

---

# AN EVALUATION FRAMEWORK FOR COMPARING CAUSAL INFERENCE MODELS

---

A PREPRINT

 **Niki Kiriakidou\***

Department of Informatics and Telematics,  
Harokopio University of Athens,  
Athens, GR177 78.  
kiriakidou@hua.gr

 **Christos Diou**

Department of Informatics and Telematics,  
Harokopio University of Athens,  
Athens, GR177 78.  
cdiou@hua.gr

September 2, 2022

## ABSTRACT

Estimation of causal effects is the core objective of many scientific disciplines. However, it remains a challenging task, especially when the effects are estimated from observational data. Recently, several promising machine learning models have been proposed for causal effect estimation. The evaluation of these models has been based on the mean values of the error of the Average Treatment Effect (ATE) as well as of the Precision in Estimation of Heterogeneous Effect (PEHE). In this paper, we propose to complement the evaluation of causal inference models using concrete statistical evidence, including the performance profiles of Dolan and Moré, as well as non-parametric and post-hoc statistical tests. The main motivation behind this approach is the elimination of the influence of a small number of instances or simulation on the benchmarking process, which in some cases dominate the results. We use the proposed evaluation methodology to compare several state-of-the-art causal effect estimation models.

**Keywords** Causal inference · treatment effects · performance profiles · non-parametric tests · post-hoc tests

## 1 Introduction

Causal inference is a fundamental problem in many scientific areas such as medicine (Höfler, 2005), education (Cordero et al., 2018) and economy (Hoover, 2012). The most effective way to infer the causal effect of a treatment (i.e. intervention) to an outcome is through a randomized controlled trial. However, in many cases, conducting a randomized controlled trial is not possible, due to financial, ethical or other constraints. Therefore, researchers must determine the effect of treatments by relying on observational data.

Observational (as opposed to experimental) data refer to data obtained without any control over independent variables. Researchers simply observe and record the data and do not affect sampling and treatment assignment. In today's big data world, observational data are abundant; nevertheless, using such data to infer causal effects still remains a challenge.

The main obstacle is that for every subject only the factual outcome is observed, i.e. the treatment/outcome combination, which actually took place. The counterfactual outcome is what would have happened, in case we have chosen a different treatment and kept everything else constant. In addition to this problem, treatment assignment is not completely at random, and depends on other factors (e.g., age, socio-economic status or other medical conditions, in the health domain). This results in significant differences between the population that received the treatment (treatment group) and the rest (control group).

---

\*Corresponding author.

In this paper, we consider the case of a binary treatment  $T$ . Let  $\mathbf{X}^{n \times m}$  be the design matrix, where  $n$  is the number of instances and  $m$  is the number of features,  $T \in \{0, 1\}^n$  is the treatment assignment and  $Y$  is a random variable, with  $Y_0$  denoting the outcome for a sample when  $T = 0$  and  $Y_1$  the outcome when  $T = 1$ .

Let's assume that  $\mathbf{x}_i \in \mathbf{X}$ ,  $t_i \in [0, 1]$  and  $y(\mathbf{x}_i, t_i)$  is the outcome of a unit  $i = 1, 2, \dots, n$ . The average treatment effect is defined as follows:

$$ATE = \frac{1}{n} \sum_{i=1}^n [y(\mathbf{x}_i, 1) - y(\mathbf{x}_i, 0)] \quad (1)$$

It is worth mentioning that, for every unit  $i$  in the dataset, only the factual outcome is observed, i.e.  $y(\mathbf{x}_i, 1)$  or  $y(\mathbf{x}_i, 0)$ , which stands for the outcome of each unit in treatment and control group, respectively. As a result, we need to estimate the counterfactual outcomes. The estimated outcomes are denoted as  $\hat{y}(\mathbf{x}_i, 1)$  and  $\hat{y}(\mathbf{x}_i, 0)$ , for treatment and control group, respectively. Then, the estimation of average treatment effect is defined as:

$$\hat{ATE} = \frac{1}{n} \sum_{i=1}^n [\hat{y}(\mathbf{x}_i, 1) - \hat{y}(\mathbf{x}_i, 0)] \quad (2)$$

For the population causal effect we report the absolute error on the average treatment effect,

$$|\epsilon_{ATE}| = \left| \frac{1}{n} \sum_{i=1}^n [y(\mathbf{x}_i, 1) - y(\mathbf{x}_i, 0)] - \frac{1}{n} \sum_{i=1}^n [\hat{y}(\mathbf{x}_i, 1) - \hat{y}(\mathbf{x}_i, 0)] \right| \quad (3)$$

In order to measure the accuracy of the individual treatment effect estimation, we use the Precision in Estimation of Heterogeneous Effect (PEHE), which is defined as:

$$\epsilon_{PEHE} = \frac{1}{n} \sum_{i=1}^n [(y(\mathbf{x}_i, 1) - y(\mathbf{x}_i, 0)) - (\hat{y}(\mathbf{x}_i, 1) - \hat{y}(\mathbf{x}_i, 0))]^2 \quad (4)$$

In the literature, the traditional approach to evaluate the performance of causal inference models is to compare the average value of  $|\epsilon_{ATE}|$  and/or  $\epsilon_{PEHE}$  (Louizos et al., 2017; Shi et al., 2019; Shalit et al., 2017). However, this approach may provide us misleading results, since all simulations are equally considered which implies that the difficulty of each simulation of the benchmarking process is not taken into account. As a result, a small number of the most difficult problems can tend to dominate these results (Dolan and Moré, 2002; Livieris and Pintelas, 2020). Additionally, another major drawback is that this approach does not provide us any information whether the performance of two or more models is equivalent neither quantify the difference between their performance.

In this work, we propose a comprehensive evaluation framework for comparing models for treatment effects and PEHE. The proposed framework is based on evaluating the performance of causal effect estimation models as estimators and as predictors. For providing concrete and empirical evidence, we utilize performance profiles of Dolan and Moré (2002), as well as non-parametric and post-hoc tests (Derrac et al., 2011). The rationale behind our approach is that when evaluating causal effect estimation models using Eqs. (3) and (4) in multiple experimental evaluations, it is common for a small number of simulations to dominate the results. In these cases, reporting only the average value of the errors may be misleading.

