

# CausalMTA: Eliminating the User Confounding Bias for Causal Multi-touch Attribution

Di Yao\*

Institute of Computing Technology,  
Chinese Academy of Sciences  
China  
yaodi@ict.ac.cn

Chang Gong

Institute of Computing Technology,  
Chinese Academy of Sciences  
University of Chinese Academy of  
Sciences  
China  
gongchang21z@ict.ac.cn

Lei Zhang

Strategic Data Solutions (SDS) Group,  
Alibaba Inc  
China  
zl165646@alibaba-inc.com

Sheng Chen

Strategic Data Solutions (SDS) Group,  
Alibaba Inc  
China  
chensheng.cs@alibaba-inc.com

Jingping Bi\*

Institute of Computing Technology,  
Chinese Academy of Sciences  
China  
bjp@ict.ac.cn

## ABSTRACT

Multi-touch attribution (MTA), aiming to estimate the contribution of each advertisement touchpoint in conversion journeys, is essential for budget allocation and automatically advertising. Existing methods first train a model to predict the conversion probability of the advertisement journeys with historical data and calculate the attribution of each touchpoint by using the results counterfactual predictions. An assumption of these works is the conversion prediction model is unbiased. It can give accurate predictions on any randomly assigned journey, including both the factual and counterfactual ones. Nevertheless, this assumption does not always hold as the user preferences act as the common cause for both ad generation and user conversion, involving the confounding bias and leading to an out-of-distribution (OOD) problem in the counterfactual prediction. In this paper, we define the causal MTA task and propose CausalMTA to solve this problem. It systemically eliminates the confounding bias from both static and dynamic perspectives and learn an unbiased conversion prediction model using historical data. We also provide a theoretical analysis to prove the effectiveness of CausalMTA with sufficient ad journeys. Extensive experiments on both synthetic and real data in Alibaba advertising platform show that CausalMTA can not only achieve better prediction performance than the state-of-the-art method but also generate meaningful attribution credits across different advertising channels.

## CCS CONCEPTS

• **Applied computing** → **Electronic commerce**; • **Information systems** → **Computational advertising**.

\* Corresponding authors.



This work is licensed under a Creative Commons Attribution International 4.0 License.

KDD '22, August 14–18, 2022, Washington, DC, USA  
© 2022 Copyright held by the owner/author(s).  
ACM ISBN 978-1-4503-9385-0/22/08.  
<https://doi.org/10.1145/3534678.3539108>

## KEYWORDS

multi-touch attribution, counterfactual prediction, computational advertising

### ACM Reference Format:

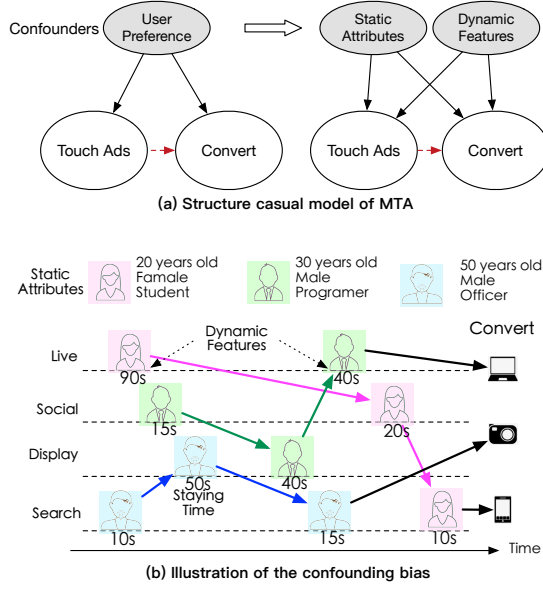
Di Yao, Chang Gong, Lei Zhang, Sheng Chen, and Jingping Bi. 2022. CausalMTA: Eliminating the User Confounding Bias for Causal Multi-touch Attribution. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '22)*, August 14–18, 2022, Washington, DC, USA. ACM, Washington, DC, USA, 11 pages. <https://doi.org/10.1145/3534678.3539108>

## 1 INTRODUCTION

Online advertising platforms have been widely deployed to help advertisers launch their advertisements (ads) across multiple marketing channels, such as social media, feed stream, and paid search. During the usage, the ad exposure sequences and conversion feedbacks of all customers are collected. Multi-touch attribution, short for MTA, aims to estimate each ad touchpoint's relative contribution in user conversion journeys. The attribution results will shed light on the budget allocation and automatically advertising.

Artificial intelligence (AI) coupled with promising machine learning (ML) techniques well known from computer science is broadly affecting many aspects of various fields including science and technology, industry, and even our day-to-day life [28]. Nowadays, instead of attributing the ad touchpoints by heuristic rules [6], data-driven methods [3, 8, 13, 21, 23, 30] which estimate the attribution credits according to the historical data have become the mainstream techniques. These methods learn a conversion prediction model with all observed historical data and then generate the counterfactual ad journeys by removing or replacing some touchpoints. The attribution credits can be estimated using the prediction results of these counterfactual journeys based on some criteria, such as the Shapley value [24]. One essential assumption of these methods is the conversion prediction model should be unbiased, which means the model can give fair predictions on any randomly assigned journeys, including the factual and counterfactual ones. Unfortunately, this assumption does not hold in online advertising.

As shown in Figure 1(a), the user preference is the common cause of both 'touch ads' and 'convert'. Specifically, the ad exposures are



**Figure 1: The motivation of CausalMTA. (a) shows the structure casual model of MTA problem. (b) illustrates the decomposition of user preference resulting in the confounding bias.**

recommended according to the user preferences while the user preference can also lead to the conversion. This common cause is a confounder in MTA, which involves spurious correlation in observed data, *i.e.*, a connection between ads and conversion that appears to be causal but not. As a result, the learned conversion prediction model is biased. The discrepancy between observed training data and counterfactual data causes an out-of-distribution (OOD) problem in counterfactual prediction, which would harm the fairness of attribution. Thus, removing the influence of confounders is critical and necessary for MTA. We define the attribution of the ad journeys with an unbiased prediction model as causal MTA.

Nevertheless, it is not trivial to eliminate the confounding bias of user preferences in MTA. The reasons are two folds: (1) **Multiple confounders**. As illustrated in Figure 1 (a), the confounders in ad exposure generation consist of the static user attributes, such as genders, ages and education background, and dynamic features, *e.g.*, previously viewed ads and favorite items. Both the static and dynamic features should be taken into account for unbiased causal MTA. Existing works either focus on the static settings [4, 14, 15, 33] using IPW and propensity score matching method for deconfounding, or are dedicated to the dynamic confounders [7, 18] learning an unbiased representation for prediction at each time step. All these works rely on the strong ignorability assumption [20], *i.e.*, no hidden confounders. In their settings, the static and dynamic features are hidden confounders mutually that disable the usage. (2) **Delay feedback**. The conversion results are observed at the end of the journey. Unlike those tasks such as [31], there is no explicit feedback available at each touchpoint. Existing sequential deconfounding methods [7, 18, 22, 29] are designed for instant feedbacks, *e.g.*, the blood pressure, which can be observed immediately after taking the hypotensor. CAMTA [16] is the most related method of our work. However, it takes the click action as the "pseudo" feedback at each touchpoint, which would involve other confounders. Above all, due

to the peculiarities of advertising, there are no existed methods that can be used for unbiased causal MTA.

In this paper, we propose a novel method, namely CausalMTA, to mitigate the effect of user preferences-related confoundedness and achieve causal MTA. It learns an unbiased counterfactual prediction model which systemically eliminates the confounding bias from both static user attributes and dynamic features. One fundamental assumption of CausalMTA is that the influence of static user attributes and dynamic features are independent. This assumption is reasonable in online advertising because user attributes usually determine their item interests, and dynamic features determine how likely the users want to buy. As shown in Figure 1 (b), twenty years old students tend to be attracted by fancy phones and cosmetics, whereas the middle age guys usually like high cost-performance phones and anti-bald goods. Dynamic features, such as previously visited ads and staying time, reflect the purchase intention. The main contributions can be summarized as follows:

- We decompose the confounding bias of user preferences into static user attributes and dynamic features, and define the causal MTA problem.
- We propose the first method CausalMTA for causal MTA, which is provable for eliminating the confounding bias of user preferences in counterfactual prediction.
- Extensive experiments on a synthetic dataset, an open-source dataset and a real-world dataset of mobile phones shops from Alibaba demonstrate CausalMTA's superiority.

## 2 RELATED WORK

Existing works can be categorized into two orthogonal groups, *i.e.* data-driven MTA and counterfactual prediction.

**Data-driven multi-touch attribution.** Previously, marketers have applied simple rules, *e.g.*, the last touch, to attribute the influence of touched ads [6], which either ignore the effects of other channels or neglect the channel difference. To overcome these drawbacks, researchers proposed data-driven methods. The data-driven MTA model was first proposed in [23], and has been combined with survival analysis [32] and hazard rate [13] to reflect the influence of ad exposure. However, the data-driven methods mentioned above neglect the customers' features and cannot directly allocate personalized attribution. Besides, the temporal dependency and dynamic interaction between channels need to be modeled. Recently, many DNN-based data-driven MTA methods have been proposed to address the issues, such as channel interaction, time dependency, user characteristics. In some studies [3, 10, 16, 21, 30], RNNs are used to model longitudinal data. DNAMTA [3] is an LSTM based deep sequential model which captures the touchpoint contextual dependency via attention mechanism and incorporates user context information and survival time-decay functions. DARNN [21] is a dual attention model that combines post-view and post-click attribution patterns for final conversion estimation.

