

A Theoretical Results

A.1 τ -BALD

Theorem 1. *Under the following assumptions:*

1. *Unconfoundedness* $(Y^0, Y^1) \perp\!\!\!\perp T \mid \mathbf{X}$;
2. *Consistency* $Y \mid T = Y^t$;
3. Y^1 and Y^0 , when conditioned on realizations \mathbf{x} of the r.v. \mathbf{X} and t of the r.v. T , are independent-normally distributed or joint-normally distributed r.v.s.
4. $\hat{\mu}_\omega(\mathbf{x}, t)$ is a consistent estimator of $\mathbb{E}[Y \mid T = t, \mathbf{X} = \mathbf{x}]$

the information gain for Ω if we could observe a label for the difference in potential outcomes $Y^1 - Y^0$ given measured covariates \mathbf{x} , treatment t and a dataset of observations $\mathcal{D}_{\text{train}} = \{\mathbf{x}_i, t_i, y_i\}_{i=1}^n$ is approximated as

$$I(Y^1 - Y^0; \Omega \mid \mathbf{x}, t, \mathcal{D}_{\text{train}}) \approx \text{Var}_{\omega \sim p(\Omega \mid \mathcal{D}_{\text{train}})} (\hat{\mu}_\omega(\mathbf{x}, 1) - \hat{\mu}_\omega(\mathbf{x}, 0)) \quad (10)$$

Proof.

$$I(Y^1 - Y^0; \Omega \mid \mathbf{x}, \mathcal{D}_{\text{train}}) = H(Y^1 - Y^0 \mid \mathbf{x}, \mathcal{D}_{\text{train}}) - \mathbb{E}_{p(\Omega \mid \mathcal{D}_{\text{train}})} [H(Y^1 - Y^0 \mid \mathbf{x}, \omega)] \quad (11a)$$

$$\approx \text{Var}(Y^1 - Y^0 \mid \mathbf{x}, \mathcal{D}_{\text{train}}) - \mathbb{E}_{p(\Omega \mid \mathcal{D}_{\text{train}})} [\text{Var}(Y^1 - Y^0 \mid \mathbf{x}, \omega)] \quad (11b)$$

$$= \text{Var}_{p(\Omega \mid \mathcal{D}_{\text{train}})} (\mathbb{E}[Y^1 - Y^0 \mid \mathbf{x}, \omega]) \quad (11c)$$

$$= \text{Var}_{p(\Omega \mid \mathcal{D}_{\text{train}})} (\hat{\mu}_\omega(\mathbf{x}, 1) - \hat{\mu}_\omega(\mathbf{x}, 0)) \quad (11d)$$

In (11a) we adapt the result of Houlsby et al. [18] and express the information gain as the mutual information between the observable difference in potential outcomes $Y^1 - Y^0$ and the parameters Ω ; given observed covariates \mathbf{x} , treatment t , and training data $\mathcal{D}_{\text{train}} = \{(\mathbf{x}_i, t_i, y_i)\}_{i=1}^{n_{\text{train}}}$. In (11b) we apply lemma 1.1 to the r.h.s terms of (11a). We then use the result in Jesson et al. [20] and move from (11b) to (11c) by application of the law of total variance. Finally, under the consistency and unconfoundedness assumptions we express the information gain in terms of the identifiable difference in expected outcomes $\hat{\mu}_\omega(\mathbf{x}, 1) - \hat{\mu}_\omega(\mathbf{x}, 0)$. \square

Lemma 1.1. *Under the following assumptions:*

1. Y^1, Y^0 are independent-normally distributed or joint-normally distributed r.v.s;
2. With $A = \text{Var}(Y^1 - Y^0)$: let $|A - 1| \leq 1$ and $A \neq 0$. That is to say, the predictive variance must be greater than 0 and less than or equal to 2;

$$H(Y^1 - Y^0) \approx \text{Var}(Y^1 - Y^0) \quad (12)$$

Proof. By assumption 1, $Y^1 - Y^0$ is also a normally distributed random variable. By corollary 1.1,

$$H(Y^1 - Y^0) = \frac{1}{2} + \frac{1}{2} \log(2\pi \text{Var}(Y^1 - Y^0)) \quad (13)$$

So given assumption 2, the first order Taylor polynomial of $H(Y^1 - Y^0)$ is

$$\begin{aligned} \frac{1}{2} + \frac{1}{2} \log(2\pi \text{Var}(Y^1 - Y^0)) &\approx \frac{1}{2} + \frac{1}{2} (2\pi \text{Var}(Y^1 - Y^0) - 1) \\ &= \frac{1}{2} + \pi \text{Var}(Y^1 - Y^0) - \frac{1}{2} \\ &= \pi \text{Var}(Y^1 - Y^0) \\ &\propto \text{Var}(Y^1 - Y^0) \end{aligned} \quad (14)$$

\square

Corollary 1.1. *The entropy of a normally distributed random variable with variance σ^2 is $\frac{1}{2} + \frac{1}{2} \log(2\pi\sigma^2)$*

