

# Topology-Aware Robust Representation Balancing for Estimating Causal Effects

Anonymous Authors<sup>1</sup>

## Abstract

Representation learning in high-dimensional spaces faces significant robustness challenges with noisy inputs, particularly with heavy-tailed noise. Arguing that topological data analysis (TDA) offers a solution, we leverage TDA to enhance representation stability in neural networks. Our theoretical analysis establishes conditions under which incorporating topological summaries improves robustness to input noise, especially for heavy-tailed distributions. Extending these results to representation-balancing methods used in causal inference, we propose the *Topology-Aware Treatment Effect Estimation* (TATEE) framework, through which we demonstrate how topological awareness can lead to learning more robust representations. A key advantage of this approach is that it requires no ground-truth or validation data, making it suitable for observational settings common in causal inference. The method remains computationally efficient with overhead scaling linearly with data size while staying constant in input dimension. Through extensive experiments with  $\alpha$ -stable noise distributions, we validate our theoretical results, demonstrating that TATEE consistently outperforms existing methods across noise regimes. This work extends stability properties of topological summaries to representation learning via a tractable framework scalable for high-dimensional inputs, providing insights into how it can enhance robustness, with applications extending to domains facing challenges with noisy data, such as causal inference.

## 1. Introduction

Robust representation learning is critical across domains, including those involving high-dimensional data, yet developing methods that handle noisy inputs remains challenging,

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

especially for heavy-tailed noise—prevalent in finance and signal processing (Barkat & Stanković, 2004; Kim et al., 2008; Stoyanov et al., 2011; Simsekli et al., 2019; Gorbunov et al., 2020; Yang et al., 2022). Techniques for improving robustness without requiring ground-truth are particularly valuable for observational settings where validation data is unavailable, typical in causal inference applications. Topological Data Analysis (TDA) offers a principled solution through the stability properties of topological summaries, enabling a purely structural approach to robustness. We leverage these properties to enhance the robustness of representation learning, demonstrating their effectiveness through representation-balancing neural networks for treatment effect estimation (Johansson et al., 2016; Shalit et al., 2017; Kazemi & Ester, 2024; Wang et al., 2024). Our approach addresses the challenge of representations’ robustness to noise, and building on our theoretical results, we introduce a topology-aware framework for treatment effect estimation that demonstrates these stability benefits while maintaining computational scalability with high-dimensional inputs.

Machine learning methods have shown promise for causal inference (Morgan & Winship, 2015; Louizos et al., 2017; Yoon et al., 2018; Shi et al., 2019; Cui et al., 2020; Shi et al., 2021; Ghosh et al., 2023), with representation-balancing approaches gaining popularity for treatment effect estimation (Johansson et al., 2016; Shalit et al., 2017). However, noisy observations present a key challenge (Wickens, 1972; Kuroki & Pearl, 2014), and despite progress (Kallus et al., 2018; Shu & Yi, 2020), techniques for handling noise remain limited, with significant gaps in addressing non-Gaussian noise. Our work addresses this gap by integrating topological summaries in treatment effect estimation, enhancing robustness in a scalable fashion. Topological Data Analysis (TDA) offers a compelling approach to robust representation learning by characterizing the shape of data through its global topology, which remains stable under local geometric perturbations (Carlsson, 2009). A standard tool for this characterization is persistent homology, which effectively summarizes topological invariants (Edelsbrunner et al., 2002; Edelsbrunner & Harer, 2008), with well-established stability theorems demonstrating how the resulting topological signatures can enhance robustness in data analysis pipelines (Cohen-Steiner et al., 2005; 2007; 2010). Recent work suggests that incorporating these topological

summaries into deep learning frameworks can improve their resilience to noise (Gabrielsson et al., 2020; Southern et al., 2023), motivating our work and making it the first to bridge TDA and treatment effect estimation.

We leverage the stability properties of topological summaries to enhance the robustness of representations against noise through an approach that is both purely structural and computationally scalable—critical advantages for observational and high-dimensional settings. In Section 3, we derive conditions under which topological summaries enhance representations’ robustness by improving metric stability, especially with heavy-tailed noise. This underpins our *Topology-Aware Treatment Effect Estimation (TATEE)* framework proposed in Section 4, which improves the robustness of counterfactual regression (CFR) (Shalit et al., 2017) by imposing both topological and distributional similarities between the representations of treatment and control groups. TATEE’s implementation is scalable, with complexity linear in data size and constant in input dimension, once the dimensionality of the representations is fixed—making it suitable for high-dimensional applications. Our experiments confirm that TATEE consistently outperforms existing methods across a range of  $\alpha$ -stable noise distributions, including Gaussian and heavy-tailed cases. In conclusion, we establish how topological properties enhance representations’ robustness under noise, provide theoretical conditions for this improvement, and demonstrate the benefits in practice.

### Main Contributions:

**A) Conceptual.** We argue that incorporating topological awareness into representation learning offers a principled path to robust deep learning. Our approach is computationally scalable, making it suitable for high-dimensional data, and requires no ground-truth or validation data—a critical advantage for observational settings. To our knowledge, this is the first work to integrate TDA with representation-balancing neural networks for causal inference, demonstrating its effectiveness across diverse noise regimes.

**B) Theoretical.** We establish new stability-type results for representations learned from noisy data, extending foundational stability theorems from TDA to deep learning frameworks. We identify regimes where persistent homology enhances metric stability in representation learning, particularly under heavy-tailed noise, providing a rigorous foundation for the robustness gains achieved by our method.

**C) Methodological.** We introduce *Topology-Aware Treatment Effect Estimation (TATEE)*, a scalable framework that integrates persistence diagrams into representation balancing to improve counterfactual regression’s robustness to noise. Extensive experiments across  $\alpha$ -stable noise distributions, including Gaussian and heavy-tailed, validate TATEE’s ability to meet the conditions for robustness in practical settings, outperforming existing methods.

## 2. Preliminaries and Related Work

### Topological Data Analysis and Persistent Homology.

Topological Data Analysis (TDA) utilizes algebraic topology to extract shape-based features from data across scales (Carlsson, 2009). Persistent homology, a central tool in TDA, captures topological features, such as connected components and holes (Edelsbrunner & Harer, 2008). These features are characterized by their lifespan through varying scales, represented as points in a persistence diagram, where each point corresponds to a feature’s birth and death in a filtration of the input—a series of nested simplicial complexes determined by the image of a filtration function (Ghrist, 2008). Figure 1 visualizes this concept, with further details in Appendix A. Resulting from the *stability theorem* and fundamental to our work, stability of persistence diagrams is a key property ensuring robustness of these topological summaries to perturbations in data (Cohen-Steiner et al., 2007). Consider a triangulable compact metric space  $(\mathcal{Z}, d)$  for some metric  $d$ , and let  $f, g : \mathcal{Z} \rightarrow \mathbb{R}$  be Lipschitz filtration functions. The following generalization of the stability theorem holds under conditions in Appendix A.3.

**Theorem 2.1** (Cohen-Steiner et al. (2010)). *For some constants  $k \geq 1$  and  $C$ , we have for all  $l$ ,*

$$W_p(\mu_f^l, \mu_g^l) \leq C^{\frac{1}{p}} \|f - g\|_{\infty}^{1 - \frac{k}{p}}, \quad (1)$$

where  $\mu_f^l$  and  $\mu_g^l$  are measures on the space of persistence diagrams corresponding to the  $l$ -dimensional homology classes (Mileyko et al., 2011),  $W_p$  denotes the Wasserstein- $p$  distance, and constants  $C$  and  $k$  are described in Appendix A.3. Since this inequality holds for all  $l$ , we shall drop the superscript  $l$  from here on for ease of notation.

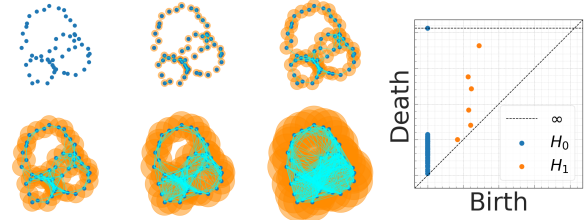


Figure 1: A visualization of the Vietoris-Rips filtration of a pointcloud and the corresponding persistence diagram. The zeroth and first homology groups ( $H_0$  and  $H_1$ ) correspond to the connected components and 1-dimensional holes.

**Causal Inference.** Causal inference aims to determine the effect of a treatment on an outcome  $Y$  given covariates  $X$ . This effect can be quantified by the conditional average treatment effect (CATE), when conditioned on features. For individuals with covariates  $x$ , CATE is given by

$$\tau(x) = \mathbb{E}[Y(1) - Y(0) | X = x], \quad (2)$$

where  $Y(1)$  and  $Y(0)$  denote the potential outcomes. The challenge in estimating  $\tau(x)$  arises from the unobservability of *counterfactual* outcomes, leading to *the fundamental problem of causal inference*.

**Related Work.** Causal inference and deep learning have been integrated in various contexts (Cui et al., 2020; Luo et al., 2020; Schölkopf et al., 2021), with representation-balancing frameworks showing effectiveness in treatment effect estimation (Louizos et al., 2017; Shalit et al., 2017) and counterfactual reasoning (Johansson et al., 2016; Pawlowski et al., 2020). Prior work on robust causal effect estimation (Kallus et al., 2018; Shu & Yi, 2020) provides important advances but assumes finite variance noise and often requires large observation counts or validation data. Despite recent progress (Lagemann et al., 2023; Pöllänen & Marttinen, 2023), methods for robust treatment effect estimation with heavy-tailed noise or limited data remain underdeveloped. TDA, which captures the intrinsic shape of data (Carlsson, 2009; Chazal & Michel, 2021), has been integrated into machine learning to improve robustness (Bronstein et al., 2017; Gabrielsson & Carlsson, 2019; Papamarkou et al., 2024). We extend these applications to representation balancing for causal inference, introducing a framework that enhances robustness without ground-truth, clean, or large datasets, and accommodates noise beyond finite-variance distributions.

### 3. Learning Robust Representations via Persistence Diagrams

How can we improve the robustness of representations learned by neural networks under noise—particularly heavy-tailed noise? Here, we introduce a new stability result characterizing when incorporating persistent homology into learning improves metric stability. Our theorem establishes a condition on neural networks’ Lipschitz constants that leads to persistence diagrams of representations being more robust to noise than the raw representations themselves.

**Problem Setup.** Let  $X$  and  $E$  be random variables representing features and noise, respectively, and define the noise-corrupted features as  $\tilde{X} := X + E$ . Denote by  $\mathbf{X}, \mathbf{E}, \tilde{\mathbf{X}}$  the corresponding finite-sample matrices. We are interested in the stability of the representations  $\varphi(X)$  and  $\varphi(\tilde{X})$  where  $\varphi : \mathcal{X} \rightarrow \mathcal{Z}$ , is a neural network mapping to a representation space  $\mathcal{Z}$ . Observe that  $\varphi(X)$  and  $\varphi(\tilde{X})$  induce measures on  $\mathcal{Z}$ , which we denote by  $\mu$  and  $\tilde{\mu}$ . Furthermore, suppose that there exist filtration functions  $f$  and  $\tilde{f}$ , which satisfy  $\tilde{f}(\varphi(X)) = f(\varphi(\tilde{X}))$ , yielding persistence diagrams for both clean and noisy representations and allowing us to invoke the stability theorem. The explicit construction of such filtration functions is provided in Appendix B. While the filtration is typically fixed once the data is given, the network  $\varphi$  remains trainable. We therefore seek conditions on  $\varphi$  that improve robustness of the persistence diagrams.

Importantly, we exploit the Lipschitz continuity of standard neural networks (Virmaux & Scaman, 2018; Gouk et al., 2021) to derive the intended condition on  $\varphi$ .

**Finite-Sample Stability and Noise Distribution.** Let  $\varphi$  be a  $K_\varphi$ -Lipschitz network. Then Lipschitz continuity arguments yield a finite-sample upper bound  $\hat{M}$  on the Wassertein distance  $W_p(\mu, \tilde{\mu})$  between the representations of clean and noisy inputs. Similarly, under standard assumptions on the filtration functions, Theorem 2.1 provides a bound  $\hat{K}_{\text{topo}}$  on the Wasserstein distance between their persistence diagrams,  $W_p(\mu_f, \mu_{\tilde{f}})$ . These bounds, detailed in Appendix B.1, lead to the following.