**Counterfactual Prediction.** Positioned as a causal estimation problem by [8], the calculation of attribution credits is actually based on counterfactual estimation [10, 25, 27, 32]. A limitation of the models mentioned above is the lack of exogenous variation in user exposure to advertising, which hazards the reliability of the attribution results as the training data of the counterfactual predictor is biased by confounders. Albeit extant papers [5, 19] mitigate

the issue with full or quasi-randomization, the cost and complexity of such randomization restrict the number of users and ad-types. The idea of calibrating the conversion prediction model in MTA by removing confounding bias is inspired by works in counterfactual prediction. There is a large number of methods for counterfactual prediction using observational data in the static setting, involving utilizing propensity score matching [4], learning unbiased representation for prediction [14, 15, 33], conducting propensity-aware hyperparameter tuning [1, 2]. For estimating the effects of time-varying treatments in the area such as epidemiology where the treatments have instant feedback, many approaches [7, 18, 22, 29] addressing the longitudinal setting are proposed. Because of the gap between the longitudinal data in epidemiology and ad journeys in MTA, those methods cannot be directly used in our task.

### 3 PRELIMINARY

#### 3.1 Problem Definition.

We consider an ad exposure dataset  $\mathcal{D}$  which consists of  $N$  conversion journeys of  $U$  users. Each journey can be formulated as a triplet, *i.e.*,  $(\mathbf{u}^i, \mathbf{J}^i, y^i)$ .  $\mathbf{u}^i$  stands for the static user attributes, which is unlikely to be changed during the user conversion journey.  $\mathbf{J}^i$  is a sequence of touchpoints, *i.e.*,  $\{\mathbf{p}_t^i\}_{t=1}^{T^i}$ . Each touchpoint  $\mathbf{p}_t^i = (\mathbf{c}_t^i, \mathbf{f}_t^i)$  contains a channel index  $\mathbf{c}_t^i$  and a feature vector  $\mathbf{z}_t^i$ . Specifically,  $\mathbf{c}_t^i \in \{\mathbf{c}_1, \dots, \mathbf{c}_k, \dots, \mathbf{c}_K\}$  indicates the exposed channels, where  $K$  is the number of ad channels. The feature vector  $\mathbf{f}_t^i$  includes the dynamic side information of this touchpoint, *e.g.*, advertising form and staying time.  $y^i$  is a binary indicator that records whether the journey leads to a conversion event or not. The goal of MTA is to model the sequential pattern and assign the attribution credits to all the touchpoints  $\mathbf{p}_t^i$  according to the whole information in  $\mathcal{D}$ . Nevertheless, historical data in  $\mathcal{D}$  often exhibits confounding bias due to user preferences, which could be a fatal challenge for estimating the attribution credits. The choice of the channel at a touch-point is likely to be influenced by multiple factors like user attributes and previously visited goods. Causal multi-touch attribution aims to estimate unbiased attribution credits  $\{\mathbf{p}_t^i\}_{t=1}^{T^i}$  of all touchpoints.

#### 3.2 Method Overview.

As shown in Figure 2, CausalMTA is a novel model-agnostic framework consisting of two key modules, *i.e.*, journey reweighting and causal conversion prediction, which mitigate the confounding bias of user static attributes and dynamic features respectively. In journey reweighting, we employ the Variational Recurrent Auto-encoders (VRAE) to learn the generation probabilities of pure channels journeys, and conduct user demographic-based density estimation to calculate the likelihoods of the channels being randomly assigned that is used for weights computation. For causal conversion prediction, CausalMTA utilizes RNNs to model the dynamic features of journeys. A gradient reverse layer is built upon the outputs of each time step to ensure the model is unable to predict the next ad channel. It derives balancing representation, which removes the association between dynamic features and the ad exposure. The last hidden output is trained to estimate the conversion probability using the learnt weights of journey reweighting. After that, we can obtain an unbiased prediction model. Lastly, with the constructed counterfactual journeys, the attribution credits can be estimated under Shapley value measure.

## 4 METHODOLOGY

In this section, we first specify the journey reweighting and causal conversion prediction respectively. After that, the calculation of attribution credits is detailed. In the end, we provide the theoretical analysis of CausalMTA.

### 4.1 Journey Reweighting

To mitigate the bias of static user attributes, the journey reweighting module takes pure channel sequences in  $\mathcal{D}$  as the input and estimates the sample weights of the prediction model according to how likely the journey be generated randomly. It consists of two procedures, *i.e.*, generation model for channel sequences and weights estimation of journeys.

**Generation Model for Channel Sequences.** We utilize VRAE (Variational Recurrent Auto-encoders) [11] to model the generation of channel sequences. When there is enough training data, the distribution of pure channel sequences tends to be random, regardless of user preferences. In this setting, the learned VRAE is capable of generating unbiased predictions of observed channel sequences.

For each ad journey  $(\mathbf{u}, \mathbf{J}, y)$  in  $\mathcal{D}$ , we only concern with the channel information in this procedure and extract the pure channel sequence  $\mathbf{C} = \{\mathbf{c}_t\}_{t=1}^T$ . Taking  $\mathbf{C}$  as the input, CausalMTA employs the channel embedding affiliated with LSTM as the encoder and utilizes the final hidden state to generate the distribution over latent representation:

$$\begin{aligned} \{\mathbf{h}_t\}_{t=1}^T &= \text{LSTM}_{\text{enc}}(\mathbf{C}, \mathbf{h}_0), \\ \mu_z &= W_\mu \mathbf{h}_T + \mathbf{b}_\mu, \\ \log(\sigma_z) &= W_\sigma \mathbf{h}_T + \mathbf{b}_\sigma, \end{aligned}$$

where  $\mathbf{h}_0$  is the initial hidden state of the encoder. Leveraging the reparametrization trick, we sample a vector  $\mathbf{z}$  from the distribution to initialize the hidden state of the decoder:

$$\begin{aligned} \mathbf{h}'_0 &= \tanh(W_z^T \mathbf{z} + \mathbf{b}_z), \\ \{\mathbf{h}'_t\}_{t=1}^T &= \text{LSTM}_{\text{dec}}(\mathbf{C}_{\text{out}}, \mathbf{h}'_0), \\ \mathbf{c}'_t &= \text{softmax}(W_o \mathbf{h}'_t + \mathbf{b}_o), \end{aligned}$$

where  $\mathbf{h}'_0$  is the initial hidden state of the decoder;  $\mathbf{C}_{\text{out}}$  is the feed previous input which takes the output of previous step as the input;  $\mathbf{c}'_t$  is the decoded channel sequence.

The loss function is composed of two parts: 1) the reconstruction loss which is defined as the cross-entropy between  $\mathbf{c}_t$  and  $\mathbf{c}'_t$ . 2) the KL divergence between the posterior and prior distribution over the latent variable:

$$\mathcal{L}_w = \alpha \sum_{i=1}^N \sum_{t=1}^{T^i} CE(\mathbf{c}_t, \mathbf{c}'_t) + \beta D_{KL}(q_\phi(\mathbf{z}) || p_\theta(\mathbf{z})), \quad (1)$$

where  $p_\theta(\mathbf{z})$  is the prior distribution usually assumed to be a standard normal distribution  $\mathcal{N}(0, \mathbf{I})$ ;  $q_\phi(\mathbf{z} | \mathbf{c}^i)$  is the posterior approximation  $(\mathcal{N}(\mu^i, (\sigma^i)^2))$ ;  $\alpha$  and  $\beta$  are hyperparameters that control the importance each parts.

**Weights Estimation for Ad Journeys.** To eliminate the bias of user static features, we estimate the weights of observed journeys. The journeys approximating to randomly assigned have higher weights in conversion prediction training than those are severely affected by user preferences. Formally, the learned weights  $W_T(\mathbf{u}, \mathbf{C})$

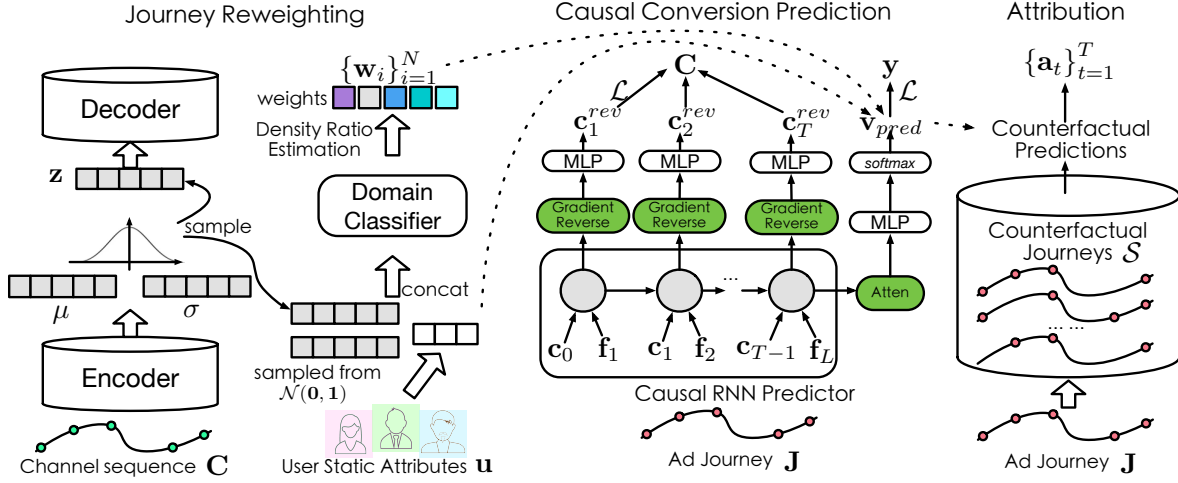


Figure 2: Architecture of CausalMTA.

should be subject to  $W_T(\mathbf{u}, \mathbf{C}) = p(\mathbf{C})/p(\mathbf{C}|\mathbf{u})$  [12, 33]. When we learn a variational distribution  $q_\phi(\mathbf{z}|\mathbf{c})$ , the variational sample weights can be computed as follows:

$$\mathbf{w}^i = W_T(\mathbf{u}^i, \mathbf{c}^i) = \{\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{c}^i)} [\frac{1}{W_z(\mathbf{u}^i, \mathbf{z})}]\}^{-1}, \quad (2)$$

where  $W_z(\mathbf{u}, \mathbf{z})$  can be viewed as the density ratio estimation to decorrelate  $\mathbf{u}$  and  $\mathbf{z}$  for points in space  $\square \times \mathcal{Z}$ . The detailed proof can be found in the appendix.

