
CAUSAL KNOWLEDGE GUIDED SOCIETAL EVENT FORECASTING

Songgaojun Deng
Stevens Institute of Technology
Hoboken, NJ
sdeng4@stevens.edu

Huzefa Rangwala
George Mason University
Fairfax, VA
rangwala@gmu.edu

Yue Ning
Stevens Institute of Technology
Hoboken, NJ
yue.ning@stevens.edu

ABSTRACT

Data-driven societal event forecasting methods exploit relevant historical information to predict future events. These methods rely on historical labeled data and cannot accurately predict events when data are limited or of poor quality. Studying causal effects between events goes beyond correlation analysis and can contribute to a more robust prediction of events. However, incorporating causality analysis in data-driven event forecasting is challenging due to several factors: (i) Events occur in a complex and dynamic social environment. Many unobserved variables, i.e., hidden confounders, affect both potential causes and outcomes. (ii) Given spatiotemporal non-independent and identically distributed (non-IID) data, modeling hidden confounders for accurate causal effect estimation is not trivial. In this work, we introduce a deep learning framework that integrates causal effect estimation into event forecasting. We first study the problem of Individual Treatment Effect (ITE) estimation from observational event data with spatiotemporal attributes and present a novel causal inference model to estimate ITEs. We then incorporate the learned event-related causal information into event prediction as prior knowledge. Two robust learning modules, including a feature reweighting module and an approximate constraint loss, are introduced to enable prior knowledge injection. We evaluate the proposed causal inference model on real-world event datasets and validate the effectiveness of proposed robust learning modules in event prediction by feeding learned causal information into different deep learning methods. Experimental results demonstrate the strengths of the proposed causal inference model for ITE estimation in societal events and showcase the beneficial properties of robust learning modules in societal event forecasting.

Keywords Event Forecasting · Representation Learning · Causal Inference

1 Introduction

Predicting large-scale societal events, such as disease outbreaks, organized crime, and civil unrest movements, from social media streams, web logs, and news media is of great significance for decision-making and resource allocation. Previous methods mainly focus on improving the predictive accuracy of a given event type or multiple event types using historical event data [1, 2]. Recently, to enhance model explainability, many approaches identify salient features or supporting evidence, such as precursor documents [3], relationships represented as graphs [4], and major actors participating in the events [5]. However, existing work explains the occurrence of events based on correlation-based indicators.

Attempts to study causality in event analysis and prediction have focused on extracting pairs of causal events from unstructured text [6], or using human-defined causally related historical events to predict events of interest [7]. Causal effect learning has shown advantages in improving predictions in various machine learning problems, such as recommender systems [8], disease diagnosis prediction [9], and computer vision tasks [10]. This suggests the potential of causal effect learning for better prediction of societal events. Leveraging causal effects can presumably provide new insights into causal-level interpretation and improve the robustness of event prediction, e.g., less susceptible to noise in data. In this study, we explore societal event forecasting methods with the help of causal effect learning.

Traditionally, learning causal effects (aka treatment effects) from observational data involves estimating causal effects of a treatment variable (e.g., medication) on an outcome variable (e.g., recovery) given observable covariates (e.g., gender). In practice, there are also unobserved covariates, i.e., *hidden confounders*, that affect both treatment and outcome variables. For instance, consider a study to evaluate the effectiveness of a medication. Gender as a covariate affects whether a patient chooses to take the medication and the corresponding outcome. The patient’s living habits can be hidden confounders that affect both the patient’s medication and outcome. Exploring hidden confounders allows for more accurate estimations of treatment effects [11, 12, 13].

In this work, we formulate the problem of estimating treatment effects in the context of societal events. Societal events can be classified into different types. Given a time window, we look at multiple types of events (e.g., “appeal”, “investigation”) at a location and define treatment variables to be the detection of increased counts of these events compared to the previous time window. If the sudden and frequent occurrence of such events triggers some event of interest, the implied causal effect can be used to guide and interpret event predictions. We define the outcome as the occurrence of an event of interest (e.g., “protest”) at a future time. Both treatment and outcome variables can be affected by hidden social factors (i.e., hidden confounders) that are difficult to explicitly capture due to complex dependencies. Intuitively, exploring hidden confounders can allow us to estimate causal effects more accurately. To this end, we formulate our main research question as: *can we build a robust event predictive model by incorporating treatment effect estimation with hidden confounder learning?* There are some challenges in solving this problem:

- Societal events have geographical characteristics and exhibit a high degree of temporal dependency [2, 3, 5]. Modeling spatiotemporal information requires an in-depth investigation of the dynamic spatial dependencies of societal events. However, few studies have focused on modeling spatiotemporal dependencies in causal effect learning, which poses a challenge for learning causal effects from societal events.
- Events occur in a complex and evolving social environment. Many unknown social factors increase the difficulty of accurately estimating causal effects of events. Moreover, events are often caused by a variety of factors rather than a single determinant. Utilizing causal effects to assist in event prediction is a new challenge.

