# THE BEST OF BOTH WORLDS: IMPROVED OUTCOME PREDICTION USING CAUSAL STRUCTURE LEARNING

Anonymous authors

Paper under double-blind review

# ABSTRACT

In limited data settings as in the medical domain, causal structure learning can be a powerful tool for understanding the relationships between variables and achieving out of sample generalisation for the prediction of a specific target variable. Most methods that learn causal structure from observational data rely on strong assumptions, such as the absence of unmeasured confounders, that are not valid in real world scenarios. In addition, due to evolving conditions and treatment approaches, causal relationships between the variables change over time. Moreover in a clinical setting, symptoms often need to be managed before finding the root cause of a problem, which puts the emphasis on accurate outcome prediction. Consequently, prediction of a specific target variable from retrospective observational data based on causal relationships alone will not be sufficient for generalisation to prospective data. To overcome these limitations, we opt for *the best of both worlds* in this work by learning a shared representation between causal structure learning and outcome prediction. We provide extensive empirical evidence to show that this would not only facilitate out-of-sample generalisation in outcome prediction but also enhance robust causal discovery for the outcome variable. We also highlight the strengths of our model in terms of time efficiency and interpretability. Code is available at:

- 029 030 031

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028

# 1 INTRODUCTION

Personalised medicine is a branch of medicine, which aims at providing individualised therapy based 032 on patient's phenotype. This is a closed loop process involving analyses of treatment response or 033 outcomes, and treatment adjustment. The treatment response and patient outcomes are influenced by 034 various factors such as disease characteristics, patient traits and the environment (Ambrosone et al., 2006). Randomised controlled trials (RCTs) enable the prospective evaluation of treatment response in randomised groups of patients under controlled conditions. Therefore, RCTs provide a reliable 037 means to assess the cause-effect relationships between treatments and outcomes by eliminating con-038 founding bias (Hariton & Locascio, 2018; Bhide et al., 2018). However, it is not always feasible to conduct RCTs as they can be time-consuming, expensive and suffer from post-randomisation biases (Fernainy et al., 2024). 040

041 In recent times, machine learning-based methods have been successfully used for the prediction of 042 patient outcomes based on observational data (Lee et al., 2021; 2024; Babaei Rikan et al., 2024; Alaa 043 et al., 2017). Most modern machine learning methods find linear or non-linear associations between 044 observational data and outcome. As the associations are learnt on a sample of the data, larger sample sizes increase the generalisability of the associations to unseen samples of the data (Chekroud et al., 2024). Ideally, these methods are evaluated using observed outcomes or expert annotations, which 046 are both time expensive. Consequently, data in the medical domain is limited, unstructured and 047 incomplete. This in turn makes generalisation to out-of-sample data more difficult for machine 048 learning based outcome prediction methods (Goetz et al., 2024). The lack of transparency in some 049 of these machine learning methods makes them difficult to interpret, compounding the challenges. 050

Causal structure learning is concerned with learning causal relationships from observational data.
 Popular techniques employ machine learning methods to model the causal relationships between
 the variables of observational data by imposing certain topological constraints (Zheng et al., 2018;
 Yu et al., 2019; Ng et al., 2019). Causal structure learning methods have the potential to improve

interpretability in the medical domain by finding causal relationships between observed variables
and the outcome for various downstream analyses (Feuerriegel et al., 2024; Piccininni et al., 2020).
Consequently, they can bridge the gap between observational studies and RCTs. However, most of
these methods make strong assumptions about the data, which might not be valid in a real world
setting (Montagna et al., 2024).

One such assumption is the absence of unmeasured confounders (Kalisch & Bühlman, 2007; Shimizu et al., 2006). This is not realisable without domain knowledge or time-expensive expert intervention (Bica et al., 2021). Moreover, patient outcomes are also influenced by evolving knowledge, treatment approaches and the environment (Futoma et al., 2020; Petzschner, 2024). Consequently, relying solely on causal relationships to predict outcomes presents a significant challenge for generalizing methods to prospective data, as these associations are derived from retrospectively observed data. This limitation underscores the difficulty of ensuring that findings translate effectively to future scenarios.

We overcome these limitations in this work by opting for *the best of both worlds* — causal structure
 learning and machine learning-based outcome prediction. Our contributions are as follows.

• We designed our approach to learn outcome prediction and causal structure simultaneously. In our framework, causal structure learning functions as an auxiliary task to support outcome prediction, sharing representations of the input through the hidden layers of our network architecture and employing task-specific heads for refined predictions.

We provide empirical evidence to show that this learning strategy enables (i) interpretability by visualisation of the learnt causal graph (ii) out-of-sample generalisation for outcome prediction. The primary focus of our work is to improve generalisation for outcome prediction in the medical domain. Causal structure learning functions as an auxiliary task to support outcome prediction. Despite this, we also provide evidence demonstrating the benefits of our approach in robust causal discovery for the outcome variable.

We provide a case study by applying the method to survival analysis. We show that the proposed framework improves interpretability of the model and generalisability to unseen data in real world scenarios. We also comment on the clinical relevance of the results.