A.2 μ -BALD

Theorem 2. *Under the following assumptions:*

1. *Unconfoundedness $(Y^0, Y^1) \perp\!\!\!\perp T \mid \mathbf{X}$,*
2. *Consistency $Y \mid T = Y^t$,*
3. *Y conditioned on \mathbf{x} and t is a normally distributed random variable,*
4. *$\hat{\mu}_\omega(\mathbf{x}, t)$ is a consistent estimator of $\mathbb{E}[Y \mid T = t, \mathbf{X} = \mathbf{x}]$,*

the information gain for Ω when we observe a label for the potential outcome Y^t given measured covariates \mathbf{x} , treatment t and a dataset of observations $\mathcal{D}_{\text{train}} = \{\mathbf{x}_i, t_i, y_i\}_{i=1}^n$ can be approximated as is

$$I(Y^t; \Omega \mid \mathbf{x}, t, \mathcal{D}_{\text{train}}) \approx \frac{1}{2} \log \left(\frac{\text{Var}(Y \mid \mathbf{x}, t, \mathcal{D}_{\text{train}})}{\mathbb{E}_\omega[\text{Var}(Y \mid \mathbf{x}, t, \omega)]} \right), \quad (15)$$

or

$$I(Y^t; \Omega \mid \mathbf{x}, t, \mathcal{D}_{\text{train}}) \approx \text{Var}_{\omega \sim p(\Omega \mid \mathcal{D}_{\text{train}})}(\hat{\mu}_\omega(\mathbf{x}, t)). \quad (16)$$

Equation (15) expresses the information gain as the logarithm of a ratio between predictive and aleatoric uncertainty in the outcome. Whereas, equation (16) expresses the information gain as a direct estimate of the epistemic uncertainty.

Proof.

$$I(Y^t; \Omega \mid \mathbf{x}, t, \mathcal{D}_{\text{train}}) = H(Y \mid \mathbf{x}, t, \mathcal{D}_{\text{train}}) - \mathbb{E}_{p(\Omega \mid \mathcal{D}_{\text{train}})}[H(Y \mid \mathbf{x}, t, \omega)] \quad (17a)$$

$$= \frac{1}{2} \log(2\pi \text{Var}(Y \mid \mathbf{x}, t, \mathcal{D}_{\text{train}})) - \mathbb{E}_{p(\Omega \mid \mathcal{D}_{\text{train}})} \left[\frac{1}{2} \log(2\pi \text{Var}(Y \mid \omega, \mathbf{x}, t)) \right] \quad (17b)$$

$$\geq \frac{1}{2} \log(2\pi \text{Var}(Y \mid \mathbf{x}, t, \mathcal{D}_{\text{train}})) - \frac{1}{2} \log \left(2\pi \mathbb{E}_{p(\Omega \mid \mathcal{D}_{\text{train}})} \text{Var}(Y \mid \omega, \mathbf{x}, t) \right) \quad (17c)$$

$$= \frac{1}{2} \log \left(\frac{\text{Var}(Y \mid \mathbf{x}, t, \mathcal{D}_{\text{train}})}{\mathbb{E}_\omega[\text{Var}(Y \mid \mathbf{x}, t, \omega)]} \right) \quad (17d)$$

$$I(Y^t; \Omega \mid \mathbf{x}, t, \mathcal{D}_{\text{train}}) = H(Y \mid \mathbf{x}, t, \mathcal{D}_{\text{train}}) - \mathbb{E}_{p(\Omega \mid \mathcal{D}_{\text{train}})}[H(Y \mid \mathbf{x}, t, \omega)] \quad (18a)$$

$$\approx \text{Var}[Y \mid \mathbf{x}, t, \mathcal{D}_{\text{train}}] - \mathbb{E}_{p(\Omega \mid \mathcal{D}_{\text{train}})}[\text{Var}[Y \mid \mathbf{x}, t, \omega]] \quad (18b)$$

$$= \text{Var}_{\omega \sim p(\Omega \mid \mathcal{D}_{\text{train}})}(\hat{\mu}_\omega(\mathbf{x}, t)) \quad (18c)$$

In (18a) we express the information gain as the mutual information between the observed potential outcome Y^t and the parameters Ω ; given observed covariates \mathbf{x} , treatment t , and training data $\mathcal{D}_{\text{train}}$. By consistency, we can drop the superscript on the potential outcome. In (18b) we approximate the r.h.s terms of (18a) by application of Lemma 1.1. Finally, we can move from (18b) to (18c) by application of the law of total variance. \square

Note that for discrete or categorical Y , it is straightforward to evaluate Equation (18a) directly.

A.3 ρ -BALD

Theorem 3. *Under the following assumptions*

1. $\{\hat{\mu}_\omega(\mathbf{x}, t) : t \in \{0, 1\}\}$ are instances of the independent-normally distributed or joint-normally distributed random variables $\{\hat{\mu}_\Omega^t = \mathbb{E}[Y \mid \Omega, T = t, \mathbf{x}] : t \in \{0, 1\}\}$,
2. $\text{Var}_{\omega \sim p(\Omega \mid \mathcal{D}_{\text{train}})}(\hat{\mu}_\omega(\mathbf{x}, t')) > 0$.

Let $\hat{\tau}_\omega(\mathbf{x})$ be a realization of the random variable $\hat{\tau}_\Omega = \hat{\mu}_\Omega^1 - \hat{\mu}_\Omega^0$. The information gain for $\hat{\tau}_\Omega$ if we observe the label for the potential outcome Y^t given measured covariates \mathbf{x} , treatment t and a dataset of observations $\mathcal{D}_{\text{train}} = \{\mathbf{x}_i, t_i, y_i\}_{i=1}^n$ is approximately

$$\begin{aligned} I(Y^t; \hat{\tau}_\Omega \mid \mathbf{x}, t, \mathcal{D}_{\text{train}}) &\approx \frac{\text{Var}_\omega(\hat{\tau}_\omega(\mathbf{x}))}{\text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, t'))}, \\ &= \frac{\text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, t)) - 2\text{Cov}_\omega(\hat{\mu}_\omega(\mathbf{x}, t), \hat{\mu}_\omega(\mathbf{x}, t'))}{\text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, t'))} + 1, \end{aligned} \quad (19)$$

where for binary $T = t$, $t' = (1 - t)$.

