

Stable Prediction across Unknown Environments

Kun Kuang*, Ruoxuan Xiong*, Peng Cui, Susan Athey, and Bo Li

Abstract—In many important machine learning applications, the training distribution used to learn a probabilistic classifier differs from the testing distribution on which the classifier will be used to make predictions. Traditional methods correct the distribution shift by reweighting the training data with the ratio of the density between test and training data. In many applications training takes place without prior knowledge of the testing distribution on which the algorithm will be applied in the future. Recently, methods have been proposed to address the shift by learning causal structure, but those methods rely on the diversity of multiple training data to a good performance, and have complexity limitations in high dimensions. In this paper, we propose a novel Deep Global Balancing Regression (DGBR) algorithm to jointly optimize a deep auto-encoder model for feature selection and a global balancing model for stable prediction across unknown environments. The global balancing model constructs balancing weights that facilitate estimating of partial effects of features (holding fixed all other features), a problem that is challenging in high dimensions, and thus helps to identify stable, causal relationships between features and outcomes. The deep auto-encoder model is designed to reduce the dimensionality of the feature space, thus making global balancing easier. We show, both theoretically and with empirical experiments, that our algorithm can make stable predictions across unknown environments. Our experiments on both synthetic and real world datasets demonstrate that our DGBR algorithm outperforms the state-of-the-art methods for stable prediction across unknown environments.

Index Terms—Stability, Stable Prediction, Unknown Environments, Confounder Balancing.

1 INTRODUCTION

Predicting unknown outcome values based on their observed features using a model estimated on a training data set is a common statistical problem. Many machine learning and data mining methods have been proposed and shown to be successful when the test data and training data come from the same distribution. However, the best-performing models for a given distribution of training data typically exploit subtle statistical relationships among features, making them potentially more prone to prediction error when applied to test data sets where, for example, the joint distribution of features differs from that in the training data. Therefore, it can be useful to develop predictive algorithms that are robust to shifts in the environment, particularly in application areas where models can not be retrained as quickly as the environment changes, i.e., online prediction.

Recently, many methods [1], [2], [3], [4], [5], [6] have been proposed to address this problem. The main idea of these methods is to reweight training data with a density ratio, so that its distribution can become more closely aligned with the distribution of test data. The methods have achieved good performance for correcting for distribution shift, but they require prior knowledge of the test distribution when estimating the density ratio.

For the case of unknown test data, some researchers have proposed learning methods where training takes place across multiple training datasets. By exploring the invariance across multiple datasets, Peters et al. [7] proposed an algorithm to identify causal features, and Rojas-Carulla et al. [8] proposed a causal transform framework to learn invariant structure. Similarly, domain generalization methods [9] try to learn an invariant representation of data. The performance of these methods relies on the diversity of

their multiple training data, and they cannot address distribution shifts which do not appear in their training data. Moreover, most of these methods are highly complex, with training complexity growing exponentially with the dimension of the feature space in the worst case, which is not acceptable in high dimensional settings.

In this paper, we focus on an environment where the expected value of the outcome conditional on all covariates is stable across environments. Further, covariates fall into one of two categories: for the first category, the conditional expectation has a non-zero dependence on the covariates; we call these “causal” variables, although in some applications they might better be described as variables that have a structural relationship with the outcome. For example, ears, noses, and whiskers are structural features of cats that are stable across different environments where images of animals may be taken. A second category of variable are termed “noisy variables,” which are variables that are correlated with either the causal variables, the outcome, or both, but do not themselves have a causal effect on the outcome; conditional on the full set of causal variables, they do not affect expected outcomes. Further, we consider a setting where the analyst may not know a prior which variables fall into each category. Finally, we assume that there are no unobserved confounders, so that it is possible to estimate the causal effect of each causal variable with a very large dataset when all covariates are adequately controlled for. We focus on settings when there are many features and perhaps limited data.

One way to improve the stability of prediction algorithms in such a setting is to isolate the impact of each individual feature. If the expectation of the outcome conditional on covariates is stable across environments, and variability in the joint distribution of features is the source of instability, then the stable prediction problem can be solved by estimating the conditional expectation function accurately. With a small number of discrete features and a large enough dataset, simple estimation methods such as ordinary least squares can accomplish this goal. If there is a larger number of features but only a few matter for the conditional expectation

* Equal Contribution

- Kun Kuang, Peng Cui and Bo Li are with Tsinghua University, Beijing 100084, China. Email: kkun2010@gmail.com, cuip@tsinghua.edu.cn, libo@sem.tsinghua.edu.cn.
- Ruoxuan Xiong and Susan Athey are with Stanford University, CA 94305, USA. Email: rxiong@stanford.edu, athey@stanford.edu.

We are grateful for helpful comments from Vitor Hadad.

(that is, the true outcome model is sparse), regularized regression can be applied to consistently estimate the conditional expectation function. However, with a larger set of causal features relative to the number of observations, regularized regression will no longer consistently estimate partial effects. For example, LASSO will omit many variables from the regression, while the coefficients on included variables depend on the covariance of the outcome with the omitted variables as well as on the covariance between the omitted and included variables. This results in instability: if the covariance among features differs across environments, then prediction based on such a model will be unstable across environments. In such high-dimensional cases, alternative approaches are required.

Here, we use an approach motivated by the literature on causal inference, where variable balancing strategies are used for estimating the average effect of changing a single binary covariate (the treatment). Causal inference methods optimize a different objective than prediction-based methods; they prioritize consistent estimation of treatment effects over prediction in a given training data set. The methods are designed for a scenario where the analyst has domain knowledge about which variable has a causal effect, so that the focus of the analysis is on estimating the effect of the treatment in the presence of other features which are known to be confounders (variables that affect both treatment assignment and potential outcomes). Indeed, only after controlling for confounders can the difference in the expectation of the outcome between treatment and control groups be interpreted as a treatment effect. One approach to estimating treatment effects in the presence of confounders is to use variable balancing methods, which attempt to construct weights that balance the distribution of covariates between a treatment and a control group. They either employ propensity scores [10], [11], [12], [13], [14], or optimize balancing weights directly [15], [16], [17], [18]. These methods provide an efficient approach to estimate causal effects with a small number of treatment variables in observational studies, but most of them can not handle well settings where there may be many causal variables and the analyst does not know which ones are causal; as such, existing covariate balancing methods do not immediately extend to the general stable prediction problem.

Inspired by balancing methods from the causal inference literature, we propose a Deep Global Balancing Regression (DGBR) algorithm for stable prediction. The framework is illustrated in Figure 2, which consists of three (jointly optimized) sub-models: (i) a deep auto-encoder to reduce the dimensionality of the features, (ii) construction of balancing weights that enable the effect of each covariate to be isolated, and (iii) estimation of a predictive model using the encoded features and balancing weights. As this algorithm explicitly prioritizes covariate balancing (at the expense of a singular focus on predictive accuracy in a given training dataset), it is able to achieve greater stability than a purely predictive model. Using both empirical experiments and theoretical analysis, we establish that our algorithm achieves stability in prediction across unknown environments. The experimental results on both synthetic and real world datasets demonstrate that our algorithm outperforms all the baselines for the stable prediction problem.

In summary, the contributions of this paper are listed as follows:

- We investigate the problem of stable prediction across unknown environments, where the distribution of agnostic test data might be very different with the training data.

- We propose a novel DGBR algorithm to jointly optimize deep auto-encoder for dimension reduction and global balancing for estimation of causal effects, and simultaneously address the stable prediction problem.
- We give theoretical analysis on our proposed algorithm and prove that our algorithm can make a stable prediction across unknown environments by global balancing.
- The advantages of our DGBR algorithm are demonstrated on both synthetic and real world datasets.

The rest of the paper is organized as follows. Section 2 reviews the related work. In Section 3, we give problem formulation and introduce our DGBR algorithm. Section 4 gives the optimization and discussion on our algorithm. Section 5 gives the theoretical analysis on our algorithm. Section 6 gives the experimental results. Finally, Section 7 concludes.

2 RELATED WORK

In this section, we investigate the previous related work, including literatures on covariate shift, variable balancing, and invariant learning.

The covariate shift literature [1] focuses on settings where the data distribution for training is different than the data distribution for testing. To correct for the differences, [1] introduced the idea of reweighting samples in training data by the ratio of the density in the testing data to the density in the training data. A variety of techniques have been proposed to estimate the density ratio, including discriminative estimation [2], Kullback-Leibler importance estimation [3], kernel mean matching [4] [19], maximum entropy methods [5], minimax optimization [20], and robust bias-aware approach [6]. These methods achieved good performance for correcting for covariate shifts, but most of them require prior knowledge of testing distribution to estimate the density ratio. In contrast, we focus on the stable prediction across unknown environments in this paper.

Adjusting for confounders is a key challenge for estimating causal effects in observational studies. To precisely estimate causal effects in the presence of many confounders, covariate balancing methods have been proposed [10], [13], [15], [16], [17], [18], [21]. In a seminal paper, Rosenbaum and Rubin [10] proposed to achieve variable balancing by reweighting observations by the inverse of propensity score. Kuang et al. [13] proposed a data-driven variable decomposition method for variable balancing. Li et al. [22] balanced the variables by matching on their nonlinear representation. Hainmueller [15] introduced entropy balancing method for variable balancing across a range of statistical tasks. Athey et al. [17] proposed approximate residual balancing algorithm, which, motivated by doubly robust approaches, combines outcome modeling using the LASSO with balancing weights constructed to approximately balance covariates between treatment and control groups. Kuang et al. [18] proposed a differentiated variable balancing algorithm by jointly optimizing sample weights and variable weights. These methods provide an effective way to estimate causal effects in observational studies, but they are limited to estimate causal effect of one variable, and are not designed for the case with many causal variables; further, the methods assume that the analyst has prior knowledge of which covariates have a causal effect and which do not.

Recently, some methods have been proposed to make prediction on agnostic test data using the method of invariant learning. Peters et al. [7] proposed an algorithm to identify causal predictors

TABLE 1: Symbols and definitions.

| Symbols | Definitions |
|--|-----------------------|
| n | Sample size |
| p | Dimension of features |
| $\mathbf{X} = \{\mathbf{S}, \mathbf{V}\} \in \{0, 1\}^p$ | Features |
| $\mathbf{S} \in \{0, 1\}^{p_s}$ | Stable features |
| $\mathbf{V} \in \{0, 1\}^{p_v}$ | Noisy features |
| $Y \in \{0, 1\}$ | Outcome |
| $W \in \mathbb{R}^{+n \times 1}$ | Global sample weights |
| $\phi(\cdot)$ | Embedding function |

by exploring the invariance of the conditional distribution of the outcome across multiple training datasets. Rojas-Carulla et al. [8] proposed a causal transfer framework to identify invariant predictors across multiple datasets and then use them for prediction. Similarly, domain generalization [9] methods estimate an invariant representation of data by minimizing the dissimilarity across training domains. Invariant learning methods can be used to estimate a model that will in principle perform well for an unknown test dataset, but the performance of these methods relies on the diversity of their multiple training data, and they cannot address the distribution shift which does not appear in their training data.

3 PROBLEM AND OUR ALGORITHM

In this section, we first give problem formulation, then introduce the details of our deep global balancing regression algorithm. Finally, we give theoretical analysis about our proposed algorithm.

3.1 Problem Formulation

Let \mathcal{X} denote the space of observed features and \mathcal{Y} denote the outcome space. For simplicity, we consider the case where the features have finite support, which without loss of generality can be represented as a set of binary features: $\mathcal{X} = \{0, 1\}^p$. We also focus on the case where the outcome space is binary: $\mathcal{Y} = \{0, 1\}$. We define an **environment** to be a joint distribution P_{XY} on $\mathcal{X} \times \mathcal{Y}$, and let \mathcal{E} denote the set of all environments. In each environment $e \in \mathcal{E}$, we have dataset $D^e = (\mathbf{X}^e, Y^e)$, where $\mathbf{X}^e \in \mathcal{X}$ are predictor variables and $Y^e \in \mathcal{Y}$ is a response variable. The joint distribution of features and outcomes on (\mathbf{X}, Y) can vary across environments: $P_{XY}^e \neq P_{XY}^{e'}$ for $e, e' \in \mathcal{E}$, and $e \neq e'$.

In this paper, our goal is to learn a predictive model, which can make a stable prediction across unknown environments. Before giving problem formulation, we first define *Average_Error* and *Stability_Error* across environments of a predictive model as:

$$\text{Average_Error} = \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \text{Error}(D^e), \quad (1)$$

$$\text{Stability_Error} = \sqrt{\frac{1}{|\mathcal{E}|-1} \sum_{e \in \mathcal{E}} (\text{Error}(D^e) - \text{Average_Error})^2} \quad (2)$$

where $|\mathcal{E}|$ refers to the number of environments, and $\text{Error}(D^e)$ represents the predictive error on dataset D^e from environment e .

In this paper, we define Stability [23] by *Stability_Error*. The smaller *Stability_Error*, the better a model is ranked in terms of Stability. Then, we define the stable prediction problem as follow:

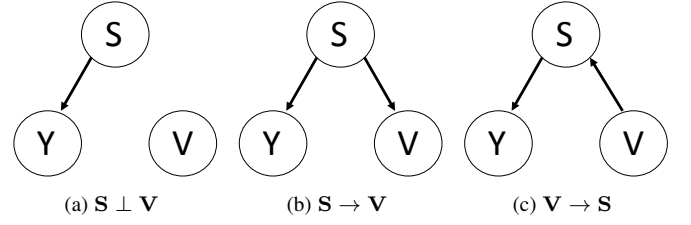


Fig. 1: Three diagrams for stable features \mathbf{S} , noisy features \mathbf{V} , and response variable Y .

Problem 1 (Stable Prediction). *Given one training environment $e \in \mathcal{E}$ with dataset $D^e = (\mathbf{X}^e, Y^e)$, the task is to **learn** a predictive model to predict across unknown environment \mathcal{E} with not only small *Average_Error* but also small *Stability_Error*.*

Suppose $\mathbf{X} = \{\mathbf{S}, \mathbf{V}\}$. We define \mathbf{S} as *stable features*, and refer to the other features $\mathbf{V} = \mathbf{X} \setminus \mathbf{S}$ as *noisy features*, where the following assumption gives their defining properties:

Assumption 1. *There exists a probability mass function $P(y|s)$ such that for all environments $e \in \mathcal{E}$, $Pr(Y^e = y | \mathbf{S}^e = s, \mathbf{V}^e = v) = Pr(Y^e = y | \mathbf{S}^e = s) = P(y|s)$.*

With Assumption 1, we can address the stable prediction problem by building a model that learns the stable function $P(y|s)$. To understand the content of Assumption 1, without loss of generality we can write a generative model for the outcome unit i in environment e with stable features s , where $h(\cdot)$ is a known function to account for discreteness of Y :

$$Y_i^e(s) = h(g(s) + \epsilon_{s,i}^e), \text{ and } Y_i^e = Y_i^e(\mathbf{S}_i) = h(g(\mathbf{S}_i) + \epsilon_{\mathbf{S}_i,i}^e).$$

$Y_i^e(s)$ is the outcome that would occur for unit i in environment e if the input is equal to s . If we allow $\epsilon_{s,i}^e$ to be correlated with the unit's features \mathbf{X}_i in arbitrary ways, Assumption 1 may fail, for example if \mathbf{V}_i^e is positively correlated with $\epsilon_{s,i}^e$ then units with higher values of \mathbf{V}_i^e would have higher than average values of Y_i^e , so that \mathbf{V}_i^e would be a useful predictor in a given environment, but that relationship might vary across environments, leading to instability. If we first impose the condition that for each s , $\epsilon_{s,i}^e$ is independent of \mathbf{V}_i^e conditional on \mathbf{S}_i^e , then given the model specification, \mathbf{V}_i^e is no longer needed as a predictor for outcomes conditional on \mathbf{S}_i^e . If we second impose the condition that for each s , $\epsilon_{s,i}^e$ is independent of \mathbf{S}_i^e conditional on \mathbf{V}_i^e , then instability in the distribution of $\epsilon_{s,i}^e$ across environments will not affect $Pr(Y^e = y | \mathbf{S}^e = s, \mathbf{V}^e = v)$. Maintaining the first condition, the second condition is sufficient not only for Assumption 1 but also to enable consistent estimation of $g(\cdot)$ using techniques from the causal inference literature in a setting with sufficient sample size and when the analyst has prior knowledge of the set of stable features; we propose a method that will estimate g without prior knowledge of which features are stable. We also observe that a stronger but simpler condition can replace the second condition to guarantee Assumption 1, namely that the distribution of $\epsilon_{s,i}^e$ does not vary with $\{e, s\}$. Fig. 1 illustrates three relationships between predictor variables $\mathbf{X}^e = \{\mathbf{S}^e, \mathbf{V}^e\}$ and response variable Y^e consistent with the conditions, including $\mathbf{S} \perp \mathbf{V}$, $\mathbf{S} \rightarrow \mathbf{V}$, and $\mathbf{V} \rightarrow \mathbf{S}$.