083 084

# 2 RELATED WORK

085

Most causal structure learning methods have been developed to learn causal relationships from ob-087 servational data based on the foundations of causal graphical model (Pearl, 2009). These methods 088 can be classified into three broad categories: (i) constraint based methods that use conditional inde-089 pendence tests to infer the direction of causal relationship between variables (Kalisch & Bühlman, 090 2007), (ii) methods that use functional causal models to identify the causal structure by making as-091 sumptions about the data distribution (Shimizu et al., 2006), (iii) score-based methods which either 092 adopt greedy search algorithms to determine the causal structure (Chickering, 2002) or impose topo-093 logical constraints to learn the causal structure (Zheng et al., 2018). Most of these methods make strong assumptions about the data (Montagna et al., 2024) which might not be realisable in a real 094 world setting. 095

096 Recent works (Kyono et al., 2020; Ge et al., 2023) use causal structure learning to improve general-097 isation in supervised learning. Ge et al. (2023) build upon (Zheng et al., 2018) to learn robust causal 098 structures that are invariant to the data environments by getting rid of spurious correlations arising from the data environment. This is contradictory to our aim of banking on the rich information from 099 evolving conditions to predict the target. Kyono et al. (2020) introduce a causal structure learning 100 based regularizer, CASTLE, for improving generalisation in supervised learning. They add a super-101 vised loss term to the non-linear framework from Zheng et al. (2018) to learn the target variable. 102 CASTLE (Kyono et al., 2020) is one of the revolutionary works which demonstrated the superior 103 performance of causality based regularaisation over commonly used regularisation techniques for 104 deep learning such as L1-norm, L2-norm, dropout and early stopping (Tibshirani, 1996; Hoerl & 105 Kennard, 1970; Goodfellow et al., 2016). 106

107 However, we observe several unsolved challenges of this work: (i) the feed-forward architecture used by CASTLE does not scale with the feature variables, (ii) CASTLE treats the target variable

108 reconstructed as a part of causal structure learning as the final output which hinders not only causal structure learning but also outcome prediction. We address these research gaps by (i): adopting a 110 graph autoencoder-based causal structure learning method (Ng et al., 2019), which builds a single 111 graph for all the variables and scales with the number of features, (ii) we introduce an additional task-112 specific head for outcome prediction, which exploits the representation shared with causal structure learning to reliably predict the outcome and generalise well to unseen data. 113

#### 3 CAUSAL STRUCTURE LEARNING

117 Given observational data  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , consisting of n i.i.d. samples of the random vector X =118  $(X_1, X_2, \dots, X_d)$ , score based methods learn an optimal causal directed acyclic graph (DAG),  $\mathcal{G}(\mathbf{W})$ 119 on d nodes from a discrete space of DAGs  $\mathbb{D}$  for the joint distribution  $\mathbb{P}(X)$  (Spirtes et al., 2001). 120 Here, X is modelled by considering the data generating process in a linear structural equation model (SEM) defined by the weighted adjacency matrix  $\mathbf{W} \in \mathbb{R}^{d \times d}$  as in (Hoyer et al., 2008), 121

$$X_j := W_j^T \mathbf{X} + Z_j$$

for j = 1, 2, ..., d;  $Z = (Z_1, Z_2, ..., Z_d)$  is a random noise vector. Zheng et al. (2018) impose smooth acylicity constraint on W and convert the combinatorial optimisation problem of finding  $\mathcal{G}(\mathbf{W}) \in \mathbb{D}$  to a continuous one:

$$\min_{\mathbf{W}} \frac{1}{2n} \sum_{i=1}^{n} \left\| X^{(i)} - \mathbf{W}^{T} X^{(i)} \right\|_{F}^{2} + \lambda \left\| \mathbf{W} \right\|_{1}$$
(1)

subject to  $\operatorname{tr}(e^{\mathbf{W} \odot \mathbf{W}}) - d = 0$ ,

where  $e^{\mathbf{M}}$  denotes the matrix exponential of  $\mathbf{M}$ ,  $\odot$  denotes the Hadamard product and n is the sample size. The L1 regularization term  $\|\mathbf{W}\|_1$  encourages sparsity in the learnt DAGs.

#### THE BEST OF BOTH WORLDS - PARADIGM 4

(Ng et al., 2019) generalises the formulation in (1) to the non-linear case and draws parallels to the graph autoencoder (GAE) framework (Cen et al., 2019). For the linear case, we can rewrite  $\mathbf{W}^T X^{(i)}$ in (1) as  $\mathbf{W}^T X^{(i)} = f(X^{(i)}, \mathbf{W})$ , where f is the data generating model with parameters  $\Theta$ . Ng et al. 142 (2019) extends this to the non-linear case by considering: 143

$$f(X^{(i)}, \mathbf{W}) = g_2(\mathbf{W}^T g_1(X^{(i)})),$$
(2)

where each variable  $X^{(i)}$  is vector valued, i.e.,  $X^{(i)} \in \mathbb{R}^l$ ;  $g_1 : \mathbb{R}^l \to \mathbb{R}^l$  and  $g_2 : \mathbb{R}^l \to \mathbb{R}^l$ 146 are Multilayer Perceptrons (MLPs) with shared weights across all variables  $X_i$ . The formulation 147 in (2) is considered similar to the GAE framework, if we view  $g_1$  and  $g_2$  as variable-wise encoder 148 and decoder modules and  $\mathbf{W}^T g_1(X^{(i)})$  as a linear transformation of the latent representation. The 149 dimension of the latent representation can be adjusted based on the intrinsic dimension of  $\mathbf{X}$ . We 150 refer to this framework of Ng et al. (2019) as CausalGAE framework. Let  $\hat{X}^{(i)} = q_2(\mathbf{W}^T q_1(X^{(i)}))$ 151 be the reconstructed output and  $\Theta_1$ ,  $\Theta_2$  be the parameters of  $g_1$  and  $g_2$  respectively. Then, the 152 framework optimises the following objective function: 153

