# Learning Disentangled Representations for Counterfactual Regression via Mutual Information Minimization

### Mingyuan Cheng
wanyu.cmy@alibaba-inc.com
Alibaba Group
Beijing, China

### Xinru Liao
xinru.lxr@alibaba-inc.com
Alibaba Group
Hangzhou, China

### Quan Liu[*]
lq204691@alibaba-inc.com
Alibaba Group
Hangzhou, China

### Bin Ma
mabin.mb@alibaba-inc.com
Alibaba Group
Hangzhou, China

### Jian Xu
xiyu.xj@alibaba-inc.com
Alibaba Group
Beijing, China

### Bo Zheng
bozheng@alibaba-inc.com
Alibaba Group
Beijing, China

## ABSTRACT

Learning individual-level treatment effect is a fundamental problem in causal inference and has received increasing attention in many areas, especially in the user growth area which concerns many internet companies. Recently, disentangled representation learning methods that decompose covariates into three latent factors, including instrumental, confounding and adjustment factors, have witnessed great success in treatment effect estimation. However, it remains an open problem how to learn the underlying disentangled factors precisely. Specifically, previous methods fail to obtain independent disentangled factors, which is a necessary condition for identifying treatment effect. In this paper, we propose Disentangled Representations for Counterfactual Regression via Mutual Information Minimization (MIM-DRCFR), which uses a multi-task learning framework to share information when learning the latent factors and incorporates MI minimization learning criteria to ensure the independence of these factors. Extensive experiments including public benchmarks and real-world industrial user growth datasets demonstrate that our method performs much better than state-of-the-art methods.

## CCS CONCEPTS

• **Information systems → Computational advertising**.

## KEYWORDS

Causal Inference, Disentangled Representations, Mutual Information Minimization, Multi-task Learning

[*]Corresponding author.

## 1 INTRODUCTION

Estimating treatment effect is one of the most important topics in many domains, such as policy making [1, 20], medicine prediction [23], advertisement [4, 11, 24], recommendation [25, 29] and user growth [9]. It often needs to answer counterfactual problems [21] like *"Would this patient have low blood sugar had she received a medication?"* or *"Would the customer buy the product had he got a 70% discount?"*. Specifically, in the user growth area, companies may take many activities such as sending coupons and pushing messages to increase user acquisition or retention, where the counterfactual problem becomes *"Would the user act more actively on the platform had he received the coupon or message?"*.

One golden standard approach to learn causal effect is to perform Randomized Controlled Trial [21], where the treatment is randomly assigned to individuals. However, this is often expensive, unethical or even infeasible, thus we usually focus on estimating treatment effect from observational data. In the observational study, the treatment often depends on some attributes of the individual $x$, which causes the problem of **selection bias** [16] (i.e., $p(t|x) \neq p(t)$). Taking the medicine scenario for example, the economic status affects both the medications and the patient's recovery rate. And it is vital to find all such confounding variables (i.e., affecting both the treatment and outcome) and control them to make precise predictions. This means *unconfoundedness* assumption often needs to be satisfied in the observational study to make the treatment effect identifiable [21].

Even though we already have all confounders in our variables, we still face a difficult problem of identifying them and further balancing them with the *backdoor criterion* [21]. Existing methods achieve balancing either by propensity score weighting methods [2] or representation learning methods which reduce the discrepancy between the treated and control group (e.g., BNN [17] and CFRnet [23]) while ignoring identification of other latent factors. Recently, disentangled representation learning methods, $D^2$VD [19], DR-CFR [14] and TEDVAE [28] have been proposed to learn disentangled factors $\{\Gamma, \Upsilon, \Delta\}$, respectively representing factor that