In CausalMTA, we design a binary domain classifier to help estimate  $W_z(\mathbf{u}, \mathbf{z})$ . Training data of the classifier is generated cooperating with the encoder of VRAE. We label static user attributes with real latent representation  $\{(\mathbf{u}^i, \mathbf{z})\}_{1 \leq i \leq N}$ ,  $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{c}^i)$  as positive ones, whereas samples with latent representation sampled from standard normal distribution  $\{(\mathbf{u}^i, \mathbf{z})\}_{1 \leq i \leq N}$ ,  $\mathbf{z} \sim p_\theta(\mathbf{z})$  as negative ones. We first embed the user attributes into latent vectors and train a domain classifier to fit these samples:

$$\mathbf{e}_u = \text{Embedding}(\mathbf{u}),$$

$$\mathbf{x} = \text{concat}(\mathbf{e}_u, \mathbf{z}),$$

$$p_{\theta_d}(L|\mathbf{u}, \mathbf{z}) = \text{sigmoid}(\text{MLP}(\mathbf{x})).$$

Now that we have  $p(L=0) = p(L=1)$ , the density ratio estimation  $W_z(\mathbf{u}, \mathbf{z})$  can be conducted as follows:

$$W_z(\mathbf{u}, \mathbf{z}) = \frac{p(\mathbf{u}, \mathbf{z}|L=0)}{p(\mathbf{u}, \mathbf{z}|L=1)} = \frac{p(L=0|\mathbf{u}, \mathbf{z})}{p(L=1|\mathbf{u}, \mathbf{z})}, \quad (3)$$

Using this formula, we can obtain the weights of all journeys, i.e.,  $\{\mathbf{w}_i\}_{i=1}^N$ .

## 4.2 Causal Conversion Prediction

After the generation of sample weights, CausalMTA utilizes them to train a trustworthy conversion prediction model. Besides the static attributes, the biases of prediction are also caused by dynamic user features. To mitigate them, we borrow the idea from CRN [7] and involve a gradient reverse layer to learn balancing representation.

Due to the delay feedback problem, we refine the structure of CRN to make it suitable for MTA.

Formally, we first reorganize the dataset. For each journey  $(\mathbf{u}, \mathbf{J}, \mathbf{y})$ , we adopt one step offset on the channel sequence and fill the

blank position with a unified placeholder, i.e.,  $\mathbf{C}_+ = \{\mathbf{c}_0, \mathbf{c}_t\}_{t=1}^{T-1}$ . CausalMTA takes  $\mathbf{C}_+$  along with other dynamic features  $\mathbf{F} = \{\mathbf{f}_t\}_{t=1}^T$  as the input and employs LSTM with the attention mechanism to obtain the trustworthy prediction:

$$\mathbf{e}_u, \mathbf{e}_{\mathbf{C}_+}, \mathbf{e}_F = \text{Embedding}(\mathbf{u}, \mathbf{C}_+, \mathbf{F}),$$

$$\mathbf{v}_{in} = \text{concat}(\mathbf{e}_{\mathbf{C}_+}, \mathbf{e}_F),$$

$$\{\mathbf{out}_t\}_{t=1}^T = \text{LSTM}_{\text{pred}}(\mathbf{v}_{in}, \mathbf{h}_0),$$

where  $\mathbf{e}_{\mathbf{C}_+}, \mathbf{e}_F, \mathbf{v}_{in}$  are the sequences of latent vectors. Once the output vectors are generated, we adopt them for two parallel processes. One for eliminating the bias of dynamic features, and the other for conversion prediction.

$$\{\mathbf{v}_t^{\text{rev}}\}_{t=1}^T = \text{MLP}(\text{GRL}(\{\mathbf{out}_t\}_{t=1}^T)),$$

$$\{\mathbf{c}_t^{\text{rev}}\}_{t=1}^T = \text{softmax}(\{\mathbf{v}_t^{\text{rev}}\}_{t=1}^T),$$

$$\mathbf{v}_{\text{attn}} = \text{Attention}(\mathbf{out}_T, \{\mathbf{out}_t\}_{t=1}^T),$$

$$\mathbf{v}_{\text{pred}} = \text{softmax}(\text{MLP}(\mathbf{v}_{\text{attn}})),$$

where GRL is the gradient reverse layer that ensures  $\mathbf{out}_t$  can not predict  $\mathbf{c}_t$ . The loss function of causal conversion prediction consists of two parts, i.e., reverse channel prediction and conversion prediction:

$$\mathcal{L}_p = \gamma \sum_{i=1}^N \sum_{t=1}^{T_i} \text{CE}(\mathbf{c}_t^{\text{rev}}, \mathbf{c}_t) + \delta \sum_{i=1}^N \mathbf{w}_i \cdot \text{CE}(\mathbf{v}_{\text{pred}}^i, \mathbf{y}^i), \quad (4)$$

where  $\text{CE}$  is the cross-entropy loss;  $\gamma$  and  $\delta$  are hyperparameters;  $\mathbf{w}_i$  is the learned journey weights. With the well-trained conversion prediction model, we can calculate the attribution credits of each touchpoint by constructing some counterfactual journeys.

## 4.3 Attribution Credits Calculation

CausalMTA computes Shapley values [24] for ad credits allocation. Based on assessing the marginal contribution of each player in the game, the Shapley value method is a general credit distribution method, and it has been widely used in MTA [27, 30] due to its advantage of having an axiomatic foundation and catering to fairness consideration.

**Algorithm 1** Learning procedure of CausalMTA.**Input:**

The ad exposure dataset  $\mathcal{D}$ , i.e.,  $\{(\mathbf{u}^i, \mathbf{J}^i, y^i)\}_{i=1}^N$ ;

**Output:**

Attribution credits  $\{\alpha_t^i\}_{t=1}^{T_i}$  for touchpoints  $\{\mathbf{a}_t^i\}_{t=1}^{T_i}$ .

```

1: # Generation model for channel sequences
2: for each journey  $(\mathbf{u}^i, \mathbf{J}^i, y^i)$  in  $\mathcal{D}$  do
3:   Evaluate  $\mathcal{L}_w$  according to Eq.(1) and update the parameters
     of VRAE model.
4: end for
5: # Weights estimation for ad journeys
6: for each ad journey  $(\mathbf{u}^i, \mathbf{J}^i, y^i)$  in  $\mathcal{D}$  do
7:   Generate latent representation for positive and negative sam-
     ples respectively.
8:   Optimize the parameters of domain classifier.
9: end for
10: Conduct density ratio estimation and calculate sample weights
      $\mathbf{w}^i$  according to Eq.(2) and Eq.(3).
11: # Causal conversion prediction
12: for each ad journey  $(\mathbf{u}^i, \mathbf{J}^i, y^i)$  in  $\mathcal{D}$  do
13:   Evaluate  $\mathcal{L}_p$  according to Eq.(4) and update the parameters
     of conversion prediction model.
14: end for
15: # Calculation of Attribution Credits
16: for each ad journey  $(\mathbf{u}^i, \mathbf{J}^i, y^i)$  in  $\mathcal{D}$  do
17:   for each touchpoint  $\mathbf{c}_t^i$  in channel sequence  $\mathbf{C}^i$  do
18:     Compute  $SV_t^i$  according to its definition.
19:     Calculate attribution credits  $\mathbf{a}_t^i$  according to Eq.(5).
20:   end for
21: end for
22: return  $\{\{\mathbf{a}_t^i\}_{t=1}^{T_i}\}_{i=1}^N$ 

```

Formally, let  $\mathbf{J}^i \setminus \{\mathbf{p}_t^i\}$  denote the counterfactual ad journey by removing touchpoint  $\mathbf{p}_t^i$ .  $S$  can be viewed as a subsequence of the counterfactual ad journey  $\mathbf{J}^i \setminus \{\mathbf{p}_t^i\}$ . If we denote the result of causal conversion prediction for channel sequence  $\mathbf{J}^i$  as  $p(\mathbf{J}^i)$ , the Shapley values for ad exposure  $\{\mathbf{c}_t^i\}$  can be defined as  $SV_t^i = \sum_{S \subseteq \mathbf{J}^i \setminus \{\mathbf{p}_t^i\}} \frac{|\mathbf{S}|!(|\mathbf{J}^i| - |\mathbf{S}| - 1)!}{|\mathbf{J}^i|!} [p(\mathbf{S} \cup \{\mathbf{p}_t^i\}) - p(\mathbf{S})]$  where  $|\mathbf{C}^i|$ ,  $|\mathbf{S}|$  are the cardinalities of these sets. If  $SV_t^i$  is negative, we set it zero. Then we normalize all incremental scores for each ad exposure as follows,

$$\mathbf{a}_t^i = \sigma(SV_t^i) / \sum_{t=1}^{T_i} \sigma(SV_t^i), \quad (5)$$

where  $\sigma(x) = \max(0, x)$  and  $\mathbf{a}_t^i$  are the attribution credits of the corresponding ad exposures. The pseudo-code of CausalMTA is shown in Algorithm 1.