Proof.

$$I(Y^t; \hat{\tau}_\Omega \mid \mathbf{x}, t, \mathcal{D}) = H(\hat{\tau}_\Omega \mid \mathbf{x}, t, \mathcal{D}) - H(\hat{\tau}_\Omega \mid Y^t, \mathbf{x}, t, \mathcal{D}) \quad (20a)$$

$$= H(\hat{\tau}_\Omega \mid \mathbf{x}, t, \mathcal{D}) - \mathbb{E}_{y^t \sim p(Y^t \mid \mathbf{x}, t, \mathcal{D})} H(\hat{\tau}_\Omega \mid y^t, \mathbf{x}, t) \quad (20b)$$

$$= \frac{1}{2} \log(2\pi \text{Var}(\hat{\tau}_\Omega)) - \mathbb{E}_{y^t \sim p(Y^t \mid \mathbf{x}, t, \mathcal{D})} \left[\frac{1}{2} \log(2\pi \text{Var}(\hat{\tau}_\Omega \mid y^t)) \right] \quad (20c)$$

$$\geq \frac{1}{2} \log(2\pi \text{Var}(\hat{\tau}_\Omega)) - \frac{1}{2} \log(2\pi \mathbb{E}[\text{Var}(\hat{\tau}_\Omega \mid y^t)]) \quad (20d)$$

$$= \frac{1}{2} \log \left(\frac{\text{Var}(\hat{\tau}_\Omega)}{\mathbb{E}[\text{Var}(\hat{\tau}_\Omega \mid y^t)]} \right), \quad (20e)$$

and we can further expand the fraction to

$$\frac{\text{Var}(\hat{\tau}_\Omega \mid \mathbf{x}, t, \mathcal{D})}{\mathbb{E}[\text{Var}(\hat{\tau}_\Omega \mid y^t)]} = \frac{\text{Var}(\hat{\tau}_\Omega \mid \mathbf{x}, t, \mathcal{D})}{\text{Var}_{\omega \sim p(\Omega \mid \mathcal{D})}(\hat{\mu}_\omega(\mathbf{x}, t'))} \quad (20f)$$

$$= \frac{\text{Var}_{\omega \sim p(\Omega \mid \mathcal{D})}(\hat{\tau}_\omega(\mathbf{x}) \mid t)}{\text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, t'))} \quad (20g)$$

$$= \frac{\text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, 1) - \hat{\mu}_\omega(\mathbf{x}, 0) \mid t)}{\text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, t'))} \quad (20h)$$

$$= \frac{\text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, t) - \hat{\mu}_\omega(\mathbf{x}, t'))}{\text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, t'))} \quad (20i)$$

$$= \frac{\text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, t)) + \text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, t')) - 2\text{Cov}_\omega(\hat{\mu}_\omega(\mathbf{x}, t), \hat{\mu}_\omega(\mathbf{x}, t'))}{\text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, t'))} \quad (20j)$$

$$= \frac{\text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, t)) - 2\text{Cov}_\omega(\hat{\mu}_\omega(\mathbf{x}, t), \hat{\mu}_\omega(\mathbf{x}, t'))}{\text{Var}_\omega(\hat{\mu}_\omega(\mathbf{x}, t'))} + 1, \quad (20k)$$

where (20a) by definition of mutual information; (20a)-(20b) from the result of Houlby et al. [18]; (20b)-(20c) by Assumption 1. and Corollary 1.1; (20c)-(20d) by Jensen's inequality; (20d)-(20e) by the logarithmic quotient identity; (20f) by Lemma 3.1; (20f)-(20g) by definition of the variance. (20g)-(20h) by definition of $\hat{\tau}_\omega$; (20h)-(20i) by symmetry of the variance of the difference of two random variables; (20i)-(20j) by the definition of the variance of the difference of two random variables; and (20j)-(20k) by cancelling terms. \square

Lemma 3.1. *Under the following assumptions*

1. *Consistency* $Y \mid T = Y^t$;
2. *Unconfoundedness* $(Y^0, Y^1) \perp\!\!\!\perp T \mid \mathbf{X}$;

$$\mathbb{E}_{y^t \sim p(Y^t | \mathbf{x}, t, \mathcal{D})} [\text{Var}(\hat{\tau}_\Omega \mid y^t)] \approx \mathbb{E}_{y^t \sim p(Y^t | \mathbf{x}, t, \mathcal{D})} \left[\text{Var}_{\omega \sim p(\Omega | \mathcal{D}_{\text{train}})} (\hat{\mu}_\omega(\mathbf{x}, t')) \right], \quad (21)$$

where for binary $T = t$, $t' = (1 - t)$.

Proof.

$$\mathbb{E}_{y^t \sim p(Y^t | \mathbf{x}, t, \mathcal{D})} [\text{Var}(\hat{\tau}_\Omega \mid y^t)] = \mathbb{E}_{p(y^t)} \left[\mathbb{E}_{p(\omega)} \left[\left(\hat{\tau}_\omega - \mathbb{E}_{p(\omega)} [\hat{\tau}_\omega \mid y^t] \right)^2 \mid y^t \right] \right], \quad (22a)$$

$$= \mathbb{E}_{p(y^t)} \left[\mathbb{E}_{p(\omega)} \left[\left(\mathbb{E}[Y^1 - Y^0 \mid \mathbf{x}, \omega] - \mathbb{E}_{p(\omega)} [\mathbb{E}[Y^1 - Y^0 \mid \mathbf{x}, \omega] \mid y^t] \right)^2 \mid y^t \right] \right], \quad (22b)$$

$$= \mathbb{E}_{p(y^t)} \left[\mathbb{E}_{p(\omega)} \left[\left(\mathbb{E}[Y^1 \mid \mathbf{x}, \omega] - \mathbb{E}[Y^0 \mid \mathbf{x}, \omega] - \mathbb{E}_{p(\omega)} [\mathbb{E}[Y^1 \mid \mathbf{x}, \omega] \mid y^t] + \mathbb{E}_{p(\omega)} [\mathbb{E}[Y^0 \mid \mathbf{x}, \omega] \mid y^t] \right)^2 \mid y^t \right] \right], \quad (22c)$$