The remainder of this paper is organised as follows: Section 2 presents a review of neural network-based and tree-based models for the estimation of treatment effects. Section 3 provides a detailed description of the theoretical framework of performance profiles as well as a complete presentation of statistical multiple comparison analysis, in order to evaluate the performance of causal inference models. Section 4 provides information about the two datasets we used in this work. Section 5 presents a detailed experimental analysis, focusing on the evaluation of the proposed model and an extended multiple comparisons statistical analysis of models Section 6 summarizes the main findings and concludes this paper by identifying interesting directions for future work.

## 2 Causal inference models

In the literature, several models have been proposed for the estimation of causal effects. Here, we focus our attention to the most widely used and efficient causal inference models, which can be divided into two categories; neural network-based models and tree-based models.

## 2.1 Neural network-based models

Shalit et al. (2017) proposed the Counterfactual Regression (CFR) framework for estimating individual treatment effects. CFR aims at learning a balanced representation using a prediction model, so that the control and the treatment group distributions look similar. The authors used two integral probability metrics: Wasserstein distance (Wass) (Villani, 2009) and Maximum Mean Discrepancy (MMD) (Gretton et al., 2012), in order to measure the distances between two distributions. Furthermore, they introduced a generalization bound for the estimation of individual treatment effect where every individual is only identified by its features. The authors also proposed the Treatment Agnostic Representation Network (TARNet), which is a variant without balance regularization.

Shi et al. (2019) proposed Dragonnet, which consists of a neural network model for estimating treatment effects from observational data. Dragonnet aims at improving the estimations of average treatment effect and individual treatment effect through the propensity score (i.e., the probability  $Pr(T = 1 | X = \mathbf{x})$  that a particular sample with covariates  $\mathbf{x}$  has received the treatment). Furthermore, the authors proposed a procedure to induce bias based on non-parametric estimation theory to further improve treatment effect estimation. The authors also evaluated Dragonnet to a multi-stage procedure, named NEDnet, which is a neural network with similar architecture with Dragonnet. More specifically, NEDnet first trained using a pure treatment prediction objective. Then, the last layer is removed, and replaced with an outcome-prediction neural network matching the one used by Dragonnet.

## 2.2 Tree-based models

Chipman et al. (2010) developed a Bayesian “sum-of-trees” model named Bayesian Additive Regression Trees (BART). BART model is based on a Bayesian non-parametric approach, which fits a parameter rich model by utilizing a strongly influential prior distribution. More specifically, it uses the sum of trees to approximate the average value of the outcomes given a set of covariates,  $E[Y|\mathbf{x}]$ . The main idea of BART model is to impose a prior, which regularizes the fit by keeping the individual tree effects small in order to elaborate the sum-of-trees model. Additionally, for fitting the sum-of-trees model, BART uses a tailored version of Bayesian backfitting Markov Chain Monte Carlo (Hastie and Tibshirani, 2000).

Künzel et al. (2019) proposed a new methodology for predicting treatment effects. The main idea is to estimate the outcome by using all of the features together with the treatment indicator, without proving any special role to the treatment indicator. In more simple words, the treatment indicator is included and treated by the based learner like any other feature. In the literature, a variety of base learners such as Linear Regression (Montgomery et al., 2021), Random-Forest (Breiman, 2001) and  $k$ -Nearest Neighbors (Aha, 2013) were used providing some interesting results. However, this approach has two disadvantages; (i) in case the treatment and control groups are very different in covariates, a single model is probably not sufficient to encode the different relevant dimensions and smoothness of features for both groups (Alaa and Schaar, 2018); (ii) a tree-based base learner may completely ignore the treatment assignment by not choosing/splitting on it (Künzel et al., 2019).

Wager and Athey (2018) developed a forest-based method for estimating heterogeneous treatment effect. This method consists of an extension of the efficient and widely used Random-Forest algorithm (Breiman, 2001). In more detail, their proposed method, named Causal Forest (C-Forest) is composed by a number of causal trees, which estimates the effect of the treatment at the leaves of the trees. A significant advantage of C-Forest over the traditional approaches to non-parametric estimation of heterogeneous treatment effects is that the performance of the former is not degraded as the number of covariates is increasing (Wager and Athey, 2018).

## 3 Proposed evaluation framework

In this section, we present a comprehensive framework for the performance evaluation of the neural network-based and tree-based models. In the literature, the traditional way for evaluating the performance of causal inference models is through the average value of the  $|\epsilon_{ATE}|$  and  $\epsilon_{PEHE}$ . However, this approach gives us misleading results, since all the number of problems, even the difficult ones, are equally considered for the evaluation of the model. For this purpose, we use the *performance profiles* of Dolan and Moré (Dolan and Moré, 2002) and a *statistical multiple comparison analysis* for evaluating the performance of the models in order to overcome the problem of the domination of a small number of the most difficult problems over the results.

Next, we provide a detailed presentation of the tools which compose the proposed framework for evaluating the performance of models to infer causality.

### 3.1 Performance profiles

Dolan and Moré (2002) presented a methodology for evaluating the effectiveness and robustness of the set of models  $M$  on a test set  $S$ . More specifically, the authors proposed the performance profiles, which provides a wealth of information such as model’s efficiency, robustness and probability of success in compact form (Livieris et al., 2018; Livieris, 2020).

In more detail, suppose that there are  $n_p$  problems and a set of  $n_M$  models for every simulation  $s$ . Furthermore, let  $a_{s,m}$  be  $\epsilon_{ATE}$  or  $\epsilon_{PEHE}$  by model  $m$  for simulation  $s$ . It is compared the performance on simulation  $s$  by model  $m$  with the best performance by any model on this problem, utilizing the performance ratio defined as follows

$$r_{s,m} = \frac{a_{s,m}}{\min[a_{s,m} : m \in M]} \quad (5)$$

It is also required to obtain an overall assessment of performance, except of the performance of the model  $m$  on a given simulation. Therefore, we are calculating the cumulative distribution function for the performance ratio

$$p_m(a) = \frac{1}{n_s} \text{size}[s \in S : r_{s,m} \leq a], \quad (6)$$

where  $a \in \mathbb{R}$  is the factor of the best possible ratio.

Notice that the *performance profile*  $p_m(a) : \mathbb{R} \rightarrow [0, 1]$  for a model is a non-decreasing, piecewise constant function, continuous from the right at each breakpoint (Dolan and Moré, 2002). In particular, the performance profile plots the fraction  $P$  of problems for which any given model is within a factor  $a$  of the best model. This means that the model which outperforms the rest of the models, is the one whose performance profile plot lies on top right.