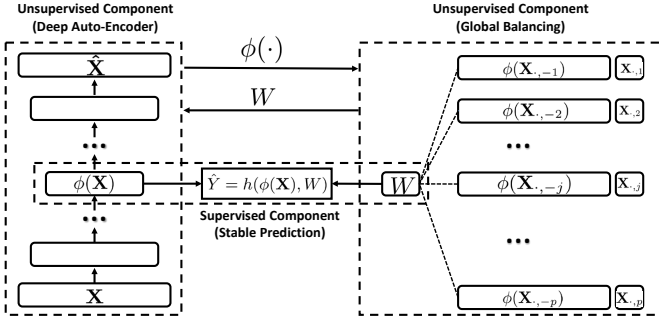


Fig. 2: The framework of our proposed DGBR model.

3.2 The Model

3.2.1 Framework

We propose a Deep Global Balancing Regression (DGBR) algorithm to identify stable features and capture non-linear structure for stable prediction. Its framework is shown in Figure 2. To identify the stable features, we propose a global balancing model, where we learn global sample weights which can be used to estimate the effect of each feature while controlling for the other features and thus identify stable features. To capture the non-linear structure between stable features and response variable, we employ a deep auto-encoder model, which is composed of multiple non-linear mapping functions to map the input data to a non-linear and low dimensional space. Balancing in a low dimensional space simplifies the problem of global balancing, since for each covariate j , the weights balance the constructed covariates from the dimension reduction $\phi(\mathbf{X}_{\cdot,-j})$ across realizations of $\mathbf{X}_{\cdot,j}$. Finally, weighting observations with the global sample weights, we learn a predictive model for outcomes as a function of the low-dimensional representation of covariates using regularized regression. All three components of the model are jointly optimized in the algorithm.

3.2.2 Global Balancing Regression Algorithm

In this section, we develop the construction of global balancing weights. To be self-contained, we briefly revisit the key idea of variable balancing technique. Variable balancing techniques are often used for **causal effect estimation** in observational studies, where the distributions of covariates are different **between treated and control groups because of non-random treatment assignment**, but treatment assignment is independent of potential outcomes conditional on covariates. To consistently estimate causal effects in such a setting, one has to **balance the distribution of covariates between treatment and control**. Most variables balancing approaches exploit moments to characterize distributions, and balance them between treated and control groups by adjusting sample weights W as following:

$$W = \arg \min_W \left\| \frac{\sum_{i:T_i=1} W_i \cdot \mathbf{X}_i}{\sum_{i:T_i=1} W_i} - \frac{\sum_{i:T_i=0} W_i \cdot \mathbf{X}_i}{\sum_{i:T_i=0} W_i} \right\|_2^2. \quad (3)$$

Given a treatment variable T , the $\frac{\sum_{i:T_i=1} W_i \cdot \mathbf{X}_i}{\sum_{i:T_i=1} W_i}$ and $\frac{\sum_{i:T_i=0} W_i \cdot \mathbf{X}_i}{\sum_{i:T_i=0} W_i}$ represent the first-order moments of variables \mathbf{X} on treated ($T = 1$) and control ($T = 0$) groups, respectively. By sample reweighting with W learnt from Eq. (3), one can estimate the causal effect of treatment variable on response variable by comparing the average difference of Y between treated and control

groups. In high-dimensional problems, approximate balancing can be used for consistent estimation under some additional assumptions [17], where to control variance of estimates the sum of squared weights is also penalized in the minimization.

In low dimensions, the same approach could be employed to estimate $Pr(Y = y | \mathbf{X} = x)$ for different values of x . However, **when p is large, there may not be sufficient data to do so**, and so approximate balancing techniques generalized to the case where \mathbf{X} is a vector of indicator variables may perform well in practice, and also help identify stable features from the larger vector \mathbf{X} . We propose a global balancing regularizer, where we successively **regard each variable as treatment variable** and balance all of them together via learning global sample weights by minimizing:

$$\sum_{j=1}^p \left\| \frac{\mathbf{X}_{\cdot,-j}^T \cdot (W \odot \mathbf{X}_{\cdot,j})}{W^T \cdot \mathbf{X}_{\cdot,j}} - \frac{\mathbf{X}_{\cdot,-j}^T \cdot (W \odot (1 - \mathbf{X}_{\cdot,j}))}{W^T \cdot (1 - \mathbf{X}_{\cdot,j})} \right\|_2^2, \quad (4)$$

where W is global sample weights, $\mathbf{X}_{\cdot,j}$ is the j^{th} variable in \mathbf{X} , and $\mathbf{X}_{\cdot,-j} = \mathbf{X} \setminus \{\mathbf{X}_{\cdot,j}\}$ means all the remaining variables by removing the j^{th} variable in \mathbf{X} ¹. The summand represents the loss from covariate imbalance when setting variable $\mathbf{X}_{\cdot,j}$ as the treatment variable, and \odot refers to Hadamard product. Note that only first-order moment is considered in Eq. (4), but higher order moments can be easily incorporated by including interaction features of \mathbf{X} .

By sample reweighting with W learnt from Eq. (4), we can identify stable features S by checking if there is any correlation between Y and X covariate by covariate, because, as we show below, only stable features are correlated with Y after sample reweighting by W .

With the global balancing regularizer in Eq. (4), we propose a Global Balancing Regression (GBR) algorithm to jointly optimize global sample weights W and regression coefficients β for stable prediction based on traditional logistical regression as:

$$\begin{aligned} \min \quad & \sum_{i=1}^n W_i \cdot \log(1 + \exp((1 - 2Y_i) \cdot (\mathbf{X}_i \beta))), \\ \text{s.t.} \quad & \sum_{j=1}^p \left\| \frac{\mathbf{X}_{\cdot,-j}^T \cdot (W \odot \mathbf{X}_{\cdot,j})}{W^T \cdot \mathbf{X}_{\cdot,j}} - \frac{\mathbf{X}_{\cdot,-j}^T \cdot (W \odot (1 - \mathbf{X}_{\cdot,j}))}{W^T \cdot (1 - \mathbf{X}_{\cdot,j})} \right\|_2^2 \leq \lambda_1, \quad W \succeq 0, \\ & \|W\|_2^2 \leq \lambda_2, \quad \|\beta\|_2^2 \leq \lambda_3, \quad \|\beta\|_1 \leq \lambda_4, \quad (\sum_{k=1}^n W_k - 1)^2 \leq \lambda_5 \end{aligned} \quad (5)$$

where \mathbf{X}_i is the i^{th} row / sample in \mathbf{X} , and $\sum_{i=1}^n W_i \cdot \log(1 + \exp((1 - 2Y_i) \cdot (\mathbf{X}_i \beta)))$ is the weighted loss of logistic regression and the loss is defined as the minus log likelihood. The terms $W \succeq 0$ constrain each of sample weights to be non-negative. With norm $\|W\|_2^2 \leq \lambda_2$, we can reduce the variance of the sample weights. Elastic net constraints $\|\beta\|_2^2 \leq \lambda_3$ and $\|\beta\|_1 \leq \lambda_4$ help to avoid overfitting. The formula $(\sum_{k=1}^n W_k - 1)^2 \leq \lambda_5$ avoids all the sample weights to be zero.

3.2.3 Deep Global Balancing Regression Algorithm

The proposed GBR algorithm in Eq. (5) can help to identify stable features and make a stable prediction, but with many features relative to observations, it may be difficult to estimate the effects of all the features as well as their interactions, and it might also be challenging for GBR to learn global sample weights.

To address these challenges, we propose a Deep Global Balancing Regression (DGBR) algorithm by jointly optimizing Deep auto-encoder and Global Balancing Regression. Following standard approaches [24], the deep auto-encoder consists of multiple non-linear mapping functions to map the input data to a

1. We obtain $\mathbf{X}_{\cdot,-j}$ in experiment by setting the value of j^{th} variable in \mathbf{X} as zero.

low dimensional space while capturing the underlying features interactions. Deep auto-encoder is an unsupervised model which is composed of two parts, the encoder and decoder. The encoder maps the input data to low-dimensional representations, while the decoder reconstructs the original input space from the representations. Given the input \mathbf{X}_i , the hidden representations for each layer are shown as follows:

$$\begin{aligned}\phi(\mathbf{X}_i)^{(1)} &= \sigma(\mathbf{A}^{(1)}\mathbf{X}_i + b^{(1)}) \\ \phi(\mathbf{X}_i)^{(k)} &= \sigma(\mathbf{A}^{(k)}\phi(\mathbf{X}_i)^{(k-1)} + b^{(k)}), k = 2, \dots, K\end{aligned}$$

where K is the number of layer. $\mathbf{A}^{(k)}$ and $b^{(k)}$ are weight matrix and bias on k^{th} layer. $\sigma(\cdot)$ represents non-linear activation function.²

After obtaining the representation $\phi(\mathbf{X}_i)^{(K)}$, we can obtain the reconstruction $\hat{\mathbf{X}}_i$ by reversing the calculation process of encoder with parameters $\hat{\mathbf{A}}^{(k)}$ and $\hat{b}^{(k)}$. The goal of deep auto-encoder is to minimize the reconstruction error between the input \mathbf{X}_i and the reconstruction $\hat{\mathbf{X}}_i$ with the following loss function.

$$\mathcal{L} = \sum_{i=1}^n \|\mathbf{X}_i - \hat{\mathbf{X}}_i\|_2^2. \quad (6)$$

By combining the loss functions of deep auto-encoder in Eq. (6) and GBR algorithm in Eq. (5), we give the objective function of our Deep Global Balancing Regression algorithm as:

$$\begin{aligned}\min \quad & \sum_{i=1}^n W_i \cdot \log(1 + \exp((1 - 2Y_i) \cdot (\phi(\mathbf{X}_i)\beta))), \quad (7) \\ \text{s.t.} \quad & \sum_{j=1}^p \left\| \frac{\phi(\mathbf{X}_{\cdot,-j})^T \cdot (W \odot \mathbf{X}_{\cdot,j})}{W^T \mathbf{X}_{\cdot,j}} - \frac{\phi(\mathbf{X}_{\cdot,-j})^T \cdot (W \odot (1 - \mathbf{X}_{\cdot,j}))}{W^T (1 - \mathbf{X}_{\cdot,j})} \right\|_2^2 \leq \lambda_1, \text{ and therefore, we have following equation with probability 1:} \\ & \|(W \cdot \mathbf{1}) \odot (X - \hat{X})\|_F^2 \leq \lambda_2, \quad W \succeq 0, \quad \|W\|_2^2 \leq \lambda_3, \quad \lim_{n \rightarrow \infty} \left(\frac{\mathbf{x}_{\cdot,k}^T (W^* \odot \mathbf{x}_{\cdot,j})}{W^{*T} \mathbf{x}_{\cdot,j}} - \frac{\mathbf{x}_{\cdot,k}^T (W^* \odot (1 - \mathbf{x}_{\cdot,j}))}{W^{*T} (1 - \mathbf{x}_{\cdot,j})} \right) = \frac{2^{p-2}}{2^{p-1}} - \frac{2^{p-2}}{2^{p-1}} = 0. \\ & \|\beta\|_2^2 \leq \lambda_4, \quad \|\beta\|_1 \leq \lambda_5, \quad (\sum_{k=1}^n W_k - 1)^2 \leq \lambda_6 \\ & \sum_{k=1}^K (\|\mathbf{A}^{(k)}\|_F^2 + \|\hat{\mathbf{A}}^{(k)}\|_F^2) \leq \lambda_7,\end{aligned}$$

where $\phi(\cdot) = \phi(\cdot)^{(K)}$ for brevity. $\|(W \cdot \mathbf{1}) \odot (X - \hat{X})\|_F^2$ represents the reconstruction error between input \mathbf{X} and reconstruction $\hat{\mathbf{X}}$ with global sample weights W . The term $\sum_{k=1}^K (\|\mathbf{A}^{(k)}\|_F^2 + \|\hat{\mathbf{A}}^{(k)}\|_F^2) \leq \lambda_7$ regularizes the coefficients of the deep auto-encoder model.

4 THEORETICAL ANALYSIS

In this section, we give theoretical analysis about our algorithm. We prove it can make a stable prediction across unknown environments with sufficient data, and analyze the upper bound about our proposed algorithm.

4.1 Analysis on Stable Prediction

A key requirement for the method to work is the overlap assumption, which is a common assumption in the literature of treatment effect estimation [17]. We suppress the notation for the environment e in the first part of this section.

Assumption 2 (Overlap). For any variable $\mathbf{X}_{\cdot,j}$ when setting it as the treatment variable, it has $\forall j, 0 < P(\mathbf{X}_{\cdot,j} = 1 | \mathbf{X}_{\cdot,-j}) < 1$.

Then, we have following Lemma and Theorem:

Lemma 1. If $\forall j, 0 < P(\mathbf{X}_{\cdot,j} = 1 | \mathbf{X}_{\cdot,-j}) < 1$, and \mathbf{X} are binary, then $\forall i, 0 < P(\mathbf{X}_i = x) < 1$, where \mathbf{X}_i is i^{th} row in \mathbf{X} .

Proof. See Appendix A. \square

2. We use sigmoid function $\sigma(x) = \frac{1}{1 + \exp(-x)}$ as non-linear activation function.

Theorem 1. Let $\mathbf{X} \in \mathbb{R}^{n \times p}$. Under Lemma 1, if number of covariates p is finite, then $\exists W$ such that

$$\lim_{n \rightarrow \infty} \sum_{j=1}^p \left\| \frac{\mathbf{x}_{\cdot,j}^T (W \odot \mathbf{x}_{\cdot,j})}{W^T \mathbf{x}_{\cdot,j}} - \frac{\mathbf{x}_{\cdot,j}^T (W \odot (1 - \mathbf{x}_{\cdot,j}))}{W^T (1 - \mathbf{x}_{\cdot,j})} \right\|_2^2 = 0 \quad (8)$$

with probability 1. In particular, a W that satisfies (8) is $W_i^* = \frac{1}{P(\mathbf{X}_i = x)}$.

