

# Deep Propensity Network using a Sparse Autoencoder for Estimation of Treatment Effects

Shantanu Ghosh,<sup>1</sup> Jiang Bian,<sup>2</sup> Yi Guo,<sup>2</sup> Mattia Prosperi<sup>3</sup>

<sup>1</sup>Computer and Information Science and Engineering

<sup>2</sup>Department of Health Outcomes and Biomedical Informatics, College of Medicine

<sup>3</sup>Department of Epidemiology, College of Public Health and Health Professions & College of Medicine  
University of Florida, Florida 32611, USA

{shantanughosh, bianjiang@ufl.edu, yiguo@ufl.edu, m.prosperi}@ufl.edu

## Abstract

Understanding causality is critical, especially in biomedical and social science; and the gold-standard solution is to perform randomized controlled experiments, enacting interventional probabilities as opposed to conditional probabilities. However, randomized experiments are not always feasible due to operational or ethical constraints. On the other hand, drawing causal estimates from pre-existing, observational data is problematic, because datasets are often littered with underlying bias. Identifying the true causal effects however is important to evaluate what-if scenarios, i.e. counterfactuals. In fact, a prediction model based only on conditional probabilities, even with a perfect accuracy, is neither guaranteed to estimate the correct causal effect, in terms of individual or average treatment effects (ITE, ATE), nor to calculate the correct counterfactuals. Propensity score matching (PSM) is a popular statistical approach for observational data that attempts to estimate the effect of an intervention, e.g. a medical treatment, by taking into account other factors that may bias the chance to undertake the intervention itself, e.g. social discrimination. PSM is typically implemented with logistic regression, but its performance can be limited due to linearity, high-dimensionality, and residual confounding in the feature space. Recently, deep counterfactual neural networks with propensity dropout (DCN-PD) have been introduced, enabling nonlinear PSM with advantage over classical methods in terms of treatment assignment error and ATE. In this work, we propose a deep propensity network using a sparse autoencoder (DPN-SA), a novel deep learning architecture for PSM to tackle the problems of high-dimensionality and residual confounding. Tests performed on real-world observational data showed that the DPN-SA outperforms logistic regression-, LASSO-, and DCN-PD in estimating treatment assignment probability, ITE and ATE. The code is available under the MIT license on Github at: <https://github.com/Shantanu48114860/>

## 1 Introduction

In many research fields, such as biomedical and social sciences, preexisting (i.e. observational) data may contain underlying bias, arising in various steps of the data generation or collation process, for which datasets cannot be used

seamlessly to draw causal claims (Prosperi et al. 2020). For instance, one could be interested in studying the effectiveness of certain medical treatments or interventions in the population, but the way in which people access the healthcare system could be different (e.g. due to social inequality or systemic racism). Therefore, due to such bias, the causal effects of the treatments could not be estimated properly. One solution would be to force the interventions to be non-discriminated, performing randomized controlled experiments or trials (RCTs) (Sibbald and Roland 1998) (Concato, Shah, and Horwitz 2000). In a RCT, individuals are assigned to different treatment (or control, no treatment) groups at random, regardless their background characteristics (i.e. the covariate or feature space). Because of the randomization process, it leads to strong ignorability of individuals' pre-treatment characteristics, and thus the causal effect of the treatment versus control can be evaluated (ROSENBAUM and RUBIN 1983a). The mean difference between the observed treatment outcomes of two different groups is called the average treatment effect (ATE). Note that the individual treatment effect (ITE) is a missing data problem (Pearl, Glymour, and Jewell 2016) (Hernan and Robins 2019), because only one factual outcome can be observed (a person cannot be assigned to both treatment and control groups).

Since RCTs are not always feasible due to ethical or operational constraints, e.g. conducting a RCT to ask individuals to smoke and then assess the effect of smoking toward the development of lung cancer, observational data are often used in attempts to draw these causal conclusions. Nevertheless, when using observational data, one must account for the possible types of underlying bias, including confounders and colliders (Greenland and Morgenstern 2001). Propensity score matching (PSM) is a popular statistical approach for observational data that attempts to estimate the causal effect of a treatment variable with respect to an outcome, taking into account possible confounding bias from other pre-treatment characteristics (Rosenbaum and Rubin 1983b) (Rosenbaum and Rubin 1984). The propensity score is a scalar estimate  $\pi(\mathbf{x})$  representing the conditional probability of receiving a certain treatment  $T = 1$ , versus the control group or no treatment  $T = 0$ , given a set of mea-

sured pre-treatment covariates  $\mathbf{X}$ , denoted as

$$\pi(\mathbf{x}) = P(T = 1 | \mathbf{X} = \mathbf{x}), \quad (1)$$

Hence, PSM balances the pre-treatment potential confounders by achieving a quasi-randomization of the different treatment group assignments, allowing a better estimation of the treatment effect. However, traditional PSM approach accounts only for measured (and measurable) covariates, and latent bias may remain after matching (Garrido and others 2014).