#### 4.4 Theoretical Analysis of CausalMTA

Under the assumption of independence, we can decompose the overall confounding bias  $\mathcal{B}$  into the bias introduced by user static attributes  $\mathcal{B}_u$  and the bias introduced by dynamic user features  $\mathcal{B}_f$ , i.e.,  $\mathcal{B} = \mathcal{B}_u + \mathcal{B}_f$ . CausalMTA aims to obtain an unbiased prediction model and achieve  $\mathcal{B} = 0$ .

We prove that the confounding bias from static user attributes  $\mathcal{B}_u$  can be mitigated by estimating sample weights  $\mathbf{w}^i$  for ad journeys. Formally, let  $\mathcal{E}_{cf}$  denote the counterfactual prediction error, which is the target to be minimized. Unfortunately,  $\mathcal{E}_{cf}$  can not be directly measured on the observational dataset. We can derive the upper bound of  $\mathcal{E}_{cf}$ , which is given by

$$\mathcal{B}_u = \mathcal{E}_{cf} - \mathcal{E}_f^w \leq IPM_G(W_T(\mathbf{u}, \mathbf{C})p(\mathbf{u}, \mathbf{C}), p(\mathbf{u})p(\mathbf{C})),$$

where  $\mathcal{E}_f^w$  is the prediction error on the re-weighted data and  $IPM$  denotes Integral Probability Metric. When  $W_T(\mathbf{u}, \mathbf{C}) = \frac{p(\mathbf{C})}{p(\mathbf{C}|\mathbf{u})}$ , the equation  $\mathcal{E}_{cf} = \mathcal{E}_f^w$  can be proved. More details of the proof are available in the appendix.

In dynamic settings,  $\mathcal{B}_f$  equals zero if we can prove that the learned representation removes the association between dynamic features and the ad exposure. We build the representation  $\mathbf{v}_t^{rev}$  invariant across different ad channels to eliminate biases caused by dynamic user features. We achieve this by minimizing the formula  $\mathcal{L}_{rev} = \sum_{i=1}^N \sum_{t=1}^{T_i} CE(\mathbf{c}_t^{rev}, \mathbf{c}_t)$  in Eq.(4). We can prove that

$$\mathcal{L}_{rev} = K \cdot JSD(p(\mathbf{v}_t^{rev}|\mathbf{c}_1), \dots, p(\mathbf{v}_t^{rev}|\mathbf{c}_K)) - K \log K,$$

where  $K \log K$  is a constant, and  $JSD(\cdot, \dots, \cdot)$  denotes the multi-distribution Jensen-Shannon Divergence [17], which is non-negative and 0 if and only if all distributions are equal. To minimize  $\mathcal{L}_{rev}$ , we derive  $p(\mathbf{v}_t^{rev}|\mathbf{c}_1) = \dots = p(\mathbf{v}_t^{rev}|\mathbf{c}_K)$ , where  $\mathbf{v}_t^{rev}$  is the learned representation invariant across different ad channels. For details, see the appendix.

## 5 EXPERIMENT

In this section, we evaluate the performance of CausalMTA and answer the following questions:

- **Q1:** What is the performance of CausalMTA in terms of eliminating confounding bias?
- **Q2:** Does CausalMTA outperform the state-of-the-art MTA methods in conversion prediction?
- **Q3:** How does CausalMTA perform on real-world ad impression datasets?
- **Q4:** What are the capabilities of the journey re-weighting module and the causal conversion prediction module?

### 5.1 Experimental Settings

This section provides an overview of the data, experimental protocol, evaluation metrics, compared baselines, and hyperparameter settings. More details of this part can be found in the Appendix A.2. All the code and data are available in the supplementary file, and will be released after acceptance.

**5.1.1 Data Descriptions.** On three datasets, the performance of CausalMTA is evaluated. The first dataset is a synthetic dataset created to test CausalMTA's ability to solve the issue of confounding bias. The second, **Criteo**<sup>1</sup>, is a publicly available dataset on ad bidding that is commonly utilized in MTA [9, 16, 21]. We use the same experimental setup as CAMTA[16] to preprocess it. More details of the preprocessing are specified in the Appendix A.2.1. The third dataset is a real ad impression dataset from **Alibaba**, which

<sup>1</sup><http://ailab.criteo.com/criteo-attribution-modeling-bidding-dataset/>

**Table 1: Results of conversion prediction on synthetic dataset. MTA-UB is the upper-bound performance trained on unbiased data.**

Method	AUC	Log-loss	AUC	Log-loss	AUC	Log-loss
	dynamic-only		static-only		hybrid	
LR	0.6256±0.02	0.8414±0.001	0.6181±0.01	0.8781±0.002	0.5883±0.03	1.1226±0.002
SP	0.5731±0.00	1.0712±0.00	0.5514±0.00	1.3361±0.00	0.5210±0.00	1.8853±0.00
AH	0.6328±0.01	0.7942±0.001	0.6231±0.02	0.8515±0.002	0.5832±0.01	1.1136±0.002
DNAMTA	0.6497±0.02	0.6795±0.001	0.6465±0.02	0.6624±0.002	0.6147±0.03	0.6497±0.002
DeepMTA	0.6519±0.02	0.6778±0.001	0.6427±0.02	0.6678±0.002	0.6073±0.03	0.6519±0.002
CAMTA	0.6926±0.02	0.6583±0.001	0.6531±0.01	0.6927±0.002	0.6485±0.03	0.6872±0.002
CausalMTA	0.7034±0.01	0.6472±0.002	0.6882±0.02	0.6521±0.003	0.6814±0.01	0.6424±0.003
MTA-UB	0.7268±0.01	0.6454±0.002	0.7205±0.01	0.6353±0.002	0.7116±0.01	0.6285±0.002

**Table 2: The overview of the Criteo dataset**

Statistics	Raw	Processed
No. of users	6,142,256	157,331
No. of campaigns	675	10
No. of journeys	6,514,319	196,560
No. of convert journeys	435,810	19,890
No. of touchpoints	16,468,027	787,483

includes 30 days of ad impression data from mobile phone shops. These touchpoints are categorized into 40 channels, including interact, feed, display, search, live show, *etc.*

**5.1.2 Experiment Protocol.** To evaluate CausalMTA’s performance, we conduct experiments on three datasets. For synthetic dataset, we simulate various confounder settings to obtain the biased and unbiased data, and we quantify the ability of CausalMTA to eliminate confounding bias quantitatively with conversion prediction (Section 5.2). For Criteo dataset, we compare CausalMTA’s performance to that of the state-of-the-art methods under two tasks, *i.e.*, conversion prediction and data reply (Section 5.3). We also report the experimental results of CausalMTA on Alibaba advertising platforms, which contains the attribution value analysis and profit comparison (Section 5.4). Moreover, we provide the ablation studies for verifying the effectiveness of the proposed journey re-weighting and causal conversion prediction modules (Section 5.5).

**5.1.3 Evaluation Metrics.** For conversion prediction, we evaluate the performance in terms of **log-loss**, **AUC**. For fairness, the log-loss only contains the conversion prediction part of Equation 4. It can be reckoned as a standard measurement to estimate the classification performance. AUC can be a metric reflecting the pairwise ranking performance of the estimation between converted and non-converted ad impression sequences.

For data reply experiments, we follow the work in [21] and utilize the **return on investment** (ROI) as the metric to compute the overall budget allocation. Then, the historical data are selected to fit the budget. We compared the performance of different budget allocation in back evaluation with two metrics, *i.e.*, Cost per Action (CPA) and Conversion Rate (CVR). CPA is the total monetary cost normalized by the number of conversions, which measures the efficiency of advertising campaign. And CVR is the number

of converted sequences averaged by number of testing sequences, which reflects the ratio of gain for the ad exposure.

**5.1.4 Compared Methods.** In our experiments, CausalMTA is compared with 8 baseline methods which can be divided into three categories, *i.e.*, statistical learning-based methods, deep learning-based methods, and causal learning-based methods. The statistical learning-based methods consist of three methods, *i.e.*, Logistic Regression [23] (LR), Simple Probabilistic [8] (SP), and Additive Hazard [32] (AH). Deep learning-based methods contain three methods, *i.e.*, DNAMTA [3], DARNN [21], and DeepMTA [30]. The causal learning-based methods also have two works, *i.e.*, JDMTA [10] and CAMTA [16]. Besides, we also compare CausalMTA with two ablation methods, *i.e.*, CM-rw and CM-causal. Detailed descriptions of these methods are available in the Appendix A.2.2.