Proof. Since $\|\cdot\| \geq 0$, Eq. (8) can be simplified to $\forall j, \forall k \neq j$

$$\lim_{n \rightarrow \infty} \left(\frac{\sum_{i: \mathbf{X}_{i,k}=1, \mathbf{X}_{i,j}=1} W_i}{\sum_{i: \mathbf{X}_{i,j}=1} W_i} - \frac{\sum_{i: \mathbf{X}_{i,k}=1, \mathbf{X}_{i,j}=0} W_i}{\sum_{i: \mathbf{X}_{i,j}=0} W_i} \right) = 0$$

with probability 1. For W^* , from Lemma 1, $0 < P(\mathbf{X}_i = x) < 1, \forall x, \forall i, t = 1$ or 0 ,

$$\begin{aligned}\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i: \mathbf{X}_{i,j}=t} W_i^* &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{x: x_j=t} \sum_{i: \mathbf{X}_i=x} W_i^* \\ &= \lim_{n \rightarrow \infty} \sum_{x: x_j=t} \frac{1}{n} \sum_{i: \mathbf{X}_i=x} \frac{1}{P(\mathbf{X}_i=x)} \\ &= \lim_{n \rightarrow \infty} \sum_{x: x_j=t} P(\mathbf{X}_i = x) \cdot \frac{1}{P(\mathbf{X}_i=x)} = 2^{p-1}\end{aligned}$$

with probability 1 from Law of Large Number. Since features are binary,

$$\begin{aligned}\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i: \mathbf{X}_{i,k}=1, \mathbf{X}_{i,j}=1} W_i^* &= 2^{p-2} \\ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i: \mathbf{X}_{i,j}=0} W_i^* &= 2^{p-1}, \\ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i: \mathbf{X}_{i,k}=1, \mathbf{X}_{i,j}=0} W_i^* &= 2^{p-2}\end{aligned}$$

and therefore, we have following equation with probability 1:

$$\lim_{n \rightarrow \infty} \left(\frac{\mathbf{x}_{\cdot,k}^T (W^* \odot \mathbf{x}_{\cdot,j})}{W^{*T} \mathbf{x}_{\cdot,j}} - \frac{\mathbf{x}_{\cdot,k}^T (W^* \odot (1 - \mathbf{x}_{\cdot,j}))}{W^{*T} (1 - \mathbf{x}_{\cdot,j})} \right) = \frac{2^{p-2}}{2^{p-1}} - \frac{2^{p-2}}{2^{p-1}} = 0. \quad \square$$

The following result shows that if there is sufficient data such that all realizations of x appear in the data, exact balancing weights can be derived. Subsequently, we show that in this case, the components of \mathbf{X} are mutually independent in the reweighted data. This highlights that overlap is a strong assumption. In real-world data sets, when p is large the cardinality of \mathcal{X} is large, and so exactly balancing weights are not available, but the results still highlight that balancing weights will reduce the covariance among features.

Then, based on Lemma 1 and Theorem 8, we have following propositions for stable prediction.

Proposition 1. If $0 < \hat{P}(\mathbf{X}_i = x) < 1$ for all x , where $\hat{P}(\mathbf{X}_i = x) = \frac{1}{n} \sum_i \mathbb{1}(\mathbf{X}_i = x)$, there exists a solution W^* satisfies equation (4) equals 0 and variables in \mathbf{X} are independent after balancing by W^* .

Proof. See Appendix B. \square

Proposition 2. If $0 < \hat{P}(\mathbf{X}_i^e = x) < 1$ for all x in environment e , $Y^{e'}$ and $\mathbf{V}^{e'}$ are independent when the joint probability mass function of $(\mathbf{X}^{e'}, Y^{e'})$ is given by reweighting the distribution from environment e using weights W^* , so that $p^{e'}(x, y) = p^e(y|x) \cdot (1/|\mathcal{X}|)$.

Proof. It is immediate that $Pr(Y^{e'} = y | \mathbf{X}^{e'} = x) = Pr(Y^e = y | \mathbf{X}^e = x)$. Putting this together with Assumption 1, $Pr(Y^{e'} =$

$y|\mathbf{X}^{e'} = x) = Pr(Y^{e'} = y|\mathbf{S}^{e'} = s)$. From Proposition 1, $(\mathbf{S}^{e'}, \mathbf{V}^{e'})$ are mutually independent. Thus, we have

$$\begin{aligned} & Pr(Y^{e'} = y|\mathbf{V}^{e'} = v) \\ &= E_{\mathbf{S}^{e'}}[Pr(Y^{e'} = y|\mathbf{S}^{e'}, \mathbf{V}^{e'} = v)|\mathbf{V}^{e'} = v] \\ &= E_{\mathbf{S}^{e'}}[Pr(Y^{e'} = y|\mathbf{S}^{e'})|\mathbf{V}^{e'} = v] \\ &= Pr(Y^{e'} = y). \end{aligned}$$

Thus, $Y^{e'}$ and $\mathbf{V}^{e'}$ are independent. \square

Propositions 1 and 2 suggest that the GBR algorithm can make a stable prediction across environments that satisfy Assumption 1, since after reweighting, only the stable features are correlated with outcomes, and $p(y|s)$ is unchanged in the reweighted dataset. The objective function of GBR algorithm is to equivalent to maximize log-likelihood of logistic regression, which is known to be consistent. Even though the regularization constraints will cause some bias to the estimated $p(y|s)$, but the bias reduces with sample size n . Thus, with sufficient data, the GBR algorithm should learn $p(y|s)$.

Now consider the properties of the DGBR algorithm:

- 1) *Preserves the above properties of the GBR algorithm while making the overlap property easier to satisfy and reducing the variance of balancing weights.* The Johnson-Lindenstrauss (JL) lemma [25] implies that for any $0 < \epsilon < 1/2$ and $x_1, \dots, x_n \in \mathbb{R}^p$, there exists a mapping $f : \mathbb{R}^p \rightarrow \mathbb{R}^k$, with $k = O(\epsilon^{-2} \log n)$, such that $\forall i, j$ $(1 - \epsilon)\|x_i - x_j\|^2 \leq \|f(x_i) - f(x_j)\|^2 \leq (1 + \epsilon)\|x_i - x_j\|^2$, we can transform high-dimensional data into a lower suitable dimensional space while approximately preserving the original distances between points. Our DGBR algorithm reduces the feature dimension, so that the population overlap assumption is more likely to be satisfied and we are less likely to see extreme values of balancing weights, so that better balance can be attained while maintaining low variance of the weights.
- 2) *Enables more accurate estimation of $p(y|s)$, because with multiple non-linear mapping functions in our DGBR algorithm, it can more easily capture the underlying non-linear relationship between stable features and response variables even with many stable features.*

4.2 Analysis on Upper Bound

4.2.1 Notation

Theorem 1 states that if $W_i^* = 1/P(\mathbf{X}_i = x)$, the global balancing regularizer in Eq. (4) converges to 0 as number of observations n goes to infinity. Nevertheless, in finite samples, Eq. (4) may not be equal to 0 for any W . We define the maximum covariate imbalance as

$$\alpha = \max_j \left\| \frac{\sum_{i: X_{i,k}=1, X_{i,j}=1} \hat{W}_i}{\sum_{i: X_{i,j}=1} \hat{W}_i} - \frac{\sum_{i: X_{i,k}=1, X_{i,j}=0} \hat{W}_i}{\sum_{i: X_{i,j}=0} \hat{W}_i} \right\|_{\infty}, \quad (9)$$

$$= \max_j \max_{k \neq j} \left| \frac{\sum_{i: X_{i,k}=1, X_{i,j}=1} \hat{W}_i}{\sum_{i: X_{i,j}=1} \hat{W}_i} - \frac{\sum_{i: X_{i,k}=1, X_{i,j}=0} \hat{W}_i}{\sum_{i: X_{i,j}=0} \hat{W}_i} \right|. \quad (10)$$

We define m as the number of values in \mathcal{X} that do not appear in \mathbf{X} ,

$$m = |\mathcal{X}| - |\mathcal{X}_{\mathbf{X}}| = 2^p - |\mathcal{X}_{\mathbf{X}}|,$$

where $\mathcal{X}_{\mathbf{X}} = \{x|\mathbf{X}_i = x, \text{ for some } i\}$, $|\mathcal{X}|$ and $|\mathcal{X}_{\mathbf{X}}|$ are the cardinalities of \mathcal{X} and $\mathcal{X}_{\mathbf{X}}$. m is random, where the variation of m comes from sampling a finite sample \mathbf{X} from the population distribution $P_{\mathbf{X}}$ on \mathcal{X} .

We define $\mathbb{E}[\alpha]$ as the expectation of α over the random sample \mathbf{X} .

4.2.2 Upper Bound of Global Balancing Regularizer

The maximum covariate imbalance α is different under different random finite sample \mathbf{X} . The following Lemma states that α is determined by \mathbf{X} through m , the number of values in \mathcal{X} that do not appear in \mathbf{X} , as well as the number of covariates p .

Lemma 2. *Given $\mathbf{X} \in \mathbb{R}^{n \times p}$, $p \geq 2$, if m different values in \mathcal{X} do not appear in \mathbf{X} , then we have*

- 1) if $m = 0$, then $\alpha = 0$
- 2) if $0 < m \leq 2^{p-2}$, then $\alpha = \frac{2^{p-2}}{2^{p-1}-m} - \frac{1}{2}$
- 3) if $2^{p-2} < m < 2^{p-1}$, then $\alpha = 1 - \frac{2^{p-1}-m}{3 \times 2^{p-2}-m}$
- 4) if $2^{p-1} \leq m \leq 2^p - 2^3$, then $\alpha = 1$

Proof. See Appendix C. \square

m measures how severe the overlap assumption is violated in the empirical distribution of \mathbf{X} . α increases with m for a fixed p . If m is fixed, when $0 < m \leq 2^{p-2}$, α is decreasing in p ; when $2^{p-2} < m < 2^{p-1}$, α is increasing in p ; ⁴ when $2^{p-1} \leq m \leq 2^p - 1$, α does not depend on p . A more realistic scenario is that the number of observations n is fixed, if p increases, the overlap assumption is harder to be satisfied in the empirical distribution, so m will also increase, which might result in an increase of α . Although n is not explicitly expressed in α , n affects m in that m converges to 0 as n goes to infinity, and therefore α may decrease with n through m .

Theorem 2 extends Lemma 2 and states $\mathbb{E}[\alpha]$, the expected value of α over the random sample \mathbf{X} . We show that $\mathbb{E}[\alpha]$ also equals to the expected value of α over m that is determined by \mathbf{X} .

Theorem 2. *Let $\alpha = \max_j \left\| \frac{\sum_{i: X_{i,k}=1, X_{i,j}=1} \hat{W}_i}{\sum_{i: X_{i,j}=1} \hat{W}_i} - \frac{\sum_{i: X_{i,k}=1, X_{i,j}=0} \hat{W}_i}{\sum_{i: X_{i,j}=0} \hat{W}_i} \right\|_{\infty}$, we have*

$$\mathbb{E}[\alpha] = \frac{1}{\binom{n+2^{p-1}}{2^{p-1}}} \sum_{m=0}^{2^p-1} \binom{2^p}{m} \binom{n-1}{2^{p-1}-m} g(p, m) \quad (11)$$

where $g(p, m)$ is

$$g(p, m) = \begin{cases} \frac{2^{p-2}}{2^{p-1}-m} - \frac{1}{2}, & \text{if } 0 \leq m \leq 2^{p-2} \\ 1 - \frac{2^{p-1}-m}{3 \times 2^{p-2}-m} & \text{if } 2^{p-2} < m < 2^{p-1} \\ 1 & \text{if } 2^{p-1} \leq m \leq 2^p - 1 \end{cases} \quad (12)$$

Proof. See Appendix D. \square

$\mathbb{E}[\alpha]$ depends on n and p . In particular, $\mathbb{E}[\alpha]$ is decreasing in n , but increasing in p , see Fig. 3.

4.2.3 Upper Bound of the Risk in Approximate Balancing

It is possible that our DGBR algorithm may not perfectly balance all covariates in \mathbf{X} . That is, the minimum value *zero* may not be attained by any W in Eq. (4). Denote W^* as an **approximate** solution to the global balancing problem if W^* is the optimal solution to our DGBR algorithm and Eq. (4) evaluated at W^* is not *zero*. Define the **imbalance** at x in \mathbf{X} balanced by W^* as the

3. Given p , $0 \leq m \leq 2^p - 1$. When $m = 2^p - 1$, either $\frac{\sum_{i: X_{i,k}=1, X_{i,j}=1} \hat{W}_i}{\sum_{i: X_{i,j}=1} \hat{W}_i}$ or $\frac{\sum_{i: X_{i,k}=1, X_{i,j}=0} \hat{W}_i}{\sum_{i: X_{i,j}=0} \hat{W}_i}$ is $\frac{0}{0}$, which is undefined. For simplicity and completeness, we define $\frac{0}{0} = 1$, and therefore $\alpha = 1$.

4. This scenario can never happen, because given m and p , if $2^{p-2} < m < 2^{p-1}$, when p increases, it will have $m < 2^{p-2}$.

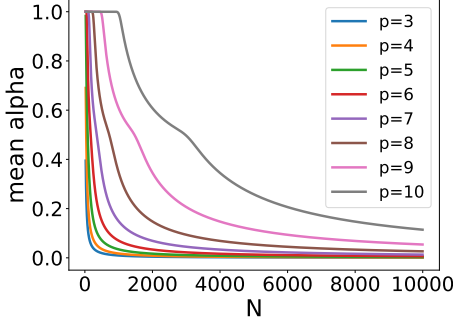


Fig. 3: $\mathbb{E}[\alpha]$ with different N and p

difference between the joint weighted probability and the product of weighted marginal probabilities, that is,

$$\epsilon_x = \tilde{p}_x - \prod_j \tilde{p}_{x_j} = \tilde{p}_x - p_x \quad (13)$$

where $\tilde{p}_x = \frac{1}{n} \sum_{i=1}^n W_i^* \mathbb{1}(\mathbf{X}_i = x)$, $\tilde{n} = \sum_{i=1}^n W_i^*$, $\tilde{p}_{x_j} = \frac{1}{\tilde{n}} \sum_{i=1}^n W_i^* \mathbb{1}(\mathbf{X}_{ij} = x_j)$ and $p_x = \prod_j \tilde{p}_{x_j}$. If Eq. (4) evaluated at W^* is not zero, ϵ_x must be nonzero at some x and ϵ_x measures how far covariates in weighted X are from independence.

In this subsection, we would like to evaluate our DGBR algorithm if optimal W^* in our DGBR algorithm is an approximate solution. That is, we would like to upper bound the expected risk of the optimal \hat{f} learned from our DGBR algorithm, where $f(\cdot)$ is a function to predict the response variable Y_i from covariates \mathbf{X}_i . The expected risk between $f(\mathbf{X})$ and Y is defined as $L_P(f) = E_P(l(f(\mathbf{X}_i), Y_i))$. In $L_P(f)$, the loss function $l(f(\mathbf{X}_i), Y_i) = \log(1 + \exp((1 - 2Y_i) \cdot (\phi(\mathbf{X}_i)\beta)))$ is the same as the objective function in our DGBR algorithm; the probability mass function $P(\mathbf{X}_i, Y_i) = P(\mathbf{X}_i)P(Y_i|\mathbf{X}_i)$ has $P(Y_i = y|\mathbf{X}_i = x) = P(Y_i = y|\mathbf{S}_i = s, \mathbf{V}_i = v) = P(y|s)$ to be the same as that in Assumption 1 and $P(\mathbf{X}_i = x) = p_x$, where p_x is defined in Eq. (13). In Empirical Risk Minimization (ERM), the population probability mass function is assumed to be fixed (but unknown). However, p_x is not fixed because it is determined by W^* and W^* is learned from our DGBR algorithm. To avoid the identification problem, we set W as fixed after we sequentially update β , W and θ for some steps; we denote this W as W^* and only update β and θ afterwards. This W^* is used to calculate p_x . We analyze the empirical risk of the following fixed weight DGBR (FWDGBR) algorithm

$$\begin{aligned} \min \quad & \frac{1}{\sum_{i=1}^n W_i^*} \sum_{i=1}^n W_i^* \cdot \log(1 + \exp((1 - 2Y_i) \cdot (\phi(\mathbf{X}_i)\beta))), \\ \text{s.t.} \quad & \|(W^* \cdot \mathbf{1}) \odot (X - \hat{X})\|_F^2 \leq \lambda_2, \\ & \|\beta\|_2^2 \leq \lambda_4, \|\beta\|_1 \leq \lambda_5, \\ & \sum_{k=1}^K (\|A^{(k)}\|_F^2 + \|\hat{A}^{(k)}\|_F^2) \leq \lambda_7, \|b^{(k)}\|_2 \leq M^{(k)}, \\ & \text{for } k = 1, 2, \dots, K. \end{aligned}$$

The empirical risk is defined as $\hat{L}(f) = \frac{1}{n} \sum_{i=1}^n W_i^* l(f(\mathbf{X}_i), Y_i)$ and $\tilde{n} = \sum_{i=1}^n W_i^*$. In the FWDGBR algorithm, it does not have the constraints related to W and has an additional bias constraint $\|b^{(k)}\|_2 \leq M^{(k)}$. The additional bias constraint guarantees that $f(\cdot)$ is bounded for any $f(\cdot)$ that satisfies the constraints. All the other constraints in the FWDGBR algorithm are the same as those in the DGBR algorithm.