**Theorem 3.1.** *If  $K_\varphi < \Lambda$ , then  $\hat{K}_{\text{topo}} < \hat{M}$ , where  $\Lambda$  depends on the noise distribution. In particular,  $\Lambda$  is increasing in  $\|\mathbf{E}\|_\infty^{p/k-1} \|\bar{\mathbf{E}}\|^{-p/k}$ , where  $\|\bar{\mathbf{E}}\|$  and  $\|\mathbf{E}\|_\infty$  are the sample average and  $\infty$ -norm of the error.*

In the statement above,  $p$  is the degree of the Wasserstein distance and  $k$  is the constant in Theorem 2.1. The proof and details about  $\Lambda$  are provided in Appendix B.1. Intuitively, the condition asserts that if the Lipschitz constant of  $\varphi$  is smaller than  $\Lambda$ , the upper bound  $\hat{K}_{\text{topo}}$  is smaller than that on the representations,  $\hat{M}$ . This suggests that with the appropriate neural network, the Wasserstein space over the persistence diagrams of the representations is more robust than that over the representations themselves—a metric stability which enhances robustness to input noise. Notably, the condition that determines if a neural network is an ‘appropriate’ one depends on the distribution of the error, particularly, its tail. Theorem 2.1 leads to the promised stability properties for  $p > k$ . In this case, the ratio  $\|\mathbf{E}\|_\infty^{p/k-1} \|\bar{\mathbf{E}}\|^{-p/k}$  is larger for heavy-tailed distributions. Since  $\Lambda$  is linear in this ratio, a slow-decaying tail of the empirical distribution of the noise corresponds to a more easily achievable neural network that satisfies the condition in Theorem 3.1.

**Implications.** Theorem 3.1 establishes that topological summaries can enhance representations’ robustness, particularly with heavy-tailed noise. This insight directly informs our approach to robust causal inference through TATEE. Our experiments in Section 5 confirm that neural networks in TATEE can be trained to satisfy the theoretical conditions and achieve the predicted robustness benefits in practice.

### 4. Topology-Aware Treatment Effect Estimation

Building on the stability results from Section 3, we introduce *Topology-Aware Treatment Effect Estimation* (TATEE), which incorporates topological awareness into representation-balancing neural networks for estimating causal effects. TATEE enhances robustness to input noise by leveraging the stability properties of persistence diagrams in a scalable fashion, showcasing the benefits of topologi-



cal summaries for robust representation learning. Here, we discuss the main aspects of TATEE’s design, analysis, and implications. More details are included in Appendix C.

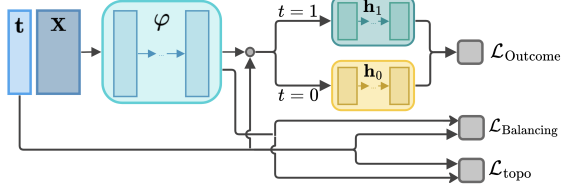


Figure 2: The neural network architecture and loss terms used in TATEE, which adopts the two-headed network, outcome loss, and balancing loss from CFR (Shalit et al., 2017). The topological signature is incorporated as the regularization term  $\mathcal{L}_{\text{topo}}$ , based on the output of  $\varphi$ .

**Architecture and Training.** TATEE incorporates topological awareness in representation-balancing neural networks for counterfactual regression (CFR) (Shalit et al., 2017) through a regularization term in the training objective, using the Wasserstein distance between persistence diagrams of treatment and control representations. The architecture follows CFR’s two-headed design with a shared encoder  $\varphi$  that maps inputs to a representation space, from which two separate heads  $h_1$  and  $h_0$  estimate the potential outcomes for the treatment and control groups. In addition to prediction accuracy, the training objective  $\mathcal{L}_{\text{TATEE}} = \mathcal{L}_{\text{Outcome}} + \lambda \mathcal{L}_{\text{Balance}} + \lambda_{\text{topo}} \mathcal{L}_{\text{topo}}$  encourages distributional and topological similarities between the inputs of  $h_1$  and  $h_0$ . Complete implementation details are provided in Appendix C.2. This implementation ensures computational scalability by computing persistence diagrams on mini-batches and applying topological regularization to representations rather than raw inputs, resulting in linear overhead with respect to data volume and constant overhead with input dimensionality, as long as the dimensionality of the representations is fixed. Figure 3 illustrates this via a simulation, where overhead elapsed time remains nearly constant as input dimensionality increases from 16 to 256 (with a fixed representation dimensionality), and scales linearly as data volume varies from 800 to 12800 samples.

**Robustness of TATEE.** Using the stability results presented before, we show that TATEE can improve the robustness of counterfactual regression under the conditions of Theorem 3.1. In particular, this theorem implies that TATEE’s training objective is more stable under input noise than the original CFR’s. This is shown by Proposition C.1 in Appendix C, which is informally stated below.

**Proposition 4.1 (Informal).** *If the neural network  $\varphi$  satisfies the constraint from Theorem 3.1, then the upper bound of the noise-induced change in TATEE’s training objective is smaller than CFR’s.*

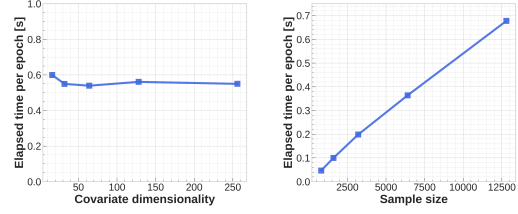


Figure 3: The average overhead elapsed time for computing  $\mathcal{L}_{\text{topo}}$  per epoch of training TATEE. As the dimensionality of the input (left) and input volume (right) increases, the overhead cost remains constant and scales linearly, respectively. The elapsed times are averaged over 20 runs.

As in Theorem 3.1, the constraint on the neural network for achieving this robustness becomes more permissive when the noise distribution has a heavier tail. Figure 4 demonstrates TATEE’s enhanced robustness in a simple example where treatment and control groups have distinct topologies (line vs. circle). While CFR enforces distributional similarity but allows topological divergence, TATEE enforces both distributional and topological similarities. Notably, noise considerably impacts CFR’s ability to enforce distributional similarity, while TATEE achieves its objective equally well in both noisy and clean environments.

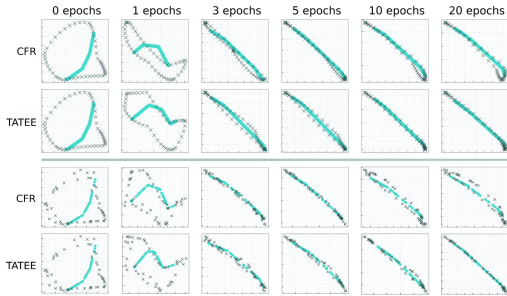


Figure 4: The representations learned by  $\varphi$  in CFR and TATEE throughout 20 epochs of training for the control (gray) and treatment (turquoise) groups, starting without input noise (top rows) and with Gaussian noise (bottom rows). Epoch 0 shows the representations before training.

## 5. Experimental Results

We evaluate TATEE’s capability to enhance robustness in neural networks for causal effect estimation, as indicated by our theoretical analysis. Our experiments across standard causal inference benchmarks confirm that incorporating topological awareness consistently improves robustness to input noise compared to CFR and other deep learning methods for treatment effect estimation. Details of the experimental setup, thorough discussion of the results, and additional experiments are included in Appendix D.

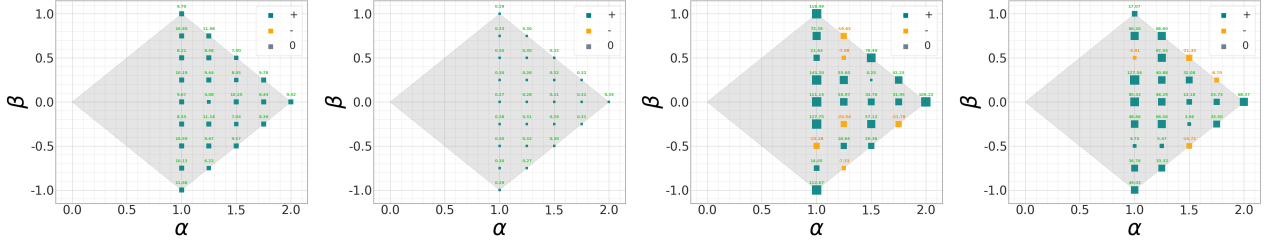


Figure 5: Evaluation of TATEE’s robustness when estimating CATE on noise-corrupted features from the IHDP, Twins, Jobs, and ACIC datasets (from left to right) for a range of  $(\alpha, \beta)$  parameters of  $\alpha$ -stable distributed noise. The Feller-Takayasu diamond (shaded) marks valid  $(\alpha, \beta)$  values. The relative gain in robustness from using TATEE, quantified by  $\rho_{\text{TATEE}} - \rho_{\text{CFR}}$ . The green and orange squares mark positive/negative gain; the size of each square is proportional to the magnitude.

**Experimental Setup.** We evaluate TATEE’s robustness to various input noises sampled from the family of  $\alpha$ -stable distributions, characterized with tail and skewness parameters— $\alpha$  and  $\beta$ . We use standard causal inference benchmarks (IHDP, Twins, Jobs, and ACIC) and measure both performance and robustness. The performance is evaluated via CATE estimation error, which is quantified by the *Precision in Estimation of Heterogeneous Effect* (PEHE), denoted  $\epsilon_{\text{PEHE}}$ , following the conventions in causal inference (e.g., Hill (2011); Louizos et al. (2017); Shi et al. (2019)). We evaluate robustness using  $\rho = 1 - (\epsilon_{\text{PEHE}} \text{ with noise} / \epsilon_{\text{PEHE}} \text{ without noise})$ , quantifying resistance to performance degradation under noise. Complete details, an elaborate discussion of these results, and additional results are provided in Appendix D.

**Main Results.** Our experiments confirm that TATEE achieves superior robustness to input noise while maintaining comparable performance without noise. Figure 5 shows this for our main comparison against TATEE’s counterpart without topological awareness—the original CFR—with  $\alpha$ -stable noise distributions across 25 valid  $(\alpha, \beta)$  values on four datasets. The gains are most pronounced with heavier-tailed noise ( $\alpha$  closer to 1), reaching near 12, 0.3, 51.0, and 68.6 percents for the IHDP, Twins, Jobs, and ACIC datasets (respectively). Moreover, pairwise comparisons against seven other causal inference methods further confirm TATEE’s superior robustness across various noise conditions and datasets. This is demonstrated via the matrix visualized in Figure 6, whose  $(i, j)$  entry shows the proportion of dataset-parameter pairs in which the model corresponding to row  $i$  is more robust than the one for column  $j$ , evaluated by having a weakly larger  $\rho$ .<sup>1</sup> Complete performance metrics and detailed comparisons are provided in Appendix D.

## 6. Discussion

We proposed incorporating topological awareness into representation learning to enhance robustness against input noise, particularly with heavy-tailed distributions. By lever-

<sup>1</sup>We consider a model to be weakly more robust than itself in all cases, hence the value 1.0 on the diagonals.

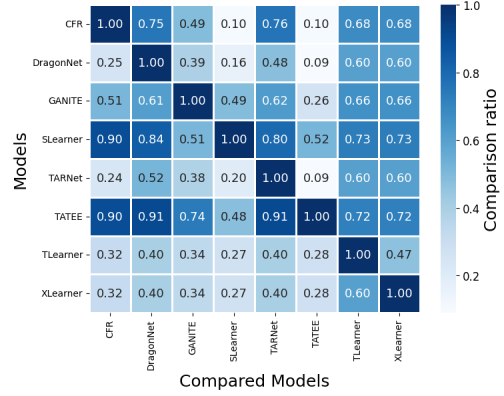


Figure 6: Comparing the robustness of TATEE with 7 other methods. Each entry shows the proportion of noise parameter and dataset cases (over 25  $(\alpha, \beta)$  values for 4 datasets) where the  $\rho$  for the model in the row is at least as large as the  $\rho$  for the model in the column. A larger value in row  $i$  and column  $j$  means method  $i$  is more robust than method  $j$  in a larger proportion of the cases. TATEE achieves more large values in its row than the other methods, with the largest row average of 0.80, indicating its superior robustness.

aging the stability properties of persistence diagrams, we showed that using topological summaries can improve representation stability in a scalable fashion without requiring ground-truth or validation data—critical advantages for observational settings. This concept is rooted in foundational stability theorems from TDA and extended through our work to neural networks and the TATEE framework we introduced for treatment effect estimation. Our theoretical analysis establishes conditions for topological awareness to enhance metric stability, especially with heavy-tailed noise, and our experiments validate that TATEE meets these conditions in practice, consistently outperforming existing methods across noise regimes. While we demonstrated our arguments through a causal inference framework, the theoretical results underpinning TATEE have broader implications for robust deep learning. Future work can explore additional topological features and investigate topology-aware methods across a wider range of representation learning problems.