154

114 115

116

122

123

124

125

132 133 134

135

136 137

138 139

140

141

144 145

$$\min_{\mathbf{W},\Theta_{1},\Theta_{2}} \frac{1}{2n} \sum_{i=1}^{n} \left\| X^{(i)} - \hat{X}^{(i)} \right\|_{F}^{2} + \lambda \left\| \mathbf{W} \right\|_{1}$$
(3)

subject to  $\operatorname{tr}(e^{\mathbf{W} \odot \mathbf{W}}) - d = 0$ ,

We build upon this work to derive a formulation for simultaneous causal structure learning and 161 outcome prediction. We consider outcome prediction as a supervised learning task that is concerned with predicting Y from  $\tilde{\mathbf{X}} := (X_1, X_2, ..., X_{d-1}) \in \mathbb{R}^{n \times d}$  variables. We consider  $\mathbf{X}_{(d-1)}$  to be the outcome variable Y. The formulation in (2) would then restrict approximation of Y to a nonlinear function of its causal parents. Outcome prediction is a complex task that is dependent on dynamically changing variables and environment. This will lead to relationships in the data that are not explained by the causal structure alone but are necessary to predict the outcome. Therefore, we hypothesize that a non-linear function of the target variable's causal parents alone is not sufficient to approximate Y and propose the following:

$$\hat{Y}^{(i)} = g_3(\mathbf{W}^T g_1(X^{(i)})), \tag{4}$$

171 172 where  $g_3$  is a variable-wise non-linear function with parameters  $\Theta_3$ . In the case of classification,  $g_3$ 173 is a projection layer  $g_3 : \mathbb{R}^d \to \mathbb{R}^c$ , where c is the number of classes. For simplicity, we consider 174 variable  $X^{(i)} \in \mathbb{R}^l$  to be scalar valued, i.e., l = 1. To summarise, we use  $g_1$ , a variable-wise encoder 175 to learn a latent representation of the data, H. Next, we perform a linear transformation of H using 176 the weighted adjacency matrix  $\mathbf{W}$  to produce  $\hat{H}$ . We then feed  $\hat{H}$  to task specific variable-wise 177 decoders  $g_2$  and  $g_3$  to provide reconstructed output  $\hat{X}$  and  $\hat{Y}$  respectively. The same is explained in 178 the following:

180	$H^{(i)} = g_1(X^{(i)})$
181	$\hat{H}^{(i)} = \mathbf{W}^T \boldsymbol{a}_i (\mathbf{Y}^{(i)})$
182	$\widehat{f}_{1}(i) = \bigvee \widehat{f}_{1}(X^{(i)})$
183	$X^{(i)} = g_2(\mathbf{W}^I g_1(X^{(i)}))$
184	$\hat{Y}^{(i)} = g_3(\mathbf{W}^T g_1(X^{(i)}))$

We learn the parameters of the shared encoder and target specific decoders jointly by optimising the following objective function:

$$\min_{\mathbf{W},\Theta_{1},\Theta_{2},\Theta_{3}} \frac{(1-\kappa)}{2n} \sum_{i=1}^{n} \left\| X^{(i)} - \hat{X}^{(i)} \right\|_{F}^{2} + \lambda \left\| \mathbf{W} \right\|_{1} + \frac{\kappa}{n} \sum_{i=1}^{n} \left\| Y^{(i)} - \hat{Y}^{(i)} \right\|_{F}^{2}$$
(5)

subject to 
$$\operatorname{tr}(e^{\mathbf{W} \odot \mathbf{W}}) - d = 0$$
,

where  $\kappa$  is a hyperparameter that can be tuned depending on the dataset and the outcome prediction task. A sensitivity analysis for the same can be found in Appendix C. We use Augmented Lagrangian method to optimise the constrained optimisation problem in (5) (Appendix A). The simplified form of the loss is as follows:

$$\mathcal{L}_{\rho}(\mathbf{W},\Theta_{1},\Theta_{2},\Theta_{3},\alpha) = (1-\kappa)\mathcal{L}_{DAG}(\mathbf{W},\Theta_{1},\Theta_{2}) + \kappa\mathcal{L}_{sup}(\mathbf{W},\Theta_{1},\Theta_{3}), \tag{6}$$

where  $\alpha$  is the Lagrange multiplier. For classification, we use cross entropy loss as the supervised loss.

# 5 Results

169 170

179

186

187

195 196

197

199

200 201 202

203

204 205

206 207

We present empirical evidence of our model's generalization performance through a series of experiments on both synthetic and real-world datasets, as detailed below. Furthermore, we demonstrate the model's ability to learn robust representations while enhancing interpretability and scalability. Through a case study, we also show the clinical relevance of our model.

Experimental setup. We perform the experiments by splitting the data into 90% training and 10% test datasets. The training dataset is used in a 10-fold cross validation setting to train the models.
We choose CASTLE network (Kyono et al., 2020), which has outperformed various regularisation methods like like dropout, data augmentation and batch normalisation, as our primary baseline. We also compare our method with Multilayer Perceptron (MLP) (Pedregosa et al., 2011) and its

regularised variants by employing L2-norm with early stopping (L2+ES) based on training loss, and
 L2-norm with early stopping based on validation score (ES).