In ERM, if empirical average converges to the expectation

(equivalent to exact balancing), the difference between the expectation of each function and the empirical average of the function ($|L_P(f) - \hat{L}(f)|$) can be bounded in terms of the Rademacher complexity of the model class and an error term depending on the confidence parameter and sample size. Note that the Rademacher complexity of the FWDGBR algorithm is determined by its constraints. In other words, the Rademacher complexity decreases with λ_2 , λ_4 , λ_5 and λ_7 . Furthermore, if empirical average converges to the expectation, $|L_P(\hat{f}) - \hat{L}(f^*)|$ can also be bounded in terms of the Rademacher complexity of the model class and an error term depending on the confidence parameter and sample size, where \hat{f} minimizes $\hat{L}(\cdot)$, f^* minimizes $L_P(\cdot)$ and both \hat{f} and f^* are in the model class. Note that f^* also needs to satisfy the constraints in the FWDGBR algorithm.

However, in approximate balancing, the empirical average $\hat{L}(f)$ for some f does not converge to expectation $L_P(f)$ because the joint distribution does not equal to the product of the marginals. In this case, $|L_P(\hat{f}) - \hat{L}(f^*)|$ can still be bounded, but with an extra term measuring covariates' imbalance, that is, ϵ_x .

Theorem 3. Let $B_k = \sqrt{\lambda_7 + (M^{(k)})^2}$, l_k be the size of $[\phi(\mathbf{X}_i)^{(k)}; 1]$ for $k = 1, 2, \dots, K$ and $l_0 = p + 1$. With probability at least $\geq 1 - \delta$,

$$\begin{aligned} L_P(\hat{f}) \leq & L_P(f^*) + 2^{K+3} \sqrt{\frac{2\log(2p)}{n}} \min(\sqrt{\lambda_4 l_K}, \lambda_5) \prod_{k=1}^K B_k (l_{k-1})^{1/2} \\ & + 3\sqrt{\frac{\log(2/\delta)}{2n}} + 2 \max_{x,f} \mathbb{E}[l(f(x), y)|x] \sum_x |\epsilon_x|, \end{aligned} \quad (14)$$

where $\hat{f} = \arg \min_f \hat{L}(f)$ and $f^* = \arg \min_f L_P(f)$.

Proof. See Appendix E. \square

Theorem 3 states that in approximate balancing, the upper bound for $L_P(\hat{f})$ consists of four components: 1. The minimum expected risk $L_P(f^*)$; 2. The Rademacher complexity, $2^{K+3} \sqrt{\frac{2\log(2p)}{n}} \min(\sqrt{\lambda_4 l_K}, \lambda_5) \prod_{k=1}^K B_k (l_{k-1})^{1/2}$, which depends on number of layers, number of hidden units and the constraints in the FWDGBR algorithm; 3. An error term depending on the confidence parameter δ and sample size n , that is, $3\sqrt{\frac{\log(2/\delta)}{2n}}$; 4. The risk from covariates' imbalance in approximate balancing, $2 \max_{x,f} \mathbb{E}[l(f(x), y)|x] \sum_x |\epsilon_x|$. Note that $\max_{x,f} \mathbb{E}[l(f(x), y)|x]$ is bounded because x and y are binary and all weights in $f(\cdot)$ are bounded.

The term related to Rademacher complexity in Theorem 3, $RC = 2^{K+3} \sqrt{\frac{2\log(2p)}{n}} \min(\sqrt{\lambda_4 l_K}, \lambda_5) \prod_{k=1}^K B_k (l_{k-1})^{1/2}$, is an upper bound of the size of the model class in the FWDGBR algorithm. This is derived from the upper bound of the size of the model class in the reduced fixed weight DGBR (RFWDGBR) algorithm

$$\begin{aligned} \min \quad & \frac{1}{\sum_{i=1}^n W_i^*} \sum_{i=1}^n W_i^* \cdot \log(1 + \exp((1 - 2Y_i) \cdot (\phi(\mathbf{X}_i)\beta))), \\ \text{s.t.} \quad & \|\beta\|_2^2 \leq \lambda_4, \|\beta\|_1 \leq \lambda_5 \\ & \|A^{(k)}\|_F^2 \leq \lambda_7, \text{ for } k = 1, 2, \dots, K. \end{aligned}$$

This algorithm does not have the auto-encoder and -decoder constraint $\|(W^* \cdot \mathbf{1}) \odot (X - \hat{X})\|_F^2 \leq \lambda_2$. The weight constraints of the auto-encoder and -decoder are relaxed. The model class described by the RFWDGBR algorithm is larger than the model class of the fixed weight DGBR algorithm. Thus, the upper bound of the Rademacher complexity of the RFWDGBR algorithm is also an upper bound of the Rademacher complexity of the FWDGBR algorithm.

The derivation of the term RC is closely related to deriving the Rademacher complexity with dropouts in Neural Networks in [26] and [27]. Compare with [26] and [27], the RFWDGBR algorithm does not have dropouts, but it has both ℓ_1 and ℓ_2 constraints. The model class with both ℓ_1 and ℓ_2 constraints is no larger than the minimum of the model class with either ℓ_1 or ℓ_2 constraint, which explains the component, $\min(\sqrt{\lambda_4 l_K}, \lambda_5)$, in the term RC.

Moreover, the upper bound for $L_P(\hat{f})$ depends on how close the approximate balancing is to the exact balancing, measured by ϵ_x . If the approximate balancing gets very close to the exact balancing, the term $2 \max_{x,f} \mathbb{E}[l(f(x), y)|x] \sum_x |\epsilon_x|$ approaches 0. The upper bound for $L_P(\hat{f})$ will mainly depend on three other terms in the right-hand side of Inequality (14).

From this upper bound of $L_P(\hat{f})$, we have two main conclusions: 1. If the feasible set of the auto-encoder and the weights in the logistic regression is smaller, implying simpler models, the upper bound for $L_P(\hat{f})$ is smaller; 2. If covariates' imbalance is smaller, the upper bound for $L_P(\hat{f})$ is smaller.

5 OPTIMIZATION AND DISCUSSION

5.1 Optimization

To optimize the aforementioned DGBR model in Eq. (7), we need to minimize following \mathcal{L}_{mix} as a function of parameter W , β , and $\theta = \{\mathbf{A}^{(k)}, \hat{\mathbf{A}}^{(k)}, b^{(k)}, \hat{b}^{(k)}\}$.

$$\mathcal{L}_{mix} = \mathcal{L}_{Pre} + \lambda_1 \mathcal{L}_{Bal} + \lambda_2 \mathcal{L}_{AE} + \mathcal{L}_{Reg}, \quad (15)$$

s.t. $W \succeq 0$,

where

$$\mathcal{L}_{Pre} = \sum_{i=1}^n W_i \cdot \log(1 + \exp((1 - 2Y_i) \cdot (\phi(\mathbf{X}_i) \cdot \beta))), \quad (16)$$

$$\mathcal{L}_{Bal} = \sum_{j=1}^p \left\| \frac{\phi(\mathbf{X}_{:, -j})^T \cdot (W \odot \mathbf{X}_{:, j})}{W^T \cdot \mathbf{X}_{:, j}} - \frac{\phi(\mathbf{X}_{:, -j})^T \cdot (W \odot (1 - \mathbf{X}_{:, j}))}{W^T \cdot (1 - \mathbf{X}_{:, j})} \right\|_2^2, \quad (17)$$

$$\mathcal{L}_{AE} = \|(W \cdot \mathbf{1}) \odot (X - \hat{X})\|_F^2, \quad (18)$$

$$\mathcal{L}_{Reg} = \lambda_3 \|W\|_2^2 + \lambda_4 \|\beta\|_2^2 + \lambda_5 \|\beta\|_1 + \lambda_6 (\sum_{i=1}^n W_i - 1)^2 + \lambda_7 \sum_{k=1}^K (\|\mathbf{A}^{(k)}\|_F^2 + \|\hat{\mathbf{A}}^{(k)}\|_F^2). \quad (19)$$

Here, we propose an iterative method to minimize the above objective function in Eq. (15). Starting from some random initialization on parameters W , β and θ , we update each of them alternatively with the other two parameters as fixed at each iteration until convergence. These steps are described below:

Update β : When fixing W and θ , the problem (15) is equivalent to optimize following objective function:

$$\mathcal{L}_{mix}(\beta) = \sum_{i=1}^n W_i \cdot \log(1 + \exp((1 - 2Y_i) \cdot (\phi(\mathbf{X}_i) \cdot \beta))) + \lambda_4 \|\beta\|_2^2 + \lambda_5 \|\beta\|_1, \quad (20)$$

which is a standard ℓ_1 norm regularized least squares problem and can be easily solved by any LASSO (or elastic net) solver.

Update W : By fixing β and θ , the key step for updating W is to calculate the partial derivative of $\frac{\partial \mathcal{L}_{mix}}{\partial W}$. The detailed mathematical form of the partial derivative is shown as following:

$$\frac{\partial \mathcal{L}_{mix}}{\partial W} = \frac{\partial \mathcal{L}_{Pre}}{\partial W} + \frac{\partial \mathcal{L}_{Bal}}{\partial W} + \frac{\partial \mathcal{L}_{AE}}{\partial W} + \frac{\partial \mathcal{L}_{Reg}}{\partial W}, \quad (21)$$

where

$$\frac{\partial \mathcal{L}_{Pre}}{\partial W} = \log(1 + \exp((1 - 2Y) \cdot (\phi(\mathbf{X}) \cdot \beta))), \quad (22)$$

$$\frac{\partial \mathcal{L}_{Bal}}{\partial W} = 2 \sum_{j=1}^p \mathcal{L}_{Bal_j} \cdot \frac{\partial \mathcal{L}_{Bal_j}}{\partial W}, \quad (23)$$

$$\frac{\partial \mathcal{L}_{AE}}{\partial W} = 2((W \cdot \mathbf{1}) \odot (\mathbf{X} - \hat{\mathbf{X}})) \odot (\mathbf{X} - \hat{\mathbf{X}}) \cdot \mathbf{1}^T, \quad (24)$$

$$\frac{\partial \mathcal{L}_{Reg}}{\partial W} = 2\lambda_3 W. \quad (25)$$

where $\mathcal{L}_{Bal_j} = \frac{\phi(\mathbf{X}_{:, -j})^T \cdot (W \odot \mathbf{X}_{:, j})}{W^T \cdot \mathbf{X}_{:, j}} - \frac{\phi(\mathbf{X}_{:, -j})^T \cdot (W \odot (1 - \mathbf{X}_{:, j}))}{W^T \cdot (1 - \mathbf{X}_{:, j})}$ and

$$\begin{aligned} \frac{\partial \mathcal{L}_{Bal_j}}{\partial W} &= \frac{\phi(\mathbf{X}_{:, -j})^T \odot (\mathbf{X}_{:, j} \cdot \mathbf{1}^T) \cdot (W^T \cdot \mathbf{X}_{:, j})}{(W^T \cdot \mathbf{X}_{:, j})^2} \\ &\quad - \frac{\phi(\mathbf{X}_{:, -j})^T \odot ((1 - \mathbf{X}_{:, j}) \cdot \mathbf{1}^T) \cdot (W^T \cdot (1 - \mathbf{X}_{:, j}))}{(W^T \cdot (1 - \mathbf{X}_{:, j}))^2}. \end{aligned}$$

For ensuring the non-negative of W with constraint $W \succeq 0$, we let $W = \omega \odot \omega$, where $\omega \in \mathbb{R}^{n \times 1}$. Then we update W by updating ω with following partial derivative.

$$\frac{\partial \mathcal{L}_{mix}}{\partial \omega} = \frac{\partial \mathcal{L}_{mix}}{\partial W} \cdot \frac{\partial W}{\partial \omega} \quad (26)$$

Update θ : By fixing W and β , the key step for updating θ is to calculate the partial derivative of $\frac{\partial \mathcal{L}_{mix}}{\partial A^{(k)}}$ and $\frac{\partial \mathcal{L}_{mix}}{\partial \hat{A}^{(k)}}$. The detailed mathematical form of the partial derivative is shown as following:

$$\frac{\partial \mathcal{L}_{mix}}{\partial A^{(k)}} = \frac{\partial \mathcal{L}_{Pre}}{\partial A^{(k)}} + \lambda_1 \frac{\partial \mathcal{L}_{Bal}}{\partial A^{(k)}} + \lambda_2 \frac{\partial \mathcal{L}_{AE}}{\partial A^{(k)}} + \frac{\partial \mathcal{L}_{Reg}}{\partial A^{(k)}}, \quad (27)$$

$$\frac{\partial \mathcal{L}_{mix}}{\partial \hat{A}^{(k)}} = \lambda_2 \frac{\partial \mathcal{L}_{AE}}{\partial \hat{A}^{(k)}} + \frac{\partial \mathcal{L}_{Reg}}{\partial \hat{A}^{(k)}}, \quad k = 1, \dots, K \quad (28)$$

First we look at the first term $\frac{\partial \mathcal{L}_{Pre}}{\partial A^{(k)}}$ in $\frac{\partial \mathcal{L}_{mix}}{\partial A^{(k)}}$, which can be rephrased as follows:

$$\frac{\partial \mathcal{L}_{Pre}}{\partial A^{(k)}} = \frac{\partial \mathcal{L}_{Pre}}{\partial \phi(X)} \cdot \frac{\partial \phi(X)}{\partial A^{(k)}}, \quad (29)$$

According to Eq. 16, we can obtain $\frac{\partial \mathcal{L}_{Pre}}{\partial \phi(X)}$. The calculation of the second term $\frac{\partial \phi(X)}{\partial A^{(k)}}$ is easy since $\phi(X) = \sigma(\phi(X)^{(K-1)} A^{(K)} + b^{(K)})$. Then $\frac{\partial \mathcal{L}_{Pre}}{\partial A^{(K)}}$ is accessible. Based on the back-propagation, we can iteratively obtain $\frac{\partial \mathcal{L}_{Pre}}{\partial A^{(k)}}$, $k = 1, \dots, K - 1$. Now the calculation of the partial derivative of \mathcal{L}_{Pre} is finished.

Similarly, by using back-propagation we can finish the calculation of \mathcal{L}_{Bal} and \mathcal{L}_{AE} . Finally we can obtain the $\frac{\partial \mathcal{L}_{mix}}{\partial A^{(k)}}$ and $\frac{\partial \mathcal{L}_{mix}}{\partial \hat{A}^{(k)}}$ for $k = 1, \dots, K$, and update our parameter θ .

We update W , β and $\theta = \{\mathbf{A}^{(k)}, \hat{\mathbf{A}}^{(k)}, b^{(k)}, \hat{b}^{(k)}\}$ iteratively until the objective function converges. The whole algorithm is summarized in Algorithm 1.

Finally, with the optimized regression coefficient β and deep auto-encoder parameters θ by our DGBR algorithm, we can make a stable prediction on various agnostic test datasets.

5.2 Complexity Analysis

During the procedure of optimization, the main time cost is to calculate the loss function, update parameters W , β and θ . For calculating the loss function, its complexity is $O(npd)$, where n is the sample size, p is the dimension of observed variables and d is the maximum dimension of the hidden layer in deep auto-encoder model. For updating parameter W , its complexity is dominated by the step of calculating the partial gradients of loss function with respect to variable W . Its complexity is also $O(npd)$. For

Algorithm 1 Deep Global Balancing Regression algorithm

Input: Observed Variables Matrix \mathbf{X} and Response Variable Y .

Output: Updated Parameters W, β, θ .

- 1: Initialize parameters $W^{(0)}, \beta^{(0)}$ and $\theta^{(0)}$,
 - 2: Calculate the current value of $\mathcal{L}_{mix}^{(0)} = \mathcal{J}(W^{(0)}, \beta^{(0)}, \theta^{(0)})$ with Equation (15),
 - 3: Initialize the iteration variable $t \leftarrow 0$,
 - 4: **repeat**
 - 5: $t \leftarrow t + 1$,
 - 6: Update $W^{(t)}$ based on Eq. (21),
 - 7: Update $\beta^{(t)}$ by solving $\mathcal{L}_{mix}(\beta^{(t-1)})$ in Equation (20),
 - 8: Based on Eq. (27) and (28), use $\frac{\partial \mathcal{L}_{mix}}{\partial \theta}$ to back-propagate through the entire deep network to get updated parameters θ ,
 - 9: Calculate $\mathcal{L}_{mix}^{(t)} = \mathcal{J}(W^{(t)}, \beta^{(t)}, \theta^{(t)})$,
 - 10: **until** $\mathcal{L}_{mix}^{(t)}$ converges or max iteration is reached.
 - 11: **return** W, β, θ .
-

updating parameter β , it is a standard LASSO problem and its complexity is $O(nd)$. For updating θ , its complexity is $O(npd)$.