## References

- Athey, S. and Wager, S. Estimating treatment effects with causal forests: An application. Observational Studies, 5 (2):37–51, 2019.
- Barkat, B. and Stanković, L. Analysis of polynomial fm signals corrupted by heavy-tailed noise. Signal Processing, 84(1):69–75, 2004.
- Bauer, U. Ripser: efficient computation of Vietoris-Rips persistence barcodes. Journal of Applied and Computational Topology, 5(3):391–423, 2021.
- Bronstein, M. M., Bruna, J., LeCun, Y., Szlam, A., and Vandergheynst, P. Geometric deep learning: going beyond euclidean data. IEEE Signal Processing Magazine, 34 (4):18–42, 2017.
- Bubenik, P. et al. Statistical topological data analysis using persistence landscapes. Journal of Machine Learning Research, 16(1):77–102, 2015.
- Carlsson, G. Topology and data. Bulletin of the American Mathematical Society, 46(2):255–308, 2009.
- Carrière, M., Cuturi, M., and Oudot, S. Sliced wasserstein kernel for persistence diagrams. In International Conference on Machine Learning, pp. 664–673. JMLR.org, 2017.
- Chazal, F. and Michel, B. An introduction to topological data analysis: fundamental and practical aspects for data scientists. Frontiers in Artificial Intelligence, 4:108, 2021.
- Chazal, F., de Silva, V., Glisse, M., and Oudot, S. Persistence stability for geometric complexes. Geometriae Dedicata, 173(1):193–214, 2014.
- Cohen-Steiner, D., Edelsbrunner, H., and Harer, J. Stability of persistence diagrams. In Proceedings of the Twenty-first Annual Symposium on Computational Geometry, pp. 263–271, 2005.
- Cohen-Steiner, D., Edelsbrunner, H., and Harer, J. Stability of persistence diagrams. Discrete & Computational Geometry, 37(1):103–120, 2007.
- Cohen-Steiner, D., Edelsbrunner, H., Harer, J., and Milroy, Y. Lipschitz functions have  $l_p$ -stable persistence. Foundations of Computational Mathematics, 10(2):127–139, 2010.
- Cui, P., Shen, Z., Li, S., Yao, L., Li, Y., Chu, Z., and Gao, J. Causal inference meets machine learning. In Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pp. 3527–3528, 2020.
- Dehejia, R. H. and Wahba, S. Causal effects in non-experimental studies: Reevaluating the evaluation of training programs. Journal of the American Statistical Association, 94(448):1053–1062, 1999.
- Edelsbrunner, Letscher, and Zomorodian. Topological persistence and simplification. Discrete & Computational Geometry, 28:511–533, 2002.
- Edelsbrunner, H. and Harer, J. Persistent homology—a survey. Contemporary Mathematics, 453:257–282, 2008.
- Feydy, J., Séjourné, T., Vialard, F.-X., Amari, S.-i., Trounev, A., and Peyré, G. Interpolating between optimal transport and mmd using sinkhorn divergences. In International Conference on Artificial Intelligence and Statistics, pp. 2681–2690. PMLR, 2019.
- Gabrielsson, R. B. and Carlsson, G. Topology of deep neural networks. Journal of Machine Learning Research, 20(196):1–40, 2019.
- Gabrielsson, R. B., Nelson, B. J., Dwarknath, A., and Skraba, P. A topology layer for machine learning. In International Conference on Artificial Intelligence and Statistics, pp. 1553–1563. PMLR, 2020.
- Ghosh, S., Feng, Z., Bian, J., Butler, K., and Prospero, M. Dr-vidal-doubly robust variational information-theoretic deep adversarial learning for counterfactual prediction and treatment effect estimation on real world data. In AMIA Annual Symposium Proceedings, volume 2022, pp. 485, 2023.
- Ghrist, R. Barcodes: The persistent topology of data. Bulletin of the American Mathematical Society, 45(1): 61–75, 2008.
- Gnedenko, B. and Kolmogorov, A. Limit Distributions for Sums of Independent Random Variables. Creative Media Partners, LLC, 2021. English Translation from the Russian edition (1954).
- Gorbunov, E., Danilova, M., and Gasnikov, A. Stochastic optimization with heavy-tailed noise via accelerated gradient clipping. In Advances in Neural Information Processing Systems, volume 33, pp. 15042–15053, 2020.
- Gouk, H., Frank, E., Pfahringer, B., and Cree, M. J. Regularisation of neural networks by enforcing lipschitz continuity. Machine Learning, 110:393–416, 2021.
- Hill, J. L. Bayesian nonparametric modeling for causal inference. Journal of Computational and Graphical Statistics, 20(1):217–240, 2011.
- Johansson, F., Shalit, U., and Sontag, D. Learning representations for counterfactual inference. In International

- Conference on Machine Learning, pp. 3020–3029. PMLR, 2016.
- Kallus, N., Mao, X., and Udell, M. Causal inference with noisy and missing covariates via matrix factorization. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- Kazemi, A. and Ester, M. Adversarially balanced representation for continuous treatment effect estimation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 13085–13093, 2024.
- Kim, Y. S., Rachev, S. T., Bianchi, M. L., and Fabozzi, F. J. Financial market models with lévy processes and time-varying volatility. *Journal of Banking & Finance*, 32(7):1363–1378, 2008.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- Künzel, S. R., Sekhon, J. S., Bickel, P. J., and Yu, B. Metalearners for estimating heterogeneous treatment effects using machine learning. *Proceedings of the National Academy of Sciences*, 116(10):4156–4165, 2019.
- Kuroki, M. and Pearl, J. Measurement bias and effect restoration in causal inference. *Biometrika*, 101(2):423–437, 2014.
- Lagemann, K., Lagemann, C., Taschler, B., and Mukherjee, S. Deep learning of causal structures in high dimensions under data limitations. *Nature Machine Intelligence*, 5(11):1306–1316, 2023.
- LaLonde, R. J. Evaluating the econometric evaluations of training programs with experimental data. *The American Economic Review*, pp. 604–620, 1986.
- Louizos, C., Shalit, U., Mooij, J. M., Sontag, D., Zemel, R., and Welling, M. Causal effect inference with deep latent-variable models. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Luo, Y., Peng, J., and Ma, J. When causal inference meets deep learning. *Nature Machine Intelligence*, 2(8):426–427, 2020.
- Mainardi, F. Lévy stable distributions in the theory of probability. *Lecture Notes on Mathematical Physics*, 2007.
- Mileyko, Y., Mukherjee, S., and Harer, J. Probability measures on the space of persistence diagrams. *Inverse Problems*, 27(12):124007, 2011.
- Morgan, S. L. and Winship, C. *Counterfactuals and Causal Inference*. Cambridge University Press, 2015.
- Papamarkou, T., Birdal, T., Bronstein, M. M., Carlsson, G. E., Curry, J., Gao, Y., Hajj, M., Kwitt, R., Lio, P., Di Lorenzo, P., et al. Position: Topological deep learning is the new frontier for relational learning. In *International Conference on Machine Learning*. PMLR, 2024.
- Pawlowski, N., Coelho de Castro, D., and Glocker, B. Deep structural causal models for tractable counterfactual inference. In *Advances in Neural Information Processing Systems*, volume 33, pp. 857–869, 2020.
- Pitera, M., Chechkin, A., and Wyłomańska, A. Goodness-of-fit test for  $\alpha$ -stable distribution based on the quantile conditional variance statistics. *Statistical Methods & Applications*, 31(2):387–424, 2022.
- Pöllänen, A. and Marttinen, P. Identifiable causal inference with noisy treatment and no side information. *arXiv preprint arXiv:2306.10614*, 2023.
- Schölkopf, B., Locatello, F., Bauer, S., Ke, N. R., Kalchbrenner, N., Goyal, A., and Bengio, Y. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.
- Séjourné, T., Feydy, J., Vialard, F.-X., Trouvé, A., and Peyré, G. Sinkhorn divergences for unbalanced optimal transport. *arXiv preprint arXiv:1910.12958*, 2019.
- Shalit, U., Johansson, F. D., and Sontag, D. Estimating individual treatment effect: generalization bounds and algorithms. In *International Conference on Machine Learning*, pp. 3076–3085. PMLR, 2017.
- Shi, C., Blei, D., and Veitch, V. Adapting neural networks for the estimation of treatment effects. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Shi, C., Veitch, V., and Blei, D. M. Invariant representation learning for treatment effect estimation. In *Uncertainty in Artificial Intelligence*, pp. 1546–1555. PMLR, 2021.
- Shu, D. and Yi, G. Y. Causal inference with noisy data: Bias analysis and estimation approaches to simultaneously addressing missingness and misclassification in binary outcomes. *Statistics in Medicine*, 39(4):456–468, 2020.
- Simsekli, U., Sagun, L., and Gurbuzbalaban, M. A tail-index analysis of stochastic gradient noise in deep neural networks. In *International Conference on Machine Learning*, pp. 5827–5837. PMLR, 2019.
- Southern, J., Wayland, J., Bronstein, M. M., and Rieck, B. Curvature filtrations for graph generative model evaluation. In *Advances on Neural Information Processing Systems*, 2023.

- Stoyanov, S. V., Rachev, S. T., Racheva-Yotova, B., and Fabozzi, F. J. Fat-tailed models for risk estimation. The Journal of Portfolio Management, 37(2):107–117, 2011.
- Tralie, C., Saul, N., and Bar-On, R. Ripser.py: A lean persistent homology library for python. The Journal of Open Source Software, 3(29):925, Sep 2018.
- Virmaux, A. and Scaman, K. Lipschitz regularity of deep neural networks: analysis and efficient estimation. In Advances in Neural Information Processing Systems, volume 31, 2018.
- Wang, H., Fan, J., Chen, Z., Li, H., Liu, W., Liu, T., Dai, Q., Wang, Y., Dong, Z., and Tang, R. Optimal transport for treatment effect estimation. In Advances in Neural Information Processing Systems, volume 36, 2024.
- Wickens, M. R. A note on the use of proxy variables. Econometrica: Journal of the Econometric Society, pp. 759–761, 1972.
- Yang, H., Qiu, P., and Liu, J. Taming fat-tailed (“heavier-tailed” with potentially infinite variance) noise in federated learning. In Advances in Neural Information Processing Systems, volume 35, pp. 17017–17029, 2022.
- Yeager, D. S., Romero, C., Paunesku, D., Hulleman, C. S., Schneider, B., Hinojosa, C., Lee, H. Y., O’Brien, J., Flint, K., Roberts, A., et al. Using design thinking to improve psychological interventions: The case of the growth mindset during the transition to high school. Journal of Educational Psychology, 108(3):374, 2016.
- Yoon, J., Jordon, J., and Van Der Schaar, M. Ganite: Estimation of individualized treatment effects using generative adversarial nets. In International Conference on Learning Representations, 2018.
- Zomorodian, A. and Carlsson, G. Computing persistent homology. Discrete & Computational Geometry, 33(2): 249–274, 2005.



## A. Topological Data Analysis, Persistent Homology, and Stability Theorems

In this appendix we provide further details on topological data analysis (TDA) and in particular the key concept of *persistent homology*, which crucially helps with understanding complex data structures (Edelsbrunner & Harer, 2008). Before doing so, we provide a brief review of the basics of TDA. While a more concise review of the preliminaries on TDA was provided in Section 2, here we repeat that review with additional elaborations in order to better familiarize interested readers.

### A.1. Review of the Preliminaries

Topological Data Analysis (TDA) applies the principles of algebraic topology to extract informative features from data. It is particularly adept at uncovering invariants such as the shape and connectivity of data across multiple scales (Carlsson, 2009; Papamarkou et al., 2024). Persistent homology, a central tool in TDA, provides a multiscale representation of topological features (Edelsbrunner & Harer, 2008). Applications of persistent homology span from feature extraction in computer vision to the analysis of complex datasets in machine learning, which benefit from its intrinsic metric and coordinate-free approach (Carrière et al., 2017).

Persistent homology captures the persistence of topological features such as connected components and holes as a scale parameter varies. These features are represented in a persistence diagram, a collection of points in the plane, each point corresponding to a feature’s birth and death in a filtration of the input (Ghrist, 2008). The construction of a filtration, a series of nested simplicial complexes determined by the image of a filtration function, is the first step in applying persistent homology to a dataset. The persistent homology of this filtration is then computed, yielding a persistence diagram, which, again, it represents the lifespan of topological features as points marking their birth and death (Zomorodian & Carlsson, 2005). Figure 1 visualizes this in a simple example. We delve into further details on persistent homology in the following subsection.

Resulting from *the stability theorem*, stability of the persistence diagram is a key property that ensures the robustness of these topological summaries to perturbations in the data (Cohen-Steiner et al., 2007). This theorem states a bound on the bottleneck distance between two persistence diagrams obtained via filtration functions  $f$  and  $g$ .<sup>2</sup> Perhaps regarded as the most central theorem underlying the applicability of persistent homology, the stability theorem provides a guarantee that the persistence diagram is stable under small perturbations, making it particularly conducive to analyzing noisy data (Chazal et al., 2014). A generalization of the stability theorem used in our analysis is stated in Section 2, and further details are discussed in the Appendix section A.3.

### A.2. Persistent Homology

Persistent homology, a fundamental tool in TDA, helps in quantifying the topological features of data. This section aims to define and elucidate key concepts related to filtration functions and persistence diagrams, providing a background for discussing the stability of these constructs and the assumptions behind our theoretical results. We begin with filtration, which yields a multi-scale representation of data, essential for understanding the evolution of topological features.