PSM has been implemented historically through logistic regression, which calculates the probability of treatment assignment given the pre-treatment covariates (Kurth et al. 2005). In presence of high-dimensional datasets, e.g., compiled from large electronic health record (EHR) databases (Prosperi et al. 2018), different feature selection methods within PSM have been employed, such as the high-dimensional propensity score (hdPS) (Schneeweiss et al. 2009) or L1-logistic regression (Tian, Schuemie, and Suchard 2018). However, logistic regression is limited because it calculates a mere linear combination of input variables, thus it is not able to capture complex relationships between the pre-treatment covariates and the treatment assignment. This is particularly true in high-dimensional settings, where it is difficult to explicitly define variable-to-variable interactions, e.g. as higher-order terms in the logistic function, and computationally burdensome to scan all of them.

An artificial neural network is a universal approximator and can smooth polynomial functions regardless of the order of the polynomial or the number of interaction terms (Barron 1994) (Bishop and others 1995) (Mhaskar 1996). In addition, it does not require *a priori* knowledge of what interactions and functional forms are likely to be relevant among covariates. Therefore, it is suited to overcome the issues in logistic regression-based PSM approach. A neural network can be built to provide the estimation of treatment group assignment probability (e.g. through *softmax*), and consequently of ATE and ITE.

Alaa, Weisz and van Der Schaar designed a multitask deep counterfactual network with propensity dropout (DCN-PD) (Alaa, Weisz, and Van Der Schaar 2017), where a feed-forward network with a set of shared layers is used to calculate the counterfactual outcomes, but being regularized by another network that calculates propensity scores and a dropout probability for training examples (Gal and Ghahramani 2016). Through alternating training phases (treatment vs. control groups), both the shared and the outcome layers' weights are updated. The DCN-PD led to lower mean squared error (MSE) in treatment assignments and more accurate estimates of ATE than other classic PSM methods.

In this work, we propose a novel deep neural architecture –the deep propensity network using a sparse autoencoder (DPN-SA)– that addresses the problem of high-dimensional PSM, yet maintains the MSE advantage of the DCN-PD approach as compared to others. The DPN-SA estimates the propensity score using a sparse autoencoder (Ng and others 2011), which at the same time learns a nonlinear feature representation and reduces the dimensionality of the pre-treatment covariate space. The sparse autoencoder is learnt

in an end-to-end manner, while the counterfactual prediction follows the original DCN-PD network layout.

## 2 Methodology

### 2.1 Problem Formulation

Let assume a population sample (independent and identically distributed) of  $N$  ( $1 \dots n$ ) individuals, given background set of pre-treatment covariates  $\mathbf{X}$ , a treatment  $T$  (binary, for simplicity) and a health outcome  $Y$ . Each subject  $i$  is represented by the tuple  $\{\mathbf{X}_i, T_i, Y_i\}$ . Let  $Y_i^0$  and  $Y_i^1$  the potential outcome for individual  $i$  under treatment  $T_i = 0$  and  $T_i = 1$ , respectively (Rubin 1974) (Rubin 2005). Given  $\mathbf{X}_i = \mathbf{x}$ , the ITE  $\tau(\mathbf{x})$  is defined as the difference in the mean potential outcomes for the individual  $i$  under both treatments, conditional on the observed covariate vector  $\mathbf{x}$

$$\tau(\mathbf{x}) = \mathbb{E}[Y_i^1 - Y_i^0 | \mathbf{X}_i = \mathbf{x}] \quad (2)$$

The ITE formulation as  $\tau(\mathbf{x})$  –called the counterfactual framework– is usually incalculable in reality, since people cannot be assigned two different treatments at the same time. However, under the assumption of strongly ignorable treatment assignment (SITA), the potential outcomes are independent of treatment conditional on background variables, i.e.  $\{Y_i^1 Y_i^0\} \perp T | \mathbf{X}$  (Imbens 2000) (Lechner 2001) (Pearl, Glymour, and Jewell 2016) (Peters, Janzing, and Schölkopf 2017). Under the assumption of SITA, the ITE can then be calculated as  $\tau(\mathbf{x}) = \mathbb{E}[Y^1 | T = 1, \mathbf{X} = \mathbf{x}] - \mathbb{E}[Y^0 | T = 0, \mathbf{X} = \mathbf{x}] = \mathbb{E}[Y | T = 1, \mathbf{X} = \mathbf{x}] - \mathbb{E}[Y | T = 0, \mathbf{X} = \mathbf{x}]$ .

Further, under SITA and by averaging over the distribution of  $\mathbf{X}$ , the ATE  $\tau_{01}$  can be calculated as

$$\tau_{01} = \mathbb{E}[\tau(\mathbf{X})] = \mathbb{E}[Y | T = 1] - \mathbb{E}[Y | T = 0] \quad (3)$$

However, by assuming SITA, ITE and ATE can be calculated only with  $\mathbf{x}$  being the same in treatment groups, which becomes quickly unfeasible when the dimension of  $\mathbf{x}$  grows. PSM, through the conditional probability  $\pi(\mathbf{x})$  (see eq. 1) attempts at balancing the probability of receiving  $T$  given  $\mathbf{X} = \mathbf{x}$ . Once propensity scores are obtained for a population sample, the individuals in the treatment group must be matched those in the control group to make sure that they are balanced with respect to the background covariates; this is a second problem for which a number of (approximate) solutions can be used, including  $k$ -nearest neighbor, Caliper matching (Austin 2014), and propensity weighing (Abadie and Imbens 2016) (on which the DCN-PD was based).