affects only the treatment, only the outcome, and both the treatment and the outcome (*aka* instrumental, adjustment and confounding factors). Although disentangled representation learning methods greatly help achieve exact identification of the latent factors, the above methods still face the following limitations: $D^2VD$ only decomposes features into two factors $\{\Upsilon, \Delta\}$, DR-CFR cannot effectively distinguish the difference between the $\Delta$ and $\{\Gamma, \Upsilon\}$ and TEDVAE uses the generative model which might greatly increase model training complexity. Besides, these methods cannot obtain **independent disentangled representations**, which is a necessary condition for identifying treatment effect. To solve the problem of independence of disentangled representations, we propose to use MI minimization [5, 6] method, which has obtained increasing attention in domain adaptation [10] and style transfer [18] recently. It is typically utilized as a learning criterion in loss function to ensure the independence between variables. Specifically, [6] proposes a MI upper bound called Contrastive Log-Ratio Upper Bound (CLUB) to deal with the MI minimization task and various experiments have demonstrated the effectiveness of this method.

In this paper, we propose an easy-handling and well-identifying model to deal with the problems of disentanglement in treatment effect estimation. We incorporate the multi-task learning framework such as shared-bottom structure and MI minimization criteria to learn the disentangled factors. And our main contributions are:

- We propose a multi-task learning structure represented by the disentangled representation layer to share information across these latent factors, instead of three independent representation networks which is commonly used by previous work.
- We introduce the MI minimization method into causal inference to learn the latent factors, which uses CLUB as MI upper bound to obtain ideally independent disentangled representations.
- We carry out extensive experiments on both public benchmarks and industrial datasets of user growth (e.g., message pushing and coupon sending), which demonstrate the superiority of our method.

## 2 THE PROPOSED METHOD

In this section, we introduce our model architecture as shown in Figure 1. We first present the basic definition and assumption in 2.1. Then, we explain the multi-task disentangled representations learning framework in 2.2. Last, we illustrate MI minimization regularizer for causal inference in 2.3.

### 2.1 Preliminary

We first present some notations in our context. Given the observational dataset $\mathcal{D} = \{(x_i, t_i, y_i^{t_i}(x_i, t_i))\}_{i=1}^n$, where $n$ is the number of data samples, $x_i \in \mathcal{X}$ is the input features referring individual context information, $y_i^{t_i}(x_i, t_i) \in \mathcal{Y}$ is observed factual outcome and counterfactual outcome $y_i^{1-t_i}(x_i, t_i)$ is missing here, $t_i \in \mathcal{T}$ refers to potential interventions (e.g., for binary treatment $t \in \{0, 1\}$). Mathematically, we define our goal in this paper is to learn a function $\mathcal{F} : \mathcal{X} \times \mathcal{T} \to \mathcal{Y}$ to predict the potential outcomes and then
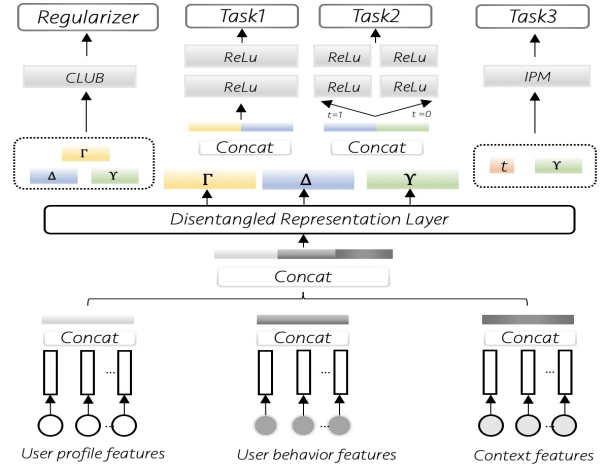


**Figure 1: The proposed model architecture of MIM-DRCFR in our industrial datasets. The input variable contains users' profile, behavior and context features. The disentangled representation layer consists of a shared-bottom structure and three factor-specific layers and then outputs three latent factors $\{\Gamma, \Delta, \Upsilon\}$. The top of the figure shows multi-task objectives, including treatment prediction, outcome prediction, IPM constraint and the MI minimization regularizer.**

estimate the *individual treatment effect (ITE)*[1] and the *average treatment effect (ATE)*:

**Definition 1.** *The individual treatment effect is formulated as*:

$$\tau_i = y_i^1(x_i, t_i) - y_i^0(x_i, t_i) \tag{1}$$

**Definition 2.** *The average treatment effect is formulated as*:

$$ATE = \frac{1}{n} \sum_{i=1}^{n} \tau_i \tag{2}$$

The following fundamental assumptions [22] need to be satisfied in individual treatment effect estimation:

**Assumption 1.** (**SUTVA**) *The Stable Unit Treatment Value Assumption requires that the response of a unit depends only on the treatment to which he himself was assigned and not affected by others.*

**Assumption 2.** (**Unconfoundedness**) *The treatment assignment mechanism is independent of the potential outcome when conditioning on the observed variables, Formally*: $Y_0, Y_1 \perp t \mid x$.

**Assumption 3.** (**Positivity**) *Each unit has a non-zero probability to be assigned to each treatment when given the observed contexts, i.e.,* $0 < P(t = 1|x) < 1$.

### 2.2 Multi-task Disentangled Representation Learning

Without loss of generality, we assume the dataset $\mathcal{D}$ is generated from three underlying factors $\{\Gamma, \Delta, \Upsilon\}$ [14]. In our user growth scene, $x_i$ consists of user's profiles, behavior and context features. $t_i$ can be sending coupon or pushing message to user. $y_i$ can be user's login rate or click-through rate. Then we aim to encode

---

[1]The individual treatment effect (ITE), *aka* conditional average treatment effect (CATE).

the input features $x_i$ into three separate embedding parts through disentangled representation layer (*DRL*), which is formulated as:

$$\Gamma, \Delta, \Upsilon = DRL(x),\tag{3}$$

where $x \subseteq \mathbb{R}^{1\times d}$ and $\Gamma, \Delta, \Upsilon \subseteq \mathbb{R}^{1\times d}$, $d$ refers to feature dim.

DR-CFR directly uses three separate representation networks to learn these factors, while our experiments show that this method cannot effectively distinguish the difference between the $\Delta$ and $\{\Gamma, \Upsilon\}$. Inspired by multi-task learning, we use a shared-bottom structure to learn the feature embedding from input variables and then use three factor-specific layers to decode the embedding into latent factors $\{\Gamma, \Delta, \Upsilon\}$ (*aka* SFD layer[2]). Then we learn the latent factors by following tasks:

**Task1.** Predict the treatment from $\Omega = \text{CONCAT}(\Gamma, \Delta)$ and define the loss $\mathcal{L}_{\text{treat}} = \mathcal{L}[t_i, \pi(\Omega(x_i))]$. $\pi$ is a classifier. Minimizing the loss $\mathcal{L}_{\text{treat}}$ ensures that the information of treatment is captured in the union of $\Gamma$ and $\Delta$.

**Task2.** Predict the outcome from $\Phi = \text{CONCAT}(\Upsilon, \Delta)$ and define the loss $\mathcal{L}_{\text{pred}} = \mathcal{L}[y_i, h^{t_i}(\Phi(x_i))]$. $h^{t_i}$ is regression network for each treatment arm. We ensure the information of outcome is captured in the union of $\Upsilon$ and $\Delta$ by minimizing $\mathcal{L}_{\text{pred}}$.

**Task3.** Restrict discrepancy distance and define the loss as $\mathcal{L}_{\text{disc}} = \text{IPM}\left(\{\Upsilon(x_i)\}_{i:t_i=0}, \{\Upsilon(x_i)\}_{i:t_i=1}\right)$. We ensure the latent factor $\Upsilon$ is irrelevant to treatment by minimizing $\mathcal{L}_{\text{disc}}$.

We expect that all confounding factors are captured in $\Delta$ when we can completely distinguish $\Delta$ from $\{\Upsilon, \Gamma\}$ by the following multi-objective function:

$$\mathcal{L}_{\text{MAIN}} = \mathcal{L}_{\text{pred}} + \alpha \cdot \mathcal{L}_{\text{treat}} + \beta \cdot \mathcal{L}_{\text{disc}}\tag{4}$$

where $\alpha$ and $\beta$ are weights for each task, and we use *Wasserstein distance* as our integral probability metric in this paper.