#### A.2.1. FILTRATION

A filtration is a nested collection of subspaces  $\{\mathcal{Z}_a\}_{a \in \mathbb{R}}$  of a topological space  $\mathcal{Z}$ , such that  $\mathcal{Z}_a \subseteq \mathcal{Z}_b$  whenever  $a \leq b$ . This process can be driven by a real-valued function  $f : \mathcal{Z} \rightarrow \mathbb{R}$ , referred to as the *filtration function*, which assigns a real number to each point in  $\mathcal{Z}$ . Considering a bounded continuous function  $f : \mathcal{Z} \rightarrow \mathbb{R}$ , the sublevel sets  $f^{-1}(-\infty, a_i]$  at various thresholds  $a_0 \leq a_1 \leq \dots \leq a_n$  give rise to a filtration. The filtration captures the evolution of the topological structure of  $\mathcal{Z}$  as the threshold varies, revealing critical values where topological features appear or disappear.

Given a filtration function, for every pair of threshold values  $a \leq b$ , the inclusion relationship between their corresponding subspaces,  $\mathcal{Z}_a \subset \mathcal{Z}_b$ , induces homomorphisms of the  $l$  dimensional homology groups  $H_l(\mathcal{Z}_a)$  and  $H_l(\mathcal{Z}_b)$ . If there exists a dimension  $l$ , a threshold value  $c \in \mathbb{R}$ , and a value  $\delta > 0$ , such that for all  $\epsilon \in (0, \delta)$  the homomorphism induced by  $H_l(\mathcal{Z}_{c-\epsilon}) \subset H_l(\mathcal{Z}_{c+\epsilon})$  is not an isomorphism,  $c$  is called a *homological critical value*. These critical values mark the levels where the homology of the sublevel sets changes (Cohen-Steiner et al., 2005).

**Definition A.1** (Cohen-Steiner et al. (2005)). A filtration function,  $f : \mathcal{Z} \rightarrow \mathbb{R}$ , is *tame*, if it only has a finite number of homological critical values, and if for all threshold values  $a \in \mathbb{R}$  and dimensions  $l$ , the homology groups  $H_l(f^{-1}(-\infty, a])$

<sup>2</sup>The bottleneck distance between two diagrams is the cost of the optimal matching between their points.

are finite dimensional.

The filtration function could be any mapping to a meaningful real-valued representation of the data. Studying how topology changes through the filtration gives insight into the structure of data at different scales. The stability theorem assumes the filtration function is tame.

### A.2.2. PERSISTENCE DIAGRAM

Given a filtration, the persistence diagram compactly represents the lifespan of homological features through their birth and death thresholds. Homology classes are born at critical threshold values where new features appear in the filtration. Subsequently, some classes die at larger thresholds. The lifespan of a class that is born at threshold  $a$  and dies entering threshold  $b$  is characterized by the *persistence* value of the corresponding point,  $x$ , in the persistence diagram, defined as  $\text{pers}(x) := b - a$  (Edelsbrunner et al., 2002). Classes with larger persistence values are considered to be more prominent features. The birth and death of homology groups can be represented in a 2-dimensional *persistence diagram* as follows (Cohen-Steiner et al., 2007): Points  $(a, b)$  denoting classes born at threshold value  $a$  and dying at threshold value  $b$ , and points of the form  $(a, \infty)$  representing essential homology classes that never die. This low-dimensional representation allows us to easily interpret and analyze the topological features of the data over different scales. Figure 7 visualizes an example of a Vietoris-Rips filtration (Carlsson, 2009) and the corresponding persistence diagram of 0- and 1-dimensional homology classes. As the figure demonstrates, the persistence diagram shows the birth and death of topological features of the data with the points farther from the diagonal marking more persistent features.

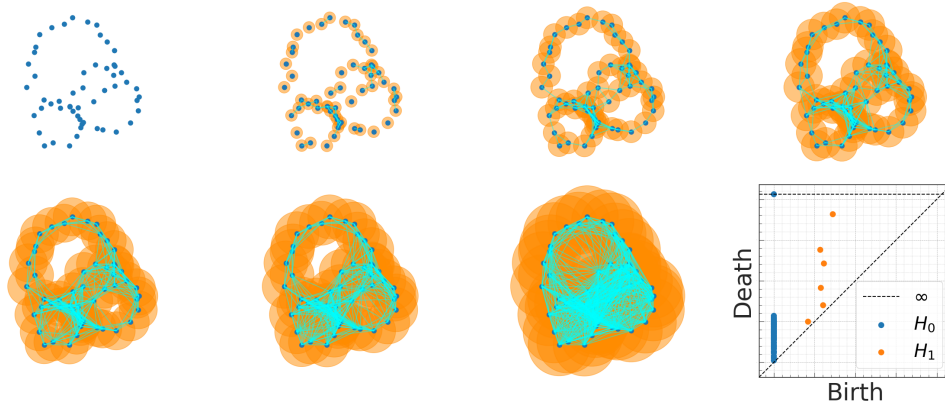


Figure 7: A visualization of the Vietoris-Rips filtration of an example pointcloud and the corresponding persistence diagram.

Before discussing key stability properties of persistence diagrams, we state two definitions related to properties of the filtration function. These definitions relate to *Degree- $k$  total persistence* of a persistence diagram corresponding to a filtration function  $f$ . Degree- $k$  total persistence, which sets one of the assumptions for the stability theorem, is defined as the sum of the  $k^{\text{th}}$  powers of the persistence values of all points in the persistence diagram of  $f$ .

**Definition A.2** (Cohen-Steiner et al. (2010)). Given a filtration function  $f : \mathcal{Z} \rightarrow \mathbb{R}$ , let  $\text{Dgm}_f$  denote the corresponding persistence diagram. The degree- $k$  total persistence is given by

$$\text{Pers}_k(f) = \sum_{x \in \text{Dgm}_f} \text{pers}(x)^k.$$

**Definition A.3** (Cohen-Steiner et al. (2010)). We say that a space  $\mathcal{Z}$  implies *bounded degree- $k$  total persistence*, if there exists a constant  $C_{\mathcal{Z}}$  that depends only on  $\mathcal{Z}$ , such that for every tame Lipschitz function  $f : \mathcal{Z} \rightarrow \mathbb{R}$  with Lipschitz constant  $K_f$ , we have

$$\text{Pers}_k(f) \leq C_{\mathcal{Z}} K_f^k. \quad (3)$$

### A.3. Stability Theorems

The stability properties of persistence diagrams are crucial, as they imply that the topological signature captured by these diagrams is robust to small perturbations and noise in the data. As such, they serve as the main justification for using

persistence diagrams in order to improve the robustness of deep learning frameworks to input noise.

The *stability theorem*—the main theoretical result regarding the stability of persistence diagrams— bounds the bottleneck distance between two diagrams as follows.

**Theorem A.4** (Cohen-Steiner et al. (2005)). *Let  $\mathcal{Z}$  be a triangulable space with continuous tame filtration functions  $f, g : \mathcal{Z} \rightarrow \mathbb{R}$ . Then the corresponding persistence diagrams satisfy*

$$d_B(\mu_f, \mu_g) \leq \|f - g\|_\infty,$$

where  $d_B$  is the bottleneck distance, and  $\mu_f$  and  $\mu_g$  are the probability measures induced by  $\text{Dgm}_f$  and  $\text{Dgm}_g$ .

Stating that the distance between the persistence diagrams is controlled by the  $L_\infty$  distance between the corresponding filtration functions, this theorem ultimately suggests that the persistence diagram is more stable than the geometry of the data it represents.

In a generalization of this statement, the *Wasserstein Stability Theorem*, extends this result to the Wasserstein- $p$  distance between the diagrams for  $p \geq k \geq 1$ , when  $\mathcal{Z}$  implies bounded degree- $k$  total persistence.

**Theorem A.5** (Cohen-Steiner et al. (2010)). *Let  $\mathcal{Z}$  be a triangulable, compact metric space that implies bounded degree- $k$  total persistence, for some  $k \geq 1$ . Let  $f, g : \mathcal{Z} \rightarrow \mathbb{R}$  be two tame Lipschitz filtration functions with Lipschitz constants  $K_f$  and  $K_g$ . Then for all  $l$ -dimensional homology classes and all  $p \geq k$  we have*

$$W_p(\mu_f^l, \mu_g^l) \leq C^{\frac{1}{p}} \|f - g\|_\infty^{1 - \frac{k}{p}}.$$

The constant  $C$  is given by  $C = C_{\mathcal{Z}} \max\{K_f^k, K_g^k\}$ , where  $C_{\mathcal{Z}}$  is the constant in Equation 3.

Note that as  $p \rightarrow \infty$ , this generalized formulation gives the statement in Theorem A.4. To conclude, these stability results show persistence diagrams are robust topological summaries for analyzing complex data.

## B. Learning Robust Representations via Persistence Diagrams

Here we provide additional details regarding the theoretical results discussed in Section 3. In particular, while providing the proofs and the background leading to the theorems, we detail the conditions and elaborate on the formulations of the bounds stated in the theorems. Theorem 3.1 offer insight into the conditions under which using persistence diagrams corresponds to enhanced robustness of the representations, depending on the properties of the neural network. In particular, this theorem constrains the Lipschitz constant of the neural network by a term that contains information on the empirical distribution of the noise, including the sample average and the infinity norm of the sample noise. Importantly, the combination of sample average and infinity norm establishes a connection with the tail of the noise distribution, enabling our analysis of TATEE’s robustness through Proposition C.1 in Section C.5, which in turn provides theoretical grounds for TATEE’s advantage, and validated in empirical settings by the experimental results observed in Section 5. Below, we provide the proof and additional details about the bounds in Theorem 3.1.

**Notation.** Following the notation in Section 3,  $X$ ,  $E$ , and  $\tilde{X} := X + E$  denote the random vectors of features, noise, and noise-corrupted features on a sample space  $\Omega$ ; and we use  $\mathbf{X}$ ,  $\tilde{\mathbf{X}}$ , and  $\mathbf{E}$  to refer to their corresponding sample matrices.  $\varphi : \Omega \rightarrow \mathcal{Z}$  is the neural network mapping  $X$  and  $\tilde{X}$  to  $\varphi(X)$  and  $\varphi(\tilde{X})$  in the representation space  $\mathcal{Z}$ ,  $\mu$  and  $\tilde{\mu}$  are the measures induced by these representations, and  $f, \tilde{f} : \mathcal{Z} \rightarrow \mathbb{R}$  are tame and Lipschitz filtration functions used to compute the persistence diagrams inducing the measures  $\mu_f$  and  $\mu_{\tilde{f}}$  on the space of persistence diagrams. For ease of notation, we denote the distances  $W_p(\mu, \tilde{\mu})$  and  $W_p(\mu_f, \mu_{\tilde{f}})$  by  $\Delta$  and  $\Delta_{\text{topo}}$ , and their corresponding sample equivalents by  $\hat{\Delta}$  and  $\hat{\Delta}_{\text{topo}}$ .

**Main Assumptions.** The main assumptions for Theorem 3.1 are related to Lipschitz continuity of the functions involved and the bounds of their co-domains. In particular, following the standard assumptions for stability theorems, we assume the filtration functions are Lipschitz. We also assume that  $\mathcal{Z}$  implies bounded degree- $k$  total persistence, as defined in Appendix A; another assumption that is made for the stability theorems. Additionally, we assume the neural network  $\varphi$  is Lipschitz, which, as we discussed in Section 3, is not a restrictive assumption and is satisfied in many standard scenarios. The restrictive assumption for the proofs of our theorem, which is only technical and for simplicity, is due to the formulation

of  $\tilde{f}$  and its Lipschitz continuity. Recall that we consider  $f$  and  $\tilde{f}$  which satisfy  $\tilde{f}(\varphi(X)) = f(\varphi(\tilde{X}))$ . For simplicity and to avoid auxiliary constructions, we require  $\tilde{f} = f \circ \varphi \circ S \circ \varphi^{-1}$ , where  $S(X) := X + E$ , so that  $\tilde{f}(\varphi(X)) = f(\varphi(\tilde{X}))$  holds exactly.

We remind the reader that these assumptions are technical and serve to enable the derivation of the proofs through simple expressions of the quantities of interest. With that in mind, for clarity, we discuss the assumption that allows  $\tilde{f} = f \circ \varphi \circ S \circ \varphi^{-1}$  to hold exactly. Strictly speaking, this is restrictive since it requires  $\tilde{f}$  to be Lipschitz continuous and  $\varphi$  to have a Lipschitz inverse, which is typically violated by commonly-used neural networks. However, the specific formulation of  $\tilde{f}$  that requires this assumption is only to simplify the derivation and expression of the bounds on  $\|f - \tilde{f}\|_\infty$ . That is, the assumption allows us to avoid tedious constructions and keep the formulation of the variables used and the steps of the proof simple, concise, clear, and focused, and otherwise could be replaced with a less restrictive assumptions, to derive effectively similar results. While this relaxation is not trivial to the authors' knowledge and hence beyond the scope of this discussion, one such relaxation could be done by assuming that  $\phi(\cdot)$  is only locally invertible on the input sample and adding an approximation error term with the  $\tilde{f}$  that is defined for the local Lipschitz pseudo-inverse. Moreover, our experimental results indicate that while the assumptions facilitate the steps for our theoretical analysis in a simplified setting, the overall intuition and implications hold more broadly in practical scenarios. With this in mind, we state the theorem, which provide significant insights, validated and confirmed by our experiments.