### 2.2 Proposed Approach

The DCN-PD method (Alaa, Weisz, and Van Der Schaar 2017) offered a novel and effective way to calculate nonlinear propensity scores as well as to match individuals, providing better ITE and ATE estimation over classical PSM methods. However, DCN-PD might be affected by curse of dimensionality, with associated residual confounding in study settings where the covariate space has potentially very high cardinality, e.g. in target trial designs that use EHR databases (Zhang et al. 2018) (Li et al. 2020). Feature selection and shrinkage has been proven effective for PSM by logistic

regression (Schneeweiss et al. 2009) (Tian, Schuemie, and Suchard 2018), and in the same way dimensionality reduction can be implemented in deep neural architectures (Hinton and Salakhutdinov 2006). An autoencoder is a neural network that learns to copy its input to its output (encoder-decoder), but the input is coded into a lower dimension within the hidden layer(s) (Hinton and Zemel 1993). In its simplest form, i.e. a single layer with linear/sigmoid activations, the autoencoder is closely related to principal component analysis (PCA), while highly nonlinear codes can be achieved by augmenting the layer architecture, e.g. deep beliefs networks (Zhou et al. 2014). Autoencoders have been employed in several applications, from machine translation to drug discovery (Cho et al. 2014) (Zhavoronkov et al. 2019). The sparse autoencoder is an approach that includes extra units (more than inputs) in the hidden layer, but only a small number of them units are allowed to be activated depending on the input (Ng and others 2011) (Makhzani and Frey 2013). It has been also broadly applied in biomedical studies, including imaging and -omics datasets (Yu et al. 2017) (Mao et al. 2018) (Praveen et al. 2018) (Lemsara, Quadfel, and Fröhlich 2020).

Our DPN-SA exploits a deep stacked sparse autoencoder to encode the covariate space  $\mathbf{X}$  into a lower dimensional, nonlinear feature representation. After training, the decoder is replaced by a *softmax* classifier which calculates the probability of treatment  $T$  assignment, thus estimating the propensity score  $\pi(\mathbf{X})$ . The matching procedure uses the propensity dropout component of the DCN-PD, which is also used for ITE and ATE calculations. Different training (end-to-end vs greedy stacked) procedures and layer (multiple vs single) architectures were tested on real and synthetic high-dimensional data, as detailed in the next sections.

### 2.3 Architecture and Training of DPN-SA

A deep stacked sparse autoencoder consists of multiple encoders and decoders to learn the identity function of the feature vectors, stacked to each other, with the setup of sparsity constraints in the hidden layers, whose neurons can be activated or not depending on the input, as shown in Figure 1.

Let  $\mathbf{X} = \{x(1), x(2), \dots, x(N)\}$  represent the population of  $x(i)$  individuals ( $j = 1 \dots N$  i.i.d samples). Each  $x(i) \in \mathbb{R}^d$  is a  $d$ -dimensional vector, where  $d$  represents the number of pre-treatment covariates, i.e. the number of features. The DPN-SA uses up to  $l$  encoder and  $l$  decoder layers stacked together one after the other, with an additional linear layer at the end of the last decoder. The activation function of each encoder and decoder is the hyperbolic tangent  $\tanh$ . As  $x(i) \in \mathbb{R}^d$ , the reconstructed output also needs to be  $x'(i) \in \mathbb{R}^d$ . Supposing we set  $l = 2$ , the sample input covariate vector  $x(i)$  is reconstructed in the forward propa-

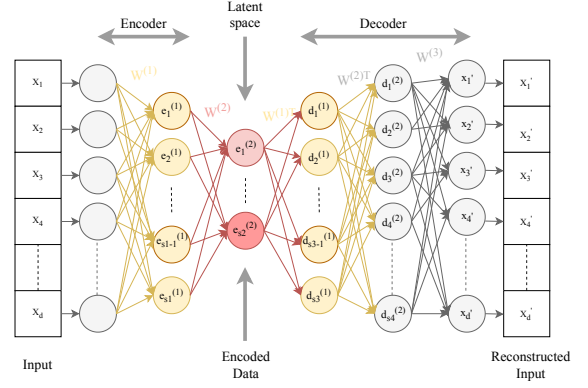


Figure 1: Architecture of a deep stacked sparse autoencoder with the sparsity constraint (brighter neurons are active) enforced in the  $2^{nd}$  layer