$$= \mathbb{E}_{p(y^t)} \left[\mathbb{E}_{p(\omega)} \left[\left(\left(\mathbb{E}[Y^1 \mid \mathbf{x}, \omega] - \mathbb{E}_{p(\omega)} [\mathbb{E}[Y^1 \mid \mathbf{x}, \omega] \mid y^t] \right) - \left(\mathbb{E}[Y^0 \mid \mathbf{x}, \omega] - \mathbb{E}_{p(\omega)} [\mathbb{E}[Y^0 \mid \mathbf{x}, \omega] \mid y^t] \right) \right)^2 \mid y^t \right] \right], \quad (22d)$$

$$= \mathbb{E}_{p(y^t)} \left[\mathbb{E}_{p(\omega)} \left[\left(\left(\mathbb{E}[Y^t \mid \mathbf{x}, \omega] - \mathbb{E}_{p(\omega)} [\mathbb{E}[Y^t \mid \mathbf{x}, \omega] \mid y^t] \right) - \left(\mathbb{E}[Y^{t'} \mid \mathbf{x}, \omega] - \mathbb{E}_{p(\omega)} [\mathbb{E}[Y^{t'} \mid \mathbf{x}, \omega] \mid y^t] \right) \right)^2 \mid y^t \right] \right], \quad (22e)$$

$$= \mathbb{E}_{p(y^t)} \left[\mathbb{E}_{p(\omega | y^t)} \left[\left(\left(\mathbb{E}_{p(y^t | \mathbf{x}, \omega)} [y^t] - \mathbb{E}_{p(\omega | y^t)} \left[\mathbb{E}_{p(y^t | \mathbf{x}, \omega)} [y^t] \right] \right) - \left(\mathbb{E}_{p(y^{t'} | \mathbf{x}, \omega)} [y^{t'}] - \mathbb{E}_{p(\omega | y^t)} \left[\mathbb{E}_{p(y^{t'} | \mathbf{x}, \omega)} [y^{t'}] \right] \right) \right)^2 \right] \right], \quad (22f)$$

$$= \mathbb{E}_{p(y^t)} \left[\mathbb{E}_{p(\omega | y^t)} \left[\left(\underbrace{\left(\mathbb{E}_{p(y^t | \mathbf{x}, \omega)} [y^t] - \mathbb{E}_{p(\omega | y^t)} \left[\mathbb{E}_{p(y^t | \mathbf{x}, \omega)} [y^t] \right] \right)}_{\approx 0} - \left(\mathbb{E}_{p(y^{t'} | \mathbf{x}, \omega)} [y^{t'}] - \mathbb{E}_{p(\omega)} \left[\mathbb{E}_{p(y^{t'} | \mathbf{x}, \omega)} [y^{t'}] \right] \right) \right)^2 \right] \right], \quad (22g)$$

$$\approx \mathbb{E}_{p(y^t)} \left[\mathbb{E}_{p(\omega | y^t)} \left[\left(\mathbb{E}_{p(y^{t'} | \mathbf{x}, \omega)} [y^{t'}] - \mathbb{E}_{p(\omega)} \left[\mathbb{E}_{p(y^{t'} | \mathbf{x}, \omega)} [y^{t'}] \right] \right)^2 \right] \right], \quad (22h)$$

$$= \mathbb{E}_{p(y^t)} \left[\mathbb{E}_{p(\omega | y^t)} \left[\left(\hat{\mu}_\omega(\mathbf{x}, t') - \mathbb{E}_{p(\omega)} [\hat{\mu}_\omega(\mathbf{x}, t')] \right)^2 \right] \right], \quad (22i)$$

$$= \mathbb{E}_{y^t \sim p(Y^t | \mathbf{x}, t, \mathcal{D})} \left[\text{Var}_{\omega \sim p(\Omega | \mathcal{D}_{\text{train}})} (\hat{\mu}_\omega(\mathbf{x}, t')) \right], \quad (22j)$$

where (22a) by definition of variance; (22a)-(22b) by definition of $\hat{\tau}_\omega$; (22b)-(22c) by linearity of expectations; (22c)-(22d) by grouping terms; (22d)-(22e) by symmetry of the square; (22e)-(22f) by rewriting expectations in terms of densities; (22f)-(22g) **the observed potential outcome does not have an effect on the expectation of the model for the counterfactual outcome**; (22g)-(22h) **we drop the term as an approximation as we cannot estimate here how much the expected outcome is going to change—the conservative assumption is that will not change**; (22h)-(22i) by definition of $\hat{\mu}_\omega$; (22i)-(22j) by definition of variance; \square

B Baselines

B.1 S-type error Information Gain

In their work, Sundin et al. [45] assume that the underlying model is a Gaussian Process (GP) and also that they have access to the counterfactual outcome. Although GPs are suitable for uncertainty estimation, they do not scale up to high dimensional datasets (e.g. images). We propose to use Deep Ensembles and DUE for alleviating the capabilities issues and we modified the objective to be more suitable for our architecture.

Following the formulation from Houlsby et al. [18], the acquisition strategy becomes $\arg \max_x \mathbb{H}[\gamma|x, D] - \mathbb{E}_{\mathbb{H}[p(\theta|D)]}[\gamma|x, \theta]$, where $\gamma(x) = \text{probit}^{-1}\left(-\frac{|\mathbf{E}_{p(\tau|x, \mathcal{D}_{\text{train}})}[\tau]|}{\sqrt{\text{Var}(\tau|x, \mathcal{D}_{\text{train}})}}\right)$, $\text{probit}^{-1}(\cdot)$ is the cumulative distribution function of normal distribution and $p(\gamma|x, D) = \text{Bernoulli}(\gamma)$. With DUE (Deep Kernel Learning method) Deep Ensembles (samples from $p(\theta|D)$) we can compute those terms similarly to how we implemented our BALD objectives.