**5.1.5 Parameter Settings.** For LSTMs in CausalMTA, we stack three 3 layer LSTMs as the encoder, decoder, and the predictor respectively. MLP models in CausalMTA are composed of 4 fully connected layers with *ELU* as the activate function. CausalMTA has 4 hyperparameters *i.e.*,  $\alpha$ ,  $\beta$ ,  $\gamma$  and  $\delta$ , we empirically set  $\alpha = \beta = 0.5$ , and  $\gamma = \delta = 0.5$ . All the experiments are conducted on a high-end server with 2× NVIDIA GTX3090 GPUs. All the compared baselines are trained in 30 epochs and the best model is chosen to report.

## 5.2 Experiments on Synthetic Data

To answer Q1, we conduct the synthetic experiment. Next, we introduce the generation process of synthetic data and report the experimental results respectively.

**5.2.1 Data Generation.** The data generation procedure is composed of the ad exposure policy and the user conversion module. For the ad exposure policy, we first generate the sequence of exposure events with a Poisson process. Then, a stochastic function is designed to assign the ad types for the events. The parameters of both Poisson process and the stochastic function are related to the user preference. As for user conversions, we follow the work in [26] which set the conversion probabilities of all ad types as a function of user demographics. The conversion probability of a specified ad journey can be calculated by aggregating the probabilities of all related ads. The detailed data generation procedure can be found in the Appendix A.3.



**Table 3: Results of conversion prediction on Criteo dataset. SL, DL, CL in the left column indicate the statistical learning-based methods, deep learning-based methods and causal learning-based methods, respectively.**

	Method	AUC	Log-loss
SL	LR	0.8370±0.03	0.191±0.01
	SP	0.7616±0	0.317±0
	AH	0.7264±0.03	0.286±0.01
DL	DNAMTA	0.9127±0.02	0.1360±0.013
	DARNN	0.8726±0.02	0.165±0.006
	DeepMTA	0.9104±0.03	0.112±0.012
CL	JDMTA	0.9127±0.01	0.0838±0.007
	CAMTA	0.9347±0.02	0.0715±0.007
Ours	CausalMTA	<b>0.9659±0.01</b>	<b>0.0517±0.003</b>
	CM-RW	0.9539±0.01	0.0560±0.003
	CM-CAUSAL	0.9517±0.01	0.0534±0.002

Following the data generation process, we construct three sub-datasets, *i.e.*, dynamic-only, static-only and hybrid, which reflect different kinds of confounding bias. For dynamic-only, the user preference only influence the generation of exposure sequences and the ad types are randomly assigned. For static-only, the user preference only influence the ad types and the sequence is generated by a fixed Poisson distribution. For hybrid, the user preference influence both. Moreover, we also construct unbiased training datasets to train prediction models and treat them as the performance upper-bound.

**5.2.2 Performance comparison.** As shown in Table 1, CausalMTA significantly outperforms the compared method in three confounder settings. The performances are very close to the performance upper-bound, which indicates CausalMTA is able to alleviate the confounding bias in MTA. By considering the casual relationship, CAMTA achieves the best performance among all the competitors, but is also inferior to CausalMTA. This phenomenon verifies the importance of modeling the causal mechanism for MTA.

Comparing the results of different settings, we observe CausalMTA achieves the best performance in hybrid confounders. It proves the effectiveness of the proposed techniques. The performance of CausalMTA in dynamic-only and static-only settings does not degrade significantly, indicating that it’s applicable to a wide range.

### 5.3 Experiments on Criteo Dataset

To answer Q2, we employ the most widely used public dataset Criteo to evaluate the performance of CausalMTA. In this section, we first briefly describe the tasks and then specify the results.

**5.3.1 Task Description.** Two tasks, *i.e.* conversion prediction and data replay, are conducted to test the performance of the unbiased prediction model and the attribution weights, respectively. For conversion prediction, we directly train CausalMTA and baselines under the same setting, and compare the prediction accuracy. For data replay, we follow the experiments in [21] and utilize the attribution credits to compute the **return on investment (ROI)** and budget allocation on different ads. Based on the result, we conduct back evaluation for budget allocation, and measure the CPA and CVR of them.

**5.3.2 Performance of Conversion Prediction.** The detailed evaluation results of conversion prediction on different baselines are presented in Table 3. As shown, CausalMTA continuously outperforms all the compared baselines, which proves the validity of eliminating the confounding bias on static user attributes and dynamic features. CAMTA is the strongest baseline but also inferior to CausalMTA. It utilizes click labels as the auxiliary information, which probably involves additional confounding bias. Moreover, CAMTA does not consider the difference between static and dynamic confounders, which would also harm the performance. One interesting phenomenon is CausalMTA has a more stable confidence interval compared to other baselines. It indicates that the parameters in CausalMTA tend to converge to similar values with different initialization. To a certain extent, CausalMTA is more robust than other baselines.

Comparing the performance of different categories of methods, we can observe that SL methods are inferior to the other two categories. SL methods either use statistical laws or employ logistic regression to predict the conversion probability, which can not well model the conversion process. DL methods perform better than the SL methods but are also inferior to the CL methods. It proves that the deep learning techniques are more suitable for conversion prediction due to their large parameter space and high ability to model complex tasks. However, these methods have poor performance compared to the CL methods. It is because deep learning methods directly use the observed data to train the prediction models, which are incapable of handling confounding bias and would suffer from the out-of-distribution problem. CL methods outperform other baselines with a large margin, which demonstrates the prediction performance highly increased by eliminating the confounding bias.

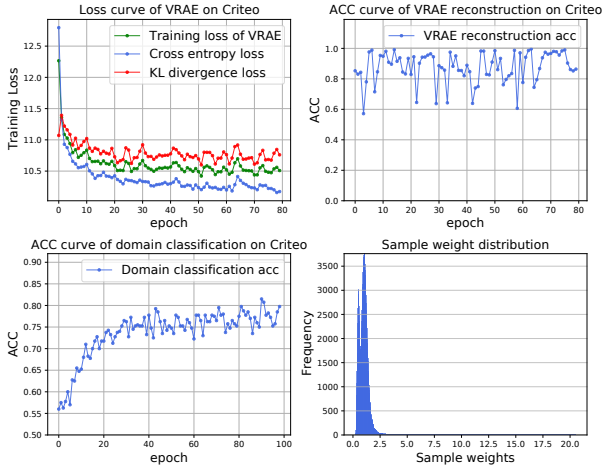
**5.3.3 Performance of Data Replay.** In this experiment, we evaluate CausalMTA and the compared baselines under four proportions of total budgets, *i.e.*, 1/2, 1/4, 1/8, 1/16. Following the setting of [16], we use value  $V(y^i) = 1$  for computing  $ROI_{ck}$  and scale the cost by 1000 to highlight the difference of CPA in the comparison. The detailed results are shown in Table 4. We can observe that: (1) In most cases, CausalMTA achieves the best performance, which verifies the effectiveness of the proposed methods. (2) The conversion number of DeepMTA is better than our method in 1/8 and 1/16 budget, but inferior to CausalMTA in 1/2 and 1/4 budget. This phenomenon shows that CausalMTA captures the intrinsic characters of attribution and performs better in large budget. (3) Among all the baselines, CAMTA is the strongest competitor. It indicates the performance can be improved by involving deconfounding mechanism. (4) Deep learning methods significantly outperform traditional methods. Benefiting from the high complexity, deep learning techniques are more suitable for modelling the user conversion.

### 5.4 Empirical Applications in Alibaba

To answer Q3, we evaluate the performance of CausalMTA on Alibaba advertising platform. Channel attributions of each shop are more meaningful to guide the budget allocation. We train the attribution models in the first 15 days and use the last data for testing. We choose all of its converted journeys in the test set for one specific shop and compute the mean credits of 40 channels. After that, we employ two experiments, *i.e.*, attribution improvement and offline data replay, to evaluate the performance of CausalMTA.

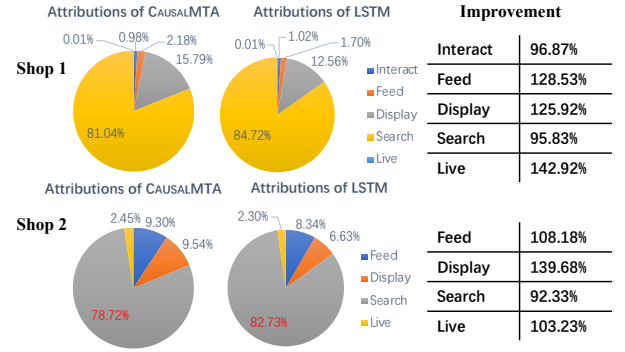
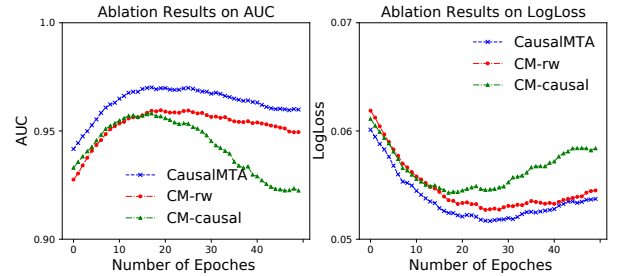
**Table 4: The results of data replay experiment on Criteo dataset. CPA: the lower the better; CVR: the higher the better.**