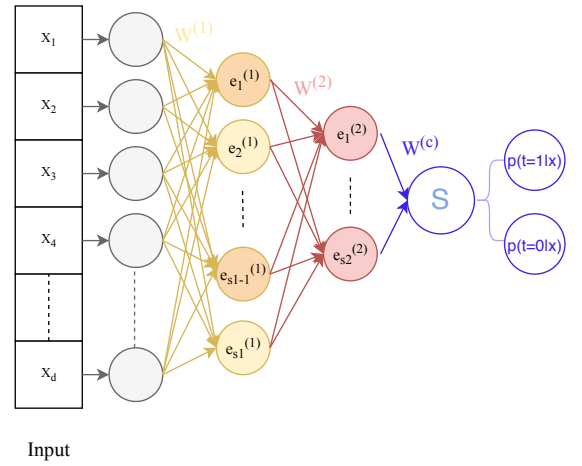


Figure 2: The deep propensity network, stacking the multi-layer sparse autoencoder next to the *softmax* classifier

gation as follows:

$$e^{(1)}(i) = f(W^{(1)}x(i) + b^{(1)}) \quad (4)$$

$$e^{(2)}(i) = f(W^{(2)}e^{(1)}(i) + b^{(2)}) \quad (5)$$

$$d^{(1)}(i) = f(W^{(2)T}e^{(2)}(i) + b^{(3)}) \quad (6)$$

$$d^{(2)}(i) = f(W^{(1)T}d^{(1)}(i) + b^{(4)}) \quad (7)$$

$$x'(i) = W^{(3)}d^{(2)}(i) + b^{(5)} \quad (8)$$

where,  $e^{(l)}(i)$ ,  $d^{(l)}(i)$ ,  $x'(i)$  are the activations of the encoder, decoder and the reconstructed input of  $l^{th}$  layer for the  $i^{th}$  sample, respectively.  $W^{(l)}$  and  $W^{(2)}$  are the weight matrices of encoder 1 and encoder 2, having sizes  $s1 \times d$  and  $s2 \times s1$ ,  $W^{(2)T}$  and  $W^{(1)T}$  are the weight matrices of decoder 1 and decoder 2 having sizes  $s1 \times s2$  and  $d \times s1$ .  $s(l)$  denotes size or number of neurons of  $l^{th}$  layer and  $f(\cdot)$  is the activation function ( $\tanh$ ).

After reconstructing the a sample input feature  $x(i)$ , the objective function of the sparse autoencoder  $J_{sparse}(W, b)$

has to be minimized, as described in (Ng and others 2011):

$$J(W, b) = \left[ \frac{1}{2N} \sum_{i=1}^N \|x'_{W,b}(i) - x(i)\|^2 \right] + \frac{\lambda}{2N} \sum_{l=1}^L \|W^{(l)}\|_F^2 \quad (9)$$

$$J_{sparse}(W, b) = J(W, b) + \beta \sum_{j=1}^{s2} KL(\rho || \rho_j) \quad (10)$$

where  $\lambda$  is the regularization parameter,  $W^{(l)}$  is the weight matrix corresponding to  $l^{th}$  layer of the network,  $s2$  is the number of neurons in the  $2^{nd}$  hidden layer,  $\beta$  is the weight of the sparsity penalty,  $KL$  is the Kullback-Leibler divergence, and the subscript  $F$  is the Frobenius norm (equivalent to the squared norm of the weight matrix). The  $J_{sparse}(W, b)$  is minimized using backpropagation, for  $K_{sa}$  number of epochs.

The DPN-SA is trained in two phases. In the first phase the sparse autoencoder is trained and optimized ( $K_{sa}$  epochs). The parameters  $(W, b)$  are updated by the Adam optimizer in each iteration. After the sparse autoencoder has learnt the latent representation of the covariates, the decoder part is removed and a *softmax* classifier is attached to the end of the last encoder layer, as shown in 2. The *softmax* classifier is trained for  $K_c$  number of epochs. The final network gives the estimation of the propensity score  $\pi(x)$ . Algorithm 1 describes the two-phase procedures to obtain the final DPN-SA.

## 2.4 Experimental Setup

**Datasets.** We used the Infant Health and Development Program (IHDP) dataset, a multi-site, longitudinal, RCT designed to evaluate the efficacy of comprehensive early intervention in enhancing the outcomes of low birth weight, prematurely born infants in the United States. The original IHDP dataset was resampled by throwing away a nonrandom subset of the treatment group (based on the ethnicity variable), thus inducing treatment imbalance, and then counterfactual outcomes were simulated using either a linear or nonlinear/nonparallel surface, thus knowing the true ITE/ATE (Hill 2011). The nonlinear surface dataset, which we chose as benchmark for this study, consisted of 747 subjects (139 in the treatment group and 608 controls), with 25 associated covariates, describing characteristics of the infants and their mothers (excluding the ethnicity).

In addition to IHDP, we also created a synthetic dataset with larger covariate space, first by doubling the original IHDP feature set through the creation of 25 random variables (shuffling the original ones), and then triplicating it with another set of 25 covariates partially correlated to the original ones (apprx.  $\rho=0.4$ ), using a Gaussian noise addition  $\mathcal{N}(0, 2\sigma(x_k))$  to each original variable  $x_k$ . The factual and counterfactual outcomes matched those of IHDP.