### 2.3 MI Minimization Regularizer

To obtain independent disentangled factors, we propose to minimize the MI among the three factors to ensure independence. MI is a fundamental measure of the dependence between two random variables. Mathematically, the definition of MI between variables $x$ and $y$ is:

$$\text{I}(x, y) = \mathbb{E}_{p(x,y)}\left[\log\frac{p(x,y)}{p(x)p(y)}\right]\tag{5}$$

Following [6], we introduce using CLUB as MI upper bound to accomplish MI minimization among latent factors, and CLUB is defined as $\text{I}_{\text{CLUB}}(x, y) = \mathbb{E}_{p(x,y)}\left[\log p(y|x)\right] - \mathbb{E}_{p(x)}\mathbb{E}_{p(y)}\left[\log p(y|x)\right]$ when the conditional distribution $p(y|x)$ is known. Unfortunately, the conditional relations between variables are unavailable in our task, and therefore we use a variational distribution $q_\theta(y|x)$ to approximate $p(y|x)$ and further extend the CLUB estimator into vCLUB, which is defined as $\text{I}_{\text{vCLUB}}(x, y) = \mathbb{E}_{p(x,y)}\left[\log q_\theta(y|x)\right] - \mathbb{E}_{p(x)}\mathbb{E}_{p(y)}\left[\log q_\theta(y|x)\right]$. $\text{I}_{\text{vCLUB}}(x, y)$ remains a MI upper bound when we have good variational approximation $q_\theta(y|x)$.

Then we use the following MI minimization regularizer to obtain independent disentangled representations for ITE estimation:

$$\mathcal{L}_{\text{CLUB}} = \text{I}_{\text{vCLUB}}(\Gamma, \Delta) + \text{I}_{\text{vCLUB}}(\Delta, \Upsilon) + \text{I}_{\text{vCLUB}}(\Upsilon, \Gamma)\tag{6}$$

---
[2]named from **S**hared-bottom and **F**actor-specific **D**isentangled representation layer.

We summarize the total objective function $\mathcal{L}_{MIM\text{-}DRCFR}$ as:

$$\mathcal{L}_{MIM\text{-}DRCFR} = \mathcal{L}_{\text{pred}} + \alpha \cdot \mathcal{L}_{\text{treat}} + \beta \cdot \mathcal{L}_{\text{disc}} + \gamma \cdot \mathcal{L}_{\text{CLUB}} + \lambda \cdot \mathcal{L}_{\text{REG}}\tag{7}$$

where $\mathcal{L}_{\text{REG}}$ is used to penalize the model complexity and $\alpha$, $\beta$, $\gamma$ and $\lambda$ are weights for these objectives.

Besides, inspired by the orthogonal regularizer in D$^2$VD [19], we introduce a criterion called **R**epresentation **L**ayer **O**rthogonality (RLO), which is an intuitive method to obtain disjoint factors:

$$\mathcal{L}_{\text{RLO}} = \bar{W}_\Gamma^T \cdot \bar{W}_\Delta + \bar{W}_\Delta^T \cdot \bar{W}_\Upsilon + \bar{W}_\Upsilon^T \cdot \bar{W}_\Gamma\tag{8}$$

where $W \subseteq \mathbb{R}^{d\times d}$ refers to products of the *DRL*, then we calculate average of $W$ ($\bar{W} \subseteq \mathbb{R}^{d\times 1}$) to represent the contribution of input variables on disentangled factors. By minimizing $\mathcal{L}_{\text{RLO}}$, we expect each dimension of $x$ is only embedded in one of $\{\Gamma, \Delta, \Upsilon\}$. $\mathcal{L}_{RLO\text{-}DRCFR}$ is obtained by replacing $\mathcal{L}_{\text{CLUB}}$ with $\mathcal{L}_{\text{RLO}}$.

## 3 EXPERIMENT

### 3.1 Benchmark Evaluation

A fundamental problem in causal inference is that we cannot observe factual outcome and counterfactual outcome simultaneously. One common used solution is to synthesize datasets where the outcomes of all possible treatments are available or synthesize outcomes from real-world covariates.