In total, the complexity of each iteration in Algorithm 1 is $O(npd)$.

5.3 Parameter Tuning

To tune the parameters for our algorithm and baselines, we need multiple validation datasets whose distributions are diverse from each other and different with the training data. In our experiments, we generate such validation datasets \mathcal{E} by non-random data resampling on training data. We calculate the *Average_Error* and *Stability_Error* of all algorithms on validation datasets by choosing *RMSE* as *Error* metrics in Eq. (1) and (2). In this paper, we tune all the parameters for our algorithm and baselines by minimizing $Average_Error + \lambda \cdot Stability_Error$ on validation datasets with cross validation by grid searching. We set $\lambda = 5$ in our experiments.

Construction of Validation Data. The key point in construction of validation data is to construct datasets where the joint distribution of the covariates changes across environments, particularly when this might create bias if we don't control for all of the stable features. However, we do not have prior knowledge about which features are noisy features. Fortunately, our estimation approach can identify noisy features as those that do not have a large estimated effect after balancing. Using the empirically identified noisy features, we can generate validation datasets that change the distribution of noisy features and use these for parameter tuning.

6 EXPERIMENTS

In this section, we evaluate our algorithm on both synthetic and real world dataset, comparing with the state-of-the-art methods.

6.1 Baselines

We implement following baselines for comparison.

- *Logistic Regression (LR)* [28]
- *Deep Logistic Regression (DLR)* [29]: Combines a deep auto-encoder and logistic regression.
- *Global Balancing Regression (GBR)*: Combines a global balancing regularizer and logistic regression as shown in Eq (5).

Since our proposed algorithm is based on logistic regression, so we compare our algorithm with only logistic regression methods. For

other predictive methods, we can propose corresponding global balancing algorithm based on them, and compare with them.

6.2 Experiments on Synthetic Data

In this section, we describe the synthetic datasets and demonstrate the effectiveness of our proposed algorithm.

6.2.1 Dataset

As shown in Fig. 1, there are three relationships between $\mathbf{X} = \{\mathbf{S}, \mathbf{V}\}$ and Y , including $\mathbf{S} \perp \mathbf{V}$, $\mathbf{S} \rightarrow \mathbf{V}$, and $\mathbf{V} \rightarrow \mathbf{S}$.

$\mathbf{S} \perp \mathbf{V}$: In this setting, \mathbf{S} and \mathbf{V} are independent. Recalling Fig. 1, we generate predictor $\mathbf{X} = \{\mathbf{S}_{\cdot,1}, \dots, \mathbf{S}_{\cdot,p_s}, \mathbf{V}_{\cdot,1}, \dots, \mathbf{V}_{\cdot,p_v}\}$ with independent Gaussian distributions as:

$$\tilde{\mathbf{S}}_{\cdot,1}, \dots, \tilde{\mathbf{S}}_{\cdot,p_s}, \tilde{\mathbf{V}}_{\cdot,1}, \dots, \tilde{\mathbf{V}}_{\cdot,p_v} \stackrel{iid}{\sim} \mathcal{N}(0, 1),$$

where $p_s = 0.4 * p$, and $\mathbf{S}_{\cdot,j}$ represents the j^{th} variable in \mathbf{S} . To make \mathbf{X} binary, we let $\mathbf{X}_{\cdot,j} = 1$ if $\tilde{\mathbf{X}}_{\cdot,j} \geq 0$, otherwise $\mathbf{X}_{\cdot,j} = 0$. **$\mathbf{S} \rightarrow \mathbf{V}$:** In this setting, the stable features \mathbf{S} are the causes of noisy features \mathbf{V} . We first generate the stable features $\tilde{\mathbf{S}}$ with independent Gaussian distributions, and let $\mathbf{S}_{\cdot,j} = 1$ if $\tilde{\mathbf{S}}_{\cdot,j} \geq 0$, otherwise $\mathbf{S}_{\cdot,j} = 0$. Then, we generate noisy features $\tilde{\mathbf{V}} = \{\tilde{\mathbf{V}}_{\cdot,1}, \dots, \tilde{\mathbf{V}}_{\cdot,p_v}\}$ based on $\tilde{\mathbf{S}}$:

$$\tilde{\mathbf{V}}_{\cdot,j} = \tilde{\mathbf{S}}_{\cdot,j} + \tilde{\mathbf{S}}_{\cdot,j+1} + \mathcal{N}(0, 2),$$

and let $\mathbf{V}_{\cdot,j} = 1$ if $\tilde{\mathbf{V}}_{\cdot,j} > 1$, otherwise $\mathbf{V}_{\cdot,j} = 0$.

$\mathbf{V} \rightarrow \mathbf{S}$: In this setting, the noisy features \mathbf{V} are the causes of stable features \mathbf{S} . We first generate the noisy features $\tilde{\mathbf{V}}$ with independent Gaussian distribution, and let $\mathbf{V}_{\cdot,j} = 1$ if $\tilde{\mathbf{V}}_{\cdot,j} \geq 0$, otherwise $\mathbf{V}_{\cdot,j} = 0$. Then, we generate stable features $\tilde{\mathbf{S}} = \{\tilde{\mathbf{S}}_{\cdot,1}, \dots, \tilde{\mathbf{S}}_{\cdot,p_s}\}$ based on $\tilde{\mathbf{V}}$:

$$\tilde{\mathbf{S}}_{\cdot,j} = \tilde{\mathbf{V}}_{\cdot,j} + \tilde{\mathbf{V}}_{\cdot,j+1} + \mathcal{N}(0, 2),$$

and let $\mathbf{S}_{\cdot,j} = 1$ if $\tilde{\mathbf{S}}_{\cdot,j} > 1$, otherwise $\mathbf{S}_{\cdot,j} = 0$.

Finally, we generate the response variable Y for all above three settings with the same function g as following:

$$Y = 1 / (1 + \exp(-\sum_{\mathbf{X}_{\cdot,i} \in \mathbf{S}_l} \alpha_i \cdot \mathbf{X}_{\cdot,i} - \sum_{\mathbf{X}_{\cdot,j} \in \mathbf{S}_n} \beta_j \cdot \mathbf{X}_{\cdot,j} \cdot \mathbf{X}_{\cdot,j+1})) + \mathcal{N}(0, 0.2),$$

where we separate the stable features \mathbf{S} into two parts, linear part \mathbf{S}_l and non-linear part \mathbf{S}_n . And $\alpha_i = (-1)^i \cdot (i \% 3 + 1) \cdot p/3$ and $\beta_j = p/2$. To make Y binary, we set $Y = 1$ when $Y \geq 0.5$, otherwise $Y = 0$.

To test the stability of all algorithms, we need to generate a set environments e , each with a distinct joint distribution. Under Assumption 1, instability in prediction arises because $P(Y|\mathbf{V})$ or $P(\mathbf{V}|\mathbf{S})$ varies across environments. Therefore, we generate different environments in our experiments by varying $P(Y|\mathbf{V})$ and $P(\mathbf{V}|\mathbf{S})$.

6.2.2 Experiments by Varying $P(Y|\mathbf{V})$ and Results

Specifically, we vary $P(Y|\mathbf{V})$ via biased sample selection with a bias rate $r \in (0, 1)$. For each sample, we select it with probability r if its noisy features equal to response variable, that is $\mathbf{V} = Y$; otherwise we select it with probability $1 - r$, where $r > .5$ corresponds to positive correlation between Y and \mathbf{V} . From the generation of Y , we know that, given stable features \mathbf{S} , noisy features \mathbf{V} are independent of Y . But after biased sample selection, \mathbf{V} could be correlated with response variable Y conditional on \mathbf{S} due to selection bias. However, since \mathbf{S} is

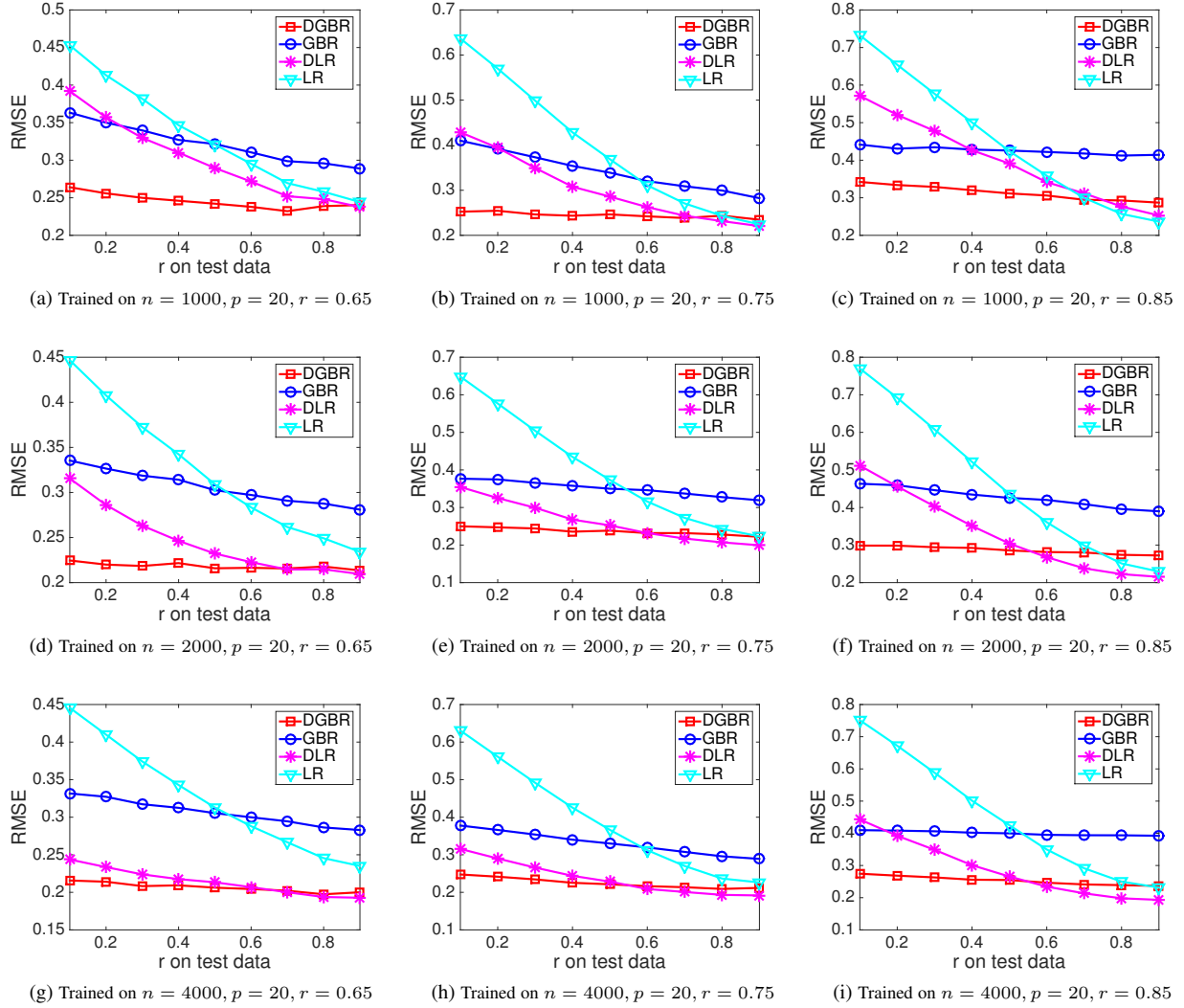


Fig. 4: **Setting $\mathbf{S} \perp \mathbf{V}$** : RMSE of outcome prediction on various test datasets by varying sample size n (vertical) and bias rate r (horizontal) on training dataset. The r of the X-axis in each figure represents the bias rate on test data.

an important factor in determining Y and thus whether a unit is selected when its noisy features are high, controlling for \mathbf{S} when estimating the correlation between Y and \mathbf{V} reduces that correlation. Note that this data-generating environment violates Assumption 1 after the sample selection is introduced. This creates an environment that is challenging for our algorithm, since we do not have strict guarantees that it will work, and also serves to illustrate that even when the strong assumptions of theory fail, the algorithm can still lead to substantial improvements.

To comprehensively and systematically evaluate stability of predictive models, we generate different synthetic data by varying sample size $n = \{1000, 2000, 4000\}$, dimensions of variables $p = \{20, 40, 80\}$, and bias rate $r = \{0.65, 0.75, 0.85\}$. We report the results of setting $\mathbf{S} \perp \mathbf{V}$ in Figure 4 & 5, and report the results of setting $\mathbf{S} \rightarrow \mathbf{V}$ and $\mathbf{V} \rightarrow \mathbf{S}$ in Figure 6 and 7, respectively.

From the results, we have following observations and analysis:

- The methods LR and DLR can not address the stable prediction problem in all settings. Since they can not remove the spurious correlation between noisy features and the response variable during model training, they often predict

large effects of the noisy features, which leads to instability across environments.

- Comparing with baselines, our method achieves a more stable prediction in different settings. The GBR method is more stable than LR, and our DGBR algorithm is more stable than DLR. The main reason is that the global balancing regularizer used in our models helps ensure accurate estimation of the effect of the stable features, and reduces the estimates of the effect of the noisy features.
- Our DGBR model makes a more precise and stable prediction than GBR model across environments. The deep embedding model in DGBR algorithm makes global balancing weights less noisy and simplifies estimates of the effect of stable features.
- By varying the sample size n , dimension of variables p and training bias rate r , the RMSE of our DGBR algorithm is consistently stable and small across environments. Another important observation is that comparing with baselines, our algorithm makes more and more significantly improvement on prediction performance when n is small relative to p and r .

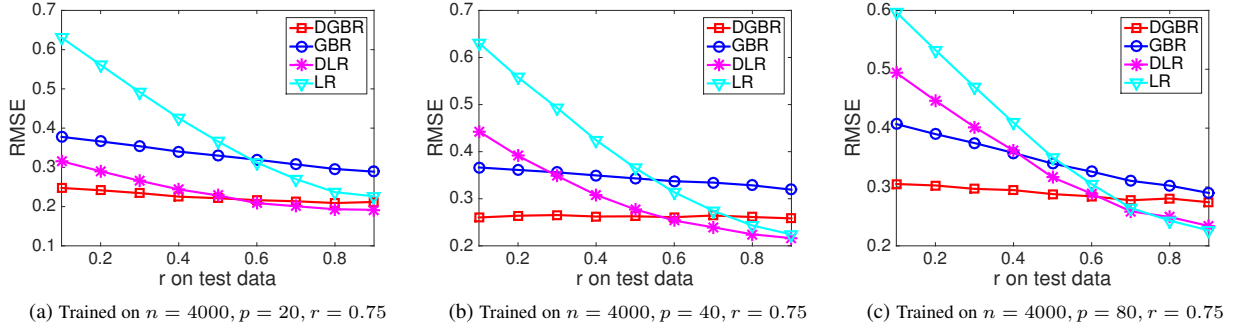


Fig. 5: **Setting $S \perp V$** : RMSE of outcome prediction on various test datasets by varying variables' dimension p on training dataset.