**DPN-SA and Other Comparison Methods.** Four different configurations of the DPN-SA were tested –on both the

---

### Algorithm 1 Training DPN-SA

---

**Input:** Batch  $B$  of random samples with assigned treatment  $T$ , training set  $X_{train}$ , number of epochs  $K_s$ ,  $K_c$ , learning rates  $lr_{SA}$ ,  $lr_c$  for the sparse autoencoder (SA) and the *softmax* classifier, respectively.

**Output:**  $PS_{out}$  consisting of propensity score for each sample  $i$ , where  $i$  ranges from 1 to  $N$ .

```

1: procedure DPN-SA TRAINING
2:   Initialise  $W^{(l)}, b^{(l)}$  of the SA.
3:   for  $epochs = 1, 2, \dots, K_s$  do
4:     for  $batches = 1, 2, \dots, B$  in  $X_{batch}$  do
5:       Compute  $x'_{batch}$  using forward propagation.
6:       Compute  $J_{sparse}$ .
7:       Compute Gradient:  $\nabla(J_{sparse})$ 
8:        $(W^{(l)}, b^{(l)}) \leftarrow \text{Adam}(X_{batch}, W^{(l)}, b^{(l)})$ 
9:     end for
10:  end for
11:  Remove the decoder from the SA
12:  Attach a softmax classifier to the last encoder
13:  Initialize  $W^{(c)}, b^{(c)}$  of the classifier
14:   $PS_{out} \leftarrow \text{Empty}$ 
15:  for  $epochs = 1, 2, \dots, K_c$  do
16:    for  $batches = 1, 2, \dots, B$  in  $X_{batch}$  do
17:       $(W^{(l)}, b^{(l)}, W^{(c)}, b^{(c)}) \leftarrow \text{Adam}(X_{batch},$ 
18:         $W^{(l)}, b^{(l)}, W^{(c)}, b^{(c)})$ 
19:      Get the propensity score  $\pi(T|X_{batch})$  and
20:      add to  $PS_{out}$ 
21:    end for
22:  end for
23: end procedure

```

---

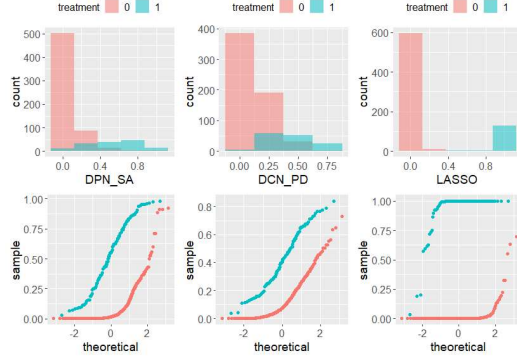


Figure 3: Histograms and quantile-quantile plots of the propensity score distributions (stratified by treatment group) for the DPN-SA, DCN-PD, and LASSO.

real and synthetic datasets: (i) 25-20-10-20-25; (ii) 25-10-25, which is similar to PCA; (iii) 25-1-25, which is similar to regularized logistic regression. For all four, both end-to-end and stacked greedy training were executed. In the greedy layer wise training, we employed two strategies - 1) we trained only the newly added layer while keeping the other layers fixed and stacked them one after the other; 2) we trained the newly added layer along with the previous layers and stacked them one after the other. The learning rates for the sparse autoencoder and for the *softmax* were 0.001 and 0.01, respectively. We set weight decay ( $\lambda$ ), sparsity parameter ( $\rho$ ), and sparse penalty ( $\beta$ ) to 0.0003, 0.8, and 0.1, respectively, with a batch size of 32. The *softmax* classifier was ran for 50 epochs ( $K_c$ ) with a batch size of 32. The DPN-SA was implemented using the *Pytorch* framework (<https://pytorch.org/>).

In addition to the DPN-SA, we ran and compared: (i) the original DCN-PD, (ii) a standard logistic regression, and (iii) logistic regression with LASSO regularization.

**Testing.** The models were trained on 80% of the data and validated on the remaining 20%, repeating the procedure for 100 times, calculating MSE, ITE and ATE in the same way as (Alaa, Weisz, and Van Der Schaar 2017). The error distributions were compared by means of a t-test with sample overlap adjustment (Nadeau and Bengio 2003).

### 3 Results

Figures 3 and 4 show the histogram, quantile-quantile, and scatterplots of the  $\pi(X)$  distribution (stratified by treatment group) of the PDN-SA and the other models (real dataset). The propensity scores among all DPN-SA were well-correlated, similarly to the correlation between DPN-SA and DCN-PD. Logistic regression and LASSO were highly correlated. The propensity scores of the DPN-SA exhibited a more polarized distribution toward the extremes as compared to DCN-PD, but not as marked as logistic regression and LASSO.