**IHDP Benchmark**. Similar to [13, 23, 28], we use a semi-synthetic dataset based on IHDP as our benchmark which was first introduced by [15]. The covariates come from a randomized experiment studying the effects of home visits by specialist on future cognitive test scores. The selection bias was introduced by removing a biased subset of the treated population and it comprises 747 instances (139 treated, 608 control) with 25 covariates measuring different attributes of children and their mothers. The simulated outcomes are implemented as both setting "A" and setting "B" in the NPCI package and follow *linear* and *nonlinear* relationship respectively.

**Performance Metrics**. Given a synthetic dataset that includes both factual and counterfactual outcomes, we evaluate treatment effect estimation methods through two performance measures. The individual-based metric is $\epsilon_{PEHE} = \frac{1}{n}\sum_{i=1}^n (\hat{\tau}_i - \tau_i)^2$, where $\hat{\tau}_i = \hat{y}_i^1 - \hat{y}_i^0$ is the predicted individual treatment effect and $\tau_i = y_i^1 - y_i^0$ is the actual effect. The population-based measure is $\epsilon_{ATE} = |\text{ATE} - \widehat{\text{ATE}}|$. ATE $= \frac{1}{n}\sum_{i=1}^n (y_i^1 - y_i^0)$ and $\widehat{\text{ATE}}$ is calculated from the estimated outcomes.

**Baselines Methods**. We compare performances of the following methods which can be divided into: *Baseline models*: **TARNET** [23], **CFR-WASS** [23], **CFR-MMD** [23], **CFR-ISW** [13]. *Disentangled models*: **DR-CFR** [14], **TEDVAE** [28], **MIM-DRCFR** (our method) and its variant **RLO-DRCFR**.

**Ablation Study.** We also conduct an ablation study to examine the contributions of different components in MIM-DRCFR.

In Table 1, we report the average results of the $\sqrt{\epsilon_{PEHE}}$ and $\epsilon_{ATE}$ metrics on IHDP-A and IHDP-B benchmarks (100 realizations with 63/27/10 proportion of train/validation/test splits). Results show that MIM-DRCFR achieves the best performance against the compared methods and its variants, which demonstrates that MIM-DRCFR is currently the most effective disentangled method in ITE estimation. The bottom part of Table 1 summarizes results of the ablation study, which demonstrate that all MIM-DRCFR variants

**Table 1: Results of different treatment effect estimation methods on IHDP Benchmark and ablation study of MIM-DRCFR**

| Dataset | IHDP-A | | IHDP-B | |
|---|---|---|---|---|
| Method | $\sqrt{\epsilon_{PEHE}}$ | $\epsilon_{ATE}$ | $\sqrt{\epsilon_{PEHE}}$ | $\epsilon_{ATE}$ |
| TARNET | 0.95 (0.38) | 0.27 (0.13) | 3.15 (0.22) | 0.42 (0.17) |
| CFR-MMD | 0.75 (0.34) | 0.30 (0.12) | 2.58 (0.18) | 0.35 (0.16) |
| CFR-WASS | 0.74 (0.35) | 0.29 (0.12) | 2.51 (0.18) | 0.34 (0.16) |
| CFR-ISW | 0.69 (0.30) | 0.23 (0.09) | 2.55 (0.16) | 0.40 (0.13) |
| DR-CFR | 0.64 (0.25) | 0.20 (0.08) | 2.33 (0.15) | 0.37 (0.10) |
| TEDVAE | 0.58 (0.22) | 0.15 (0.08) | 2.24 (0.13) | 0.31 (0.09) |
| RLO-DRCFR | 0.54 (0.16) | 0.14 (0.05) | 2.16 (0.11) | 0.31 (0.06) |
| MIM-DRCFR | **0.38 (0.09)** | **0.09 (0.01)** | **2.08 (0.09)** | **0.25 (0.04)** |
| *w/o* SFD | 0.53 (0.20) | 0.14 (0.05) | 2.28 (0.13) | 0.34 (0.08) |
| *w/o* MIM | 0.50 (0.21) | 0.13 (0.05) | 2.29 (0.12) | 0.32 (0.09) |
| *w/o* Both | 0.63 (0.25) | 0.19 (0.07) | 2.31 (0.15) | 0.37 (0.10) |

[1] The **bolded** values mean the best performance and the metric represented in the form of "mean (standard deviation)" and the result is statistically significant based on the Welch's $t$-test with $\alpha = 0.05$.