It is worth highlighting that the use of performance profiles eliminates the influence of a small number of problems on the benchmarking process, which tend to dominate the results (Livieris and Pintelas, 2020, 2019). Finally, the vertical axis provides the percentage of the simulations, which were successfully addressed by each model according to the factor of the best solver (efficiency), while the horizontal axis summarizes the percentage of the problems for which a model exhibits the best performance (robustness).

### 3.2 Statistical multiple comparison analysis

In the statistical literature, comparison among multiple models is usually carried out by means of non-parametric statistical tests, such as the Friedman test, which constitutes a well-known and widely utilized procedure for examining the differences between more than two models.

The Friedman test constitutes a non-parametric test, analogue of the parametric two-way analysis of variance (Friedman, 1940). Its main objective is to detect significant differences between the effectiveness of evaluated models, based on a sample of simulations. An attractive advantage of Friedman test is that the commensurability of the measures across different simulation is not required, since this test is non-parametric (Derrac et al., 2011). Furthermore, since it does not assume the normality of the sample means, it is robust to outliers. This non-parametric test, ranks the scores of each model and uses them in the calculation of the statistic, instead of using the scores themselves. Notice that in case of ties, average ranks are computed. It is worth mentioning that this test ranks the models from the best to the worst (Derrac et al., 2011; García et al., 2010).

Suppose that  $r_i^j$  be the rank of the  $j$ -th of  $k$  models on the  $i$ -th of  $M$  problems. The Friedman test requires the computation of the average ranks of algorithms  $R_j = \frac{1}{n} \sum_i r_i^j$ . Under the null-hypothesis  $H_0 : \{\text{all models perform similarly with non-significant differences}\}$ , the Friedman test statistic is calculated by:

$$F_j = \frac{12n}{k(k+1)} \left[ \sum_j R_j^2 - \frac{k(k+1)^2}{4} \right], \quad (7)$$

which is distributed according to a  $\chi^2$  distribution with  $k - 1$  degrees of freedom.

Once Friedman’s test rejects  $H_0$ , we may proceed with a post-hoc test in order to identify which pair of models differ significantly. In this research, we employ Bergmann-Hommel test, which was demonstrated as the most powerful test for determine the distinctive models in pairwise comparisons (Bergmann and Hommel, 1988; Garcia and Herrera, 2008).

For presenting Bergmann-Hommel’s procedure, we need the following definition.

**Definition 1 (Derrac et al., 2011).** An index set of hypotheses  $I \subseteq \{1, \dots, m\}$  is called *exhaustive* if exactly all  $H_j$  with  $j \in I$ , could be true.

Under Definition 1, the Bergmann-Hommel procedure rejects any hypothesis  $H_j$  with  $j$  is not included in the acceptance set  $A$  is the index set of null hypotheses, which are retained and it is defined by

$$A = \bigcup \{I : I \text{ exhaustive, } \min\{P_i : i \in I\} > \alpha/|I|\}. \quad (8)$$

Additionally, Wright (1992) summarized the formula for computation of each Adjusted P-Value (APV), that is

$$\text{APV}_i = \min\{v; 1\}, \quad (9)$$

where  $v = \max\{||I|| \cdot \min\{p_j, j \in I\} : I \text{ exhaustive}; i \in I\}$ .

## 4 Data

Generally, it is a challenging task to evaluate the performance of a model based on real-world data, due to the nature of the problem on causal inference, since we scarcely have access to the ground truth causal effects. However, to deal with this difficulty, we count on synthetic and semi-synthetic datasets for the empirical evaluation of causal estimation procedures.

**IHDP dataset.** The first dataset used for the estimation of individual and population causal effects is the semi-synthetic IHDP dataset, which was introduced by Hill (2011). This dataset was constructed from the Infant Health and Development Program and the outcome and treatment assignment are fully known. It comprises of 25 features regarding children and mothers and contains 747 units in which 608 belong to the control group is 608 while the rest 139 belong to the treatment group. Additionally, we studied the effect of home visits by specialists on future cognitive test scores. Finally, we utilized 1000 realizations from the NPCI package (Dorie, 2016).

**Synthetic dataset.** This dataset is a toy dataset introduced by Louizos et al. (2017) and it is generated conditioned on the hidden confounder variable  $Z$ . More analytically, the process for the generation of synthetic dataset is the following:

$$\begin{aligned} z_i &\sim \text{Bern}(0.5) \\ t_i|z_i &\sim \text{Bern}(0.75z_i + 0.25(1 - z_i)) \\ x_i|z_i &\sim \mathcal{N}(z_i, \sigma_{z_1}^2 z_i + \sigma_{z_0}^2 (1 - z_i)) \\ y_i|t_i, z_i &\sim \text{Bern}(\text{Sigmoid}(3(z_i + 2(2t_i - 1)))) \end{aligned}$$

where the treatment variable  $T$  is a mixture of Bernoulli distribution, the proxy to the confounder  $X$  is a mixture of Gaussian distribution, the outcome  $Y$  is determined as a logistic sigmoid function,  $\sigma_{z_0} = 3$  and  $\sigma_{z_1} = 5$ . This generation process introduces hidden confounding between  $T$  and  $Y$  as they both depend on the mixture assignment  $Z$ .

## 5 Experimental Analysis

In this section, we provide a detailed experimental analysis of the performance of neural network-based and tree-based models for IHDP and Synthetic datasets. For both datasets, we present the performance profiles and statistical multiple comparison analysis of the models.

The performance of each model was measured using the metrics  $|\epsilon_{ATE}|$  and  $\epsilon_{PEHE}$ , which are respectively defined by 3 and 4, respectively, as in (Shi et al., 2019; Shalit et al., 2017). It is worth highlighting, that  $|\epsilon_{ATE}|$  metric and  $\epsilon_{PEHE}$  are used to compare the evaluated neural network and tree-based models as estimators and predictors.

The implementation code was written in Python 3.7 using Keras library (Gulli and Pal, 2017) while the detailed experimental results for each model regarding both datasets can be found in [redacted for review]

Next, we evaluate the performance of:

- “Dragonnet”, which stands for Dragonnet model of Shi et al. (2019).
- “TARNet”, which stands for TARNet model of Shalit et al. (2017).
- “NEDnet”, which stands for NEDnet model of Shi et al. (2019).
- “R-Forest”, which stands for "S-learner" methodology of Künzel et al. (2019) using as base learner Random-Forest (Breiman, 2001).

- “C-Forest”, which stands for C-Forest model of Wager and Athey (2018).
- “BART”, which stands for BART model of Chipman et al. (2010).

All models used the parameters introduced in their original papers.