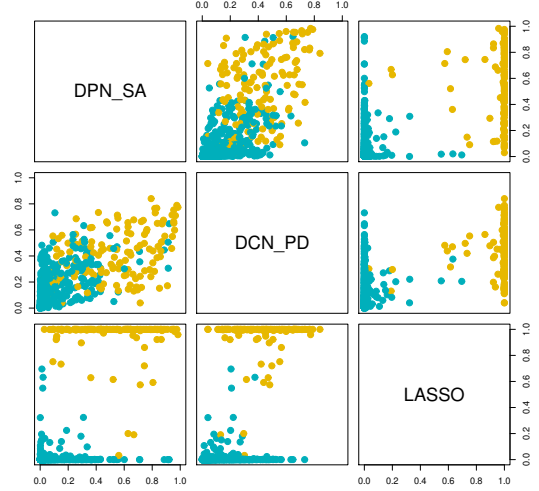


Figure 4: Scatterplots of the propensity score distributions (stratified by treatment group) comparing the DPN-SA, DCN-PD, and LASSO.

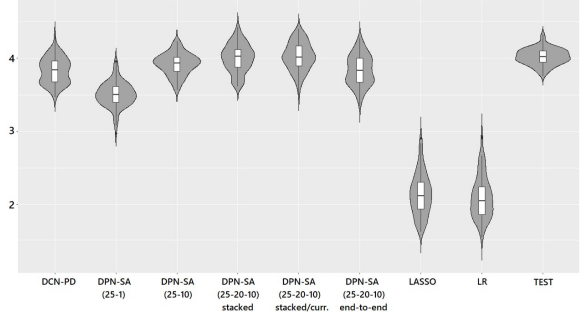


Figure 5: Violin plots of average treatment effect estimation.

Table 1 shows the MSE results for all models (real and synthetic data). On the real dataset, the DPN-SA configured with the 25-20-10-20-25 layer stacking, trained in an end-to-end manner, yielded the best performance in terms of MSE, with an improvement of 6% over the DCN-PD, and 64% over the LASSO logistic regression. The one-neuron DPN-SA 25-1-25 configuration was better than LASSO, but worse than all neural network-based classifiers, while the PCA-like DPN-SA 25-10-25 had performance comparable to all other networks. On the synthetic dataset, performance of all models decreased due to the artificial noise addition, but the regularization/sparsity constraints demonstrated to be robust against such noise and the additional correlated variables. The DPN-SA (25-20-10-20-25) with greedy stacked-current layer training was the best model, but only 1% better than the DCN-PD. Instead, the end-to-end training did not work as well as in the real dataset case. Overall, due to the limited sample size, the null hypothesis of no difference could not be rejected at the 0.05 significance level for DPN-SA

Dataset	# Covariates	Model	MSE (SD)	Bengio's p-value	Raw p-value
Real	25	<b>DPN-SA (25-20-10-) end-to-end</b>	<b>2.09 (0.22)</b>	Ref.	Ref.
		DPN-SA (25-20-10-) greedy st.	2.1 (0.2)	0.396	0.599
		DPN-SA (25-20-10-) greedy st. curr.	2.11 (0.2)	0.391	0.350
		DPN-SA (25-10-)	2.14 (0.22)	0.374	0.078
		DPN-SA (25-1-)	2.52 (0.37)	0.061	1.69e-16
		DCN-PD	2.22 (0.21)	0.295	1.54e-4
		Logistic Regression	6.02 (1.19)	0	1.01e-55
		LASSO Logistic Regression	5.89 (1.22)	0	6.03e-53
Real + Synthetic	75	DPN-SA (25-20-10-) end-to-end	3.14 (0.52)	0.096	1.13e-13
		DPN-SA (25-20-10-) greedy st.	2.7 (0.37)	0.397	0.689
		<b>DPN-SA (25-20-10-) greedy st. curr.</b>	<b>2.69 (0.35)</b>	Ref.	Ref.
		DPN-SA (25-10-)	3.52 (0.53)	0.02	2.74e-22
		DPN-SA (25-1-)	3.59 (0.54)	0.009	8.83e-26
		DCN-PD	2.71 (0.41)	0.396	0.648
		Logistic Regression	7.08 (1.27)	2.69e-10	1.56e-59
		LASSO Logistic Regression	7.35 (1.12)	1.36e-12	3.34e-65

Table 1: Performance of the models against IHDP dataset

architectures and DCN-PD (except for the simple 25-1-25 configuration).

Figure 5 shows the violin plots of each model's ATE estimation (real dataset), compared to the test set ATE. The DPN-SA (25-20-10-20-25) with greedy stacked-current configuration showed the closest resemblance to the test set ATE, followed by the end-to-end and the (25-10-25) configurations that were similar to the DCN-PD. The DPN-SA (25-1-25) significantly underestimated the test set ATE, but was closer to the neural architectures than both logistic regression and LASSO that exhibited much a lower ATE. Note that the IHDP surface is nonlinear by design, so any linear estimator would be biased.

When comparing training times, the fastest methods were logistic regression and LASSO, followed by the DCN-PD. The DPN-SA takes more time to train than the DCN-PD, with the stacked configuration being slower than the end-to-end.