[2] *w/o* SFD means using three separate representation layers instead of SFD layer. *w/o* Both means *w/o* MIM+SFD.

with some components removed witness clear performance drops when comparing to the full model on the $\sqrt{\epsilon_{PEHE}}$ metric, indicating that each of the designed components contributes to the success of MIM-DRCFR.

## 3.2 Real-world dataset Offline Evaluation

In real-world scenarios we often face the following budget constrained problem *"How to maximize global value of the population $\omega$ when we can only **intervene** on a subgroup $\lambda$ of the population $\omega$ due to the budget limit"*, which can be formulated as:

$$\max \sum_{i \in \lambda} y_i^1(x_i, t_i = 1) + \sum_{i \in \delta} y_i^0(x_i, t_i = 0)$$
$$s.t. \sum_{i \in \lambda} \mathbb{I}[t_i = 1] \leq B, \tag{9}$$

where $\lambda$ (resp., $\delta$) refers to treatment (resp., control) subgroup, $\omega = \lambda \cup \delta$ and $\lambda \cap \delta = \emptyset$. $\mathbb{I}$ denotes the indicator function and B refers to the budget (e.g., total number of treated users). We prove that this problem equals finding an optimal subgroup $\lambda^*$ that has higher non-negative *uplift value* (i.e., individual treatment effect $\tau$) than that of $\delta$. Mathematically:

$$\lambda^* = \left\{ x_i \mid \forall x_j \in \delta, \ \tau_i \geq \tau_j, \ \tau_i \geq 0 \right\} \quad and \quad |\lambda^*| \leq B \tag{10}$$

We can easily obtain $\lambda^*$ through greedy approximation algorithm [7] based on uplift value. Thus, we convert this budget constrained problem into ITE estimation problem, which has gained lots of interest in recent years under the name of **uplift modelling**. The problem consists in targeting treatment to the individuals for whom it would be the most beneficial. For instance, in marketing, one would aim to target advertisement budget to users that would be most likely to be persuadable to purchase [3].

**Message pushing Dataset.** A real-world industrial dataset with 10 million samples that was collected from a current online policy.

**Table 2: Offline AUUC on the two real-world datasets**

| Method | Coupon sending | Message pushing |
|---|---|---|
| TARNET | 1.09 | 0.56 |
| DR-CFR | 1.22 | 0.68 |
| RLO-DRCFR | 1.41 | 0.73 |
| MIM-DRCFR | **1.55** | **0.80** |

[1] We normalize the $AUUC$ value by dividing the $AUUC_\pi(1)$ .

The covariates contain users' profile, behavior and context features, treatment is defined as *"If pushing the message to user"* and the outcome is whether user log onto apps that day. In order to satisfy the *unconfoundedness* assumption, we introduce sufficient confounders (e.g, user activity features) based on our prior knowledge.

**Coupon sending Dataset.** This is similar to above with the message pushing action simply replaced by sending coupons.

**Performance Metrics.** As we cannot obtain factual and counterfactual outcomes simultaneously in the real-world datasets, we use Area Under the Uplift Curve (AUUC) [8, 12, 26] as our offline metric:

$$AUUC_\pi(k) = \left( \frac{R_\pi^{T=1}(k)}{N_\pi^{T=1}(k)} - \frac{R_\pi^{T=0}(k)}{N_\pi^{T=0}(k)} \right) (N_\pi^{T=1}(k) + N_\pi^{T=0}(k)), \tag{11}$$

where $\pi(k)$ denotes the first k proportions of population sorted in descending order of uplift value and $k \in [0, 1]$. $R_\pi^{T=1}(k)$ (resp., $R_\pi^{T=0}(k)$) are the positive outcomes (i.e., login outcome in our industrial datasets) in the treatment (resp., control) group and $N_\pi^{T=1}(k)$ (resp., $N_\pi^{T=0}(k)$) are the number of subjects in the treatment (resp., control) group from $\pi(k)$. The total AUUC is then obtained by cumulative summation [8]:

$$AUUC = \int_0^1 AUUC_\pi(\rho) d\rho \approx \sum_{k=0}^1 AUUC_\pi(k) dk \tag{12}$$