### 5.1 Results on IHDP dataset

Figure 1 presents the performance profiles of neural network-based and tree-based model, based on  $\epsilon_{ATE}$  metric. Dragonnet exhibited the best performance in terms of efficiency, slightly outperforming TARNet, R-Forest and C-Forest. More specifically, Dragonnet reported 23% of simulations with the best (lowest)  $\epsilon_{ATE}$ , while TARNet, R-Forest and C-Forest presented 20%, 20% and 19%, respectively. In contrast, BART and NEDnet exhibited the worst performance, solving only the 9% and the 8% of simulations with the lowest  $|\epsilon_{ATE}|$ , respectively. Finally, it is worth noticing that Dragonnet and TARNet demonstrated the best performance in terms of robustness, since their curves lie on top.

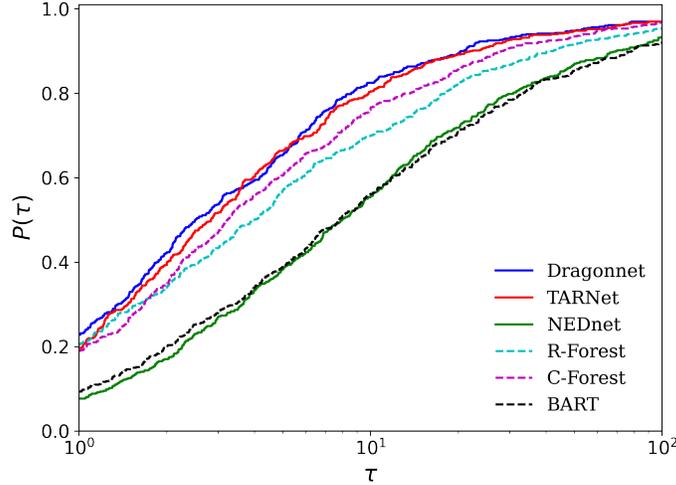


Figure 1: Log<sub>10</sub> scaled performance profiles based on  $|\epsilon_{ATE}|$

Table 1 presents the Friedman average ranking of all evaluated models, which represent the associated effectiveness of each model. Notice that the models are ordered from the best (lowest) to the worst (highest) ranking. Clearly, Dragonnet was the best performing model followed by TARNet, whereas BART and NEDnet were the worst. Furthermore, the Friedman statistic  $F_f$  with 5 degrees of freedom is equal to 314.33 while the  $p$ -value is equal to  $1.73 \cdot 10^{-10}$ , which suggests the existence of significant differences among the evaluated models.

| Algorithm | Ranking |
|-----------|---------|
| Dragonnet | 2.755   |
| TARNet    | 3.016   |
| C-Forest  | 3.227   |
| R-Forest  | 3.381   |
| BART      | 4.245   |
| NEDnet    | 4.376   |

Table 1: Friedman average rankings of evaluated models based on  $|\epsilon_{ATE}|$  for IHDP

Table 2 presents the information about the state of rejection of all the hypotheses, comparing the models, by summarizing the APVs with Bergmann Hommel’s procedure  $p_{Berg}$  with  $\alpha = 5\%$  level of significance for the 15 established comparisons. Each row contains a hypothesis if the first model (left side) outperforms the second one (right side), the corresponding  $p_{Berg}$  value and if the hypothesis is rejected or not. Notice that the hypotheses are ordered from the most to the least significant differences.

In Table 2 it is worth mentioning that Dragonnet outperforms all tree-based models as well as NEDnet, and has similar performance with TARNet, relative to  $|\epsilon_{ATE}|$ . TARNet outperforms R-Forest, BART and NEDnet and performed

equally well with C-Forest. Furthermore, C-Forest and R-Forest performed similarly and they were only statistically outperformed by Dragonnet. Finally, NEDnet and BART reported the worst performance according both Friedman’s and Bergmann’s test, since none of them outperforms any other model.

| Hypothesis            | $p_{Berg}$ |                       |
|-----------------------|------------|-----------------------|
| NEDnet vs Dragonnet   | 0          | Rejected              |
| Dragonnet vs BART     | 0          | Rejected              |
| NEDnet vs TARNet      | 0          | Rejected              |
| TARNet vs BART        | 0          | Rejected              |
| NEDnet vs C-Forest    | 0          | Rejected              |
| C-Forest vs BART      | 0          | Rejected              |
| R-Forest vs NEDnet    | 0          | Rejected              |
| R-Forest vs BART      | 0          | Rejected              |
| R-Forest vs Dragonnet | 0.000001   | Rejected              |
| Dragonnet vs C-Forest | 0.000265   | Rejected              |
| R-Forest vs TARNet    | 0.008147   | Rejected              |
| Dragonnet vs TARNet   | 0.082183   | Failed to be rejected |
| TARNet vs C-Forest    | 0.149083   | Failed to be rejected |
| R-Forest vs C-Forest  | 0.386149   | Failed to be rejected |
| NEDnet vs BART        | 0.386149   | Failed to be rejected |

Table 2: Multiple comparison test: Bergmann-Hommel’s APVs based on  $|\epsilon_{ATE}|$  for IHDP

Figure 2 summarizes the the conducted findings and conclusions of Table 2 and Table 2. More specifically,  $x$ -axis presents the models based on Friedman ranking. For each of them,  $y$ -axis collects the models which were statistically outperformed according to the Bergmann–Hommel test.