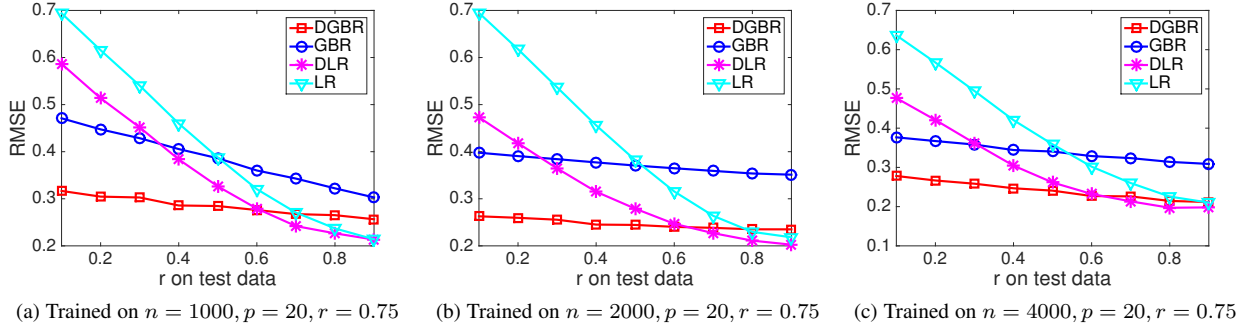


Fig. 6: **Setting $S \rightarrow V$** : RMSE of outcome prediction on various testing datasets by varying sample size n on training dataset.

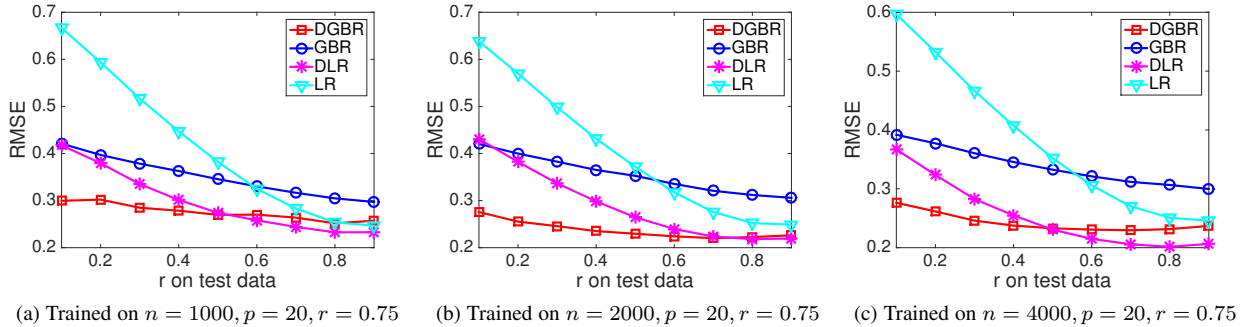


Fig. 7: **Setting $V \rightarrow S$** : RMSE of outcome prediction on various testing datasets by varying sample size n on training dataset.

6.2.3 Experiments by Varying $P(V|S)$ and Results

In this sub experiment, we generate the response variable Y with the function g as following:

$$Y = 1/(1 + \exp(-\sum_{\mathbf{x}_{:,i} \in \mathbf{S}_l} \alpha_i \cdot \mathbf{x}_{:,i} - \sum_{\mathbf{x}_{:,j} \in \mathbf{S}_n} \beta_j \cdot \mathbf{x}_{:,j} \cdot \mathbf{x}_{:,j+1})),$$

where $\alpha_i = (-1)^i$ and $\beta_j = p/2$. And to make Y binary, we set $Y = 1$ when $Y \geq 0.5$, otherwise $Y = 0$.

Here, we vary $P(V|S)$ also via biased sample selection with a bias rate $r \in (0, 1)$. Specifically, for each sample, we select it with probability r if its noisy features equal to a mediate variable \mathbf{Z} , that is $\mathbf{V}_i = \mathbf{Z}_i$; otherwise we select it with probability $1 - r$, where $\mathbf{Z}_i = \sum_{j=i}^{i+5} (-1)^j \cdot \mathbf{S}_j$, and $r > .5$ corresponds to positive correlation between \mathbf{Z}_i and \mathbf{V}_i . The same, from the generation of Y , we know that, given stable features \mathbf{S} , noisy features \mathbf{V} are independent of Y . After biased sample selection,

\mathbf{V} could be highly correlated with response variable Y , but it is still independent with Y conditional on stable features \mathbf{S} . Thus, the Assumption 1 is valid under this setting. Therefore, with identifying stable features \mathbf{S} , our algorithm can make a stable prediction across environments.

In this part experiments, we generate different synthetic data by varying sample size $n = \{1000, 2000, 4000\}$. We report the experimental results under settings $S \perp V$, $S \rightarrow V$, and $V \rightarrow S$ in Figure 8, 9 & 10, respectively. From these results, we can obtain the same observations that (i) The traditional classification methods LR and DLR can not address the stable prediction problem in all settings, (ii) Comparing with baselines, our method achieves a more stable prediction in different settings. The GBR method is more stable than LR, and our DGBR algorithm is more stable than DLR, and (iii) Our DGBR model makes a more precise and stable prediction than GBR model across environments.

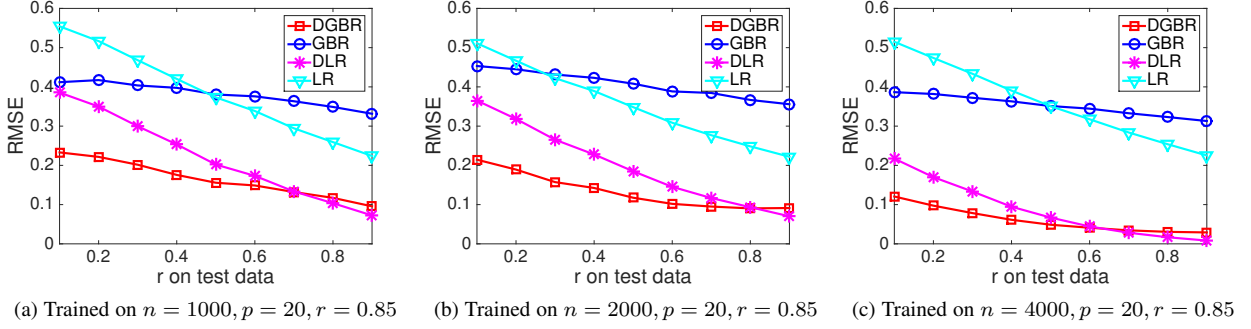


Fig. 8: **Setting $S \perp V$** : RMSE of outcome prediction on various testing datasets by varying sample size n on training dataset.

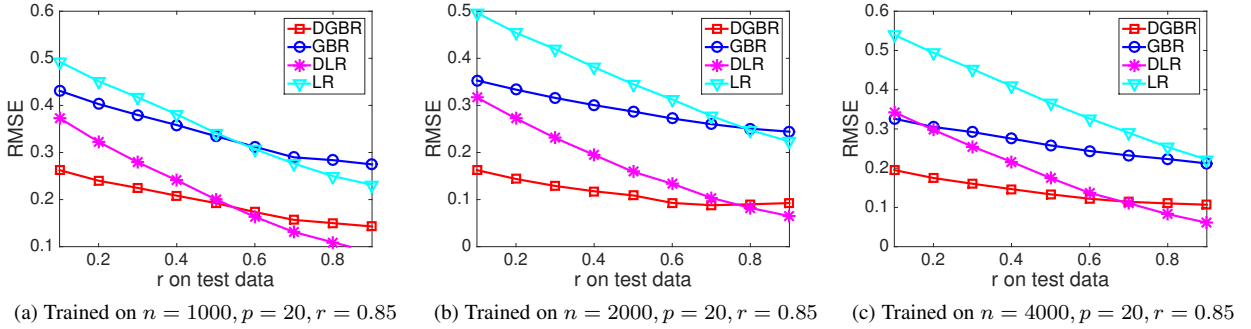


Fig. 9: **Setting $S \rightarrow V$** : RMSE of outcome prediction on various testing datasets by varying sample size n on training dataset.

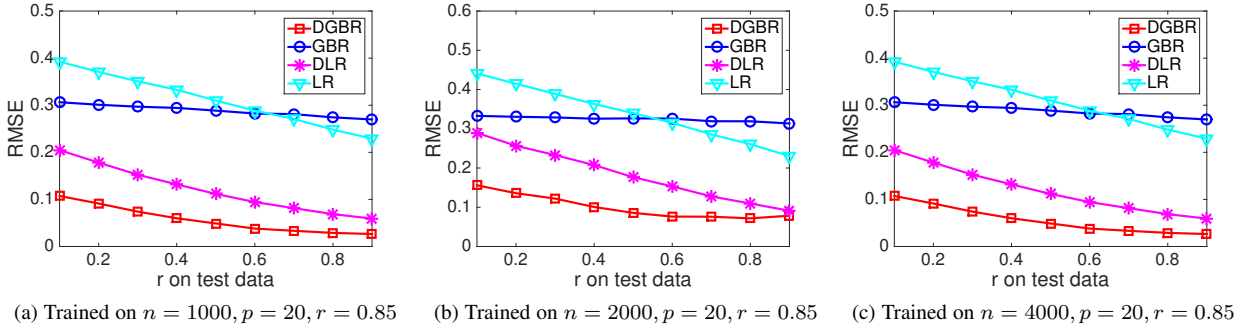


Fig. 10: **Setting $V \rightarrow S$** : RMSE of outcome prediction on various testing datasets by varying sample size n on training dataset.

6.2.4 Visualization of Embedded Features

In Figure 12, we also show that the embedded features in our DGBR algorithm have few information of noisy features V from raw space. This demonstrates that our DGBR could approximately preserve the independence between Y and V of global balancing, thus can identify stable features and make a stable prediction across unknown environments.

6.2.5 Parameter Analysis

In our DGBR algorithm, we have some hyper-parameters, such as λ_1 for constraining the error of global balancing, λ_2 constraining the loss of auto-encoder term, λ_3 constraining the variance of the global sample weights, and so on. In this section, we investigate how these hyper-parameters affect the results. We tuned these parameters in our experiments with cross validation by grid searching, based on our constructed validation data.

We report the *Average_Error*, $5 * \text{Stability_Error}$, and $\text{Average_Error} + 5 * \text{Stability_Error}$ on a synthetic dataset under setting $S \perp V$ with $n = 2000$ and $p = 20$.

Tradeoffs between prediction and covariate balancing: We first show how the hyper-parameter λ_1 affects the performance in Figure 11a. The parameter of λ_1 restrain the error of global balancing. We can see that initially the value of both *Average_Error* and *Stability_Error* decreases when the value of λ_1 increases. This is intuitive as the data could be more balanced with the increased value of λ_1 , and balanced data could help to identify stable features and remove some noise for more precise prediction. However, when the value of λ_1 increases further, the value of *Stability_Error* decreases, but the value of *Average_Error* starts to increase slowly. Large value of λ_1 makes the algorithm concentrate on global balancing component at the expense of the prediction component. Both prediction and global balancing

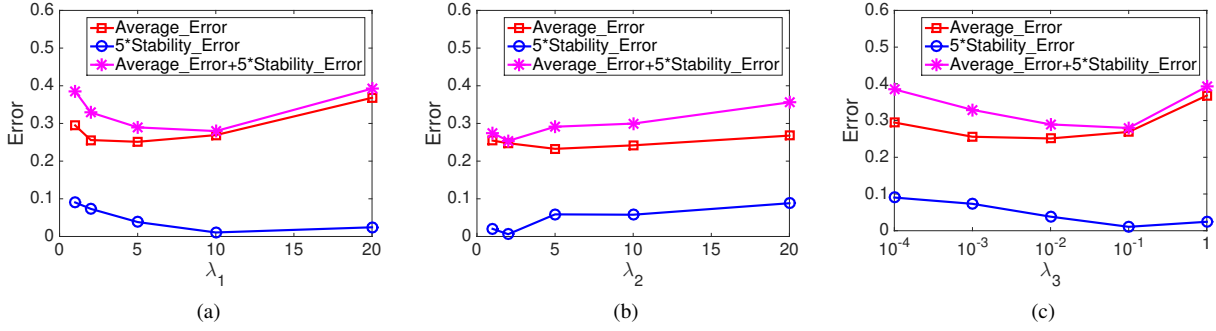


Fig. 11: The effect of hyper-parameters λ_1 , λ_2 , and λ_3 .

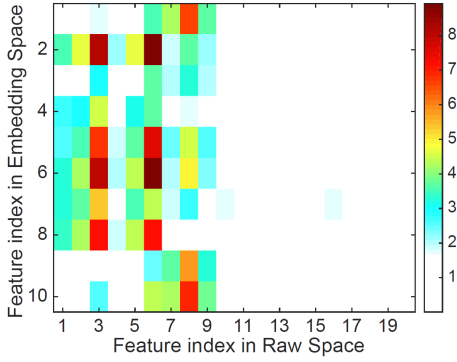


Fig. 12: Embedding weights in our DGBR algorithm, where $X_{.,1}, \dots, X_{.,9}$ are stable features \mathbf{S} and others are noisy features \mathbf{V} . It illustrates that our DGBR can achieve \mathbf{Y} and \mathbf{V} are independent, since the features in embedding space have few information of noisy feature \mathbf{V} in raw sapce.

components are essential for stable prediction.

Feature representation: Here, we show how the hyper-parameter λ_2 affects the results in Figure 11b. The value of *Average_Error* decreases with λ_2 , since a high value of λ_2 leads to more accurate prediction. Initially, *Stability_Error* decreases with λ_2 , but it starts to increase when $\lambda_2 \geq 5$. It is important to choose an appropriate value of λ_2 for learning feature representation, but our method is not very sensitive to this parameter.

The variance of global sample weights: Figure 11c shows how the value of λ_3 affect performance. Both the value of *Average_Error* and *Stability_Error* decrease when the value of λ_3 increases, since appropriate constraints on the variance of global sample weights could prevent some samples from becoming dominate in whole data, and thus help to improve the precision and robustness of prediction. However, when the value of λ_3 grows too large, those errors increase. Too large value of λ_3 could lead the learned global sample weight to fail to make appropriate tradeoffs between balancing and prediction.

6.3 Experiments on Real World Data

6.3.1 Online Advertising Dataset

The real online advertising dataset we used is collected from Tencent WeChat App⁵ during September 2015. In WeChat, each

5. <http://www.wechat.com/en/>

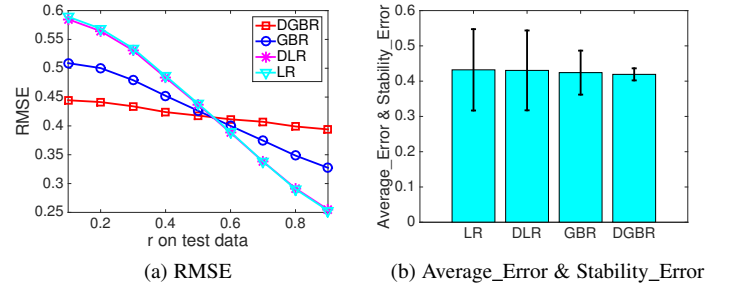


Fig. 13: Our proposed DGBR algorithm makes the most stable prediction on whether user will like or dislike an advertisement.

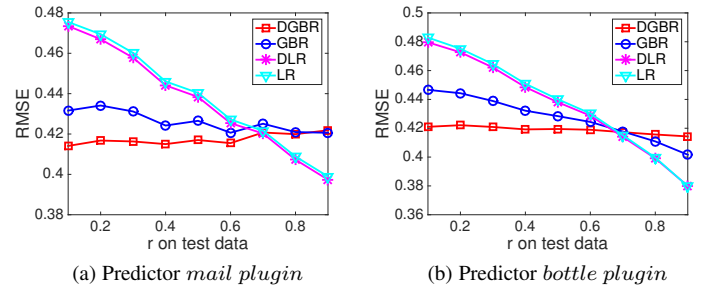


Fig. 14: RMSE of outcome prediction by varying bias rate r between one predictor and outcome.

user can share (receive) posts to (from) his/her friends as like the Twitter and Facebook. Then the advertisers could push their advertisements to users, by merging them into the list of the user's wallposts. For each advertisement, there are two types of feedbacks: "Like" and "Dislike". When the user clicks the "Like" button, his/her friends will receive the advertisements with this action.