Method	CPA(Cost pre action)				Conversion Number				CVR(Conversion rate)			
	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$
LR	58.41	56.43	55.26	55.25	913	627	342	193	0.0835	0.0879	0.0946	0.0823
SP	49.60	46.15	48.64	45.39	842	548	375	207	0.0789	0.0772	0.0914	0.0910
AH	51.66	47.30	47.67	52.65	843	594	382	135	0.0804	0.0826	0.0895	0.0775
DNAMTA	38.79	35.12	31.47	32.35	1181	778	459	264	0.1039	0.1068	0.1131	0.1148
DARNN	32.62	30.46	28.09	28.72	1286	829	480	244	0.1218	0.1237	0.1241	0.1223
DeepMTA	36.25	30.60	26.08	25.97	1372	880	<b>549</b>	<b>289</b>	0.1194	0.1202	0.1236	0.1249
CAMTA	32.61	29.73	<b>26.05</b>	26.25	1270	864	538	211	0.1127	0.1160	0.1191	0.1166
CausalMTA	<b>30.34</b>	<b>29.52</b>	26.45	<b>25.47</b>	<b>1441</b>	<b>976</b>	548	255	<b>0.1247</b>	<b>0.1265</b>	<b>0.1305</b>	<b>0.1283</b>

**Figure 3: Learning curves and sample weight distribution on Criteo dataset.****Table 5: Profit comparison of data replay.**

Method	1/2	1/4	1/8	1/16
Shop 1				
LSTM	69.8	63.7	56.2	58.3
CausalMTA	72.3	70.2	57.1	59.2
Improvement	+3.58%	+10.20%	+1.60%	+1.54%
Shop 2				
LSTM	20.9	18.3	17.5	15.2
CausalMTA	21.1	19.1	17.8	15.8
Improvement	+0.96%	+4.37%	+1.71%	+3.95%

**5.4.1 Attribution Improvement.** We compute the attribution credits of two representative cellphone shops utilizing CausalMTA and compare them with the credits calculated by an LSTM-based conversion prediction model. As shown in Figure 5, the credits on both shops decreased, indicating that the estimated contribution of search ads is reduced after eliminating the confounding bias of user preferences. This is because user tends to search the goods before paying. The attributions of the search channel are usually overestimated, and CausalMTA can mitigate this kind of bias.

**Figure 4: The comparison of attribution changes for two cellphone shops.****Figure 5: Performance comparison of ablation methods.**

**5.4.2 Performance of Data Replay.** In this experiment, we employ the attribution credits to guide the budget allocation. Based on the attribution credits, we first compute the return-on-investment (ROI) of each channel and utilize the normalized weights of ROI as budget proportion [21]. Assuming that the total cost of the test set is  $cost_t$ , we set the evaluation budgets as  $1/2, 1/4, 1/8, 1/16$  of  $cost_t$  and replay the historical data to select journeys satisfying evaluation budgets. Table 5 shows the comparison results of the profit in each evaluation budget. We can observe that the profit CausalMTA consistently outperforms the LSTM-based predictor on all evaluation budgets, which indicates that the attribution credits of CausalMTA reflect the causal relationships in advertising. It can be used to guide budget allocation and achieve better profit.



## 5.5 Ablation Studies

To explore the effectiveness of the proposed methods, we compare CausalMTA with two ablation methods on Cirteo dataset, *i.e.*, CM-rw and CM-CAUSAL, which remove the reweighting procedure and the gradient reverse layer respectively. We first show the intermediate results of journey reweighting. The reconstruction accuracy of VRAE and the classification accuracy of the domain classifier directly influence the performance of CausalMTA. As shown in Figure 3, the reconstruction accuracy is approximate to 98%, and the classification accuracy is approximate to 82%, which indicates that VRAE and domain classifier are well trained, and results of journey reweighting are significant. We also witness the reconstruction accuracy fluctuates at a high level as KL divergence dominates the training loss when the cross-entropy loss is small enough.

We summarize the metrics of ablations in Table 3 and illustrate the training procedure in Figure 5. As shown, the AUCs of CM-rw and CM-CAUSAL are inferior to CausalMTA, which proves that both the journey reweighting and causal conversion prediction help improve the performance. By removing the journey reweighting model, the performance of CausalMTA decreases from 0.9659 to 0.9539. By removing the gradient reverse layer, the performance of CausalMTA decreases from 0.9659 to 0.9617. We can observe the improvement of gradient reverse layer is more significant than journey reweighting, which indicates the confounding bias of dynamic feature are more obvious than the static user attributes. This result is consistent with our cognition. Moreover, as illustrated in Figure 5, the convergence speed of CausalMTA is faster than CM-rw and CM-CAUSAL, which shows the superiority of CausalMTA in eliminating the confounding bias of user preferences.

## 6 CONCLUSION

In this paper, we define the problem of causal MTA, which eliminates the confounding bias introduced by user preferences and assigns the attribution credits fairly over all touchpoints. We propose CausalMTA, which decomposes the confounding bias of user preferences into two independent parts, *i.e.*, the static user attributes and the dynamic features. CausalMTA employs journey reweighting and causal conversion prediction to solve these two kinds of confounding bias respectively. We prove that CausalMTA is capable of generating unbiased conversion predictions of ad journey. Extensive experiments on the public dataset and real commercial data from Alibaba show that CausalMTA outperforms all the compared baselines and works well in the real-world application.

## ACKNOWLEDGMENTS

This work has been supported by the National Natural Science Foundation of China under Grant No.: 62002343, 62077044, 61702470.

## REFERENCES

- [1] Ahmed M. Alaa and Mihaela van der Schaar. 2017. Bayesian Inference of Individualized Treatment Effects using Multi-task Gaussian Processes. In *NeurIPS'17*. 3424–3432.
- [2] Ahmed M. Alaa and Mihaela van der Schaar. 2018. Limits of Estimating Heterogeneous Treatment Effects: Guidelines for Practical Algorithm Design. In *ICML '18 (Proceedings of Machine Learning Research)*, Vol. 80. PMLR, 129–138.
- [3] Sai Kumar Arava, Chen Dong, Zhenyu Yan, Abhishek Pani, et al. 2018. Deep neural net with attention for multi-channel multi-touch attribution. *arXiv:1809.02230* (2018).
- [4] Peter C Austin. 2011. An introduction to propensity score methods for reducing the effects of confounding in observational studies. *Multivariate behavioral research* 46, 3 (2011), 399–424.
- [5] Joel Barajas, Ram Akella, Marius Holtan, and Aaron Flores. 2016. Experimental designs and estimation for online display advertising attribution in marketplaces. *Marketing Science* 35, 3 (2016), 465–483.
- [6] Ron Berman. 2018. Beyond the last touch: Attribution in online advertising. *Marketing Science* 37, 5 (2018), 771–792.
- [7] Ioana Bica, Ahmed M. Alaa, James Jordon, and Mihaela van der Schaar. 2020. Estimating counterfactual treatment outcomes over time through adversarially balanced representations. In *ICLR'20*. OpenReview.net.
- [8] Brian Dalesandro, Claudia Perlich, Ori Stitelman, and Foster Provost. 2012. Causally motivated attribution for online advertising. In *Proceedings of the sixth international workshop on data mining for online advertising and internet economy*. 1–9.
- [9] Eustache Diemert, Julien Meynet, Pierre Galland, and Damien Lefortier. 2017. Attribution Modeling Increases Efficiency of Bidding in Display Advertising. In *Proceedings of KDD Workshop ADKDD'17*. ACM, 2:1–2:6.
- [10] Ruihuan Du and etc. 2019. Causally driven incremental multi touch attribution using a recurrent neural network. *arXiv:1902.00215* (2019).
- [11] Otto Fabius, Joost R. van Amersfoort, and Diederik P. Kingma. 2015. Variational Recurrent Auto-Encoders. In *ICLR'15*.
- [12] Christian Fong, Chad Hazlett, Kosuke Imai, et al. 2018. Covariate balancing propensity score for a continuous treatment: application to the efficacy of political advertisements. *The Annals of Applied Statistics* 12, 1 (2018), 156–177.
- [13] Wendi Ji and Xiaoling Wang. 2017. Additional multi-touch attribution for online advertising. In *AAAI'17*, Vol. 31.
- [14] Fredrik Johansson, Uri Shalit, and David Sontag. 2016. Learning representations for counterfactual inference. In *ICML '16*. PMLR, 3020–3029.
- [15] Fredrik D Johansson, Nathan Kallus, Uri Shalit, and David Sontag. 2018. Learning weighted representations for generalization across designs. *arXiv:1802.08598* (2018).
- [16] Sachin Kumar, Garima Gupta, Ranjitha Prasad, Arnab Chatterjee, Lovekesh Vig, and Gautam Shroff. 2020. CAMTA: Causal Attention Model for Multi-touch Attribution. In *ICDM'20*. IEEE, 79–86.
- [17] Ya Li, Xinmei Tian, Mingming Gong, Yajing Liu, Tongliang Liu, Kun Zhang, and Dacheng Tao. 2018. Deep Domain Generalization via Conditional Invariant Adversarial Networks. In *ECCV'18 (Lecture Notes in Computer Science)*, Vol. 11219. Springer, 647–663.
- [18] Bryan Lim. 2018. Forecasting Treatment Responses Over Time Using Recurrent Marginal Structural Networks. In *NeurIPS'18*. 7494–7504.
- [19] Harikesh S Nair, Sanjog Misra, William J Hornbuckle IV, Ranjan Mishra, and Anand Acharya. 2017. Big data and marketing analytics in gaming: Combining empirical models and field experimentation. *Marketing Science* 36, 5 (2017), 699–725.
- [20] Judea Pearl. 2009. *Causality*. Cambridge university press.
- [21] Kan Ren and etc. 2018. Learning Multi-touch Conversion Attribution with Dual-attention Mechanisms for Online Advertising. In *CIKM 2018*. ACM, 1433–1442.
- [22] Jason Roy, Kirsten J Lum, and Michael J Daniels. 2017. A Bayesian nonparametric approach to marginal structural models for point treatments and a continuous or survival outcome. *Biostatistics* 18, 1 (2017), 32–47.
- [23] Xuhui Shao and Lexin Li. 2011. Data-driven multi-touch attribution models. In *KDD'11*. ACM, 258–264.
- [24] Lloyd S Shapley. 1953. A value for n-person games. *Contributions to the Theory of Games* 2, 28 (1953), 307–317.
- [25] Dinah Shender, Ali Nasiri Amini, Xinlong Bao, Mert Dikmen, Amy Richardson, and Jing Wang. 2020. A Time To Event Framework For Multi-touch Attribution. *arXiv:2009.08432* (2020).
- [26] Dinah Shender, Ali Nasiri Amini, Xinlong Bao, Mert Dikmen, Amy Richardson, and Jing Wang. 2020. A Time To Event Framework For Multi-touch Attribution. *arXiv:2009.08432* (2020).
- [27] Raghav Singal and etc. 2019. Shapley meets uniform: An axiomatic framework for attribution in online advertising. In *WWW*. 1713–1723.
- [28] Yongjun Xu, Xin Liu, Xin Cao, and etc. 2021. Artificial intelligence: A powerful paradigm for scientific research. *The Innovation* 2, 4 (2021), 100179.
- [29] Yanbo Xu, Yanxun Xu, and Suchi Saria. 2016. A Bayesian nonparametric approach for estimating individualized treatment-response curves. In *Machine Learning for Healthcare Conference*. PMLR, 282–300.
- [30] Dongdong Yang, Kevin Dyer, and Senzhang Wang. 2020. Interpretable Deep Learning Model for Online Multi-touch Attribution. *arXiv:2004.00384* (2020).
- [31] Di Yao, Chao Zhang, Jian-Hui Huang, and Jingping Bi. 2017. SERM: A Recurrent Model for Next Location Prediction in Semantic Trajectories. In *CIKM*. ACM, 2411–2414.
- [32] Ya Zhang, Yi Wei, and Jianbiao Ren. 2014. Multi-touch attribution in online advertising with survival theory. In *ICDM'14*. IEEE, 687–696.
- [33] Hao Zou, Peng Cui, Bo Li, Zheyang Shen, Jianxin Ma, Hongxia Yang, and Yue He. 2020. Counterfactual Prediction for Bundle Treatment. *NeurIPS'20* 33 (2020).