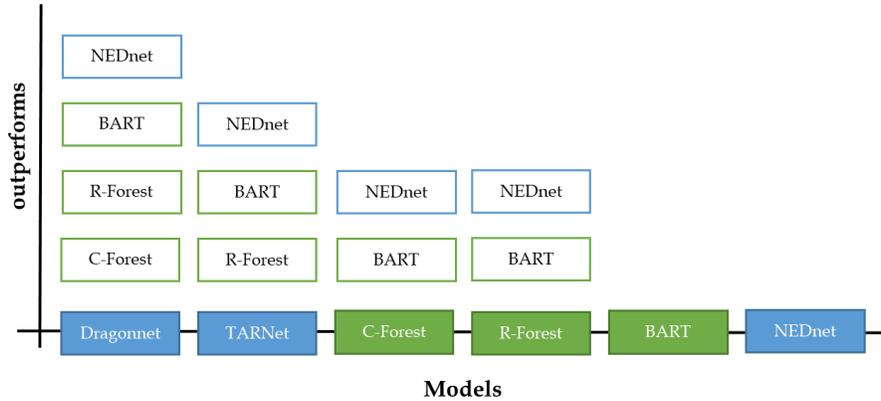
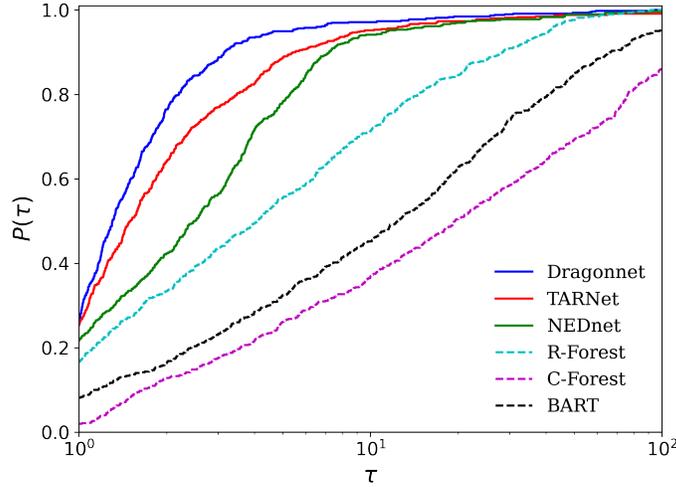


Figure 2: Conclusions for comparison based on  $|\epsilon_{ATE}|$  for IHDP

Figure 3 presents the performance profiles of the selected models, based on  $\epsilon_{PEHE}$  metric. In terms of efficiency, Dragonnet and TARNet presented the best performance, followed by NEDnet and R-Forest. In more detail, Dragonnet and TARNet reported 26% and 25% with the lowest  $\epsilon_{PEHE}$ , respectively, while NEDnet and R-Forest exhibited 22% and 17%, respectively. Additionally, BART and C-Forest presented the lowest performances solving only the 8% and 2% of simulations, respectively. In terms of robustness, Dragonnet illustrated the top curve.

Table 3 presents the Friedman average ranking of all evaluated models, which represent the associated effectiveness of each model. Dragonnet was the best performing model followed by TARNet, while C-Forest was the worst. In addition, the Friedman statistic  $F_f$  with 5 degrees of freedom is equal to 1029.74 while the  $p$ -value is equal to 0, which strongly suggests the existence of significant differences between the evaluated models.

Figure 3: Log<sub>10</sub> scaled performance profiles based on  $\epsilon_{PEHE}$ 

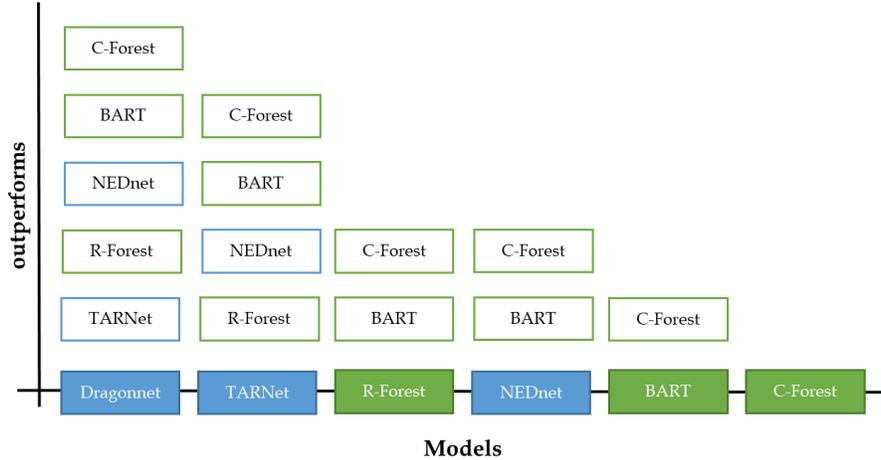
| Algorithm | Ranking |
|-----------|---------|
| Dragonnet | 2.098   |
| TARNet    | 2.782   |
| R-Forest  | 3.150   |
| NEDnet    | 3.166   |
| BART      | 4.372   |
| C-Forest  | 5.432   |

Table 3: Friedman average rankings of evaluated models based on  $\epsilon_{PEHE}$  for IHDP

Table 4 suggests that Dragonnet outperformed all other models, while TARNet outperformed NEDnet and all tree-based models. R-Forest and NEDnet had no statistically important differences in their performance and both outperformed C-Forest and BART. Eventually, BART and C-Forest exhibited the worst performance in terms of  $\epsilon_{PEHE}$ . A graphical overview of the conducted findings of the statistical analysis is presented in Figure 4.

| Hypothesis            | $p_{Berg}$ |                       |
|-----------------------|------------|-----------------------|
| Dragonnet vs C-Forest | 0          | Rejected              |
| TARNet vs C-Forest    | 0          | Rejected              |
| R-Forest vs C-Forest  | 0          | Rejected              |
| Dragonnet vs BART     | 0          | Rejected              |
| NEDnet vs C-Forest    | 0          | Rejected              |
| TARNet vs BART        | 0          | Rejected              |
| R-Forest vs BART      | 0          | Rejected              |
| NEDnet vs BART        | 0          | Rejected              |
| NEDnet vs Dragonnet   | 0          | Rejected              |
| C-Forest vs BART      | 0          | Rejected              |
| R-Forest vs Dragonnet | 0          | Rejected              |
| Dragonnet vs TARNet   | 0          | Rejected              |
| NEDnet vs TARNet      | 0.003519   | Rejected              |
| R-Forest vs TARNet    | 0.003519   | Rejected              |
| R-Forest vs NEDnet    | 0.892434   | Failed to be rejected |

Table 4: Multiple comparison test: Bergmann-Hommel's APVs based on  $\epsilon_{PEHE}$  for IHDP

Figure 4: Conclusions for comparison based on  $\epsilon_{PEHE}$  for IHDP

Summarizing, by taking into consideration both the performance profiles and statistical analysis, we are able to conclude that Dragonnet and TARNet outperformed the rest models in terms of  $|\epsilon_{ATE}|$  and  $\epsilon_{PEHE}$ . This suggests that they reported the best performance, as estimators and predictors. Finally, it is worth mentioning, that our experimental analysis illustrated that Dragonnet outperformed TARNet in terms of robustness for both metrics.