The online advertising campaign used in our paper is about the LONGCHAMP handbags for young women.⁶ This campaign contains 14,891 user feedbacks with Like and 93,108 Dislikes. For each user, we have their features including (1) demographic attributes, such as age, gender, (2) number of friends, (3) device (iOS or Android), and (4) the user settings on WeChat, for example, whether allowing strangers to see his/her album and

6. <http://en.longchamp.com/en/womens-bags>

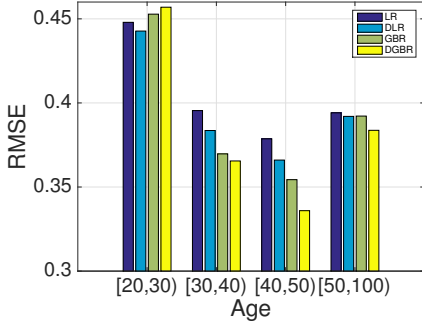


Fig. 15: Prediction across environments separated by age. The models are trained on dataset where uses' $Age \in [20, 30)$, but tested on various datasets with different users' age range.

whether installing the online payment service.

Experimental Settings. In our experiments, we set $Y_i = 1$ when user i likes the ad, otherwise $Y_i = 0$. For non-binary user features, we dichotomize them around their mean value. Considering the overlap assumption in assumption 2, we only preserve users' features which satisfied $0.2 \leq \frac{\#\{x=1\}}{\#\{x=1\} + \#\{x=0\}} \leq 0.8$. All the predictors and response variable in our experiment are binary.

In order to test the performance of our proposed model, we execute the experiments with two different settings. The first experimental setting is similar with the setting on synthetic dataset. We generate different environments by biased sample selection via bias rate r . In this setting, we choose those features which have no associations with outcome as noisy features for biased sample selection. In second experimental setting, we generate the various environments by dataset separation with users' feature. Specifically, we separate the whole dataset into 4 parts by users' age, including $Age \in [20, 30)$, $Age \in [30, 40)$, $Age \in [40, 50)$ and $Age \in [50, 100)$.

Results on Setting 1. Based on the first experimental setting, we plot the results in Figure 13 and Figure 14. Under this setting, we trained all algorithms on a dataset with bias rate $r = 0.6$ for four noisy features. Then we test the performance of our proposed algorithm and baselines on various test data with different bias rate on these four noisy features, and report the $RMSE$ in Fig. 13a. To explicitly demonstrate the advantage of our proposed algorithm, we plot the $Average_Error$ and $Stability_Error$ as defined in Eq. (1) and (2) in Fig. 13b. We further generate additional test data by varying bias rate r on other features, with results in Fig. 14. Fig. 14a and 14b show that DGBR makes the most stable prediction across test data. Overall, the results and their interpretation are very similar to the simulation experiments.

Results on Setting 2. Based on the second experimental setting, we plot the results in Figure 15, where we separate the dataset into four environments by users' age, including $Age \in [20, 30)$, $Age \in [30, 40)$, $Age \in [40, 50)$ and $Age \in [50, 100)$. We trained all algorithms on dataset where users' $Age \in [20, 30)$, then we test them on all the four environments. From Figure 15, we could find that our DGBR algorithm achieves comparable result to the baselines on test data with users' $Age \in [20, 30)$, where the distributions of variables are similar with the one on the training data. On such test data, the spurious correlation between noisy features and outcome could help baselines make a more precise prediction. While on the other three parts of test dataset, whose distributions are different with training dataset, our

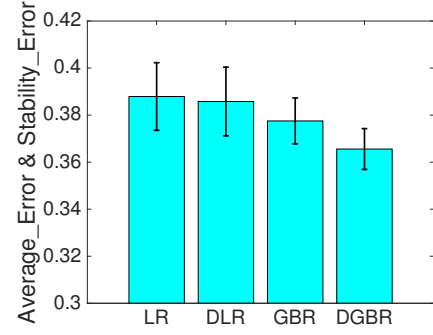


Fig. 16: $Average_Error$ and $Stability_Error$ of all algorithms across environments after fixing $P(Y)$ as the same with its value on global dataset.

DGBR algorithm obtains the best prediction performance. The main reason is that our algorithm can reduce or even remove the spurious effect of noisy features on outcome and find out the stable features for stable prediction.

We can infer that the stability of DGBR algorithm is not as good as baselines in Fig. 15; this occurs because the distribution of outcome $P(Y)$ varied across these four environments. After we fixed $P(Y)$ by data sampling on the outcome with $P(Y = 1) = \frac{14,891}{14,891 + 93,108}$ on the global dataset, we report the $Average_Error$ and $Stability_Error$ of all algorithms across four environments in Figure 16. And we find that as the $P(Y)$ is stable, DGBR outperforms baselines.

7 CONCLUSION

In this paper, we focus on how to make a stable prediction across unknown environments, where the data distribution of unknown environments might be very different with the distribution of training data. We argued that most previous methods for addressing stable prediction are deficient because either they need the distribution of test data as prior knowledge or rely on diversity of training datasets from different environments. Therefore, we propose a Deep Global Balancing Regression algorithm for stable prediction across environments by jointly optimizing the deep auto-encoder model and global balancing model. The global balancing model can identify the causal relationship between predictor variables and response variable, while the deep auto-encoder model is designed for capturing the non-linear structure among variables and making global balancing easier and less noisy. We prove that our algorithm can make a stable prediction from both theoretical analysis and empirical experiments. The experimental results on both synthetic and real world datasets show that our DGBR algorithm outperforms the baselines for stable prediction across unknown environments.

APPENDIX A PROOF OF LEMMA 1

Proof. Assume treatment variable is $T = \mathbf{X}_{i,j}$ and $\mathbf{X}_{i,-j}$ are covariates. From the propensity score is bounded away from zero and one, and $\exists(x_1^0, \dots, x_{j-1}^0, x_{j+1}^0, \dots, x_p^0)$, $P(\mathbf{X}_{i,-j} = (x_1^0, \dots, x_{j-1}^0, x_{j+1}^0, \dots, x_p^0)) > 0$, from

$$\begin{aligned}
& P(\mathbf{X}_i = (x_1^0, \dots, x_{j-1}^0, x_j, x_{j+1}^0, \dots, x_p^0)) \\
&= P(\mathbf{X}_{i,-j} = (x_1^0, \dots, x_{j-1}^0, x_{j+1}^0, \dots, x_p^0)) \cdot \\
& P(\mathbf{X}_{i,j} = x_j | \mathbf{X}_{i,-j} = (x_1^0, \dots, x_{j-1}^0, x_{j+1}^0, \dots, x_p^0))
\end{aligned}$$

we have

$$0 < P(\mathbf{X}_i = (x_1^0, \dots, x_{j-1}^0, x_j, x_{j+1}^0, \dots, x_p^0)) < 1 \quad (30)$$

for $x_j = 0$ or $x_j = 1$.

Next is to proof $\forall x$ (x is binary),

$$0 < P(\mathbf{X}_i = x) < 1$$

from inequality (30). Let $k \neq j$, from

$$\begin{aligned}
& P(\mathbf{X}_i = (x_1^0, \dots, x_{j-1}^0, x_j, x_{j+1}^0, \dots, x_p^0)) \\
&= P(\mathbf{X}_{i,-k} = (x_1^0, \dots, x_{k-1}^0, x_{k+1}^0, \dots, x_p^0)) \cdot \\
& P(\mathbf{X}_{i,k} = x_k^0 | \mathbf{X}_{i,-k} = (x_1^0, \dots, x_{k-1}^0, x_{k+1}^0, \dots, x_p^0))
\end{aligned}$$

and $0 < P(\mathbf{X}_i = (x_1^0, \dots, x_{j-1}^0, x_j, x_{j+1}^0, \dots, x_p^0)) < 1$, we have

$$P(\mathbf{X}_{i,-k} = (x_1^0, \dots, x_{k-1}^0, x_{k+1}^0, \dots, x_p^0)) > 0$$

Furthermore, $\mathbf{X}_{i,k}$ can also be viewed as the treatment variable, so

$$0 < P(\mathbf{X}_{i,k} = x_k^0 | \mathbf{X}_{i,-k} = (x_1^0, \dots, x_{k-1}^0, x_{k+1}^0, \dots, x_p^0)) < 1$$

, and therefore,

$$0 < P(\mathbf{X}_{i,k} = 1 - x_k^0 | \mathbf{X}_{i,-k} = (x_1^0, \dots, x_{k-1}^0, x_{k+1}^0, \dots, x_p^0)) < 1$$

We have (without loss of generality, we assume $k < j$), $\forall x_k, x_j$

$$0 < P(\mathbf{X}_i = (x_1^0, \dots, x_{k-1}^0, x_k, x_{k+1}^0, \dots, x_{j-1}^0, x_j, x_{j+1}^0, \dots, x_p^0)) < 1.$$

We repeat the above for all other variables one by one, we have $\forall x$,

$$0 < P(\mathbf{X}_i = x) < 1$$

□

APPENDIX B

PROOF OF PROPOSITION 1

Proof. If $0 < \hat{P}(\mathbf{X}_i = x) < 1$, from Theorem 1, $W_i^* = \frac{1}{\hat{P}(\mathbf{X}_i = x)}$ satisfies equation (4) equals 0. Next is to show all variables in \mathbf{X} are independent after balancing by this W^* . Note that

$$\begin{aligned}
\sum_{i=1}^n W_i^* &= n \sum_x \frac{1}{n} \sum_{i: \mathbf{X}_i = x} W_i^* \\
&= n \sum_x \hat{P}(\mathbf{X}_i = x) \cdot \frac{1}{\hat{P}(\mathbf{X}_i = x)} = n \cdot 2^p
\end{aligned}$$

Similarly, $\sum_{i: \mathbf{X}_{i,j}=1} W_i^* = n \cdot 2^{p-1}$ and $\sum_{i: \mathbf{X}_{i,j}=0} W_i^* = n \cdot 2^{p-1}$. Denote the probability mass function of \mathbf{X} weighted by W^* as \tilde{P} . Thus, for $x = (x_1, \dots, x_p)$,

$$\tilde{P}(\mathbf{X}_i = (x_1, \dots, x_p)) = \frac{\sum_{i: \mathbf{X}_{i,j}=x} W_i^*}{\sum_i W_i^*} = \frac{1}{2^p}$$

and $\forall j$, $\tilde{P}(\mathbf{X}_{i,j} = x_j) = \frac{\sum_{i: \mathbf{X}_{i,j}=x_j} W_i^*}{\sum_i W_i^*} = \frac{1}{2}$, so we have

$$\tilde{P}(\mathbf{X}_i = (x_1, \dots, x_p)) = \tilde{P}(\mathbf{X}_{i,1} = x_1) \cdots \tilde{P}(\mathbf{X}_{i,p} = x_p),$$

which implies that covariates in \mathbf{X} are independent after balanced by W^* . □

APPENDIX C

PROOF OF LEMMA 2

Proof. $\forall k, j, k \neq j$, it has $0 \leq \frac{\sum_{i: \mathbf{X}_{i,k}=1, \mathbf{X}_{i,j}=1} \hat{W}_i}{\sum_{i: \mathbf{X}_{i,j}=1} \hat{W}_i} \leq 1$

and $0 \leq \frac{\sum_{i: \mathbf{X}_{i,k}=1, \mathbf{X}_{i,j}=0} \hat{W}_i}{\sum_{i: \mathbf{X}_{i,j}=0} \hat{W}_i} \leq 1$. Thus, $0 \leq \alpha \leq 1, \forall m$.

Assume for k, j and $k \neq j$, $\sum_{x: x_k=1, x_j=1} \mathbb{1}(\sum_{i=1}^n \mathbb{1}(X_i = x) = 0) = m_1$, $\sum_{x: x_k=1, x_j=0} \mathbb{1}(\sum_{i=1}^n \mathbb{1}(X_i = x) = 0) = m_2$, $\sum_{x: x_k=0, x_j=1} \mathbb{1}(\sum_{i=1}^n \mathbb{1}(X_i = x) = 0) = m_3$ and $\sum_{x: x_k=0, x_j=0} \mathbb{1}(\sum_{i=1}^n \mathbb{1}(X_i = x) = 0) = m_4$.

- 1) If $m = 0, \alpha = 0$ is a direct result from Theorem 1
- 2) If $0 < m \leq 2^{p-2}$, without loss of generality, assume $m_2 \geq m_4$,

$$\begin{aligned}
\alpha_{jk} &= \left| \frac{\sum_{i: \mathbf{X}_{i,k}=1, \mathbf{X}_{i,j}=1} \hat{W}_i}{\sum_{i: \mathbf{X}_{i,j}=1} \hat{W}_i} - \frac{\sum_{i: \mathbf{X}_{i,k}=1, \mathbf{X}_{i,j}=0} \hat{W}_i}{\sum_{i: \mathbf{X}_{i,j}=0} \hat{W}_i} \right| \\
&= \left| \frac{2^{p-2} - m_1}{2^{p-1} - m_2} - \frac{2^{p-2} - m_3}{2^{p-1} - m_4} \right| \\
&\leq \frac{2^{p-2}}{2^{p-1} - m_2} - \frac{2^{p-2} - m_4}{2^{p-1} - m_4}
\end{aligned}$$

Given $m_4 = m - m_2$,

$$\frac{\partial \alpha_{jk}}{\partial m_2} = 2^{p-2} \left(\frac{1}{(2^{p-1} - m_2)^2} - \frac{1}{(2^{p-1} - m + m_2)^2} \right)$$

which is positive when $m_2 \leq m/2$ (we assume $m_2 \geq m_4$), and therefore

$$\alpha_{jk} \leq \frac{2^{p-2}}{2^{p-1} - m} - \frac{1}{2}$$

- 3) If $2^{p-2} < m < 2^{p-1}$, without loss of generality, assume $m_2 \geq m_4$, when $m_2 \leq 2^{p-2}$, from 2, we have

$$\begin{aligned}
\left| \frac{2^{p-2} - m_1}{2^{p-1} - m_2} - \frac{2^{p-2} - m_3}{2^{p-1} - m_4} \right| &\leq \frac{2^{p-2}}{2^{p-2} - m_2} - \frac{2^{p-2} - m_4}{2^{p-1} - m_4} \\
&\leq 1 - \frac{2^{p-2} - m + 2^{p-2}}{2^{p-1} - m + 2^{p-2}}
\end{aligned}$$

when $m_2 > 2^{p-1}$,

$$\begin{aligned}
\left| \frac{2^{p-2} - m_1}{2^{p-1} - m_2} - \frac{2^{p-2} - m_3}{2^{p-1} - m_4} \right| &\leq 1 - \frac{2^{p-2} - m_4}{2^{p-1} - m_4} \\
&< 1 - \frac{2^{p-2} - m + 2^{p-2}}{2^{p-1} - m + 2^{p-2}}
\end{aligned}$$

because $\frac{2^{p-2} - m_3}{2^{p-2} - m_4}$ is decreasing in m_4 . Thus

$$\alpha \leq 1 - \frac{2^{p-2} - m + 2^{p-2}}{2^{p-1} - m + 2^{p-2}} = 1 - \frac{2^{p-1} - m}{3 \times 2^{p-2} - m}$$

- 4) If $2^{p-1} \leq m$, let $m_1 = \lfloor \frac{m}{2} \rfloor - 2^{p-2}$, $m_2 = \lfloor \frac{m}{2} \rfloor$, $m_3 = 2^{p-2}$, $m_4 = \lceil \frac{m}{2} \rceil$, which satisfy $m_2 + m_4 = 1$, $m_1 \leq m_2$, and $m_3 \leq m_4$. Moreover,

$$\left| \frac{2^{p-2} - m_1}{2^{p-1} - m_2} - \frac{2^{p-2} - m_3}{2^{p-1} - m_4} \right| = 1$$

together with $\alpha \leq 1$, we have $\alpha = 1$ □

APPENDIX D

PROOF OF THEOREM 2

Proof. The probability that m different values in \mathcal{X} do not appear in \mathbf{X} equals the ratio of the number of solutions to

$$y_1 + y_2 + \cdots + y_{2^p} = n, \quad (31)$$

where m different i s have $y_i = 0$, to the total number of solution to Eq. (31) without any constraint. The denominator is $\binom{n+2^p-1}{2^p-1}$. The numerator is the number of methods to select m different i s, such that $y_i = 0$ multiplied by the number of solutions to $y_1 + y_2 + \cdots + y_{2^p-m} = n$ without any constraint, which is $\binom{2^p}{m} \binom{n-1}{2^p-1-m}$. Thus the probability that m different x s do not appear in \mathbf{X} is

$$\frac{1}{\binom{n+2^p-1}{2^p-1}} \binom{2^p}{m} \binom{n-1}{2^p-1-m}$$

With lemma 2,

$$E[\alpha] = \frac{1}{\binom{n+2^p-1}{2^p-1}} \left\{ \sum_{m=0}^{2^p-1} \binom{2^p}{m} \binom{n-1}{2^p-1-m} g(p, m) \right\},$$

where $g(p, m)$ is defined in (12). \square

APPENDIX E

PROOF OF THEOREM 3

Proof. Define $L_{\tilde{P}}(f) = E_{\tilde{P}}(l(f(\mathbf{X}), Y))$, where the probability mass function $\tilde{P}(\mathbf{X}_i, Y_i) = \tilde{P}(\mathbf{X}_i)P(Y_i|\mathbf{X}_i)$ has $P(Y_i = y|\mathbf{X}_i = x) = P(Y_i = y|\mathbf{S}_i = s, \mathbf{V}_i = v) = P(y|s)$ to be the same as that in Assumption 1 and $\tilde{P}(\mathbf{X}_i = x) = \tilde{p}_x$, where \tilde{p}_x is defined in Eq. (13) and equals $\frac{1}{n} \sum_{i=1}^n W_i^* \mathbb{1}(\mathbf{X}_i = x)$.