### B.1. Finite Sample Stability and Error Distribution

Let us denote the empirical estimators of the Wasserstein distance of  $W_p(\mu, \tilde{\mu})$  by  $\hat{\Delta}$ , and its upper bound by  $\hat{M}$ , as defined in Section 3. We denote the upper bound on  $C^{\frac{1}{p}} \|f - \tilde{f}\|_\infty^{1-\frac{k}{p}}$  by  $K_{topo}$ , which bounds  $W_p(\mu_f, \mu_{\tilde{f}})$ , as stated in Section 3. Furthermore, let us denote its finite sample equivalent from Theorem 2.1 by  $\hat{K}_{topo}$ . Under the assumptions discussed above, Theorem 3.1 states the following: Given a bounded degree- $k$  total persistence and a  $K_f$ -Lipschitz filtration function, if the Lipschitz constant of the neural network is smaller than  $\Lambda$  for a value  $\Lambda$ , then  $\hat{K}_{topo} < \hat{M}$ . Importantly, the value  $\Lambda$ , described below, is linear in  $\|\mathbf{E}\|_\infty^{p/k-1} \|\bar{\mathbf{E}}\|^{-p/k}$ , which speaks to the impact of the empirical distribution of noise on the gain in robustness through use of persistence diagrams. The proof of this theorem is presented next.

**Proof.** The upper bound  $\hat{M}$  can be derived using the Lipschitz continuity of  $\varphi$ . In the finite sample regime, this bound becomes

$$\hat{\Delta} \leq \hat{M} = K_\varphi \hat{\epsilon}, \quad (4)$$

where  $\hat{\epsilon}$  is the sample average. For  $B_{topo} = C^{\frac{1}{p}} \|f \circ \varphi - \tilde{f} \circ \varphi\|_\infty^{1-\frac{k}{p}}$ , using Lipschitz continuity of  $f$  and  $\varphi$ , we can derive the following finite sample bound,

$$\hat{\Delta}_{topo} \leq \hat{B}_{topo} \leq \hat{K}_{topo} = C^{\frac{1}{p}} (K_f K_\varphi \bar{\epsilon})^{1-\frac{k}{p}}, \quad (5)$$

where, for ease of notation, we use  $\bar{\epsilon} = \|\mathbf{E}\|_\infty$  to denote the finite sample infinity norm of the noise matrix. Additionally, note that  $C = C_{\mathcal{Z}} L$ , where  $L := \max \{K_f^k, K_{\tilde{f}}^k\}$ , as explained in Appendix Section A. Comparing the right hand sides of the inequalities 4 and 5, it follows that if

$$K_\varphi < C^{1/k} K_f^{\frac{p}{k}-1} \frac{\bar{\epsilon}^{\frac{p}{k}-1}}{\hat{\epsilon}^{\frac{p}{k}}},$$

then,  $\hat{K}_{topo} < \hat{M}$ , which completes the proof of Theorem 3.1.  $\square$

## C. Topology-Aware Treatment Effect Estimation

Building on the stability results established in Section 3, in Section 4 we introduced *Topology-Aware Treatment Effect Estimation* (TATEE)—a framework designed to improve the robustness of deep learning approaches to counterfactual regression. TATEE incorporates topological regularization by embedding the persistence diagram of the learned representations into the training objective. We briefly described the main aspects of TATEE's design, analysis, and implications in

Section 4. In this aspect, we elaborate and expand on the discussion in Section 4, delving into details of the structure of the CATE estimation network used in TATEE, along with its training procedure, scalability, and conceptual motivation and implications. We further provide a detailed theoretical analysis of TATEE’s robustness properties formally stating results mentioned in the informal Proposition 4.1, which connects the stability benefits of topological summaries to improved treatment effect estimation. Note that, as in prior works (Shalit et al., 2017; Louizos et al., 2017; Shi et al., 2019), we focus on the binary treatment setting for clarity, though TATEE is in principle extensible to multi-valued treatments.

### C.1. Representation-Balancing Neural Network for CFR

TATEE builds on the two-headed architecture introduced in the CFR framework Shalit et al. (2017), consisting of a shared encoder followed by separate branches for the treatment and control groups. The shared encoder  $\varphi$  maps inputs to a representation space where the distributions of treated and control units are approximately aligned—enabling the learning of *treatment-agnostic* representations. These shared representations are then passed to two distinct heads,  $h_0$  and  $h_1$ , which learn the potential outcomes under control and treatment, respectively. This structure corresponds to a T-Learner (Künzel et al., 2019)—the potential outcome heads for each group—to learn the factual and counterfactual outcomes separately for estimating CATE. Figure 2 visualizes this architecture, showing the shared neural network,  $\varphi$ , and two separate heads  $h_0$  and  $h_1$  which learn the potential outcomes from  $\varphi$ ’s output. As we detail in Section C.3, the weights are trained with respect to the outcome, the distribution of the representations, as well as the topological summaries of the representations as captured by the persistence diagram of the 1-dimensional homology class.

### C.2. Implementation

The architecture of the neural network in TATEE is described in Section 4, where we explain the role of each component of the model in the CFR-type architecture shown in Figure 2. We also specified the training objective in Section C.3, which, critically, incorporates topological awareness into the CFR framework and is the core distinguishing component of TATEE. In this appendix, we include the details of the implementation of the neural network and state the full algorithm.

The skeleton of the neural network in TATEE follows that of CFR, described in Shalit et al. (2017). As we mention in Appendix G, adopting the hyper-parameters used by Shalit et al. (2017), we use three fully connected layers with ELU (for *exponential linear unit*) activation functions for all three components— $\varphi$ ,  $h_0$ , and  $h_1$ . Each shared representation layer has 200 neurons, while the layers in  $h_0$  and  $h_1$  have 100 neurons each. Since the outcome in the Twins dataset takes binary values, we use a sigmoid activation function on the final layer, this affecting the loss function used as  $\mathcal{L}_{Outcome}$  in Equation 6, which is binary cross entropy for Twins dataset and mean squared error for the others. The central term in the loss function,  $\mathcal{L}_{topo}$ , uses the Wasserstein-2 distance between the persistence landscapes of the representations of a mini-batch of size 256 for the IHDP dataset, and size 128 for the others. We use the Vietoris-Rips complex (Carlsson, 2009) for computing the persistence diagrams, which are obtained using the Ripser package (Tralie et al., 2018; Bauer, 2021). The distance between the persistence landscapes for the 1-dimensional homology class are then approximated using the Sinkhorn divergence (Séjourné et al., 2019)—an efficient and differentiable approximation of the Wasserstein distance which is amenable to gradient descent for training the neural network. The GeomLoss package (Feydy et al., 2019) is used for computing the Sinkhorn divergence. The full algorithm is stated in Algorithm 1, which clarifies how the topological signature, as described in Section 4, is incorporated in the CFR framework to obtain TATEE. Note that we used a fixed number of epochs of training, hence, in our implementation, the convergence criterion in Algorithm 1 is simply completing the specified number of epochs.

### C.3. Training TATEE

The effectiveness of TATEE in enhancing the robustness of learning treatment-agnostic representations is primarily owed to incorporating the topology of the shared representations in the training process. This is achieved by adding a topological regularization term based on the Wasserstein distance between the persistence diagrams of the treatment and control groups. The CFR network architecture in Shalit et al. (2017) is predicated on the minimization of a loss function that encapsulates both the prediction accuracy and the distributional balance between treated and control groups. Incorporating the topological regularization, for a sample of size  $N$ , TATEE’s training objective is as follows:

$$\mathcal{L}_{TATEE} = \mathcal{L}_{Outcome} + \lambda \mathcal{L}_{Balance} + \lambda_{topo} \mathcal{L}_{topo}, \quad (6)$$

The three components of  $\mathcal{L}_{TATEE}$  are given by:



**Algorithm 1** TATEE Training

- 1: **Input:** Neural network composed of the components  $\varphi$  and  $h(\cdot, t_i) := (1 - t_i)h_0(\cdot) + t_i h_1(\cdot)$  with initial weights  $\theta_\varphi$  and  $\theta_h$ , sample data  $(x_1, t_1, y_1), \dots, (x_N, t_N, y_N)$ , regularization parameters  $\lambda, \lambda_{topo} > 0$ , and loss function  $\mathcal{L}$ .
- 2: Compute  $N_1 = \sum_{i=1}^N t_i$  and  $N_0 = N - N_1$ .
- 3: Compute sample weights  $w_i = \frac{N t_i}{2N_1} + \frac{N(1-t_i)}{2N_0}$  for  $i = 1 \dots n$ .
- 4: **while** not converged **do**
- 5:   Take mini-batch  $I_M := \{i_1, i_2, \dots, i_M\} \subseteq \{1, 2, \dots, N\}$
- 6:   Compute representations  $\Phi_0 := \{\varphi(x_j) : t_j = 0, j \in I_M\}$  and  $\Phi_1 := \{\varphi(x_j) : t_j = 1, j \in I_M\}$
- 7:   Compute the predicted outcomes  $\{\hat{y}_j = h(\varphi(x_j), t_j) : j \in I_M\}$
- 8:   Compute the persistence diagrams  $\text{Dgm}_1(\Phi_0)$  and  $\text{Dgm}_1(\Phi_1)$
- 9:   Compute the gradient of the empirical  $\mathcal{L}_{\text{Balancing}}$  as  $\delta_{\text{Balancing}} = \nabla_{\theta_\varphi} W_p(\Phi_0, \Phi_1)$
- 10:   Compute the gradient of the empirical  $\mathcal{L}_{\text{topo}}$  as  $\delta_{\text{topo}} = \nabla_{\theta_\varphi} W_p(\text{Dgm}_1(\Phi_0), \text{Dgm}_1(\Phi_1))$
- 11:   Compute the gradients of the empirical  $\mathcal{L}_{\text{Outcome}}$  as  
 $\delta_{\varphi, \text{Outcome}} = \nabla_{\theta_\varphi} \frac{1}{M} \sum_{j \in I_M} w_j \mathcal{L}(y_j, \hat{y}_j)$   
 $\delta_{h, \text{Outcome}} = \nabla_{\theta_h} \frac{1}{M} \sum_{j \in I_M} w_j \mathcal{L}(y_j, \hat{y}_j)$
- 12:   Determine the step size  $\eta$  using Adam
- 13:   Update weights  
 $\theta_\varphi \leftarrow \theta_\varphi - \eta [\delta_{\varphi, \text{Outcome}} + \lambda \delta_{\text{Balancing}} + \lambda_{\text{topo}} \delta_{\text{topo}}]$   
 $\theta_h \leftarrow \theta_h - \eta (\delta_{h, \text{Outcome}})$
- 14:   Check for convergence
- 15: **end while**

$$\mathcal{L}_{\text{Outcome}} = \frac{1}{N} \sum_{i=1}^N w_i \ell(\hat{y}_i, y_i), \quad \mathcal{L}_{\text{Balance}} = W_p(\Phi_0, \Phi_1), \quad \mathcal{L}_{\text{topo}} = W_p(\text{Dgm}_1(\Phi_0), \text{Dgm}_1(\Phi_1)),$$

where  $\Phi_0$  and  $\Phi_1$  are the representations of the control and treated groups, respectively;  $y_i$  the true outcome of unit  $i$ ,  $\hat{y}_i := h(\varphi(\mathbf{x}_i), t_i)$  the predicted outcome with the feature vector  $\mathbf{x}_i$  and treatment  $t_i$ ,  $\ell(\cdot)$  the outcome prediction loss function,  $\text{Dgm}_1(\cdot)$  the 1-homology class persistence diagram, and  $W_p(\cdot, \cdot)$  the Wasserstein-p distance. The weight  $w_i$  aims to deal with the imbalance in the size of the treatment and control groups, and is given by  $\frac{N t_i}{2N_1} + \frac{N(1-t_i)}{2N_0}$  where  $N_1$  and  $N_0$  are the sample sizes of the two groups. The function  $h(\cdot, t_i) := (1 - t_i)h_0(\cdot) + t_i h_1(\cdot)$  combines the two potential outcome functions. While our theoretical analysis in Equation (1) holds for any homology class, we use the first homology group in practice. This choice is motivated by simplicity and the goal of capturing holes, thereby characterizing topological features beyond connected components. We include the full algorithm in Appendix C.2, describing the implementation of TATEE in details.

**C.4. Scalability of TATEE**

TATEE incorporates topological summaries in a computationally scalable fashion, due to two key factors in our implementation: First, topological regularization operates on representations rather than raw inputs, making computational cost constant with respect to input dimensionality, once the dimensionality of the representations are fixed. Second, we compute persistence diagrams on mini-batches, hence, with a batch size of  $b$  and for  $N$  total data points, the cost of computing  $\mathcal{L}_{\text{topo}}$  for each of the  $N/b$  batches remains fixed, allowing the overhead to scale linearly with data volume. As a result, despite topological methods typically being expensive for high-dimensional or large inputs, TATEE's implementation ensures the topological component does not become a computational bottleneck. Figure 3 illustrates this scalability via a simulation, confirming that overhead elapsed time remains nearly constant as input dimensionality increases from 16 to 256 (with a fixed representation dimensionality), and scales linearly as data volume varies from 800 to 12800 samples. These results further support the fact that the topological component introduces a scalable computational overhead while providing significant robustness benefits, making TATEE practical for real-world applications.