Below is an example of how this was implemented in PyTorch:

```
tau_mu = mu1s - mu0s
tau_var = var1s + var0s + 1e-07
gammas = torch.distributions.normal.Normal(0, 1).cdf(
    -tau_mu.abs() / tau_var.sqrt()
)
gamma = gammas.mean(-1)
predictive_entropy = dist.Bernoulli(gamma).entropy()
conditional_entropy = dist.Bernoulli(gammas).entropy().mean(-1)
# it can get negative very small number
# because of numerical instabilities
scores = (predictive_entropy - conditional_entropy).clamp_min(1e-07)
```

C Datasets

C.1 Synthetic Data

We modify the synthetic dataset presented by Kallus et al. [22]. Our dataset is described by the following structural causal model (SCM):

$$\mathbf{x} := N_{\mathbf{x}}, \quad (23a)$$

$$t := N_t, \quad (23b)$$

$$y := (2t - 1)\mathbf{x} + (2t - 1) - 2\sin(2(2t - 1)\mathbf{x}) + 2(1 + 0.5\mathbf{x}) + N_y, \quad (23c)$$

where $N_{\mathbf{x}} \sim \mathcal{N}(0, 1)$, $N_t \sim \text{Bern}(\text{sigmoid}(2\mathbf{x} + 0.5))$, and $N_y \sim \mathcal{N}(0, 1)$.

Each random realization of the simulated dataset generates 10000 pool set examples, 1000 validation examples, and 1000 test examples. In the experiments we report results over 40 random realizations. The seeds for the random number generators are i , $i + 1$, and $i + 2$; $\{i \in [0, 1, \dots, 19]\}$, for the training, validation, and test sets, respectively.

C.2 IHDP Data.

Infant Health and Development Program (IHDP) is a semi-synthetic dataset [17, 42] commonly used in literature to study the performance of causal effect estimation methods. The dataset consists of 747 cases, out of which 139 are assigned in treatment group and 608 in control. Each unit is represented by 25 covariates describing different aspects of the infants and their mothers. We report results over 200 random realizations of response surface B described by Hill [17].

C.3 CMNIST Data.

Following the setup from [21], we use a simulated dataset based on MNIST [29]. CMNIST is described by the following SCM:

$$\mathbf{x} := N_{\mathbf{x}}, \quad (24a)$$

$$\phi := \left(\text{clip} \left(\frac{\mu_{N_{\mathbf{x}}} - \mu_c}{\sigma_c}; -1.4, 1.4 \right) - \text{Min}_c \right) \frac{\text{Max}_c - \text{Min}_c}{1.4 - -1.4} \quad (24b)$$

$$t := N_t, \quad (24c)$$

$$y := (2t - 1)\phi + (2t - 1) - 2 \sin(2(2t - 1)\phi) + 2(1 + 0.5\phi) + N_y, \quad (24d)$$

where N_t (swapping \mathbf{x} for ϕ), and N_y are as described in Appendix C.1. $N_{\mathbf{x}}$ is a sample of an MNIST image. The sampled image has a corresponding label $c \in [0, \dots, 9]$. $\mu_{N_{\mathbf{x}}}$ is the average intensity of the sampled image. μ_c and σ_c are the mean and standard deviation of the average image intensities over all images with label c in the MNIST training set. In other words, $\mu_c = \mathbb{E}[\mu_{N_{\mathbf{x}}} | c]$ and $\sigma_c^2 = \text{Var}[\mu_{N_{\mathbf{x}}} | c]$. To map the high dimensional images \mathbf{x} onto a one-dimensional manifold ϕ with domain $[-3, 3]$ above, we first clip the standardized average image intensity on the range $(-1.4, 1.4)$. Each digit class has its own domain in ϕ , so there is a linear transformation of the clipped value onto the range $[\text{Min}_c, \text{Max}_c]$. Finally, $\text{Min}_c = -2 + \frac{4}{10}c$, and $\text{Max}_c = -2 + \frac{4}{10}(c + 1)$.

For each random realization of the dataset, the MNIST training set is split into training ($n = 35000$) and validation ($n = 15000$) subsets using the scikit-learn function `train_test_split()`. The test set is generated using the MNIST test set ($n = 10000$). The random seeds are $\{i \in [0, 1, \dots, 19]\}$ for the 10 random realizations generated.