## A APPENDIX

### A.1 Theoretical Analysis of CausalMTA

**A.1.1 Proof of weights calculation.** In order to create a pseudo-population which debiases by means of sample re-weighting, the weights should cater to the equation  $W_T(\mathbf{U}, \mathbf{C}) = p(\mathbf{C})/p(\mathbf{C}|\mathbf{U})$  [12, 33]. When we learn a variational distribution  $q_\phi(\mathbf{z}|\mathbf{c})$  of the original channel assignment of touch-points, the variational sample weights can be computed as follows:

$$\begin{aligned} w^i &= W_T(\mathbf{u}^i, \mathbf{c}^i) = \frac{p(\mathbf{c}^i)}{p(\mathbf{c}^i|\mathbf{u}^i)} \\ &= \frac{p(\mathbf{c}^i)}{\int_{\mathbf{z}} p(\mathbf{c}^i|\mathbf{z}) p(\mathbf{z}|\mathbf{u}^i) d\mathbf{z}} = \frac{1}{\int_{\mathbf{z}} \frac{p(\mathbf{c}^i|\mathbf{z})}{p(\mathbf{c}^i)} p(\mathbf{z}|\mathbf{u}^i) d\mathbf{z}} \\ &= \frac{1}{\int_{\mathbf{z}} \frac{p(\mathbf{z}|\mathbf{c}^i)}{p(\mathbf{z})} p(\mathbf{z}|\mathbf{u}^i) d\mathbf{z}} = \frac{1}{\int_{\mathbf{z}} \frac{p(\mathbf{z}|\mathbf{u}^i)}{p(\mathbf{z})} p(\mathbf{z}|\mathbf{c}^i) d\mathbf{z}} \\ &= \frac{1}{\int_{\mathbf{z}} \frac{p(\mathbf{z}, \mathbf{u}^i)}{p(\mathbf{z}) p(\mathbf{u}^i)} p(\mathbf{z}|\mathbf{c}^i) d\mathbf{z}} = \frac{1}{\int_{\mathbf{z}} \frac{1}{W_Z(\mathbf{u}^i, \mathbf{z})} p(\mathbf{z}|\mathbf{c}^i) d\mathbf{z}} \\ &= \frac{1}{\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{c}^i)} \left[ \frac{1}{W_Z(\mathbf{u}^i, \mathbf{z})} \right]}, \end{aligned}$$

where  $W_Z(\mathbf{U}, \mathbf{Z})$  can be viewed as the density ratio estimation to decorrelate  $\mathbf{U}$  and  $\mathbf{Z}$  for points in space  $\mathcal{U} \times \mathcal{Z}$ .

**A.1.2 Proof of the journal reweighting module.** In the task of multi-touch attribution, provided with observational data, we hope to learn a hypothesis  $f_{\theta_p} : \mathcal{U} \times \mathcal{C} \mapsto \mathbb{R}$  with model parameters  $\theta_p$ , which predicts the conversion rate based on the confounders and touch-points. In this setting, the concept of counterfactual is to guarantee the learned hypothesis to predict accurate outcome when the assignment of touch-point (e.g., the channel preference) is random. For the individual  $\mathbf{U}$ , when  $\mathcal{L}(\cdot)$  denotes the error function and  $y(\cdot)$  denotes the ground-truth outcome, the prediction error can be formed as:

$$\mathcal{E}(\mathbf{U}) = \mathbb{E}_{\mathbf{C} \sim p(\mathbf{C})} \left[ \mathcal{L} \left( f_{\theta_p}(\mathbf{U}, \mathbf{C}), y(\mathbf{U}, \mathbf{C}) \right) \right].$$

The target of the counterfactual prediction error to be minimized is  $\mathcal{E}_{cf} = \mathbb{E}_{\mathbf{U} \sim p(\mathbf{U})} [\mathcal{E}(\mathbf{U})]$ . But in the observational dataset, the touch-points are assigned based on confounders, i.e.,  $\mathbf{C} \sim p(\mathbf{C}|\mathbf{U})$ . Instead of directly using supervised learning, optimizing the prediction error on the re-weighted data

$$\mathcal{E}_f^w = \mathbb{E}_{\mathbf{U}, \mathbf{C} \sim p(\mathbf{U}, \mathbf{C})} \left[ \mathcal{L} \left( f_{\theta_p}(\mathbf{U}, \mathbf{C}), y(\mathbf{U}, \mathbf{C}) \right) W_T(\mathbf{U}, \mathbf{C}) \right],$$

can lead to a more accuracy counterfactual prediction.

Assuming a family  $G$  of functions  $g : \mathcal{U} \times \mathcal{C} \mapsto \mathbb{R}$ , and we have  $\mathcal{L}(f(\mathbf{U}, \mathbf{C}), y(\mathbf{U}, \mathbf{C})) = l(\mathbf{U}, \mathbf{C}) \in G$ . We can therefore bridge the gap between the counterfactual loss and the re-weighted loss under

observational data.

$$\begin{aligned} &\mathcal{E}_{cf} - \mathcal{E}_f^w \\ &= \int_{\mathbf{U}} \int_{\mathbf{C}} (p(\mathbf{U})p(\mathbf{C}) - W_T(\mathbf{U}, \mathbf{C})p(\mathbf{U}, \mathbf{C})) \\ &\quad \cdot \mathcal{L}(f(\mathbf{U}, \mathbf{C}), y(\mathbf{U}, \mathbf{C})) d\mathbf{U} d\mathbf{C} \\ &\leq \left| \int_{\mathbf{U}} \int_{\mathbf{C}} (p(\mathbf{U})p(\mathbf{C}) - W_T(\mathbf{U}, \mathbf{C})p(\mathbf{U}, \mathbf{C})) \right. \\ &\quad \cdot \mathcal{L}(f(\mathbf{U}, \mathbf{C}), y(\mathbf{U}, \mathbf{C})) d\mathbf{U} d\mathbf{C} \left. \right| \\ &\leq \sup_{g \in G} \left| \int_{\mathbf{U}} \int_{\mathbf{C}} (p(\mathbf{U})p(\mathbf{C}) - W_T(\mathbf{U}, \mathbf{C})p(\mathbf{U}, \mathbf{C})) \right. \\ &\quad \cdot g(\mathbf{U}, \mathbf{C}) d\mathbf{U} d\mathbf{C} \left. \right| \\ &= IPM_G(W_T(\mathbf{U}, \mathbf{C})p(\mathbf{U}, \mathbf{C}), p(\mathbf{U})p(\mathbf{C})). \end{aligned}$$

When  $W_T(\mathbf{U}, \mathbf{C}) = \frac{p(\mathbf{C})}{p(\mathbf{C}|\mathbf{U})}$ , we have:

$$\begin{aligned} IPM_G(W_T(\mathbf{U}, \mathbf{C})p(\mathbf{U}, \mathbf{C}), p(\mathbf{U})p(\mathbf{C})) &= 0, \\ \mathcal{E}_f^w &= \mathcal{E}_{cf}. \end{aligned}$$