All the methods used Adam optimiser with a learning rate of 1e-3. The models were trained for 219 300 epochs with an early stopping criterion based on validation loss (except the L2-norm with early 220 stopping (ES) variant of the MLP which employed early stopping based on validation score). In 221 addition, our model and CausalGAE update Lagrange multiplier  $\alpha$  and penalty  $\rho$  over 20 iterations 222 with early stopping based on a threshold. For both models, we use the default threshold and loss 223 hyperparameters as in (Ng et al., 2019). The loss hyperparameter  $\kappa$  is set to 0.25 for our model upon 224 doing a sensitivity analysis (Appendix C). All methods used the same data splits and identical seeds. 225 For all our experiments we use a machine equipped with Intel i9-10900X processor and NVIDIA 226 RTX2080 GPUs. Our code will be publicly available at:

227 228

229

235

236 237 238

# 5.1 GENERALISATION PERFORMANCE ON SYNTHETIC DATA

We study the generalisation performance of the models on synthetic data for the regression task. We generate synthetic data as in (Ng et al., 2019). We generate a random DAG having n = 1000samples, d = 20 nodes and degree of freedom dof = 3 from Erdős–Rényi graph. The variables or features **X** are sampled from the Additive noise model (ANM) under two conditions described as follows.

**Case 1.** Non-linear causal relationships between the variables as:

$$\mathbf{X} = 2\sin(\mathbf{W}^T(\mathbf{X} + \mathbf{1}) + 0.5 \cdot \mathbf{1}) + \cos(\mathbf{W}^T(\mathbf{X} + \mathbf{1}) + 0.5 \cdot \mathbf{1}) + \mathbf{Z}$$

**Case 2.** We consider  $X_{(d-1)}$  to be the outcome variable Y. In addition to the causal relationships described in Case 1, we simulate the case where the outcome is not dependent only on its causal parents by adding an extra term to Y as:

$$Y = 2\sin(\mathbf{W}^{T}(\mathbf{X} + \mathbf{1}) + 0.5 \cdot \mathbf{1}) + \cos(\mathbf{W}^{T}(\mathbf{X} + \mathbf{1}) + 0.5 \cdot \mathbf{1}) + \mathbf{Z} + \cos((X^{non - pa(Y)_{0}} + \mathbf{1}) + 0.5 \cdot \mathbf{1}),$$

244 245 246

247

242 243

where **W** is the weighted Adjacency matrix of the random DAG, **Z** is additive noise and  $X^{non-pa(Y)_0}$  is the first non-parent of Y.

Mean squared error (MSE) is chosen as the metric to compare the performance of the models. Table shows the results on test dataset. For this experiment, we also compare our model with CausalGAE (Ng et al., 2019), the causal structure learning framework upon which our model is built. We infer the reconstructed target from the trained model and report the MSE for the test data in Table 1. Our model performs the best in both the cases illustrating its ability to learn not only from causal parents but also from other predictors of the target variable.

## 255 256 5.2 ROBUST CAUSAL DISCOVERY

257 The primary focus of our work is to improve generalisation for outcome prediction in the medical 258 domain. Causal structure learning functions as an auxiliary task to support outcome prediction. Here 259 we demonstrate the performance of our model in recovering true causal graph. We use the synthetic 260 datasets from Case 1 and Case 2, and compare the performance of our model with CausalGAE. We 261 focus our comparison on CausalGAE for two main reasons: (i) our model is built on this framework, and (ii) CausalGAE has demonstrated superior performance compared to widely-used linear and 262 graph neural network (GNN)-based causal structure learning methods, such as NOTEARS (Zheng 263 et al., 2018) and DAG-GNN (Yu et al., 2019), in synthetic data experiments, particularly in Case 1. 264

We report the false discovery rate (FDR), true positive rate (TPR), false positive rate (FPR) and
structural Hamming distance (SHD). As seen in Table 2, our model recovers the true causal graph
with SHD comparable to that of CausalGAE. In both cases, our model improves upon TPR. In
contrast to our model, CausalGAE fails to identify the causal parents of Y in both the cases as seen
in the associated causal graphs included in Appendix D. These results highlight the versatility of the
shared representations learned by our model.

Table 1: Generalisation performance on synthetic datasets. Baseline: MLP; L2+ES: MLP + L2norm with early stopping based on training loss; ES: MLP + L2 norm + early stopping based on validation MSE; Test MSE along with the gap between train and test MSE ( $\Delta :=$  Mean MSE<sub>test</sub> -Mean MSE<sub>train</sub>) are reported.

	Case 1	Case 2
Baseline L2+ES ES CASTLE CausalGAE Ours	$\begin{array}{c} 1.141 \pm 0.042 \ (0.854) \\ 1.091 \pm 0.029 \ (0.734) \\ 0.974 \pm 0.018 \ (0.233) \\ 1.073 \pm 0.108 \ (0.613) \\ 1.172 \pm 0.000 \ (0.164) \\ \textbf{0.938} \pm \textbf{0.029} \ \textbf{(0.086)} \end{array}$	$\begin{array}{c} 0.991 \pm 0.043 \; (0.704) \\ 0.882 \pm 0.037 \; (0.350) \\ 0.857 \pm 0.022 \; (0.145) \\ 0.923 \pm 0.089 \; (0.466) \\ 1.027 \pm 0.000 \; (0.148) \\ \textbf{0.815} \pm \textbf{0.029} \; \textbf{(0.002)} \end{array}$

Table 2: Performance of the models in recovering true causal graphs. FDR: false discovery rate; TPR: true positive rate; FPR: false positive rate and SHD: structural Hamming distance.