D More Results

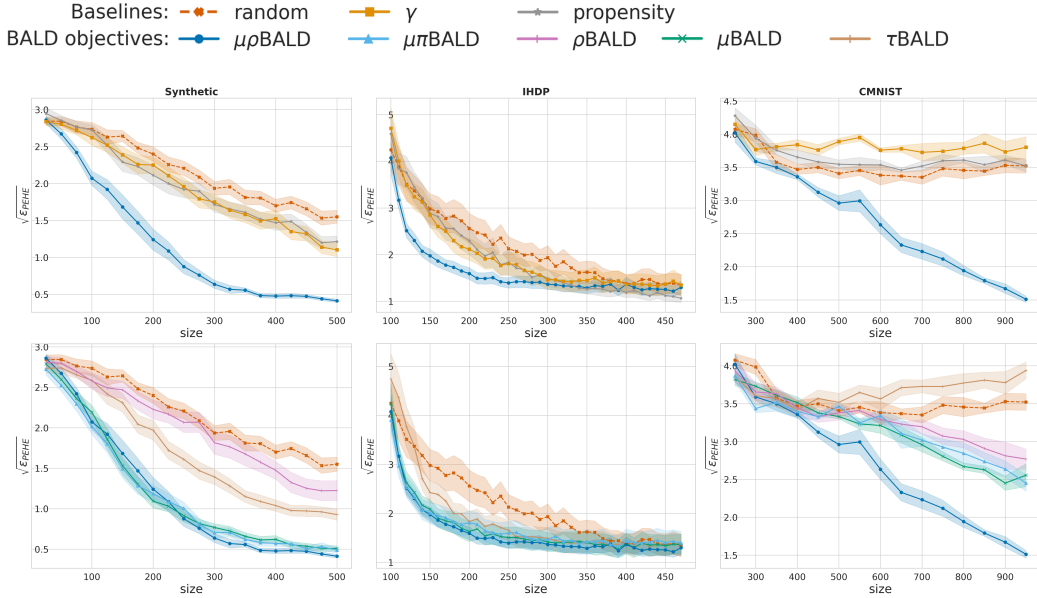


Figure 6: $\sqrt{\epsilon_{PEHE}}$ performance (shaded standard error) for Deep Ensembles based models. **(left to right) synthetic** (20 seeds), **IHDP** (50 seeds) and **CMNIST** (5 seeds) dataset results, **(top to bottom)** comparison with baselines, comparison between BALD objectives. We observe that BALD objectives outperform the **random**, **γ** and **propensity** acquisition functions significantly, suggesting that epistemic uncertainty aware methods that target reducible uncertainty can be more sample efficient.

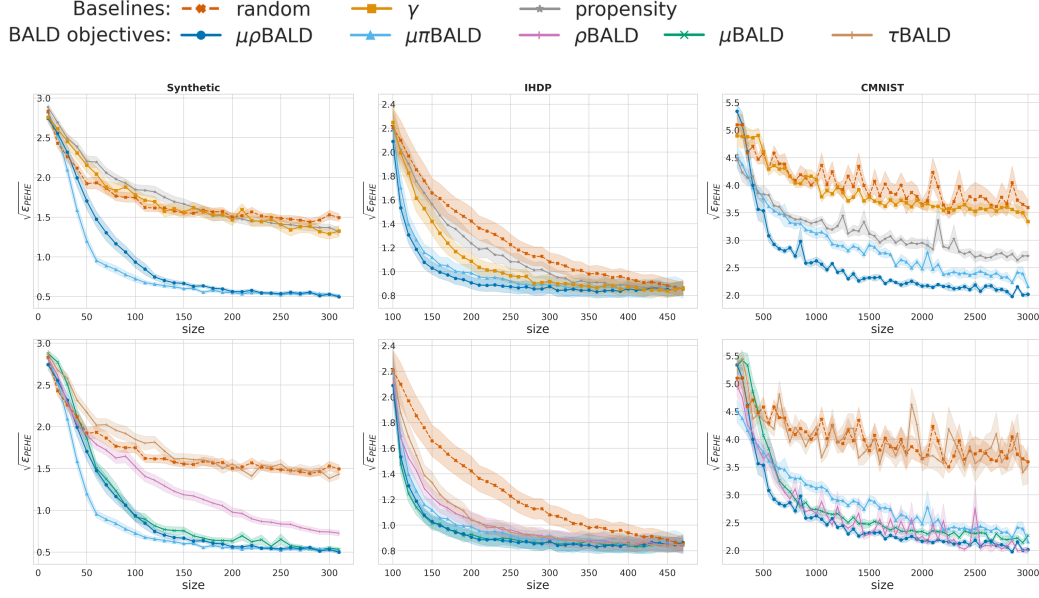


Figure 7: $\sqrt{\epsilon_{PHE}}$ performance (shaded standard error) for DUE models. **(left to right) synthetic** (40 seeds), and **IHDP** (200 seeds). We observe that BALD objectives outperform the **random**, γ and **propensity** acquisition functions significantly, suggesting that epistemic uncertainty aware methods that target reducible uncertainty can be more sample efficient.

E Compute

We used a cluster of 8 nodes with 4 GPUs each (16 RTX 2080 and 16 Titan RTX). The total GPU hours is estimated to be:

8 baselines x (.5 + 1 + 1) days per dataset x (5 ensemble components * 0.25 GPU usage + 1 DUE * 0.3 GPU usage) x 24 hours = 744 GPU hours

Code is written in python. Packages used include PyTorch [35], scikit-learn [36], Ray [32], NumPy, SciPy, and Matplotlib.

F Model Architectures

For deep ensembles, we use an ensemble of TarNETs [42]. For Due, we append the treatment variable to the features extracted, then define the GP over that input. **Synthetic Architecture**

Layer (type:depth-idx)	Output Shape
Sequential	--
NeuralNetwork: 1-1	[64, 100]
Sequential: 2-1	[64, 100]
Linear: 3-1	[64, 100]
ResidualDense: 3-2	[64, 100]
PreactivationDense: 4-1	[64, 100]
Sequential: 5-1	[64, 100]
Activation: 6-1	[64, 100]
Linear: 6-2	[64, 100]
Identity: 4-2	[64, 100]
ResidualDense: 3-3	[64, 100]
PreactivationDense: 4-3	[64, 100]
Sequential: 5-2	[64, 100]
Activation: 6-3	[64, 100]
Linear: 6-4	[64, 100]
Identity: 4-4	[64, 100]

