramanian, Rajiv Narayan, Steven M Corsello, David D Peck, Ted E Natoli, Xiaodong  
540 Lu, Joshua Gould, John F Davis, Andrew A Tubelli, Jacob K Asiedu, David L Lahr, Jodi E  
541 Hirschman, Zihan Liu, Melanie Donahue, Bina Julian, Mariya Khan, David Wadden, Ian C  
542 Smith, Daniel Lam, Arthur Liberzon, Courtney Toder, Mukta Bagul, Marek Orzechowski,  
543 Oana M Enache, Federica Piccioni, Sarah A Johnson, Nicholas J Lyons, Alice H Berger,  
544 Alykhan F Shamji, Angela N Brooks, Anita Vrcic, Corey Flynn, Jacqueline Rosains, David Y

- 545 Takeda, Roger Hu, Desiree Davison, Justin Lamb, Kristin Ardlie, Larson Hogstrom, Peyton  
546 Greenside, Nathanael S Gray, Paul A Clemons, Serena Silver, Xiaoyun Wu, Wen-Ning Zhao,  
547 Willis Read-Button, Xiaohua Wu, Stephen J Haggarty, Lucienne V Ronco, Jesse S Boehm,  
548 Stuart L Schreiber, John G Doench, Joshua A Bittker, David E Root, Bang Wong, and Todd R  
549 Golub. A next generation connectivity map: L1000 platform and the first 1,000,000 profiles.  
550 *Cell*, 171(6):1437–1452.e17, November 2017.
- 551 [71] Nicholas P Tatonetti, Patrick P Ye, Roxana Daneshjou, and Russ B Altman. Data-driven prediction  
552 of drug effects and interactions. *Science translational medicine*, 4(125):125ra31–125ra31, 2012.
- 553 [72] Sebastian Thrun and Lorien Pratt. *Learning to learn*. Springer Science & Business Media, 2012.
- 554 [73] Vladimir Vapnik. Principles of risk minimization for learning theory. *Advances in neural  
555 information processing systems*, 4, 1991.
- 556 [74] Victor Veitch, Dhanya Sridhar, and David Blei. Adapting text embeddings for causal inference.  
557 In *Conference on Uncertainty in Artificial Intelligence*, pages 919–928. PMLR, 2020.
- 558 [75] Stefan Wager. Stats 361: Causal inference, 2020.
- 559 [76] Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using  
560 random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018.
- 561 [77] Wei Wang, Vincent W Zheng, Han Yu, and Chunyan Miao. A survey of zero-shot learning:  
562 Settings, methods, and applications. *ACM Transactions on Intelligent Systems and Technology  
563 (TIST)*, 10(2):1–37, 2019.
- 564 [78] Yixin Wang and David M Blei. The blessings of multiple causes. *Journal of the American  
565 Statistical Association*, 114(528):1574–1596, 2019.
- 566 [79] Galen Weld, Peter West, Maria Glenski, David Arbour, Ryan A Rossi, and Tim Althoff.  
567 Adjusting for confounders with text: Challenges and an empirical evaluation framework for  
568 causal inference. In *Proceedings of the International AAAI Conference on Web and Social Media*,  
569 volume 16, pages 1109–1120, 2022.
- 570 [80] Yongqin Xian, Bernt Schiele, and Zeynep Akata. Zero-shot learning—the good, the bad and the  
571 ugly. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages  
572 4582–4591, 2017.
- 573 [81] Steve Yadlowsky, Scott Fleming, Nigam Shah, Emma Brunskill, and Stefan Wager. Evaluating  
574 treatment prioritization rules via rank-weighted average treatment effects. *arXiv preprint  
575 arXiv:2111.07966*, 2021.
- 576 [82] Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. QA-GNN:  
577 Reasoning with language models and knowledge graphs for question answering. In *North  
578 American Chapter of the Association for Computational Linguistics (NAACL)*, 2021.
- 579 [83] Shing-Tung Yau. Isoperimetric constants and the first eigenvalue of a compact riemannian man-  
580 ifold. In *Annales Scientifiques de l’École Normale Supérieure*, volume 8, pages 487–507, 1975.
- 581 [84] Jinsung Yoon, James Jordon, and Mihaela Van Der Schaar. Ganite: Estimation of individualized  
582 treatment effects using generative adversarial nets. In *International Conference on Learning  
583 Representations*, 2018.
- 584 [85] Long Yu, Zhicong Luo, Huanyong Liu, Deng Lin, Hongzhu Li, and Yafeng Deng. Triplere: Knowl-  
585 edge graph embeddings via tripled relation vectors. *arXiv preprint arXiv:2209.08271*, 2022.
- 586 [86] Yao Zhang, Alexis Bellot, and Mihaela Schaar. Learning overlapping representations for  
587 the estimation of individualized treatment effects. In *International Conference on Artificial  
588 Intelligence and Statistics*, pages 1005–1014. PMLR, 2020.
- 589 [87] Marinka Zitnik, Monica Agrawal, and Jure Leskovec. Modeling polypharmacy side effects with  
590 graph convolutional networks. *Bioinformatics*, 34(13):i457–i466, 2018.