	FDR	TPR	FPR	SHD
Case 1: CausalGAE	0 59	0.27	0.06	26
Case 2: CausalGAE	0.40	0.27	0.00	23
Case 1: Ours	0.59	0.62	0.14	29
Case 2: Ours	0.59	0.58	0.13	30

# 5.3 ABLATION STUDIES

To further highlight the strength of the shared representations, we present the results of our model
 with ablations of the causal structure learning and outcome prediction components. We use the
 synthetic data described in Section 5.1 and report the MSE for regression on the outcome variable Y.
 As shown in Table 3, our model achieves the lowest MSE and demonstrates superior generalization
 performance in both cases.

# **5.4 S**CALABILITY ANALYSIS

We compare the time complexity of our model with CASTLE. We use the synthetic dataset from Case 2 for the analysis. We measure the average training time across the 10 folds of both models and plot it against the number of variables *d* (Figure 1). We also report the corresponding average test MSEs of the models in Table 4. Our model efficiently scales with the number of feature variables while consistently maintaining a stable MSE score. CASTLE uses one feed forward network for each feature variable and does not scale with the number of feature variables.

# 5.5 GENERALISATION PERFORMANCE ON REAL DATA

We also study the generalisation performance of the models on publicly available datasets from The
UCI Machine Learning Repository (Markelle Kelly). We choose two binary classification datasets
Statlog Heart and Breast cancer (Wisconsin Diagnostic), and one multi-class classification dataset
Las Vegas ratings. We choose CASTLE as our primary baseline here because CASTLE also uses
causal structure learning for outcome prediction and has outperformed state-of-the-art regularisation
methods like dropout, data augmentation and batch normalisation for Statlog heart and Las Vegas ratings.

Table 3: Ablation results for causal structure learning and outcome prediction components using synthetic data. Test MSE for regression task along with the gap between train and test MSE ( $\Delta :=$ Mean  $MSE_{test}$  - Mean  $MSE_{train}$ ) are reported.

	Case 1	Case 2
causal structure learning alone	$1.172 \pm 0.000 \ (0.164)$	$1.027 \pm 0.000$
outcome prediction alone	$0.951 \pm 0.041$ (0.188)	$0.848 \pm 0.038$
Ours	$0.938 \pm 0.029 (0.086)$	$\textbf{0.815} \pm \textbf{0.029}$



Figure 1: Comparison of average training time against the number of feature variables d.

Table 5 shows the performance of the models for the classification tasks. Area Under receiver operating characteristic Curve (AUC) is chosen as the evaluation metric. We see that on fairly simple binary classification datasets, all the models perform similarly and reach an AUC greater than 0.9. However, on the relatively difficult multi-class classification task, which involves the classification of the samples into 5 classes, MLP-based models and CASTLE fail to perform well on the test data. Our model performs consistently well on all the datasets. The ES variant performs early stopping based on validation accuracy in case of classification. Therefore, the performance of ES is worse than L2+ES in the multi-class classification task where accuracy might not be a robust metric. The recovered causal graphs from our model and dataset details are included in Appendix E. 

## 5.6 CASE STUDY: APPLICATION TO SURVIVAL ANALYSIS

We choose the Worcester heart attack study dataset (Hosmer Jr et al., 2008) for our case study. The dataset was originally designed to study the trends in incidence rates and patient outcomes across multiple decades (Floyd et al., 2009). We use a publicly available subset of this dataset <sup>1</sup>, which contains 500 samples collected across three years (1997, 1999, 2001). We choose death until the length of hospital stay as our endpoint. We convert the time-to-event analysis problem to a classification problem by predicting the likelihood of death before a specific length of hospital stay = k (see Appendix F). We study two scenarios as follows. 

Scenario 1. In this scenario, we randomly split the datasets into into 90% training and 10% test data stratified according to the labels. The training data is used in a 10-fold cross validation setting to train the models. The threshold k for the endpoint is the median length of hospital stay. 

<sup>1</sup>https://web.archive.org/web/20170517071528/http://www.umass.edu/ statdata/statdata/data/whas500.txt

	CASTLE	Ours
d = 10	$0.667\pm0.027$	$0.626 \pm 0.017$
d = 20	$0.923\pm0.089$	$0.815\pm0.029$
d = 30	$1.104\pm0.100$	$0.830\pm0.023$
d = 40	$1.305\pm0.159$	$0.874\pm0.026$
d = 50	$1.018\pm0.087$	$0.693 \pm 0.022$

Table 4: Test MSE for varying number of feature variables *d*.

Table 5: Classification: Generalisation performance on real-world datasets. Baseline: MLP; L2+ES: MLP + L2-norm with early stopping based on training loss; ES: MLP + L2 norm + early stopping based on validation accuracy; Test AUC along with the gap between train and test AUC ( $\Delta :=$  Mean AUC<sub>train</sub> - Mean AUC<sub>test</sub>) are reported.

	Statlog Heart	<b>Breast cancer</b>	Las Vegas ratings
Baseline	$0.935 \pm 0.012$ (0.064)	0.991 ± 0.001 (0.008)	$0.574 \pm 0.057$ (0.380)
L2+ES	$0.936 \pm 0.011 (0.062)$	$0.996 \pm 0.000 (0.002)$	$0.571 \pm 0.056 (0.383)$
ES	$0.954 \pm 0.010 (0.067)$	$0.995 \pm 0.007$ (-0.006)	$0.336 \pm 0.044 \ (0.253)$
CASTLE	$0.928 \pm 0.020$ (0.054)	$0.997 \pm 0.003$ (-0.003)	$0.553 \pm 0.067$ (-0.042)
Ours	$0.931 \pm 0.017 (0.068)$	$0.996 \pm 0.003 \; (0.001)$	$0.658 \pm 0.049 \; (0.239)$

**Scenario 2.** Here we simulate the scenario in the real world setting by choosing the patients studied during the years 1997 and 1999 as our training data, and the patients studied during 2001 as our test data. The training data is used in a 10-fold cross validation setting to train the models. The threshold k for the endpoint is the median length of hospital stay over the training data.