### C.5. Robustness of TATEE’s Training to Noise

Using the stability results presented before, we now show that under the conditions of Theorem 3.1, TATEE can improve the robustness of counterfactual regression to input noise. In particular, we show that Theorem 3.1 implies that TATEE’s training objective in Equation 6 is more stable than the original CFR’s objective. Consider the problem setup and notation in Section 3, and let us label the variables corresponding to the treatment and control groups by superscripts  $\cdot^0$  and  $\cdot^1$ , respectively. By the triangle inequality, the terms  $\lambda(\hat{M}^1 + \hat{M}^0) + \lambda_{\text{topo}}(\hat{K}_{\text{topo}}^1 + \hat{K}_{\text{topo}}^0)$  and  $\lambda_{\text{CFR}}(\hat{M}^1 + \hat{M}^0)$  upper-bound the noise-induced change in  $\lambda\mathcal{L}_{\text{Balance}} + \lambda_{\text{topo}}\mathcal{L}_{\text{topo}}$  and  $\lambda_{\text{CFR}}\mathcal{L}_{\text{Balance}}$ , for loss coefficients  $\lambda$ ,  $\lambda_{\text{topo}}$ , and  $\lambda_{\text{CFR}}$ . Using Theorem 3.1, we derive the following result (proof in Appendix E).

**Proposition C.1.** *If the Lipschitz constant of  $\varphi$  satisfies the constraint in Theorem 3.1, for any given  $\lambda_{\text{CFR}} > 0$ , with sufficiently small choices of  $\lambda$  and  $\lambda_{\text{topo}}$ , we have*

$$\lambda(\hat{M}^1 + \hat{M}^0) + \lambda_{\text{topo}}(\hat{K}_{\text{topo}}^1 + \hat{K}_{\text{topo}}^0) \leq \lambda_{\text{CFR}}(\hat{M}^1 + \hat{M}^0).$$

The inequality in Proposition C.1 indicates that the upper bound on the sum of the balancing and topological terms in TATEE’s loss undergoes a smaller change due to input noise, compared to that on the balancing term in CFR’s loss. In other words, this proposition implies that, when  $\varphi$  in TATEE’s architecture (Figure 2) satisfies the conditions of Theorem 3.1, TATEE is trained using a loss that can be more stable under additive noise than its counterpart CFR. As in Theorem 3.1, this robustness is easier to attain when the noise is heavy-tailed—i.e., the constraint on  $\varphi$ ’s Lipschitz constant becomes more permissive. This highlights TATEE’s particular suitability for robustness under heavy-tailed perturbations. Our experiments in Section 5 provide empirical evidence that TATEE can meet the theory-indicated conditions and gain the anticipated robustness in practical settings, outperforming the original CFR from Shalit et al. (2017) and other standard causal inference baselines.

### C.6. TATEE in Action

According to our theoretical results, the topological regularization term in Equation (6) could help train  $\varphi$  such that the total loss becomes more robust to noise. To complement the theoretical analysis, we examine the intuition behind this improvement through a simple example of how the network learns the representations. As we explained earlier, the CFR framework is based on the principle of balancing the representations of treated and control groups through enforcing a distributional similarity between the two. TDA on the other hand, characterizes data using topological features which are invariant to smooth deformations. This allows TDA-based methods to capture qualitative properties of the shape of the data at a global level, while limiting their sensitivity to local geometric perturbations, hence leading to the resulting robustness. This understanding provides an explanation as to why informing a representation balancing framework with the topology of the representations could enhance robustness.

To see how this works in a simple example, we simulate features such that the 2-dimensional representations corresponding to the treated and control groups have distinct topologies: one forming a line and the other a circle—a simple difference in their 1-dimensional homology class. Figure 4 compares how  $\varphi$  learns to balance these representations in TATEE versus standard CFR, with and without noise. Visually confirming our understanding, CFR forces distributional similarity between groups while topological differences persist, even widening after 10 epochs of training without topological awareness. TATEE, however, enforces both distributional and topological similarities through the distance between the persistence diagrams, resulting in representations that converge to similar shapes with matching topologies. In other words, while the Wasserstein loss of  $\mathcal{L}_{\text{Balance}}$  may allow the representations to qualitatively diverge, Figure 4 illustrates that the topological signature captured by  $\mathcal{L}_{\text{topo}}$  effectively prevents that. Notably, noise considerably impacts CFR’s ability to enforce distributional similarity, while TATEE achieves its objective equally well in both noisy and clean settings.

## D. Experiments

Section 5 discusses our main experimental results, showing TATEE’s superior robustness. In this appendix, we provide comprehensive details on our experimental setup, elaborate on our main results, and present additional results that complement the findings described in the main paper. We begin by describing the experimental setup, including datasets, evaluation metrics, noise distributions, and implementation. We then review our main experimental results on robustness of TATEE in more detail, followed by a report on its performance. We also elaborate on the distinction between evaluating robustness and performance on noisy inputs and comment on an expected tradeoff therein. We then discuss our experiments benchmarking TATEE against other treatment effect estimation methods, showing its superior robustness. Additional details

on the implementation of these experiments including hyperparameter values, and a more detailed description of each benchmark dataset are provided in appendices G and H.

### D.1. Experimental Setup

**Models and Evaluation Metrics.** We implement TATEE with the network architecture and training algorithm described in Section 4. In this implementation, we use the Wasserstein-2 distance in the balancing and topological terms in the training objective. We also use persistence landscapes (Bubenik et al., 2015), which maintain a one-to-one correspondence with persistence diagrams while offering differentiability and better statistical tractability. Other implementation details and hyperparameters are provided in Appendices C.2 and G. Since TATEE incorporates a topological signature in CFR (Shalit et al., 2017), to evaluate the resulting improvement in robustness, we compare their performances with and without input noise. We complement our experiments by several other causal inference models included in Figure 6 and listed in Appendix G. Following the conventions in causal inference (e.g., Hill (2011); Louizos et al. (2017); Shalit et al. (2017); Shi et al. (2019)), we use the Precision in Estimation of Heterogeneous Effect (PEHE) to quantify the CATE estimation error. The empirical PEHE is given by  $\hat{\epsilon}_{\text{PEHE}} = \frac{1}{N} \sum_{i=1}^N (\tau(x_i) - \hat{\tau}(x_i))^2$ , where  $\tau(x_i)$  is the true CATE from Equation (2) and  $\hat{\tau}$  is the estimated CATE. While PEHE can be used on noisy input to assess the performance in noisy regimes, we are primarily interested in judging TATEE’s robustness. To evaluate the robustness of the models to noise, we use the following metric, denoted by  $\rho$ ,

$$\rho = 1 - \frac{\tilde{\epsilon}_{\text{PEHE}}}{\hat{\epsilon}_{\text{PEHE}}}. \quad (7)$$

$\rho$  compares PEHE on noisy and noise-free data, which we denote here by  $\tilde{\epsilon}_{\text{PEHE}}$  and  $\hat{\epsilon}_{\text{PEHE}}$ . Unless the noise itself helps the training,  $\rho$  should take negative values; more negative when less robust. In order to quantify the gain in robustness from TATEE, we compare the increase in  $\rho$  by computing  $\rho_{\text{TATEE}} - \rho_{\text{CFR}}$ , which takes positive values if TATEE is more robust than CFR. Given the main objective of our work, our evaluation of the experiments is primarily focused on robustness, while we also inspect the performance of TATEE on both clean and noisy datasets to make sure the enhanced robustness does not cost a considerable decline in performance. Note the distinction between the robustness and performance metrics: a larger (worse)  $\tilde{\epsilon}_{\text{PEHE}}$  can still yield a higher (better)  $\rho$ , when the noise-induced deterioration of performance, captured by  $\frac{\tilde{\epsilon}_{\text{PEHE}}}{\hat{\epsilon}_{\text{PEHE}}}$ , is lower. In other words,  $\rho_{\text{TATEE}} - \rho_{\text{CFR}}$  can be positive, reflecting better robustness, even if both  $\tilde{\epsilon}_{\text{PEHE}}$  and  $\hat{\epsilon}_{\text{PEHE}}$  are larger for TATEE, as long as the increase from  $\hat{\epsilon}_{\text{PEHE}}$  to  $\tilde{\epsilon}_{\text{PEHE}}$  remains smaller.

**Data.** We use four standard benchmark datasets in causal inference: IHDP, Twins, Jobs, and ACIC (Hill, 2011; Louizos et al., 2017; Shalit et al., 2017; Athey & Wager, 2019). The Twins dataset has the rare quality of having real-world values for both potential outcomes. The other datasets are semi-synthetic, with empirical values for features and simulated potential outcomes. All datasets are described in details in Appendix H. Note that using semi-synthetic data in experiments on estimating CATE is standard, and inevitable due to the *fundamental problem of causal inference*. This also dictates the use of synthetic noise as the only means to introduce noise to the treatment effect estimates. The Twins dataset stands out in this regard, allowing us to evaluate TATEE on intrinsic noise in empirical measurements.

**Noise Distribution.** We use the family of stable distributions (Mainardi, 2007) to simulate input noise for the feature matrix. This allows us to evaluate the robustness of TATEE in learning from features corrupted with various noise distributions of interest in empirical contexts (Kim et al., 2008; Stoyanov et al., 2011; Gorbunov et al., 2020; Yang et al., 2022). In particular,  $\alpha$ -stable distributions are characterized by a tail parameter  $\alpha$  and an asymmetry/skewness parameter  $\beta$ . Larger values of  $\alpha$  correspond to slower decay of the tail of the distribution, and positive/negative values of  $\beta$  correspond to positive/negative asymmetry. The valid parameter values fall within the *Feller-Takayasu diamond*, where  $\alpha \in (0, 2]$  and  $|\beta| \leq \min\{2 - \alpha, \alpha\}$  (Mainardi, 2007), giving the Gaussian distributions at  $\alpha = 2$  and  $\beta = 0$ . Appendix F contains a more detailed review of stable distributions.

**Noise Configuration.** For  $\alpha < 2$ , the noise follows a heavy-tailed distribution with infinite moments of integer order higher than 1, and for  $\alpha \leq 1$ , the mean is infinite. We therefore focus on noise distributions with finite mean, as well as  $\alpha = 1$ , which corresponds to the Cauchy distribution for  $\beta = 0$ . Varying the two parameters of this family of distributions allows us to empirically test the theoretical implications of Theorem 3.1. Here we report the average results over 30 trials. Since the Central Limit Theorem in its classical formulation does not hold for  $\alpha < 2$ , the standard sample size criteria for statistical confidence does not apply, as discussed in Appendix F.

## D.2. Main Results on Enhanced Robustness

Our experiments showcase the utility of TATEE and speak to its superior robustness, validating the applicability of our theoretical analysis in empirical settings. Reviewing our main empirical results from Section 5 with more details, we begin by confirming that TATEE has a performance comparable with CFR in the absence of synthetic noise, and then turn our attention to the main focus of this paper, showing the superior robustness of TATEE.

**Performance without Noise.** TATEE is consistently on-par with the original CFR in the noise-free regime, with  $1.977 \times 10^{-2}$ ,  $4.976 \times 10^{-4}$ , and  $2.631 \times 10^{-3}$  increase in  $\epsilon_{\text{PEHE}}$  for IHDP, Jobs, and ACIC datasets (respectively), and  $1.007 \times 10^{-2}$  decrease for Twins. Note that the theory-informed advantage of TATEE is its robustness. In fact, adding  $\mathcal{L}_{\text{topo}}$  in the loss is expected to lead to a more challenging path for optimizing the remaining terms, and hence, an increase in  $\epsilon_{\text{PEHE}}$  is not unexpected. Nevertheless, the evaluation mentioned above indicates that TATEE’s superior robustness does not come at the cost of a noticeable decline in performance. Additional experiments in Appendix D.3 indicate that TATEE also maintains a comparable or better performance on noisy data in most cases, with up to 6.9% improvement on IHDP, and 1.37%, 45.1%, and 42.3% on the Twins, Jobs, and ACIC datasets, as measured by reduction in PEHE. As clarified in our discussion on the evaluation metrics following Equation (7), a better performance in the noisy regime could still correspond to a lower robustness, if the noise-driven degradation is larger. While due to the semi-synthetic nature of all but one dataset we can only experiment with synthetic noise, the Twins dataset contains empirical treatment effects, which are computed under the assumption that a pair of twins have the same covariates. This assumption, as well as empirical measurement noise in the features, are likely to lead to real intrinsic noise in the Twins dataset, even before injecting synthetic noise. Notably, on this dataset TATEE achieves the lowest (best)  $\hat{\epsilon}_{\text{PEHE}}$  than all other benchmarks without simulated noise (see Appendix D.4) and outperforms CFR in all cases, including performance without synthetic noise, as well as both performance and robustness with simulated noise for all noise distribution parameters.