## 4 Discussion

The DPN-SA architecture conjugates the ability to calculate nonlinear propensity scores with dimension reduction, and demonstrates advantage over other methods in treatment effect estimation. In the IHDP counterfactual datasets, the response surface was made nonlinear and nonparallel across treatment groups, thus the true ATE could not be estimated by means of a single linear model.

The DPN-SA architecture allows flexibility in configurations, from the simple 1-neuron akin to LASSO, to the single-layer PCA-like, to the multi-layer setup. All the DPN-SA multi-layer configurations gave ATE estimates close to the true ones, and even the simpler configurations yielded ATE better than the linear models. However, the advantage in MSE of DPN-SA over the DCN-PD, which also employs regularization, is small and would need to be assessed on larger and more diverse datasets. The DPN-SA might be preferable because of its latent space encoding that

can be directly chosen and compared (e.g. linear vs. PCA vs. more complex nonlinear setup).

This work has some limitations. First, the choice of a *softmax* classifier as a replacement to the decoder is relatively simplistic, nonetheless, it provided lower MSE and lower-variance ATE. Other solutions for embedding the sparse autoencoder within the DCN-PD framework or within alternative approaches, such as ITE estimation with generalized adversarial networks (Yoon, Jordon, and van der Schaar 2018), could be devised. Another limitation is that the benchmark dataset has a limited sample size, and therefore the differences in average performance between models are subject to uncertainty. Finally, the distribution of propensity scores among treatment and control groups is often highly dependent on the dataset, and can be highly imbalanced, and therefore the results obtained with one experimental data set are not assured to be reproducible with others.

In conclusion, deep learning frameworks for propensity score estimation and treatment effect prediction are particularly suited for EHR-based studies, not only because such studies can include large sets of covariates, but also because there is high chance of complex heterogeneity in treatment assignments. In these cases, regularized linear propensity score methods, e.g. hdPS (Schneeweiss et al. 2009) or LASSO (Tian, Schuemie, and Suchard 2018), would not be able to provide reliable estimates and yield biased ITE/ATE. The DPN-SA provides a valid, possibly improved, alternative to DCN-PD, and to more traditional PSM methods.

## References

- Abadie, A., and Imbens, G. W. 2016. Matching on the estimated propensity score. *Econometrica* 84(2):781–807.
- Alaa, A. M.; Weisz, M.; and Van Der Schaar, M. 2017. Deep counterfactual networks with propensity-dropout. *arXiv preprint arXiv:1706.05966*.
- Austin, P. C. 2014. A comparison of 12 algorithms for