Table 6 presents the results of all the methods in both the scenarios. CASTLE, Baseline and ES fail to generalise to test data in both scenarios. L2+ES performs well when trained on data from scenario 1 but fails to generalise to test data in scenario 2. Our method performs the best in both scenarios confirming our hypotheses: (i) our method performs well in predicting outcome over evolving knowledge and treatment approaches; (ii) our method performs better than the baseline models in both scenarios emphasising the importance of proposed approach for generalisation;

Figure 2 illustrates the causal graph recovered from our models. The graph in Scenario 1 shows an association between death during hospital stay and the factors such as age, body mass index (bmi), atrial fibrillation (afb), cardiogenic shock (sho) and complete heart blockage (av3). In addition to age, bmi, sho and av3, the graph in Scenario 2 shows an association with initial heart rate (hr), order of myocardial infarction (miord) and congestive heart complications (chf). The robust performance of our model in both scenarios highlights the ability of our model to adjust to evolving real world scenarios.

424 425

378

379380381382

391

392

393

394

396 397

6 CONCLUSION

426 427

We introduce a novel paradigm that leverages the best of both worlds - causal structure learning and
 machine learning-based outcome prediction for improved outcome prediction. Through experiments
 on several synthetic and real data, we demonstrated that this strategy not only improves out-of sample generalisation for outcome prediction but also improves interpretability by learning causal
 structure. We also demonstrated advantages of our model with respect to robust causal discovery for

Table 6: Classification performance on Worcester heart attack study dataset. Baseline: MLP; L2+ES: MLP + L2-norm with early stopping based on training loss; ES: MLP + L2 norm + early stopping based on validation accuracy; Test AUC along with the gap between Train and test AUC  $(\Delta := \text{Mean AUC}_{train} - \text{Mean AUC}_{test})$  are reported. 

	Scenario 1	Scenario 2
Baseline L2+ES ES	$\begin{array}{c} 0.595 \pm 0.049 \ (0.394) \\ 0.652 \pm 0.037 \ (0.229) \\ 0.499 \pm 0.020 \ (-0.074) \\ 0.567 \pm 0.221 \ (0.095) \end{array}$	$\begin{array}{c} 0.582 \pm 0.022 \; (0.417) \\ 0.595 \pm 0.015 \; (0.351) \\ 0.426 \pm 0.009 \; (0.012) \\ 0.595 \pm 0.122 \; (0.240) \\ 0.595 \pm 0.1$
Ours	$0.567 \pm 0.221 (0.085)$ $0.834 \pm 0.127 (0.081)$	$0.588 \pm 0.122 (0.240) \\ 0.693 \pm 0.120 (0.237)$





Figure 2: Causal graph recovered from our models in both scenarios of the experiments on Worcester heart attack study dataset.

the outcome variable and time efficiency. With a case study on survival analysis, we demonstrated the potential translational value that our model provides.

# REFERENCES

- Ahmed M Alaa, Jinsung Yoon, Scott Hu, and Mihaela Van der Schaar. Personalized risk scoring for critical care prognosis using mixtures of gaussian processes. IEEE Transactions on Biomedical Engineering, 65(1):207-218, 2017.
- Christine B Ambrosone, Timothy R Rebbeck, Gareth J Morgan, Kathy S Albain, Eugenia E Calle, William E Evans, Daniel F Hayes, Lawrence H Kushi, Howard L McLeod, Julia H Rowland, et al. New developments in the epidemiology of cancer prognosis: traditional and molecular predictors of treatment response and survival. Cancer Epidemiology Biomarkers & Prevention, 15(11):2042-2046, 2006.
- Samin Babaei Rikan, Amir Sorayaie Azar, Amin Naemi, Jamshid Bagherzadeh Mohasefi, Habibol-lah Pirnejad, and Uffe Kock Wiil. Survival prediction of glioblastoma patients using modern deep learning and machine learning techniques. Scientific Reports, 14(1):2371, 2024.