**Robustness to Noise.** Sampling additive feature noise from  $\alpha$ -stable distributions for 25  $(\alpha, \beta)$  pairs in the right half of the Feller-Takayasu diamond, we compute the relative gain in  $\rho$  on the test set after training on noisy features. The results (Figure 5) strongly confirm that TATEE is more robust than CFR. Aligned with the theoretical discussion following Theorem 3.1, the experiments also show that the largest values of gain in robustness (larger  $\rho_{\text{TATEE}} - \rho_{\text{CFR}}$ ) are observed when  $\alpha$  is closer to 1, corresponding to a heavier tail of the noise distribution, reaching near 12, 0.3, 51.0, and 68.6 percents of gain for the IHDP, Twins, Jobs, and ACIC datasets (respectively), when  $\alpha > 1$ .<sup>3</sup>

## D.3. Performance on Clean and Noisy Data

While TATEE aims to improve the robustness of treatment effect estimation, our experiments confirm that TATEE also has a performance comparable with CFR both in the absence of noise and on noisy inputs. Here we report the experimental results on the performance of TATEE compared to CFR, complementing the discussion on robustness. The main performance metric measuring CATE estimation error is PEHE, which we denote by  $\epsilon_{\text{PEHE}}$  for noise-free input and by  $\tilde{\epsilon}_{\text{PEHE}}$  for input with synthetic additive noise, as described in Appendix D.1. To compare the performances of TATEE and CFR on noisy features, we assess the relative gain in reduction of PEHE due to TATEE, given by  $1 - \frac{\tilde{\epsilon}_{\text{PEHE, TATEE}}}{\tilde{\epsilon}_{\text{PEHE, CFR}}}$  and report the percentage point gains in PEHE. Figure 8 visualizes the performance of TATEE across all four datasets under various noise distributions. These results underscore the efficacy of incorporating topological awareness into representation learning, enhancing robustness against diverse noise distributions while maintaining performance on both clean and noisy data. Altogether, the performance and robustness improvements observed across all four datasets align with our theory-informed expectations and demonstrate the practical utility of the TATEE framework in robust treatment effect estimation.

**Empirical Results on Performance.** On clean data, TATEE is consistently on par with CFR, with minimal deviations in  $\hat{\epsilon}_{\text{PEHE}}$  across all four datasets, showing  $1.977 \times 10^{-2}$ ,  $4.976 \times 10^{-4}$ , and  $2.631 \times 10^{-3}$  increase for IHDP, Jobs, and ACIC (respectively), and  $1.007 \times 10^{-2}$  decrease for Twins. This indicates that the integration of topological summaries does not compromise performance in noise-free scenarios. When subjected to additive noise from  $\alpha$ -stable distributions in the inputs, TATEE demonstrates superior performance by reducing  $\hat{\epsilon}_{\text{PEHE}}$  compared to CFR in most cases. On the IHDP and Twins datasets, TATEE consistently outperforms CFR across the entire range of noise parameters. For the Jobs and ACIC datasets, while the improvements are less consistent, TATEE still offers better performance in almost as many cases as CFR, while

<sup>3</sup>Recall that the first absolute moment of the noise distribution is infinite for  $\alpha = 1$ , hence, we need to be cautious about interpreting the average of trials at  $\alpha = 1$ .

consistently exhibiting more robustness, as illustrated in Figure 5 and discussed in Section 5 and Appendix D.2.

**Distinguishing Performance and Robustness.** Note the distinction between the performance and robustness metrics: a worse (higher)  $\hat{\epsilon}_{\text{PEHE}}$  on noisy inputs can still yield a better (higher)  $\rho$  if the noise-induced deterioration,  $\tilde{\epsilon}_{\text{PEHE}}/\hat{\epsilon}_{\text{PEHE}}$ , is comparatively smaller. In other words,  $\rho_{\text{TATEE}} - \rho_{\text{CFR}}$  can be positive even if both  $\tilde{\epsilon}_{\text{PEHE}}$  and  $\hat{\epsilon}_{\text{PEHE}}$  for TATEE exceed those of CFR, as long as TATEE’s increase in error due to noise remains more modest. This subtlety underscores why TATEE might show higher  $\hat{\epsilon}_{\text{PEHE}}$  overall, yet still achieve superior robustness as reflected by a larger  $\rho$ . Consequently, we observe scenarios where TATEE exhibits a smaller gap between its clean and noisy performances, reinforcing the principle that our method’s primary goal is to mitigate noise-driven degradation.

**Expected Performance Trade-offs.** The theory-informed advantage of TATEE lies in improving robustness to noise by including  $\mathcal{L}_{\text{topo}}$  in the loss. Since this regularization term can complicate the path toward minimizing other objectives, such as  $\hat{\epsilon}_{\text{PEHE}}$  on clean inputs, making a slight performance decline is not unexpected. Nonetheless, the evaluation above affirms that TATEE’s enhanced robustness does not come at the cost of a significant decline in performance: our experiments show only minimal increases in  $\hat{\epsilon}_{\text{PEHE}}$  without noise, and in many settings, TATEE even outperforms CFR when noise is injected in terms of  $\tilde{\epsilon}_{\text{PEHE}}$ . Additionally, the Twins dataset—featuring empirical treatment effects and inherent real noise—exemplifies TATEE’s robustness, as it achieves the lowest (best)  $\hat{\epsilon}_{\text{PEHE}}$  among all methods tested, with or without synthetic noise injection. This finding implies that in contexts where empirical noise exists, TATEE can achieve lower treatment effect estimation error, aligning with our theoretical claims.

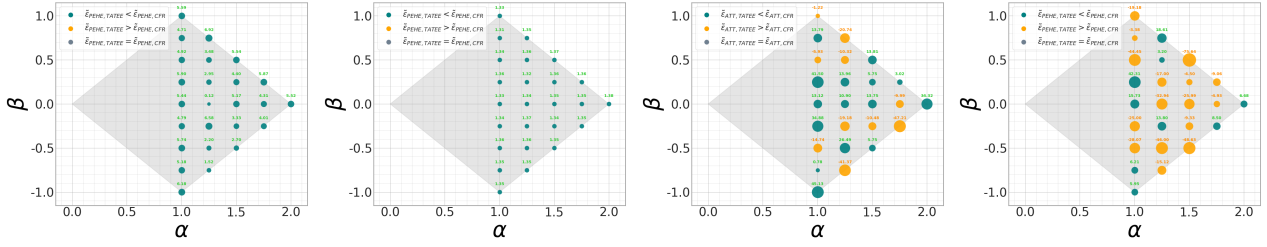


Figure 8: Evaluation of TATEE’s performance when estimating CATE on noise-corrupted features from the IHDP, Twins, Jobs, and ACIC datasets (from left to right) for a range of  $(\alpha, \beta)$  parameters of  $\alpha$ -stable distributed noise. The Feller-Takayasu diamond (shaded) marks valid  $(\alpha, \beta)$  values. The relative gain in  $\hat{\epsilon}_{\text{PEHE}}$  from using TATEE is visualized by the circles. The green and orange circles mark improvement/deterioration; the size of the circles is proportional to the magnitude of relative change.

#### D.4. Benchmark Comparisons

**Benchmark Models.** To conduct a more thorough evaluation of the performance and robustness of TATEE, we compare it against several causal inference models. These benchmark models encompass a range of methodologies, ensuring a thorough assessment across different approaches. In the main text of the paper we discussed CFR (Shalit et al., 2017), the most closely related model to TATEE, on which we base our main empirical evaluations. Here we name the other benchmark models used in our experiments described below and in Section 5. Dragonnet (Shi et al., 2019) is a neural network architecture designed to jointly model treatment assignments and potential outcomes by estimating propensity scores from the features, and GANITE (Yoon et al., 2018) leverages generative adversarial networks to generate counterfactual outcomes. S-Learner, T-Learner, and X-Learner (Künzel et al., 2019) utilize meta-learning approaches that adapt base learners to estimate treatment effects. TARNet (Shalit et al., 2017) is a CFR-type framework, learning separate representations for treated and control units, without the balancing loss term. Each of these models introduces unique mechanisms for addressing causal inference challenges. By benchmarking TATEE against these diverse models, we demonstrate its effectiveness and superior robustness across a broad spectrum of causal inference methods.

**Results.** Our experiments confirm TATEE’s enhanced robustness compared to the benchmark models in pairwise comparisons. As mentioned in Section 5, this superior robustness is shown in Figure 6, whose  $(i, j)$  entry shows the proportion of dataset-parameter pairs in which the model corresponding to row  $i$  is more robust than the one for column  $j$ , TATEE consistently demonstrates better robustness across most datasets and noise parameters. Figure 9 provides additional evidence



by showing the average robustness rank of each model (among the 8 models) for each pair of noise distribution parameters, as measured by  $\rho$ , averaged over the four datasets. A smaller number in row  $i$  and column  $j$  means, on average, the model corresponding to column  $j$  is more robust than others for the noise distribution parameters in row  $i$ . These results further confirm TATEE’s superior robustness. Meanwhile, the  $\hat{\epsilon}_{\text{PEHE}}$  values reported in Table 1 show that this gain in robustness does not come at the cost of a noticeable decline in performance. While the main purpose of TATEE is enhancing robustness, and in general, TATEE is not expected to perform better than CFR on noise-free inputs, the results indicate that TATEE has competitive performance in noise-free regimes as well, achieving the best or second best  $\hat{\epsilon}_{\text{PEHE}}$  on three out of four datasets. Notably, on the Twins dataset—the only dataset where we suspect intrinsic noise exists without injecting synthetic noise—TATEE achieves the best  $\hat{\epsilon}_{\text{PEHE}}$  without simulated noise.

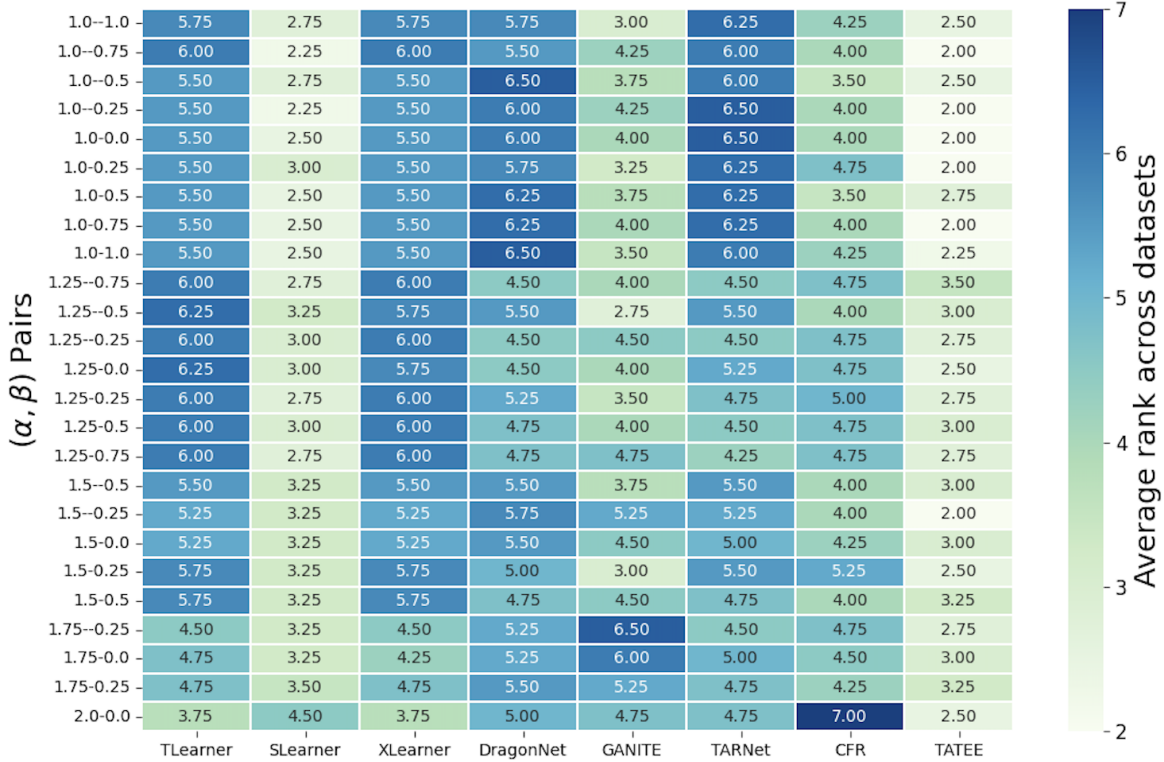


Figure 9: Comparing robustness of TATEE and the benchmark models: The quantity in column  $j$  and row  $i$  shows the average rank of the model corresponding to column  $j$  (out of 8), averaged over the model’s rank across all 4 datasets, for the noise distribution’s  $(\alpha, \beta)$  values in row  $i$ .

Table 1:  $\hat{\epsilon}_{\text{PEHE}}$  on input without synthetic noise. Lower means better performance, with best in bold and second best underlined. TATEE is best or second best (second to CFR) across all but one dataset. Notably, TATEE performs best on the Twins dataset—the only dataset which contains empirical values for ground-truth treatment effects, and hence can have real-world noise.