## 5.2 Results on Synthetic dataset

Figure 5 presents the performance profiles of tree-based models and neural network models, based on metric  $|\epsilon_{ATE}|$ . It is worth mentioning that R-Forest model, reported the best performance in terms of efficiency and robustness, illustrating the top curve. More specifically, R-Forest outperforms the rest of neural network-based models and tree-based models, exhibiting 62% of simulations with the lowest  $|\epsilon_{ATE}|$ . Neural network-based models TARNet, Dragonnet and NEDnet presented 15%, 8% and 3% respectively, with the lowest  $|\epsilon_{ATE}|$  score, while tree-based models BART and C-Forest reported poor performance, solving only 8% and 4% of the simulations, respectively.

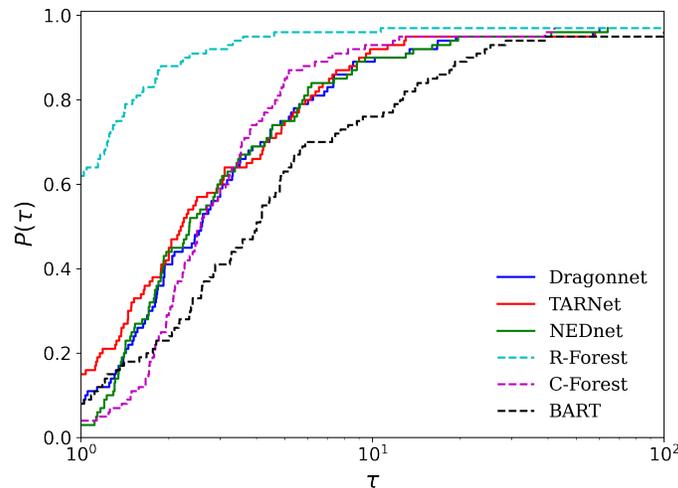
Figure 5: Log<sub>10</sub> scaled performance profiles based on  $|\epsilon_{ATE}|$

Table 5 presents the Friedman average ranking of all evaluated models, which represents the associated effectiveness of each model. Clearly, R-Forest was the best performing model, outperforming all neural network-based models, while BART and C-Forest reported the worst performance. Furthermore, the Friedman statistic  $F_f$  with 5 degrees of freedom is equal to 174.25 while the  $p$ -value is equal to  $8.6 \cdot 10^{-11}$ , which suggests the existence of differences among the evaluated models.

| Algorithm | Ranking |
|-----------|---------|
| R-Forest  | 1.665   |
| TARNet    | 3.13    |
| NEDnet    | 3.49    |
| Dragonnet | 3.685   |
| C-Forest  | 4.015   |
| BART      | 5.015   |

Table 5: Friedman average rankings of evaluated models based on  $|\epsilon_{ATE}|$  for Synthetic Dataset

Table 6 presents the information about the state of rejection of all the hypotheses, comparing the models, by summarizing the APVs with Bergmann Hommel’s procedure  $p_{Berg}$  with  $\alpha = 5\%$  level of significance for the 15 established comparisons.

More specifically, Table 6 reveals that R-Forest outperforms all of the tree-based and neural network-based models, while TARNet outperformed the rest of the neural network models as well as C-Forest and BART. Furthermore, NEDnet and Dragonnet have equally well performances, and they both outperformed C-Forest and BART. Finally, C-Forest and BART reported the worst performances. Figure 6 summarizes and presents in a compact graphical overview the conclusions extracted from Friedman and Bergmann-Hommel’s test.

| Hypothesis            | $p_{Berg}$ |                     |
|-----------------------|------------|---------------------|
| R-Forest vs C-Forest  | 0          | Rejected            |
| R-Forest vs Dragonnet | 0          | Rejected            |
| BART vs TARNet        | 0          | Rejected            |
| R-Forest vs NEDnet    | 0          | Rejected            |
| NEDnet vs BART        | 0          | Rejected            |
| R-Forest vs TARNet    | 0          | Rejected            |
| BART vs Dragonnet     | 0.000002   | Rejected            |
| BART vs C-Forest      | 0.000628   | Rejected            |
| TARNet vs C-Forest    | 0.004937   | Rejected            |
| Dragonnet vs TARNet   | 0.107794   | Fail to be rejected |
| NEDnet vs C-Forest    | 0.141663   | Fail to be rejected |
| NEDnet vs TARNet      | 0.347235   | Fail to be rejected |
| Dragonnet vs C-Forest | 0.347235   | Fail to be rejected |
| NEDnet vs Dragonnet   | 0.461104   | Fail to be rejected |

Table 6: Multiple comparison test: Bergmann-Hommel’s APVs based on  $|\epsilon_{ATE}|$  for Synthetic Dataset

In Figure 7 it is worth mentioning, that in terms of  $\epsilon_{PEHE}$  in Synthetic dataset, C-Forest is solving an extremely high percentage of simulations. More specifically, C-Forest reported 80% of the simulations. Simultaneously, BART exhibited 30% of simulations in the same situation. Additionally, the most noticeable result is that R-Forest and all of neural network-based models, were not able to solve any of the simulations with the best  $\epsilon_{PEHE}$  and performed poorly since their curves lie on the bottom. Finally, the interpretation of Figure reveals that all the neural network-based models were the worst in terms of robustness.

Table 7 presents the Friedman average ranking of all evaluated models, which represent the associated effectiveness of each model. C-Forest was the best performing model followed by BART while the neural network-based models were the worst. In addition, the Friedman statistic  $F_f$  with 5 degrees of freedom is equal to 469.41 while the  $p$ -value is equal to  $1.76 \cdot 10^{-10}$ , which suggests the existence of significant differences between the evaluated models.