495

- Amar Bhide, Prakesh S Shah, and Ganesh Acharya. A simplified guide to randomized controlled trials. *Acta obstetricia et gynecologica Scandinavica*, 97(4):380–387, 2018.
- Ioana Bica, Ahmed M Alaa, Craig Lambert, and Mihaela Van Der Schaar. From real-world patient data to individualized treatment effects using machine learning: current and future methods to address underlying challenges. *Clinical Pharmacology & Therapeutics*, 109(1):87–100, 2021.
- Keting Cen, Huawei Shen, Jinhua Gao, Qi Cao, Bingbing Xu, and Xueqi Cheng. Anae: Learning node context representation for attributed network embedding. *arXiv preprint arXiv:1906.08745*, 2019.
- Adam M Chekroud, Matt Hawrilenko, Hieronimus Loho, Julia Bondar, Ralitza Gueorguieva, Alkomiet Hasan, Joseph Kambeitz, Philip R Corlett, Nikolaos Koutsouleris, Harlan M Krumholz, et al. Illusory generalizability of clinical prediction models. *Science*, 383(6679):164–167, 2024.
- David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554, 2002.
- Pamela Fernainy, Alan A Cohen, Eleanor Murray, Elena Losina, Francois Lamontagne, and Nadia Sourial. Rethinking the pros and cons of randomized controlled trials and observational studies in the era of big data and advanced methods: a panel discussion. In *BMC proceedings*, volume 18, pp. 1. Springer, 2024.
- Stefan Feuerriegel, Dennis Frauen, Valentyn Melnychuk, Jonas Schweisthal, Konstantin Hess, Alicia Curth, Stefan Bauer, Niki Kilbertus, Isaac S Kohane, and Mihaela van der Schaar. Causal machine learning for predicting treatment outcomes. *Nature Medicine*, 30(4):958–968, 2024.
- Kevin C Floyd, Jorge Yarzebski, Frederick A Spencer, Darleen Lessard, James E Dalen, Joseph S Alpert, Joel M Gore, and Robert J Goldberg. A 30-year perspective (1975–2005) into the changing landscape of patients hospitalized with initial acute myocardial infarction: Worcester heart attack study. *Circulation: Cardiovascular Quality and Outcomes*, 2(2):88–95, 2009.
- Joseph Futoma, Morgan Simons, Trishan Panch, Finale Doshi-Velez, and Leo Anthony Celi. The
   myth of generalisability in clinical research and machine learning in health care. *The Lancet Digital Health*, 2(9):e489–e492, 2020.
- 517
  518
  519
  520
  520
  521
  521
  521
  522
  523
  524
  525
  524
  525
  526
  526
  526
  527
  528
  529
  528
  529
  529
  529
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
  520
- Lea Goetz, Nabeel Seedat, Robert Vandersluis, and Mihaela van der Schaar. Generalization—a key
   challenge for responsible ai in patient-facing clinical applications. *npj Digital Medicine*, 7(1): 126, 2024.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1.
   MIT Press, 2016.
- Eduardo Hariton and Joseph J Locascio. Randomised controlled trials—the gold standard for effectiveness research. *BJOG: an international journal of obstetrics and gynaecology*, 125(13):1716, 2018.
- Arthur E Hoerl and Robert W Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- 533 David W Hosmer Jr, Stanley Lemeshow, and Susanne May. *Applied survival analysis: regression* 534 *modeling of time-to-event data*, volume 618. John Wiley & Sons, 2008.
- Patrik Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear causal discovery with additive noise models. *Advances in neural information processing systems*, 21, 2008.
- 539 Markus Kalisch and Peter Bühlman. Estimating high-dimensional directed acyclic graphs with the pc-algorithm. *Journal of Machine Learning Research*, 8(3), 2007.

- Trent Kyono, Yao Zhang, and Mihaela van der Schaar. Castle: Regularization via auxiliary causal graph discovery. *Advances in Neural Information Processing Systems*, 33:1501–1512, 2020.
- 543 Changhee Lee, Alexander Light, Ahmed Alaa, David Thurtle, Mihaela van der Schaar, and Vin544 cent J Gnanapragasam. Application of a novel machine learning framework for predicting non545 metastatic prostate cancer-specific mortality in men using the surveillance, epidemiology, and end
  546 results (seer) database. *The Lancet Digital Health*, 3(3):e158–e165, 2021.
- 547 Kyung Hwa Lee, Gwang Hyeon Choi, Jihye Yun, Jonggi Choi, Myung Ji Goh, Dong Hyun Sinn,
  548 Young Joo Jin, Minseok Albert Kim, Su Jong Yu, Sangmi Jang, et al. Machine learning-based
  549 clinical decision support system for treatment recommendation and overall survival prediction of
  550 hepatocellular carcinoma: a multi-center study. *npj Digital Medicine*, 7(1):2, 2024.
- Kolby Nottingham Markelle Kelly, Rachel Longjohn. The UCI Machine Learning Repository.
   https://archive.ics.uci.edu.
- Francesco Montagna, Atalanti Mastakouri, Elias Eulig, Nicoletta Noceti, Lorenzo Rosasco, Dominik Janzing, Bryon Aragam, and Francesco Locatello. Assumption violations in causal discovery and the robustness of score matching. *Advances in Neural Information Processing Systems*, 36, 2024.
- Ignavier Ng, Shengyu Zhu, Zhitang Chen, and Zhuangyan Fang. A graph autoencoder approach to causal structure learning. *arXiv preprint arXiv:1911.07420*, 2019.
  - Judea Pearl. Causal inference in statistics: An overview. 2009.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Frederike H Petzschner. Practical challenges for precision medicine. *Science*, 383(6679):149–150, 2024.
- Marco Piccininni, Stefan Konigorski, Jessica L Rohmann, and Tobias Kurth. Directed acyclic graphs
   and causal thinking in clinical risk prediction modeling. *BMC medical research methodology*, 20:
   1–9, 2020.
- Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, Antti Kerminen, and Michael Jordan. A linear non-gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7(10), 2006.
  - Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, prediction, and search*. MIT press, 2001.
  - Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.
  - Yue Yu, Jie Chen, Tian Gao, and Mo Yu. Dag-gnn: Dag structure learning with graph neural networks. In *International conference on machine learning*, pp. 7154–7163. PMLR, 2019.
    - Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous optimization for structure learning. *Advances in neural information processing systems*, 31, 2018.
- 583 584 585 586