Table 2 and Figure 2 illustrate the offline AUUC and uplift curve, which demonstrates that MIM-DRCFR performs better than other methods on the real-world industrial datasets.
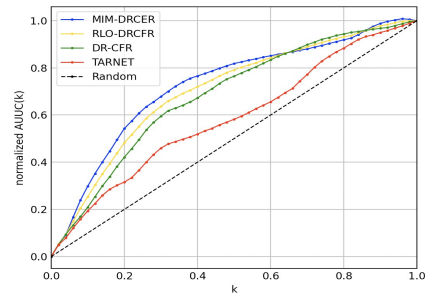


**Figure 2: Uplift curve of the Messing pushing dataset. x-axis refers to the proportion $k$ of the test dataset, y-axis denotes the normalized $AUUC_\pi(k)$ value by dividing the $AUUC_\pi(1)$.**

## 3.3 Online A/B test

We design an online A/B test to further evaluate the performance through calculating the daily login users (DLU) of the population

$\omega$ after pushing message to the estimated optimal subgroup $\hat{\lambda}^*$.

$$DLU(\omega) = L^{T=1}(\hat{\lambda}^*) + L^{T=0}(\omega - \hat{\lambda}^*), \qquad (13)$$

where $L^{T=1}(\hat{\lambda}^*)$ (resp., $L^{T=0}(\omega - \hat{\lambda}^*)$) refers to the number of login users of the treatment (resp., control) group and $\hat{\lambda}^*$ is selected based on the estimated uplift value of different models. In our online datasets, for each experiment group, $\omega$ contains 1 million users and these users are randomly sampled from the entire population who have logged onto our platform in the last 7 days. Finally we choose 60% of them to push messages.

Table 3 illustrates the online DLU result of different models during 5 days' experiment period. The average DLU result of MIM-DRCFR increased by **6.1%** compared with current online policy, which is greater than that of DR-CFR and RLO-DRCFR. This again demonstrates the improvement of MIM-DRCFR on ITE estimation.

**Table 3: Comparison of DLU of Message Pushing**

| METHOD | T | T+1 | T+2 | T+3 | T+4 | Avg |
|---|---|---|---|---|---|---|
| DR-CFR | +2.1% | +1.5% | +1.2% | +0.4% | +0.8% | +1.2% |
| RLO-DRCFR | +4.4% | +4.8% | +5.0% | +3.5% | +2.9% | +4.1% |
| MIM-DRCFR | **+6.6%** | **+8.5%** | **+6.3%** | **+5.8%** | **+3.4%** | **+6.1%** |

## 4 CONCLUSION

In this paper, we focus on disentangled representation learning for ITE estimation and propose a disentangled framework called MIM-DRCFR, which incorporates multi-task learning for the sake of information sharing during the disentangling process and MI minimization for obtaining better independence of the latent factors. Both public benchmarks and real-world industrial datasets demonstrate its superiority over state-of-the-art methods. For future work, we would like to explore more efficient disentangling framework like generative models and extend our method to multi-treatment scenarios.

## REFERENCES

[1] Athey, S. and Imbens, G. (2016). Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113:7353 − 7360.

[2] Austin, P. (2011). An introduction to propensity score methods for reducing the effects of confounding in observational studies. *Multivariate Behavioral Research*, 46:399 − 424.

[3] Betlei, A., Diemert, E., and Amini, M.-R. (2020). Treatment targeting by auuc maximization with generalization guarantees.

[4] Bottou, L., Peters, J., Candela, J. Q., Charles, D., Chickering, D. M., Portugaly, E., Ray, D., Simard, P., and Snelson, E. (2013). Counterfactual reasoning and learning systems: the example of computational advertising. *J. Mach. Learn. Res.*, 14:3207–3260.