- matching on the propensity score. *Statistics in Medicine* 33:1057–1069.
- Barron, A. R. 1994. Approximation and estimation bounds for artificial neural networks. *Machine learning* 14(1):115–133.
- Bishop, C. M., et al. 1995. *Neural networks for pattern recognition*. Oxford university press.
- Cho, K.; van Merriënboer, B.; Bahdanau, D.; and Bengio, Y. 2014. On the properties of neural machine translation: Encoder-decoder approaches.
- Concato, J.; Shah, N.; and Horwitz, R. I. 2000. Randomized, controlled trials, observational studies, and the hierarchy of research designs. *New England journal of medicine* 342(25):1887–1892.
- Gal, Y., and Ghahramani, Z. 2016. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, 1050–1059.
- Garrido, M., et al. 2014. Methods for constructing and assessing propensity scores. *Health Services Research* 49(5):1701–20.
- Greenland, S., and Morgenstern, H. 2001. Confounding in health research. *Annual review of public health* 22(1):189–212.
- Hernan, M., and Robins, J. 2019. *Causal Inference*. Chapman & Hall/CRC Monographs on Statistics & Applied Probab. Taylor & Francis.
- Hill, J. L. 2011. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics* 20(1):217–240.
- Hinton, G. E., and Salakhutdinov, R. R. 2006. Reducing the dimensionality of data with neural networks. *science* 313(5786):504–507.
- Hinton, G. E., and Zemel, R. S. 1993. Autoencoders, minimum description length and helmholtz free energy. In *Proceedings of the 6th International Conference on Neural Information Processing Systems*, NIPS’93, 3–10. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.
- Imbens, G. W. 2000. The role of the propensity score in estimating dose-response functions. *Biometrika* 87(3):706–710.
- Kurth, T.; Walker, A. M.; Glynn, R. J.; Chan, K. A.; Gaziano, J. M.; Berger, K.; and Robins, J. M. 2005. Results of Multivariable Logistic Regression, Propensity Matching, Propensity Adjustment, and Propensity-based Weighting under Conditions of Nonuniform Effect. *American Journal of Epidemiology* 163(3):262–270.
- Lechner, M. 2001. Identification and estimation of causal effects of multiple treatments under the conditional independence assumption. In *Econometric evaluation of labour market policies*. Springer. 43–58.
- Lemsara, A.; Ouadfel, S.; and Fröhlich, H. 2020. PathME: pathway based multi-modal sparse autoencoders for clustering of patient-level multi-omics data. *BMC Bioinformatics* 21(1):146.
- Li, Q.; He, Z.; Guo, Y.; Zhang, H.; George, T. J.; Hogan, W.; Charness, N.; and Bian, J. 2020. Assessing the Validity of a a priori Patient-Trial Generalizability Score using Real-world Data from a Large Clinical Data Research Network: A Colorectal Cancer Clinical Trial Case Study. *AMIA ... Annual Symposium proceedings. AMIA Symposium* 2019:1101–1110.
- Makhzani, A., and Frey, B. 2013. k-sparse autoencoders.
- Mao, K.; Tang, R.; Wang, X.; Zhang, W.; and Wu, H. 2018. Feature representation using deep autoencoder for lung nodule image classification. *Complexity* 2018.
- Mhaskar, H. N. 1996. Neural networks for optimal approximation of smooth and analytic functions. *Neural computation* 8(1):164–177.
- Nadeau, C., and Bengio, Y. 2003. Inference for the generalization error. *Mach. Learn.* 52(3):239–281.
- Ng, A., et al. 2011. Sparse autoencoder. *CS294A Lecture notes* 72(2011):1–19.
- Pearl, J.; Glymour, M.; and Jewell, N. 2016. *Causal Inference in Statistics: A Primer*. Wiley.
- Peters, J.; Janzing, D.; and Schölkopf, B. 2017. *Elements of causal inference*. The MIT Press.
- Praveen, G.; Agrawal, A.; Sundaram, P.; and Sardesai, S. 2018. Ischemic stroke lesion segmentation using stacked sparse autoencoder. *Computers in biology and medicine* 99:38–52.
- Prosperi, M.; Min, J. S.; Bian, J.; and Modave, F. 2018. Big data hurdles in precision medicine and precision public health. *BMC Medical Informatics and Decision Making* 18:139.
- Prosperi, M.; Guo, Y.; Sperrin, M.; Koopman, J. S.; Min, J. S.; He, X.; Rich, S.; Wang, M.; Buchan, I. E.; and Bian, J. 2020. Causal inference and counterfactual prediction in machine learning for actionable healthcare. *Nature Machine Intelligence* 2:369–375.
- ROSENBAUM, P. R., and RUBIN, D. B. 1983a. The central role of the propensity score in observational studies for causal effects. *Biometrika* 70(1):41–55.
- Rosenbaum, P. R., and Rubin, D. B. 1983b. The central role of the propensity score in observational studies for causal effects. *Biometrika* 70(1):41–55.
- Rosenbaum, P. R., and Rubin, D. B. 1984. Reducing bias in observational studies using subclassification on the propensity score. *Journal of the American statistical Association* 79(387):516–524.
- Rubin, D. B. 1974. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology* 66(5):688–701.
- Rubin, D. B. 2005. Causal inference using potential outcomes. *Journal of the American Statistical Association* 100(469):322–331.
- Schneeweiss, S.; Rassen, J. A.; Glynn, R. J.; Avorn, J.; Mogun, H.; and Brookhart, M. A. 2009. High-dimensional propensity score adjustment in studies of treatment effects

using health care claims data. *Epidemiology (Cambridge, Mass.)* 20(4):512.

Sibbald, B., and Roland, M. 1998. Understanding controlled trials: Why are randomised controlled trials important? *BMJ* 316(7126):201.

Tian, Y.; Schuemie, M. J.; and Suchard, M. A. 2018. Evaluating large-scale propensity score performance through real-world and synthetic data experiments. *International journal of epidemiology* 47(6):2005–2014.

Yoon, J.; Jordon, J.; and van der Schaar, M. 2018. GAN-ITE: estimation of individualized treatment effects using generative adversarial nets. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.

Yu, J.; Hong, C.; Rui, Y.; and Tao, D. 2017. Multitask autoencoder model for recovering human poses. *IEEE Transactions on Industrial Electronics* 65(6):5060–5068.

Zhang, H.; He, Z.; He, X.; Guo, Y.; Nelson, D. R.; Modave, F.; Wu, Y.; Hogan, W.; Prosperi, M.; and Bian, J. 2018. Computable Eligibility Criteria through Ontology-driven Data Access: A Case Study of Hepatitis C Virus Trials. *AMIA ... Annual Symposium proceedings. AMIA Symposium 2018*:1601–1610.

Zhavoronkov, A.; Ivanenkov, Y. A.; Aliper, A.; Veselov, M. S.; Aladinskiy, V. A.; Aladinskaya, A. V.; Terentiev, V. A.; Polykovskiy, D. A.; Kuznetsov, M. D.; Asadulaev, A.; Volkov, Y.; Zholus, A.; Shayakhmetov, R. R.; Zhebrak, A.; Minaeva, L. I.; Zagribelnyy, B. A.; Lee, L. H.; Soll, R.; Madge, D.; Xing, L.; Guo, T.; and Aspuru-Guzik, A. 2019. Deep learning enables rapid identification of potent DDR1 kinase inhibitors. *Nature Biotechnology* 37(9):1038–1040.

Zhou, Y.; Arpit, D.; Nwogu, I.; and Govindaraju, V. 2014. Is joint training better for deep auto-encoders?