**A.1.3 Proof of the causal prediction module.** The optimal prediction probabilities of ad exposure are given by

$$\mathbf{c}^{rev*} = \arg \max_{\mathbf{c}^{rev}} \sum_{k=1}^K \int_{\mathbf{v}^{rev}} p(\mathbf{v}^{rev}|\mathbf{c}^k) \log(\mathbf{c}_k^{rev}) d\mathbf{v}^{rev}.$$

By maximizing value function and leveraging Lagrange multiplies, we can derive  $\mathbf{c}^{rev*}$  by the following form

$$\arg \max_{\mathbf{c}^{rev}} \sum_{k=1}^K \left( p(\mathbf{v}^{rev}|\mathbf{c}^k) \log(\mathbf{c}_k^{rev}) \right) + \chi \left( \sum_{k=1}^K \mathbf{c}_k^{rev} - 1 \right).$$

We have  $\mathbf{c}_k^{rev} = -\frac{p(\mathbf{v}^{rev}|\mathbf{c}^k)}{\chi} = \frac{p(\mathbf{v}^{rev}|\mathbf{c}^k)}{\sum_{m=1}^K p(\mathbf{v}^{rev}|\mathbf{c}^m)}$  by setting the above derivative to zero and solving  $\mathbf{c}^{rev*}$ .

Therefore, the objective  $\min_{\mathbf{v}^{rev}} \mathcal{L}_{rev}$  for the learned representation  $\mathbf{v}^{rev}$  becomes

$$\min_{\mathbf{v}^{rev}} \sum_{k=1}^K \mathbb{E}_{\mathbf{v}^{rev} \sim p(\mathbf{v}^{rev}|\mathbf{c}_k)} \left[ \log \frac{p(\mathbf{v}^{rev}|\mathbf{c}^k)}{\sum_{m=1}^K p(\mathbf{v}^{rev}|\mathbf{c}^m)} \right].$$

We can derive that

$$\begin{aligned} &\sum_{k=1}^K \mathbb{E}_{\mathbf{v}^{rev} \sim p(\mathbf{v}^{rev}|\mathbf{c}_k)} \left[ \log \frac{p(\mathbf{v}^{rev}|\mathbf{c}^k)}{\sum_{m=1}^K p(\mathbf{v}^{rev}|\mathbf{c}^m)} \right] + K \log K \\ &= \sum_{k=1}^K \left( \mathbb{E}_{\mathbf{v}^{rev} \sim p(\mathbf{v}^{rev}|\mathbf{c}_k)} \left[ \log \frac{p(\mathbf{v}^{rev}|\mathbf{c}^k)}{\sum_{m=1}^K p(\mathbf{v}^{rev}|\mathbf{c}^m)} \right] + \log K \right) \\ &= \sum_{k=1}^K \mathbb{E}_{\mathbf{v}^{rev} \sim p(\mathbf{v}^{rev}|\mathbf{c}_k)} \left[ \log \frac{p(\mathbf{v}^{rev}|\mathbf{c}^k)}{\frac{1}{K} \sum_{m=1}^K p(\mathbf{v}^{rev}|\mathbf{c}^m)} \right] \\ &= \sum_{k=1}^K KL \left( p(\mathbf{v}^{rev}|\mathbf{c}^k) \parallel \frac{1}{K} \sum_{m=1}^K p(\mathbf{v}^{rev}|\mathbf{c}^m) \right) \\ &= K \cdot JSD \left( p(\mathbf{v}^{rev}|\mathbf{c}^1), \dots, p(\mathbf{v}^{rev}|\mathbf{c}^K) \right), \end{aligned}$$

where  $KL(\cdot|\cdot)$  is the KL divergence and  $JSD(\cdot, \dots, \cdot)$  is the Jensen-Shannon Divergence [7, 17] in the multi-distribution form. Because

$JSD(\cdot, \dots, \cdot)$  is non-negative and equals zero when all distributions are equal and  $K \log K$  is a constant, we have that  $p(v_t^{rev}|c_1) = \dots = p(v_t^{rev}|c_K)$  by minimizing  $\mathcal{L}_{rev}$ .

## A.2 Experiments

**A.2.1 Details of data preprocessing.** The attribution modeling for bidding dataset published by **Criteo** company is widely deployed in the research of modeling user behavior and ad attribution [9, 16, 21]. As a sample of 30 days of Criteo live traffic data, this dataset has more than 16 million ad impression records and 45 thousand conversions over 700 ad campaigns. Each ad impression record contains items such as timestamp, user id, ad campaign, and side information. There is also a label denotes whether a click action has occurred, and the corresponding conversion ID if this sequence of ad impressions finally leads to a conversion. We preprocess the raw Criteo dataset in the following procedures: (i) we count the top 10 campaigns with the largest number of ad impression records and filter out the ad impression records corresponding to other campaigns; (ii) we group the ad impression entries, which have the same user id and conversion id, into the same sequence and sort each sequence by timestamp; (iii) for a conversion id of -1, i.e., for a specific user, a group of ad impression records that did not cause the user to convert, the original group is divided at a time interval of 3 days based on timestamp; (iv) we filter out ad sequences that are less than 3 in length; (v) we divide it into the train set and the test set, and ensure the set of user id in the test set is a subset of user id in the train set.

**A.2.2 Details of compared baselines.** We compare CausalMTA with four kinds of baselines, i.e., statistical learning-based methods (SL), deep learning-based methods (DL), causal learning-based methods (CL) and ablations.

The statistical learning-based methods consist of three works:

- **LR** (Logistic Regression) model for ad attribution is proposed by Shao and Li [23], in which channel's attribution values are calculated as the learned coefficients.
- **SP** (Simple Probabilistic) model calculates the conversion rate taking into the conversion probability from the observed data into account. As in [8], the conversion rate is

$$p(y = 1 | \{c_j\}_{j=1}^{m_i}) = 1 - \prod_j (1 - \Pr(y = 1 | c_j = k)).$$

- **AH** (Additive Hazard) proposed by Zhang *et al.* [32] is the first user conversion estimation model based on survival analysis and additive hazard function.

The deep learning-based methods also consist of three works:

- **DNAMTA** is the Deep Neural Net with Attention Multi-touch Attribution model proposed by Arava *et al.* [3]. It leverages LSTM and attention mechanism to model the dynamic interaction between ad channels, and incorporates user-context information to reduce estimation bias.
- **DARNN** is the Dual-Attention Recurrent Neural Network proposed in [21] which uses dual-attention RNNs to combine both post-view and post-click attribution patterns together for the user conversion estimation.
- **DeepMTA** is a phased-LSTM based model [30] which combines deep neural networks and additive feature explanation model for

interpretable online multi-touch attribution. For fair comparison, we replace the phased-LSTM with vanilla LSTM.

The causal learning-based methods consist of two works:

- **JDMTA** is a causal-inspired model [10] which employs Shapley Value to compute the attribution credits for touchpoints.
- **CAMTA** is the Causal Attention Model for Multi-touch Attribution proposed by Kumar *et al.* [16]. This model leverages counterfactual recurrent network to minimize selection bias in channel assignment while conducting conversion estimation.

We also compare CausalMTA with its two ablations:

- **CM-rw** removes the journey reweighting module and treats all journeys equally. It only employs the proposed causal conversion prediction model for MTA.
- **CM-CAUSAL** replaces the causal RNN predictor with traditional RNN and only utilize the reweighting mechanism to eliminate the confounding bias of static user attributes.

## A.3 Generation of Synthetic Data

To generate synthetic data, we simulate the ad exposure policy and user conversions.

**A.3.1 Details of ad exposure policy.** Ad delivery involves two issues: serving time and advertising channels. To generate the time series data for ad exposure, we can first use a Poisson process. Then, to assign the ad types for the events, we create a stochastic function. In the dynamic-only setting, the intensity rate  $\lambda_{exp}$  is a function of the user preference while the selection of ad types is independent of user preference. In the static-only setting, user characteristics are parameters of the stochastic function. The user feature has an impact on both aspects of ad exposure in a hybrid setting. We also simulate ad sequences whose generation is invariant across different user preferences to create unbiased data.

**A.3.2 Details of user conversion module.** The user conversion behaviour in multi-touch attribution can be viewed as occurrences in an inhomogeneous Poisson process [26]. In detail, we leverage a realization of a Poisson counting process combined with a time-varying intensity function,  $\lambda(t)$ . This process can be formulated as

$$Y_i(t) - Y_i(s) \sim \text{Poisson}(\int_s^t \lambda(t) dt),$$

where  $Y_i(t)$  is the number of conversion events for customer  $i$  up until time  $t$ . As the user conversion is the effect of both ad exposures and user characteristics, we use a log-linear model for the intensity function, and allow it to depend on the impact of the previous seen ads and user features, then we have

$$\log(\lambda(t)) = \alpha_0 + \alpha_{user} + \sum_{j,k} g_k(t - t_j),$$

where  $\alpha_0$  represents the log of conversion rate before any ads are shown regardless of user preferences, and  $\alpha_{user}$  stands for the impact of user features. The item  $g_k(t - t_j)$  models the impact of an ad exposure of channel  $k$  that occurs at time  $t_j$ , which brings the jump in conversions and has an exponentially decaying effect.