	IHDP	Twins	Jobs	ACIC
TLearner	<b>0.596 ± 0.000</b>	0.454 ± 0.000	4.546 ± 0.000	0.264 ± 0.000
SLearner	0.863 ± 0.000	0.418 ± 0.000	0.004 ± 0.000	0.255 ± 0.000
XLearner	<u>0.596 ± 0.000</u>	0.453 ± 0.000	4.546 ± 0.000	0.264 ± 0.000
DragonNet	1.234 ± 0.435	0.493 ± 0.105	0.457 ± 0.409	0.334 ± 0.259
GANITE	4.184 ± 0.216	0.727 ± 0.271	0.013 ± 0.013	0.457 ± 0.229
TARNet	1.268 ± 0.456	0.504 ± 0.147	0.456 ± 0.462	0.220 ± 0.152
CFR	0.870 ± 0.089	<u>0.414 ± 0.017</u>	<b>0.002 ± 0.003</b>	<b>0.006 ± 0.005</b>
TATEE	0.888 ± 0.088	<b>0.404 ± 0.029</b>	<u>0.003 ± 0.004</u>	<u>0.009 ± 0.010</u>

## E. Analysis of TATEE's Robustness

In Appendix C.5 we present a detailed analysis of the stability of TATEE's learning objective, relying on Theorem 3.1 to derive conditions which correspond to TATEE's enhanced robustness in counterfactual regression, stated in Proposition C.1. This proposition suggests that under the assumptions of Theorem 3.1 and when the constraint stated in this theorem is satisfied, for a suitable choice of training parameters, TATEE's loss function (Equation 6) is more robust than the CFR's loss, which does not account for the topology of the representations. Here we provide a proof for this proposition. Unless stated otherwise, the derivation here follows the problem setup and the notation defined in Section 3 and Appendix B.

### Notation.

Following a similar notation as in Theorem 3.1, let  $\mu^i$  be the measure over the representations of the control ( $i = 0$ ) and treatment ( $i = 1$ ) groups and  $\mu_f^i$  the measure over the persistence diagrams with a filtration function  $f$ . Also, similar to the notation used in Section 3.1, we use  $\tilde{\cdot}$  to denote the variables corresponding to the noisy input. Let  $W^{0,1} := W_p(\mu^0, \mu^1)$ ,  $\tilde{W}^{0,1} := W_p(\tilde{\mu}^0, \tilde{\mu}^1)$ ,  $\tilde{W}^0 := W_p(\tilde{\mu}^0, \mu^0)$ ,  $\tilde{W}^1 := W_p(\tilde{\mu}^1, \mu^1)$ ,  $W_{\text{topo}}^{0,1} := W_p(\mu_f^0, \mu_f^1)$ ,  $\tilde{W}_{\text{topo}}^{0,1} := W_p(\mu_{\tilde{f}}^0, \mu_{\tilde{f}}^1)$ ,  $\tilde{W}_{\text{topo}}^0 := W_p(\mu_{\tilde{f}}^0, \mu_f^0)$ ,  $\tilde{W}_{\text{topo}}^1 := W_p(\mu_{\tilde{f}}^1, \mu_f^1)$  denote the Wasserstein- $p$  distances. Following the notation used for the upper bounds in Section 3, for  $i \in \{0, 1\}$ , we use  $M^i$  and  $K_{\text{topo}}^i$  to denote the upper bounds on  $\tilde{W}^i$  and  $\tilde{W}_{\text{topo}}^i$ , and  $\hat{M}^i$  and  $\hat{K}_{\text{topo}}^i$  for their final sample equivalents.

**Proof.** Given any positive balancing loss coefficient for the CFR loss,  $\lambda_{\text{CFR}} > 0$ , a suitable choice of  $\lambda$  and  $\lambda_{\text{topo}}$  loss term coefficients of TATEE can always lead to satisfying the inequality below.

$$1 \leq \frac{\lambda_{\text{CFR}} - \lambda}{\lambda_{\text{topo}}}. \quad (8)$$

Meanwhile, by Theorem 3.1, when  $\varphi$  satisfies the constraint stated in the theorem, we have  $\hat{K}_{\text{topo}}^0 < \hat{M}^0$  and  $\hat{K}_{\text{topo}}^1 < \hat{M}^1$ , hence, using Inequality 8, we can write

$$1 \leq \frac{\lambda_{\text{CFR}} - \lambda}{\lambda_{\text{topo}}} \frac{\hat{M}^0 + \hat{M}^1}{\hat{K}_{\text{topo}}^0 + \hat{K}_{\text{topo}}^1}.$$

Rearranging the terms in this inequality gives

$$\lambda (\hat{M}^1 + \hat{M}^0) + \lambda_{\text{topo}} (\hat{K}_{\text{topo}}^1 + \hat{K}_{\text{topo}}^0) \leq \lambda_{\text{CFR}} (\hat{M}^1 + \hat{M}^0), \quad (9)$$

completing the proof for Proposition C.1.  $\square$

Notice that the result in Inequality 9 shows that the upper bound on the noise-induced change in the sum of the balancing and topological loss terms of TATEE is smaller than the upper bound on the change in the balancing term of CFR due to noise. These upper bounds are due to the triangle and quadrilateral inequalities (or applying the reverse triangle inequality and then the triangle inequality), which yield the following.

$$\begin{aligned} \left\| \tilde{W}^{0,1} - W^{0,1} \right\| &\leq \tilde{W}^0 + \tilde{W}^1 \leq M^0 + M^1, \\ \left\| \tilde{W}_{\text{topo}}^{0,1} - W_{\text{topo}}^{0,1} \right\| &\leq \tilde{W}_{\text{topo}}^0 + \tilde{W}_{\text{topo}}^1 \leq K_{\text{topo}}^0 + K_{\text{topo}}^1, \\ \left| (\lambda \tilde{W}^{0,1} + \lambda_{\text{topo}} \tilde{W}_{\text{topo}}^{0,1}) - (\lambda W^{0,1} + \lambda_{\text{topo}} W_{\text{topo}}^{0,1}) \right| &\leq \lambda \left| \lambda \tilde{W}^{0,1} - W^{0,1} \right| \\ &\quad + \lambda_{\text{topo}} \left| \lambda_{\text{topo}} \tilde{W}_{\text{topo}}^{0,1} - W_{\text{topo}}^{0,1} \right|, \\ &\leq \lambda (M^0 + M^1) + \lambda_{\text{topo}} (K_{\text{topo}}^0 + K_{\text{topo}}^1). \end{aligned}$$

## F. $\alpha$ -Stable Distributions

$\alpha$ -stable distributions are a rich class of probability distributions that allow modeling data with heavy tails and asymmetry. Their flexible parametric form is particularly suitable for data that exhibit extreme values. Random variables with  $\alpha$ -stable distributions do not necessarily have finite mean or variance. This family of distributions facilitate a generalization of the

central limit theorem (CLT): While the normal distribution arises as the limit distribution when summing independent, thin-tailed random variables with finite second moment,  $\alpha$ -stable distributions with the tail parameter  $\alpha < 2$  arise as the limit when the random variables have infinite variance (Mainardi, 2007; Gnedenko & Kolmogorov, 2021).

The characteristic function of a random variable with an  $\alpha$ -stable distribution is as follows,

$$\varphi(t) = \begin{cases} \exp \{itm - |ct|^\alpha (1 - i\beta \operatorname{sgn}(t) \tan \frac{\pi\alpha}{2})\} & \alpha \neq 1 \\ \exp \{itm - |ct|^\alpha (1 + i\beta \frac{2}{\pi} \operatorname{sgn}(t) \log |t|)\} & \alpha = 1. \end{cases} \quad (10)$$

The parameters  $\alpha, \beta, c \geq 0$ , and  $m \in \mathbb{R}$ , which we refer to as the *characteristic exponent* or the *tail* parameter, *skewness* or *asymmetry* parameter, *scale* parameter, and *location* parameter, determine the distribution of the random variable.

Note that the main parameters of interest, determining the shape of the distribution, are the tail and skewness parameters. These parameters take values in a region of the plane where  $\alpha \in (0, 2]$  and  $|\beta| \leq \min \{\alpha, 2 - \alpha\}$ , which is referred to as the *Feller-Takayasu diamond* (Mainardi, 2007). The skewness parameter  $\beta \in [-1, 1]$  introduces asymmetry; a symmetric  $\alpha$ -stable distribution has  $\beta = 0$ . The parameter  $\alpha$  determines how fast the tail decays, with smaller values meaning heavier tails. The normal distribution is the special case when  $\alpha = 2$ , and  $\alpha = 1$  corresponds to the Cauchy distribution. Except the Gaussian case, other  $\alpha$ -stable distributions have infinite moments of order greater than or equal to  $\alpha$ , i.e., infinite variance for all  $\alpha < 2$ , and infinite absolute mean for  $\alpha \leq 1$ . Most statistical models rely on the CLT to justify using Gaussian distributions as the asymptotic distribution of the sum/mean of an independent and identically distributed (i.i.d.) sequence of random variables. However, the CLT does not hold for  $\alpha$ -stable distributions with  $\alpha < 2$ . This fact is consequential for the sample size criterion and significance level of statistical tests (Pitera et al., 2022).

## G. Implementation Details

Recall that the main purpose of our experiments is to evaluate the change in robustness to input noise due to our proposed topological regularization in TATEE. To this end, in order to assess the impact of incorporating the regularization term based on persistence diagrams into the training of the model, and not the CFR-type model architecture, we adopted the same parameters and overall configuration utilized by Shalit et al. (2017) for CFR, for all implementation and training purposes of TATEE, as well as the CFR model we compared against TATEE. Each of  $\varphi, h_0$ , and  $h_1$  have 3 fully connected layers of size 200 for  $\varphi$  and 100 for  $h_0$  and  $h_1$ . Following Shalit et al. (2017), we used the value  $\sqrt{10}$  for the  $\lambda$  regularization coefficient for the balancing term of the training objective in Equation 6, which Shalit et al. (2017) found to yield the lowest  $\epsilon_{PEHE}$  on the IHDP dataset. After a standard grid search<sup>4</sup> fine-tuning of  $\lambda_{\text{topo}}$ , we used  $\lambda_{\text{topo}} = \frac{\sqrt{10}}{4}$  for the topological regularization coefficient in Equation 6, in the case of IHDP dataset. We also performed a grid search fine tuning over batch size,  $\lambda$ , and  $\lambda_{\text{topo}}$  for the Twins and Jobs datasets, with a parameters grid of size 80, from which we obtained the values 128, 10, and  $2\sqrt{10}$  for the Twins dataset, 128, 0.1, and 1 for the Jobs dataset, and 256, 10, and 0.1 for the ACIC dataset. The learning rate is  $10^{-2}$  for mini-batch training using an Adam optimizer (Kingma & Ba, 2015) with weight decay parameter value of  $10^{-5}$ . The system specification for the computations is provided in Table 2.

Table 2: System specifications for the computations.

CPU	Intel(R) Xeon(R) CPU @ 2.20GHz
GPU	Nvidia A100 SXM4 40GB
OS	Ubuntu 22.04.3 LTS
Architecture	x86_64

## H. Datasets

**Infant Health and Development Program (IHDP).** The IHDP dataset is a semi-synthetic benchmark widely used in causal inference studies (Hill, 2011). It combines real-world covariates from a randomized experiment with simulated counterfactual outcomes, providing a ground truth for evaluating treatment effect estimation methods.

**Twins.** The Twins dataset uniquely contains both factual and counterfactual outcomes, as it includes data on twin pairs (Louizos et al., 2017). By treating one twin as treated and the other as control, the dataset provides real-world values for both

<sup>4</sup>In this case the search is in fact over a line, as we fine tuned only a single parameter.

potential outcomes, eliminating the need for synthetic counterfactuals. This characteristic allows for evaluation of treatment effects estimation against empirical ground-truth, making it the only dataset for assessing the robustness of methods in the presence of intrinsic empirical noise, without injecting synthetic noise. Notably, TATEE outperforms all benchmarks on this dataset (see Appendix D.4) in experiments without simulated noise, achieving the lowest (best)  $\epsilon_{PEHE}$ .

**Jobs.** The Jobs dataset, originally collected by LaLonde (1986) and later curated for causal inference benchmarking by Shalit et al. (2017), includes both randomized controlled trial (RCT) and observational data. The dataset combines treated units from the RCT subset with a control sample from Dehejia & Wahba (1999), ensuring that treatment assignment depends on covariates.

**Atlantic Causal Inference Conference (ACIC) 2018.** The ACIC-18 dataset (Athey & Wager, 2019) is based on data from the National Study of Learning Mindsets (Yeager et al., 2016). Assuming an RCT design, the dataset allows for the computation of the true average treatment effect on the treated (ATT). To introduce covariate shift, a mask variable sampled from a Bernoulli distribution based on a feature is applied to control units, making treatment assignment dependent on the covariates.