588 589 590

592

575

576

577

578 579

580

581

582

565

## A AUGMENTED LAGRANGIAN METHOD-BASED OPTIMISATION

We use Augmented Lagrangian method to optimise the constrained optimisation problem in (5) with the acyclicity constraint  $h(\mathbf{W}) = tr(e^{\mathbf{W} \odot \mathbf{W}}) - d$  as:

$$\mathcal{L}_{\rho}(\mathbf{W},\Theta_{1},\Theta_{2},\Theta_{3},\alpha) = \frac{(1-\kappa)}{2n} \sum_{i=1}^{n} \left\| X^{(i)} - \hat{X}^{(i)} \right\|_{F}^{2} + \lambda \left\| \mathbf{W} \right\|_{1} + \alpha h(\mathbf{W}) + \frac{\rho}{2} |h(\mathbf{W})|^{2} + \frac{\kappa}{n} \sum_{i=1}^{n} \left\| Y^{(i)} - \hat{Y}^{(i)} \right\|_{F}^{2},$$

where  $\alpha$  is the Lagrange multiplier,  $\rho > 0$  is multiplier for penalty and  $\lambda$  is the L1-norm penalty for W. We solve the following optimisation problem by using Adam optimiser at each iteration:

$$\mathbf{W}^{k+1}, \boldsymbol{\Theta_1}^{k+1}, \boldsymbol{\Theta_2}^{k+1}, \boldsymbol{\Theta_3}^{k+1} = \underset{\mathbf{W}, \boldsymbol{\Theta_1}, \boldsymbol{\Theta_2}, \boldsymbol{\Theta_3}}{\arg\min} \mathcal{L}^k_{\rho}(\mathbf{W}, \boldsymbol{\Theta_1}, \boldsymbol{\Theta_2}, \boldsymbol{\Theta_3}, \boldsymbol{\alpha}^k)$$

We then update the parameters  $\alpha$  and  $\rho$  for the next iteration as:

$$\begin{split} \alpha^{k+1} &= \alpha^k + \rho^k h(\mathbf{W}^{k+1}), \\ \rho^{k+1} &= \begin{cases} \beta \rho^k, & \text{if } |h(\mathbf{W}^{k+1})| \geq \gamma |h(\mathbf{W}^k)|, \\ \rho^k, & \text{otherwise,} \end{cases} \end{split}$$

where  $\gamma < 1$  and  $\beta > 1$  are training hyperparameters.

## **B REPRODUCIBILITY STATEMENT**

All the methods used Adam optimiser with a learning rate of 1e-3. The models were trained for 300 epochs with an early stopping criterion based on validation loss (except the L2-norm with early stopping (ES) variant of the MLP which employed early stopping based on validation score). In addition, our model and CausalGAE update Lagrange multiplier  $\alpha$  and penalty  $\rho$  over 20 iterations with early stopping based on a threshold. For both models, we use the default threshold and loss hyperparameters as in (Ng et al., 2019). All methods used the same data splits and identical seeds. For all our experiments we use a machine equipped with Intel i9-10900X processor and NVIDIA RTX2080 GPUs. We intend to make our code publicly available.

# C SENSITIVITY ANALYSIS

We perform a sensitivity analysis for the loss hyperparameter  $\kappa$  by using the synthetic datasets in Case 1:

$$\mathbf{X} = 2\sin(\mathbf{W}^{T}(\mathbf{X} + \mathbf{1}) + 0.5 \cdot \mathbf{1}) + \cos(\mathbf{W}^{T}(\mathbf{X} + \mathbf{1}) + 0.5 \cdot \mathbf{1}) + \mathbf{Z}$$

and Case 2:

$$Y = 2\sin(\mathbf{W}^{T}(\mathbf{X} + \mathbf{1}) + 0.5 \cdot \mathbf{1}) + \cos(\mathbf{W}^{T}(\mathbf{X} + \mathbf{1}) + 0.5 \cdot \mathbf{1}) + \mathbf{Z}$$
$$+ \cos((X^{non - pa(Y)_{0}} + \mathbf{1}) + 0.5 \cdot \mathbf{1})$$

This hyperparameter controls the fraction of supervised loss added to the overall loss function. Figure 3 illustrates the test MSE loss with respect to  $\kappa$ . We observe minimal variation between  $\kappa = 0.25$  and  $\kappa = 0.75$ . The MSE loss is worst at  $\kappa = 0$ , reinforcing the importance of supervised loss for generalisation. We choose  $\kappa = 0.25$  for all our experiments.

## D ROBUST CAUSAL DISCOVERY

 Figures 4 and 5 compare the true causal graph with causal graphs recovered from CausalGAE and our model for synthetic datasets from Case1 and Case 2 respectively. In contrast to our model, CausalGAE fails to identify the causal parents of *Y* in both the cases.

# E CLASSIFICATION ON REAL DATA

We show the causal graphs recovered by our model for various real datasets here (Figure ). The
Statlog heart dataset has 270 samples and 13 features. The Breast cancer (Wisconsin Diagnostic)
dataset includes 569 samples and 30 features. The Las Vegas ratings dataset includes 504 samples and 19 features.



Figure 3: Sensitivity of our model performance to the loss hyperparameter  $\kappa$ .



Figure 4: True causal graph for synthetic data Case 1 and the recovered graphs from CausalGAE and our model.

#### F **APPLICATION TO SURVIVAL ANALYSIS**

We convert the time-to-event analysis problem to classification by thresholding the continuous time and assigning ground truth labels based on the event. Specifically, a positive class label is assigned to those cases where an event (death at discharge) occurred before k days, k being the median length of hospital stay in Scenario 1 and median length of hospital stay over the training cohort in Scenario 2. A negative label is assigned for all other cases where no event occurred before k days.



Figure 5: True causal graph for synthetic data Case 2 and the recovered graphs from CausalGAE and our model.