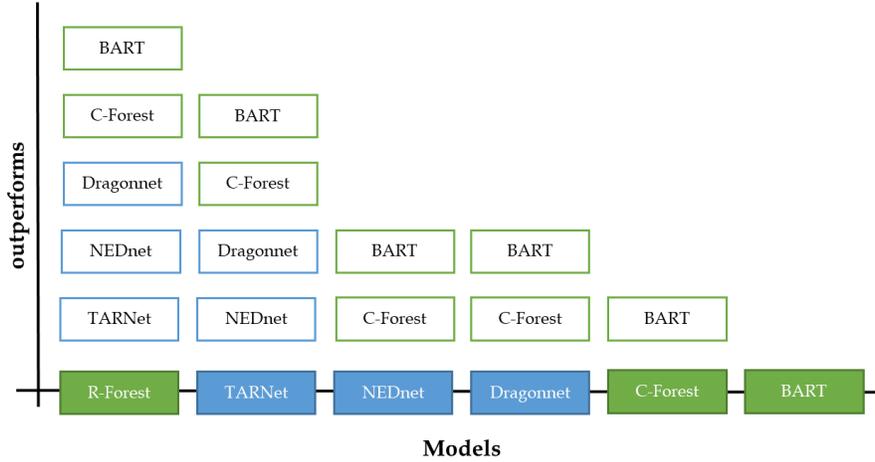


Figure 6: Conclusions for comparison based on  $|\epsilon_{ATE}|$  for Synthetic

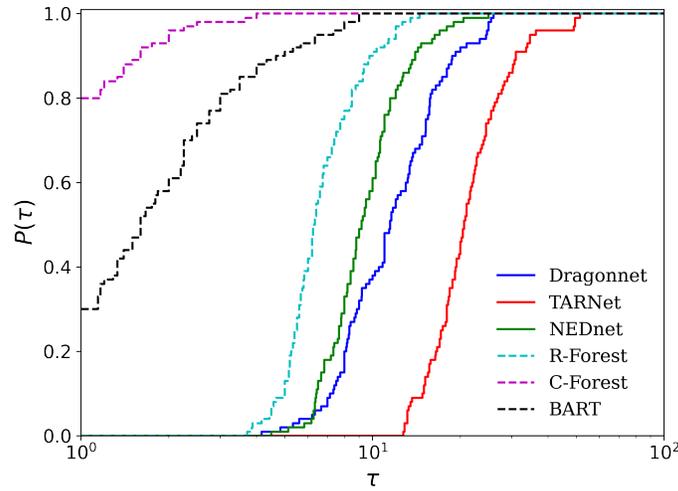


Figure 7: Log<sub>10</sub> scaled performance profiles based on  $\epsilon_{PEHE}$

| Algorithm | Ranking |
|-----------|---------|
| C-Forest  | 1.25    |
| BART      | 1.75    |
| R-Forest  | 3.12    |
| NEDnet    | 4.15    |
| Dragonnet | 4.74    |
| TARNet    | 5.99    |

Table 7: Friedman average rankings of evaluated models based on  $\epsilon_{PEHE}$  for Synthetic Dataset

Table 8 summarizes the state of rejection of all conducted hypotheses, between the evaluated models, by reporting the APVs with Bergmann Hommel’s procedure  $p_{Berg}$  with  $\alpha = 5\%$  level of significance.

In more detail, Table 8 reveals that both C-Forest and BART exhibited similar performances, simultaneously outperformed all neural network models and R-Forest. Additionally, R-Forest outperformed all neural network models, while

NEDnet and Dragonnet outperformed TARNet. The worst performance was held by TARNet. A visual outline of the conducted findings of statistical analysis is reported in Figure 8.

| Hypothesis            | $p_{Berg}$ |                       |
|-----------------------|------------|-----------------------|
| TARNet vs C-Forest    | 0          | Rejected              |
| BART vs TARNet        | 0          | Rejected              |
| Dragonnet vs C-Forest | 0          | Rejected              |
| BART vs Dragonnet     | 0          | Rejected              |
| NEDnet vs C-Forest    | 0          | Rejected              |
| R-Forest vs TARNet    | 0          | Rejected              |
| NEDnet vs BART        | 0          | Rejected              |
| R-Forest vs C-Forest  | 0          | Rejected              |
| NEDnet vs TARNet      | 0          | Rejected              |
| R-Forest vs Dragonnet | 0          | Rejected              |
| R-Forest vs BART      | 0          | Rejected              |
| Dragonnet vs TARNet   | 0.000007   | Rejected              |
| R-Forest vs NEDnet    | 0.000198   | Rejected              |
| NEDnet vs Dragonnet   | 0.051496   | Failed to be rejected |
| BART vs C-Forest      | 0.058782   | Failed to be rejected |

Table 8: Multiple comparison test: Bergmann-Hommel’s APVs based on  $\epsilon_{PEHE}$  for Synthetic Dataset

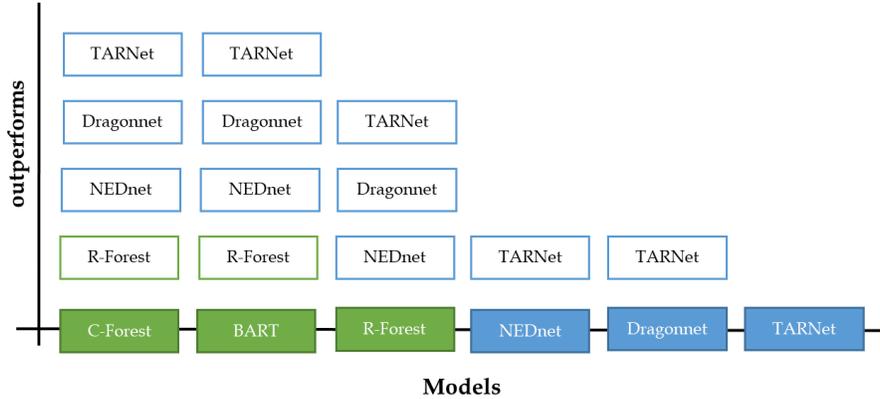


Figure 8: Conclusions for comparison based on  $\epsilon_{PEHE}$  for Synthetic

Summarizing all the above, by taking into consideration both the statistical analysis and performance profiles, it is clear that tree-based model exhibited better performance than neural network models on Synthetic dataset as regards to  $|\epsilon_{ATE}|$  and  $\epsilon_{PEHE}$ . This reveals that tree-based model exhibited the best performance as estimators and predictors.

## 6 Conclusions

In this work, we proposed a complete framework, for the evaluation of causal inference models in order to estimate treatment effects. More specifically, we presented an evaluation framework based on the performance profiles of Dolan and Moré as well as on multiple statistical analysis for providing strong and concrete statistical evidence. In the literature, the traditional approach for the evaluation of these models has been based on the mean values of the  $|\epsilon_{ATE}|$  and  $\epsilon_{PEHE}$ . The major disadvantage of this approach is that all simulations are considered equal which leads to misleading and ambiguous conclusions about the performance of each compared model. This constitutes the main motivation behind our approach focusing on eliminating the influence of a small number of simulations on the benchmarking process, which in some cases dominate the results.

In our experiments, we evaluated the performance of neural network-based and tree-based models using two different datasets (IHDP and Synthetic) with different characteristics. The experimental analysis based on the proposed evaluation

framework revealed that neural-based models presented the best performance on IHDP, while tree-based model exhibited top performance on Synthetic. Therefore, we are able to conclude that no single class of models is the dominant for both datasets. Additionally, based on the performance of each single model, we are able to conduct similar conclusions (i.e. no single model performs well on all problems.). In more detail, our analysis showed that for IHDP dataset, Dragonnet and TARNet presented the best results over all models, both as estimators and as predictors. In case of the Synthetic dataset, R-Forest performed better as estimator in terms of  $|\epsilon_{ATE}|$ , while C-Forest performed better as predictor in terms of  $\epsilon_{PEHE}$ .