	└ResidualDense: 3-4	[64, 100]
	└└PreactivationDense: 4-5	[64, 100]
	└└└Sequential: 5-3	[64, 100]
	└└└└Activation: 6-5	[64, 100]
	└└└└Linear: 6-6	[64, 100]
	└└Identity: 4-6	[64, 100]
	└Activation: 3-5	[64, 100]
	└└Sequential: 4-7	[64, 100]
	└└└Identity: 5-4	[64, 100]
	└└└LeakyReLU: 5-5	[64, 100]
	└└└Dropout: 5-6	[64, 100]
└GMM: 1-2		[64, 5]
└└Linear: 2-2		[64, 5]
└└Linear: 2-3		[64, 5]
└└Sequential: 2-4		[64, 5]
└└└Linear: 3-6		[64, 5]
└└└Softplus: 3-7		[64, 5]

=====

Total params: 32,115

IHDP Architecture

Layer (type:depth-idx)	Output Shape
Sequential	--
└TARNet: 1-1	[64, 400]
└└NeuralNetwork: 2-1	[64, 400]
└└└Sequential: 3-1	[64, 400]
└└└└Linear: 4-1	[64, 400]
└└└└ResidualDense: 4-2	[64, 400]
└└└└└PreactivationDense: 5-1	[64, 400]
└└└└└└Sequential: 6-1	[64, 400]
└└└└└└Identity: 5-2	[64, 400]
└└└└ResidualDense: 4-3	[64, 400]
└└└└└PreactivationDense: 5-3	[64, 400]
└└└└└└Sequential: 6-2	[64, 400]
└└└└└└Identity: 5-4	[64, 400]
└└Sequential: 2-2	[64, 400]
└└└ResidualDense: 3-2	[64, 400]
└└└└PreactivationDense: 4-4	[64, 400]
└└└└└Sequential: 5-5	[64, 400]
└└└└└└Activation: 6-3	[64, 401]
└└└└└└Linear: 6-4	[64, 400]
└└└└Sequential: 4-5	[64, 400]
└└└└└Dropout: 5-6	[64, 401]
└└└└└Linear: 5-7	[64, 400]
└└└ResidualDense: 3-3	[64, 400]
└└└└PreactivationDense: 4-6	[64, 400]
└└└└└Sequential: 5-8	[64, 400]
└└└└└└Activation: 6-5	[64, 400]
└└└└└└Linear: 6-6	[64, 400]
└└└└Identity: 4-7	[64, 400]
└└Activation: 3-4	[64, 400]
└└└Sequential: 4-8	[64, 400]
└└└└Identity: 5-9	[64, 400]
└└└└ELU: 5-10	[64, 400]
└└└└Dropout: 5-11	[64, 400]
└GMM: 1-2	[64, 5]
└└Linear: 2-3	[64, 5]
└└Linear: 2-4	[64, 5]
└└Sequential: 2-5	[64, 5]
└└└Linear: 3-5	[64, 5]
└└└Softplus: 3-6	[64, 5]

=====

CMNIST Architecture

Layer (type:depth-idx)	Output Shape
Sequential	--
└TARNet: 1-1	[200, 100]
└└ResNet: 2-1	[200, 48]
└└└Sequential: 3-1	[200, 48, 1, 1]
└└└└Conv2d: 4-1	[200, 12, 28, 28]
└└└└Identity: 4-2	[200, 12, 28, 28]
└└└└ResidualConv: 4-3	[200, 12, 28, 28]
└└└└└Sequential: 5-1	[200, 12, 28, 28]
└└└└└└PreactivationConv: 6-1	[200, 12, 28, 28]
└└└└└└PreactivationConv: 6-2	[200, 12, 28, 28]
└└└└└Sequential: 5-2	[200, 12, 28, 28]
└└└└└└Dropout2d: 6-3	[200, 12, 28, 28]
└└└└└└Conv2d: 6-4	[200, 12, 28, 28]
└└└└ResidualConv: 4-4	[200, 24, 14, 14]
└└└└└Sequential: 5-3	[200, 24, 14, 14]
└└└└└└PreactivationConv: 6-5	[200, 12, 28, 28]
└└└└└└PreactivationConv: 6-6	[200, 24, 14, 14]
└└└└└Sequential: 5-4	[200, 24, 14, 14]
└└└└└└Dropout2d: 6-7	[200, 12, 28, 28]
└└└└└└Conv2d: 6-8	[200, 24, 14, 14]
└└└└ResidualConv: 4-5	[200, 24, 14, 14]
└└└└└Sequential: 5-5	[200, 24, 14, 14]
└└└└└└PreactivationConv: 6-9	[200, 24, 14, 14]
└└└└└└PreactivationConv: 6-10	[200, 24, 14, 14]
└└└└└Sequential: 5-6	[200, 24, 14, 14]
└└└└└└Dropout2d: 6-11	[200, 24, 14, 14]
└└└└└└Conv2d: 6-12	[200, 24, 14, 14]
└└└└ResidualConv: 4-6	[200, 48, 7, 7]
└└└└└Sequential: 5-7	[200, 48, 7, 7]
└└└└└└PreactivationConv: 6-13	[200, 24, 14, 14]
└└└└└└PreactivationConv: 6-14	[200, 48, 7, 7]
└└└└└Sequential: 5-8	[200, 48, 7, 7]
└└└└└└Dropout2d: 6-15	[200, 24, 14, 14]
└└└└└└Conv2d: 6-16	[200, 48, 7, 7]
└└└└ResidualConv: 4-7	[200, 48, 7, 7]
└└└└└Sequential: 5-9	[200, 48, 7, 7]
└└└└└└PreactivationConv: 6-17	[200, 48, 7, 7]
└└└└└└PreactivationConv: 6-18	[200, 48, 7, 7]
└└└└└Sequential: 5-10	[200, 48, 7, 7]
└└└└└└Dropout2d: 6-19	[200, 48, 7, 7]
└└└└└└Conv2d: 6-20	[200, 48, 7, 7]
└└└└ResidualConv: 4-8	[200, 48, 7, 7]
└└└└└Sequential: 5-11	[200, 48, 7, 7]
└└└└└└PreactivationConv: 6-21	[200, 48, 7, 7]
└└└└└└PreactivationConv: 6-22	[200, 48, 7, 7]
└└└└└Sequential: 5-12	[200, 48, 7, 7]
└└└└└└Dropout2d: 6-23	[200, 48, 7, 7]
└└└└└└Conv2d: 6-24	[200, 48, 7, 7]
└└└└AdaptiveAvgPool2d: 4-9	[200, 48, 1, 1]
└└└Sequential: 2-2	[200, 100]
└└└└ResidualDense: 3-2	[200, 100]
└└└└└PreactivationDense: 4-10	[200, 100]
└└└└└└Sequential: 5-13	[200, 100]
└└└└└└└Activation: 6-25	[200, 49]
└└└└└└└Linear: 6-26	[200, 100]
└└└└└Sequential: 4-11	[200, 100]
└└└└└└Dropout: 5-14	[200, 49]
└└└└└└Linear: 5-15	[200, 100]
└└└└ResidualDense: 3-3	[200, 100]
└└└└└PreactivationDense: 4-12	[200, 100]