Let $\tilde{f}^* = \arg \min_f L_{\tilde{P}}(f)$. For all f ,

$$|L_{\tilde{P}}(f) - L_P(f)| \leq \max_x \mathbb{E}[l(f(x), y)|x] \sum_x |\epsilon_x|, \quad (32)$$

followed by

$$\begin{aligned} |L_{\tilde{P}}(f) - L_P(f)| &= |\sum_x \tilde{p}_x \mathbb{E}[l(f^*(x), y)|x] - \sum_x p_x \mathbb{E}[l(f(x), y)|x]| \\ &= |\sum_x \tilde{p}_x \mathbb{E}[l(f(x), y)|x] - \sum_x (\tilde{p}_x - \epsilon_x) \mathbb{E}[l(f(x), y)|x]| \\ &= |\sum_x \epsilon_x \mathbb{E}[l(f(x), y)|x]| \leq \max_x \mathbb{E}[l(f(x), y)|x] \sum_x |\epsilon_x|. \end{aligned}$$

$\max_x \mathbb{E}[l(\hat{f}(x), y)|x]$ is bounded because x and y are binary and all weights in $f \in \mathcal{F}$ are bounded, where \mathcal{F} is the model class defined by the constraints in FWDGBR algorithm. From Eq. (32), we have

$$L_P(\hat{f}) \leq L_{\tilde{P}}(\hat{f}) + \max_x \mathbb{E}[l(\hat{f}(x), y)|x] \sum_x |\epsilon_x|. \quad (33)$$

Next is to upper bound the difference between $L_{\tilde{P}}(\hat{f})$ and $L_{\tilde{P}}(f^*)$. Let $\mathcal{A} = \{x \mapsto l(f(x), y) : f \in \mathcal{F}\}$ to be the loss class, where $l(\cdot)$ is the cross-entropy loss and y is binary. From Lemma 3 in [26], the generalized bound of a 2-class classifier with logistic cross-entropy loss function is related empirical Rademacher complexity, with probability at least $1 - \delta$,

$$L_{\tilde{P}}(\hat{f}) \leq L_{\tilde{P}}(f^*) + 4R_n(\mathcal{A}) + 3\sqrt{\frac{\log(2/\delta)}{2n}}. \quad (34)$$

Note that the auto-encoder has K layers to construct $\phi(\mathbf{X}_i)$ and another K layers to reconstruct \mathbf{X}_i from $\phi(\mathbf{X}_i)$. Y_i is predicted by a logistic regression model on $\phi(\mathbf{X}_i)$. The Rademacher complexity depends on weight constraints $\|\beta\|_2^2 \leq \lambda_4$, $\|\beta\|_1 \leq \lambda_5$ and $\sum_{k=1}^K (\|A^{(k)}\|_F^2 + \|\hat{A}^{(k)}\|_F^2) \leq \lambda_7$ in the FWDGBR algorithm. The decoder from $\phi(\mathbf{X}_i)$ to \mathbf{X}_i is not used to predict Y_i , so the decoder does not affect the complexity $R_n(\mathcal{A})$.

Our goal is to give an upper bound on $R_n(\mathcal{A})$. Constraint $\sum_{k=1}^K (\|A^{(k)}\|_F^2 + \|\hat{A}^{(k)}\|_F^2) \leq \lambda_7$ implies that $\sum_{k=1}^K \|A^{(k)}\|_F^2 \leq \lambda_7$, and together with $\|b^{(k)}\|_2 \leq M^{(k)}$, implies that $\|A_j^{(k)}\|_2 \leq \sqrt{\lambda_7 + (M^{(k)})^2}$. Let $B_k = \sqrt{\lambda_7 + (M^{(k)})^2}$. Constraints $\sum_{k=1}^K (\|A^{(k)}\|_F^2 + \|\hat{A}^{(k)}\|_F^2) \leq \lambda_7$ and $\|b^{(k)}\|_2 \leq M^{(k)}$ imply $\|A_j^{(k)}, b_j^{(k)}\|_2 \leq B_k$.

We can employ Theorem 3.1 in [27] to obtain the empirical Rademacher complexity $R_n(\mathcal{A})$. Since \mathbf{X} is binary, $\|\mathbf{X}\|_{\max} = \max_{i,j} |X_{i,j}| = 1$. Theorem 3.1 in [27] does not have bias term in each layer. We can add constant 1 to neurons in k -th layer $\phi(\mathbf{X}_i)^{(k)}$ to fit in the framework of Theorem 3.1 in [27]. Thus, the dimension of the k -th layer is l_k for $k = 0, 1, 2, \dots, K$. The retain vector is $\theta^k = [1]^{l_k}$ in our case (corresponding to the dropout rate in each layer is 0). If $\|\beta\|_2^2 \leq \lambda_4$ is tighter than $\|\beta\|_1 \leq \lambda_5$, that is, $1/p = 1/2$ and $1/q = 1/2$ for all layers, we have

$$R_n(\mathcal{A}) \leq 2^{K+1} \sqrt{\frac{2\log(2p)}{n}} \sqrt{\lambda_4 l_K} \prod_{k=1}^K B_k (l_{k-1})^{1/2}$$

On the other hand, if $\|\beta\|_1 \leq \lambda_5$ is tighter than $\|\beta\|_2^2 \leq \lambda_4$, that is $1/p = 1$ and $1/q = 0$ for the K -th layer, so $l_K^{1/q} = 1$ and

$$R_n(\mathcal{A}) \leq 2^{K+1} \sqrt{\frac{2\log(2p)}{n}} \lambda_5 \prod_{k=1}^K B_k (l_{k-1})^{1/2}$$

We combine these two cases and have

$$R_n(\mathcal{A}) \leq 2^{K+1} \sqrt{\frac{2\log(2p)}{n}} \min(\sqrt{\lambda_4 l_K}, \lambda_5) \prod_{k=1}^K B_k (l_{k-1})^{1/2} \quad (35)$$

Plug Inequality (35) into Inequality (34), we have

$$\begin{aligned} L_{\tilde{P}}(\hat{f}) &\leq L_{\tilde{P}}(f^*) \\ &\quad + 2^{K+3} \sqrt{\frac{2\log(2p)}{n}} \min(\sqrt{\lambda_4 l_K}, \lambda_5) \prod_{k=1}^K B_k (l_{k-1})^{1/2} \\ &\quad + 3\sqrt{\frac{\log(2/\delta)}{2n}}. \end{aligned} \quad (36)$$

The last step is to bound the difference between $L_{\tilde{P}}(\tilde{f}^*)$ and $L_P(f^*)$. When $L_{\tilde{P}}(\tilde{f}^*) \geq L_P(f^*)$,

$$\begin{aligned} L_{\tilde{P}}(\tilde{f}^*) - L_P(f^*) &= \sum_x \tilde{p}_x \mathbb{E}[l(\tilde{f}^*(x), y)|x] - \sum_x p_x \mathbb{E}[l(f^*(x), y)|x] \\ &= \sum_x \tilde{p}_x \mathbb{E}[l(\tilde{f}^*(x), y)|x] - \sum_x \tilde{p}_x \mathbb{E}[l(f^*(x), y)|x] + \sum_x \epsilon_x \mathbb{E}[l(f^*(x), y)|x] \\ &\leq \sum_x \epsilon_x \mathbb{E}[l(f^*(x), y)|x] \\ &\leq \max_x \mathbb{E}[l(f^*(x), y)|x] \sum_{x:\epsilon_x > 0} \epsilon_x \end{aligned} \quad (37)$$

Eq. (37) holds followed by $\tilde{f}^*(x)$ minimizes $L_{\tilde{P}}(\tilde{f})$, so $L_{\tilde{P}}(\tilde{f}^*) \leq L_{\tilde{P}}(f^*)$, and then $\sum_x \tilde{p}_x \mathbb{E}[l(\tilde{f}^*(x), y)|x] \leq \sum_x \tilde{p}_x \mathbb{E}[l(f^*(x), y)|x]$. Thus,

$$L_{\tilde{P}}(\tilde{f}^*) \leq L_P(f^*) + \max_x \mathbb{E}[l(f^*(x), y)|x] \sum_{x:\epsilon_x > 0} \epsilon_x \quad (38)$$

always holds. From Eq. (33), (36), (38), we have

$$\begin{aligned} L_P(\hat{f}) &\leq L_P(f^*) + \max_x \mathbb{E}[l(\hat{f}(x), y)|x] \sum_x |\epsilon_x| \\ &\quad + 2^{K+3} \sqrt{\frac{2\log(2p)}{n}} \min(\sqrt{\lambda_4 l_K}, \lambda_5) \prod_{k=1}^K B_k (l_{k-1})^{1/2} \\ &\quad + 3\sqrt{\frac{\log(2/\delta)}{2n}} + \max_x \mathbb{E}[l(f^*(x), y)|x] \sum_{x:\epsilon_x > 0} \epsilon_x \\ &\leq L_P(f^*) + 2^{K+3} \sqrt{\frac{2\log(2p)}{n}} \min(\sqrt{\lambda_4 l_K}, \lambda_5) \prod_{k=1}^K B_k (l_{k-1})^{1/2} \\ &\quad + 3\sqrt{\frac{\log(2/\delta)}{2n}} + 2 \max_{x,f} \mathbb{E}[l(f(x), y)|x] \sum_x |\epsilon_x| \end{aligned}$$

with probability $\geq 1 - \delta$. \square

REFERENCES

- [1] H. Shimodaira, "Improving predictive inference under covariate shift by weighting the log-likelihood function," *Journal of statistical planning*

- and inference, vol. 90, no. 2, pp. 227–244, 2000.
- [2] S. Bickel, M. Brückner, and T. Scheffer, “Discriminative learning under covariate shift,” *Journal of Machine Learning Research*, vol. 10, no. Sep, pp. 2137–2155, 2009.
 - [3] M. Sugiyama, S. Nakajima, H. Kashima, P. V. Buenau, and M. Kawanabe, “Direct importance estimation with model selection and its application to covariate shift adaptation,” in *Advances in neural information processing systems*, 2008, pp. 1433–1440.
 - [4] J. Huang, A. Gretton, K. M. Borgwardt, B. Schölkopf, and A. J. Smola, “Correcting sample selection bias by unlabeled data,” in *Advances in neural information processing systems*, 2007, pp. 601–608.
 - [5] M. Dudík, S. J. Phillips, and R. E. Schapire, “Correcting sample selection bias in maximum entropy density estimation,” in *Advances in neural information processing systems*, 2006, pp. 323–330.
 - [6] A. Liu and B. Ziebart, “Robust classification under sample selection bias,” in *Advances in neural information processing systems*, 2014, pp. 37–45.
 - [7] J. Peters, P. Bühlmann, and N. Meinshausen, “Causal inference by using invariant prediction: identification and confidence intervals,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 78, no. 5, pp. 947–1012, 2016.
 - [8] M. Rojas-Carulla, B. Schölkopf, R. Turner, and J. Peters, “Causal transfer in machine learning,” *arXiv preprint arXiv:1507.05333*, 2015.
 - [9] K. Muandet, D. Balduzzi, and B. Schölkopf, “Domain generalization via invariant feature representation,” in *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, 2013, pp. 10–18.
 - [10] P. R. Rosenbaum and D. B. Rubin, “The central role of the propensity score in observational studies for causal effects,” *Biometrika*, vol. 70, no. 1, pp. 41–55, 1983.
 - [11] J. K. Lunceford and M. Davidian, “Stratification and weighting via the propensity score in estimation of causal treatment effects: a comparative study,” *Statistics in medicine*, vol. 23, no. 19, pp. 2937–2960, 2004.
 - [12] P. C. Austin, “An introduction to propensity score methods for reducing the effects of confounding in observational studies,” *Multivariate behavioral research*, vol. 46, no. 3, pp. 399–424, 2011.
 - [13] K. Kuang, P. Cui, B. Li, M. Jiang, S. Yang, and F. Wang, “Treatment effect estimation with data-driven variable decomposition,” in *AAAI*, 2017, pp. 140–146.
 - [14] K. Kuang, M. Jiang, P. Cui, and S. Yang, “Steering social media promotions with effective strategies,” in *Data Mining (ICDM), 2016 IEEE 16th International Conference on*. IEEE, 2016, pp. 985–990.
 - [15] J. Hainmueller, “Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies,” *Political Analysis*, vol. 20, no. 1, pp. 25–46, 2012.
 - [16] J. R. Zubizarreta, “Stable weights that balance covariates for estimation with incomplete outcome data,” *Journal of the American Statistical Association*, vol. 110, no. 511, pp. 910–922, 2015.
 - [17] S. Athey, G. W. Imbens, and S. Wager, “Approximate residual balancing: De-biased inference of average treatment effects in high dimensions,” *arXiv preprint arXiv:1604.07125*, 2016.
 - [18] K. Kuang, P. Cui, B. Li, M. Jiang, and S. Yang, “Estimating treatment effect in the wild via differentiated confounder balancing,” in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2017, pp. 265–274.
 - [19] Y. Yu and C. Szepesvári, “Analysis of kernel mean matching under covariate shift,” *arXiv preprint arXiv:1206.4650*, 2012.
 - [20] J. Wen, C.-N. Yu, and R. Greiner, “Robust learning under uncertain test distributions: Relating covariate shift to model misspecification,” in *ICML*, 2014, pp. 631–639.
 - [21] K. Kuang, M. Jiang, P. Cui, J. Sun, and S. Yang, “Effective promotional strategies selection in social media: A data-driven approach,” *IEEE Transactions on Big Data*, 2017.
 - [22] S. Li and Y. Fu, “Matching on balanced nonlinear representations for treatment effects estimation,” in *Advances in Neural Information Processing Systems*, 2017, pp. 930–940.
 - [23] B. Yu et al., “Stability,” *Bernoulli*, vol. 19, no. 4, pp. 1484–1500, 2013.
 - [24] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, “Greedy layer-wise training of deep networks,” in *Advances in neural information processing systems*, 2007, pp. 153–160.
 - [25] W. B. Johnson and J. Lindenstrauss, “Extensions of lipschitz mappings into a hilbert space,” *Contemporary mathematics*, vol. 26, no. 189–206, p. 1, 1984.
 - [26] L. Wan, M. Zeiler, S. Zhang, Y. Le Cun, and R. Fergus, “Regularization of neural networks using dropconnect,” in *International Conference on Machine Learning*, 2013, pp. 1058–1066.
 - [27] K. Zhai and H. Wang, “Adaptive dropout with rademacher complexity regularization,” in *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=Sluxsy0Z>
 - [28] S. Menard, *Applied logistic regression analysis*. Sage, 2002, vol. 106.
 - [29] Y. Chen, Z. Lin, X. Zhao, G. Wang, and Y. Gu, “Deep learning-based classification of hyperspectral data,” *IEEE Journal of Selected topics in applied earth observations and remote sensing*, vol. 7, no. 6, pp. 2094–2107, 2014.