It is worth mentioning that based on the traditional evaluation approach (i.e. comparing the mean  $|\epsilon_{ATE}|$  and  $\epsilon_{PEHE}$  over all simulations), we include similar conclusions regarding the models' ranking. However, this approach is misleading in case the variance over all simulations is high and it does not provide us any information whether the performance of two or more models is equivalent neither quantify the difference between their performance. In conclusion, we point out that the proposed evaluation framework is able to provide a deep insight to the performance of causal inference models.

In our future work, we intend to apply the proposed evaluation framework including more real-world and synthetic datasets as well as more causal inference models for the estimation of individual and average treatment effect. Our aim is to identify the models or the class of models, which is "dominated" for a class of problems with specific characteristics. This could provide more useful information about the performance as well as the advantages and disadvantages of each compared model. Finally, it is worth mentioning that our objective and expectation is that this work could be utilized as a reference for evaluating causal inference models, by offering strong statistical evidence for model's evaluation.

## Acknowledgements

The work leading to these results has received funding from the European Union's Horizon 2020 research and innovation programme under Grant Agreement No. 965231, project REBECCA (REsearch on BrEast Cancer induced chronic conditions supported by Causal Analysis of multi-source data)

## References

- Marc Höfler. Causal inference based on counterfactuals. *BMC medical research methodology*, 5(1):1–12, 2005.
- José M Cordero, Víctor Cristóbal, and Daniel Santín. Causal inference on education policies: A survey of empirical studies using pisa, timss and pirls. *Journal of Economic Surveys*, 32(3):878–915, 2018.
- Kevin D Hoover. Economic theory and causal inference. *Philosophy of economics*, 13:89–113, 2012.
- Christos Louizos, Uri Shalit, Joris M Mooij, David Sontag, Richard Zemel, and Max Welling. Causal effect inference with deep latent-variable models. *Advances in neural information processing systems*, 30, 2017.
- Claudia Shi, David Blei, and Victor Veitch. Adapting neural networks for the estimation of treatment effects. In *Advances in neural information processing systems*, volume 32, 2019.
- Uri Shalit, Fredrik D Johansson, and David Sontag. Estimating individual treatment effect: generalization bounds and algorithms. In *International Conference on Machine Learning*, pages 3076–3085. PMLR, 2017.
- Elizabeth D Dolan and Jorge J Moré. Benchmarking optimization software with performance profiles. *Mathematical programming*, 91(2):201–213, 2002.
- Ioannis E Livieris and Panagiotis Pintelas. An improved weight-constrained neural network training algorithm. *Neural Computing and Applications*, 32(9):4177–4185, 2020.
- Joaquín Derrac, Salvador García, Daniel Molina, and Francisco Herrera. A practical tutorial on the use of nonparametric statistical tests as a methodology for comparing evolutionary and swarm intelligence algorithms. *Swarm and Evolutionary Computation*, 1(1):3–18, 2011.
- Cédric Villani. *Optimal transport: old and new*, volume 338. Springer, 2009.
- Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773, 2012.
- Hugh A Chipman, Edward I George, and Robert E McCulloch. Bart: Bayesian additive regression trees. *The Annals of Applied Statistics*, 4(1):266–298, 2010.
- Trevor Hastie and Robert Tibshirani. Bayesian backfitting (with comments and a rejoinder by the authors. *Statistical Science*, 15(3):196–223, 2000.
- Sören R Künzel, Jasjeet S Sekhon, Peter J Bickel, and Bin Yu. Metalearners for estimating heterogeneous treatment effects using machine learning. *Proceedings of the national academy of sciences*, 116(10):4156–4165, 2019.

- Douglas C Montgomery, Elizabeth A Peck, and G Geoffrey Vining. *Introduction to linear regression analysis*. John Wiley & Sons, 2021.
- Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- David W Aha. *Lazy learning*. Springer Science & Business Media, 2013.
- Ahmed Alaa and Mihaela Schaar. Limits of estimating heterogeneous treatment effects: Guidelines for practical algorithm design. In *International Conference on Machine Learning*, pages 129–138. PMLR, 2018.
- Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018.
- Ioannis E Livieris, Andreas Kanavos, Gerasimos Vonitsanos, Niki Kiriakidou, Anastasios Vikatos, Konstantinos Giotopoulos, and Vassilis Tampakas. Performance evaluation of an ssl algorithm for forecasting the dow jones index stocks. In *2018 9th International Conference on Information, Intelligence, Systems and Applications (IISA)*, pages 1–8. IEEE, 2018.
- Ioannis E Livieris. An advanced active set l-bfgs algorithm for training weight-constrained neural networks. *Neural Computing and Applications*, 32(11):6669–6684, 2020.
- Ioannis E Livieris and Panagiotis Pintelas. An adaptive nonmonotone active set–weight constrained–neural network training algorithm. *Neurocomputing*, 360:294–303, 2019.
- Milton Friedman. A comparison of alternative tests of significance for the problem of m rankings. *The Annals of Mathematical Statistics*, 11(1):86–92, 1940.
- Salvador García, Alberto Fernández, Julián Luengo, and Francisco Herrera. Advanced nonparametric tests for multiple comparisons in the design of experiments in computational intelligence and data mining: Experimental analysis of power. *Information sciences*, 180(10):2044–2064, 2010.
- Beate Bergmann and Gerhard Hommel. Improvements of general multiple test procedures for redundant systems of hypotheses. In *Multiple Hypothesenprüfung/Multiple Hypotheses Testing*, pages 100–115. Springer, 1988.
- Salvador Garcia and Francisco Herrera. An extension on “statistical comparisons of classifiers over multiple data sets” for all pairwise comparisons. *Journal of machine learning research*, 9(12), 2008.
- S Paul Wright. Adjusted p-values for simultaneous inference. *Biometrics*, pages 1005–1013, 1992.
- Jennifer L Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20(1):217–240, 2011.
- Vincent Dorie. Npci: Non-parametrics for causal inference. URL: <https://github.com/vdorie/npci>, 2016.
- Antonio Gulli and Sujit Pal. *Deep learning with Keras*. Packt Publishing Ltd, 2017.