			Sequential: 5-16	[200, 100]
			Activation: 6-27	[200, 100]
			Linear: 6-28	[200, 100]
		Identity: 4-13		[200, 100]
	Activation: 3-4			[200, 100]
	Sequential: 4-14			[200, 100]
	Identity: 5-17			[200, 100]
	LeakyReLU: 5-18			[200, 100]
	Dropout: 5-19			[200, 100]
GMM: 1-2				[200, 5]
	Linear: 2-3			[200, 5]
	Linear: 2-4			[200, 5]
	Sequential: 2-5			[200, 5]
	Linear: 3-5			[200, 5]
	Softplus: 3-6			[200, 5]

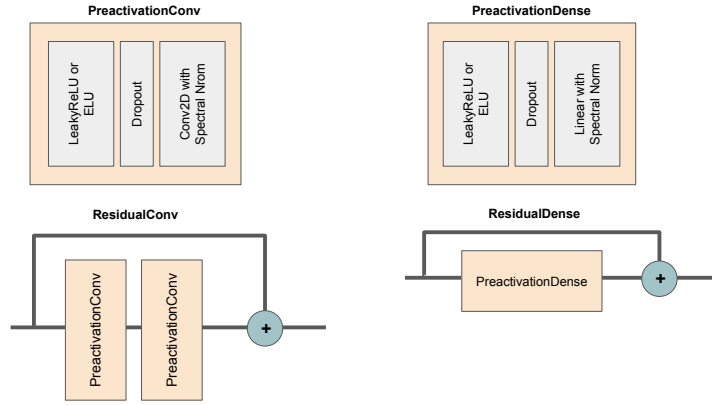


Figure 8: PreactivationConv is a convolution layer with LeakyReLU (or ELU when slope is negative) activation, dropout and spectral norm applied [15, 31]. Similarly, PreactivationDense is a dense layer with BatchNorm [19], LeakyReLU (or ELU when slope is negative) activation and spectral norm applied [15, 31]. ResidualConv is the residual convolution layer, defined as $\text{PreactivationConv}(\text{PreactivationConv}(x)) + \text{SpectralNorm}(1 \times 1 \text{Conv}(x))$ and ResidualDense are residual dense layers, defined as $\text{PreactivationDense}(x) + x$.

All experiments were trained using Adam optimizer [24].

E.1 Hyper-parameters

We use ray tune [30] with the hyperopt [6] search algorithm to optimize our network hyper-parameters. The hyper-parameter search spaces are given in Table 2 and Table 3. The hyper-parameter optimization objective for each dataset is the expected batch-wise log-likelihood of the validation data for a single dataset realization with random seed 1331. The final hyper-parameters are given in Table 4 and Table 5.

Table 2: Hyper-parameter search space for **Deep Ensemble**

Hyper-parameter	Search Space
hidden units	[100, 200, 400]
network depth	[2, 3, 4]
negative slope	[ReLU [2], 0.1, 0.2, ELU [7]]
dropout rate	[0.05, 0.1, 0.2, 0.5]
spectral norm	[None, 0.95, 1.5, 3.0]
batch size	[32, 64, 100, 200]
learning rate	[2e-4, 5e-4, 1e-3]

Table 3: Hyper-parameter search space for **DUE**

Hyper-parameter	Search Space
kernel	[RBF, Matern]
ν (Matern)	[0.5, 1.5, 2.5]
inducing points	[20, 50, 100, 200]
hidden units	[100, 200, 400]
network depth	[2, 3, 4]
negative slope	[ReLU [2], 0.1, 0.2, ELU [7]]
dropout rate	[0.05, 0.1, 0.2, 0.5]
spectral norm	[None, 0.95, 1.5, 3.0]
batch size	[32, 64, 100, 200]
learning rate	[2e-4, 5e-4, 1e-3]

Table 4: Training hyper parameters for **Deep Ensemble** experiments

Parameter	Synthetic	IHDP	CMNIST
dim hidden	100	400	100
dropout	0.0	0.15	0.1
depth	4	3	3
spectral norm	12	0.95	24
learning rate	0.001	0.001	0.001
non-linearity	ReLU	ELU	ReLU

Table 5: Training hyper parameters for **DUE** experiments

Parameter	Synthetic	IHDP	CMNIST
kernel	RBF	Matern ($\nu = 1.5$)	RBF
inducing points	20	100	100
dim hidden	100	200	200
dropout	0.2	0.1	0.05
depth	3	3	2
batch size	200	100	64
spectral norm	0.95	0.95	3.0
learning rate	0.001	0.001	0.001
non-linearity	ReLU	ELU	ELU