[5] Chen, T., Li, X., Grosse, R. B., and Duvenaud, D. (2018). Isolating sources of disentanglement in variational autoencoders. In *NeurIPS*.

[6] Cheng, P., Hao, W., Dai, S., Liu, J., Gan, Z., and Carin, L. (2020). Club: A contrastive log-ratio upper bound of mutual information. In *ICML*.

[7] Dantzig, G. B. (1957). Discrete-variable extremum problems. *Operations Research*, 5:266–288.

[8] Diemert, E. (2018). A large scale benchmark for uplift modeling.

[9] Du, S., Lee, J., and Ghaffarizadeh, F. (2019). Improve user retention with causal learning. In *The 2019 ACM SIGKDD Workshop on Causal Discovery*, pages 34–49. PMLR.

[10] Granger, C. and Lin, J.-L. (1994). Using the mutual information coefficient to identify lags in nonlinear models. *Journal of time series analysis*, 15(4):371–384.

[11] Gu, T., Kuang, K., Zhu, H., Li, J., Dong, Z., Hu, W., Li, Z., He, X., and Liu, Y. (2021). Estimating true post-click conversion via group-stratified counterfactual inference.

[12] Gutierrez, P. and Gérardy, J.-Y. (2017). Causal inference and uplift modelling: A review of the literature. In *International conference on predictive applications and APIs*, pages 1–13. PMLR.

[13] Hassanpour, N. and Greiner, R. (2019). Counterfactual regression with importance sampling weights. In *IJCAI*.

[14] Hassanpour, N. and Greiner, R. (2020). Learning disentangled representations for counterfactual regression. In *ICLR*.

[15] Hill, J. L. (2011). Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20:217 − 240.

[16] Imbens, G. and Rubin, D. (2015). Causal inference for statistics, social, and biomedical sciences: An introduction.

[17] Johansson, F. D., Shalit, U., and Sontag, D. (2016). Learning representations for counterfactual inference. *ArXiv*, abs/1605.03661.

[18] Kazemi, H., Soleymani, S., Taherkhani, F., Iranmanesh, S., and Nasrabadi, N. (2018). Unsupervised image-to-image translation using domain-specific variational information bound. *Advances in neural information processing systems*, 31.

[19] Kuang, K., Cui, P., Li, B., Jiang, M., Yang, S., and Wang, F. (2017). Treatment effect estimation with data-driven variable decomposition. In *AAAI*.

[20] LaLonde, R. J. (1986). Evaluating the econometric evaluations of training programs with experimental data. *The American economic review*, pages 604–620.

[21] Pearl, J. (2009). *Causality*. Cambridge university press.

[22] Rosenbaum, P. and Rubin, D. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70:41–55.

[23] Shalit, U., Johansson, F. D., and Sontag, D. (2017). Estimating individual treatment effect: generalization bounds and algorithms. In *ICML*.

[24] Sun, W., Wang, P., Yin, D., Yang, J., and Chang, Y. (2015). Causal inference via sparse additive models with application to online advertising. In *AAAI*.

[25] Wang, Z., Zhang, J., Xu, H., Chen, X., Zhang, Y., Zhao, W. X., and Wen, J.-R. (2021). Counterfactual data-augmented sequential recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 347–356.

[26] Zhang, W., Li, J., and Liu, L. (2021a). A unified survey of treatment effect heterogeneity modelling and uplift modelling. *ACM Computing Surveys (CSUR)*, 54(8):1–36.

[27] Zhang, W., Liu, L., and Li, J. (2020). Treatment effect estimation with disentangled latent factors. *arXiv preprint arXiv:2001.10652*.

[28] Zhang, W., Liu, L., and Li, J. (2021b). Treatment effect estimation with disentangled latent factors. In *AAAI*.

[29] Zhang, Y., Feng, F., He, X., Wei, T., Song, C., Ling, G., and Zhang, Y. (2021c). Causal intervention for leveraging popularity bias in recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 11–20.