We address the above challenges by first introducing the task of Individual Treatment Effect (ITE) estimation from societal events. ITE is defined as the expected difference between the treated outcome and control outcome, where the outcome is the occurrence of a future event (e.g., *protest*) at a specific place and time, and the treatment is a change in some event (e.g., *appeal*) in the past. We consider multiple treatments (e.g., *appeal*, *investigation*, etc.) with the motivation that the underlying causes of societal events are often complex. We model the spatiotemporal dependencies in learning the representations of hidden confounders to estimate ITEs. We then present an approach to inject the learned causal information into a data-driven predictive model to improve its predictive power. Our contributions are summarized as follows:

- We introduce a novel causal inference model for ITE estimation, which learns the representation of hidden confounders by capturing spatiotemporal dependencies of events in different locations.
- We propose two robust learning modules for event prediction that take as prior knowledge the information learned from the causal inference model. Incorporating such modules can enable event prediction models to be more robust to data noise.

We evaluate the proposed method against other state-of-the-art methods on several real-world event datasets. Through extensive experiments, we demonstrate the strengths of the proposed method in treatment effect learning and robust event prediction.

2 Related Work

2.1 Event Prediction

Event prediction focuses on forecasting future events that have not yet happened based on various social indicators, such as event occurrence rates and news reports. Related research has been conducted in various fields and applications, such as election prediction [14, 15], stock market forecasting [16], disease outbreak simulation [17, 18], and crime prediction [19]. Machine learning models such as linear regression [16] and random forests [20] were investigated to predict events of interest. Time-series methods such as autoregressive models were studied to capture the temporal evolution of event-related indicators [18]. With the increased availability of various data, more sophisticated features have been shown effective in predicting societal events such as topic-related keywords [1], document embedding [3], word graphs [4] and knowledge graphs [5, 21]. More advanced machine learning and deep learning-based models have emerged, such as multi-instance learning [3], multi-task learning [22, 23] and graph neural networks [4, 5].

Table 1: Important notations and descriptions.

Notations	Descriptions
M, T, E	sets of locations, timestamps and numbers of event types
$\mathbf{X}_i^{\leq t} \in \mathbb{R}^{E \times \Delta}$	covariates of the i -th location before time t
$\mathbf{x}_i^t \in \mathbb{R}^E$	frequency of events at time t for location i
$\mathbf{c}_i^{\leq t} \in \{0, 1\}^E$	observed treatment vector of the i -th location before time t
$c_{i(j)}^{\leq t} \in \{0, 1\}$	observed j -th treatment of the i -th location before time t
$y_{i(j)}^{t+\delta}(1), y_{i(j)}^{t+\delta}(0)$	potential outcomes for the j -th treatment of the i -th location at time $t + \delta$
$\hat{y}_{i(j)}^{t+\delta}(1), \hat{y}_{i(j)}^{t+\delta}(0)$	predicted potential outcomes for the j -th treatment of the i -th location at time $t + \delta$
$\mathbf{A} \in \mathbb{R}^{M \times M}$	connectivity of M locations
$\tau_{i(j)}^{t+\delta}$	ITE of the i -th location at time $t + \delta$ for the j -th treatment
$\mathbf{z}_{i(j)}^{\leq t}$	learned hidden confounders of the i -th location before time t when the j -th treatment is considered

Given the spatiotemporal dependencies of events, some existing research work studied spatiotemporal correlations in event prediction [24, 25, 26]. However, few studies explored the causality in event prediction. Our proposed model incorporates causal effect learning in a spatiotemporal event prediction framework. This gives us the benefit of discovering the effects of different potential causes on predicting future events.

2.2 Individual Treatment Effect Estimation

Individual treatment effect (ITE) estimation refers to estimating the causal effect of a treatment variable on its outcome. A wealth of observational data facilitates treatment effect estimation in many fields, such as health care [27], education [28], online advertising [29], and recommender systems [8]. Several methods have been studied for ITE estimation including regression and tree based model [30, 31], counterfactual inference [32], and representation learning [33]. The former approaches rely on the Ignorability assumption [34], which is often untenable in real-world studies. A deep latent variable model, CEVAE [11] learns representations of confounders through variational inference. Recent work relaxed the Ignorability assumption and studied ITE estimation from observational data with an auxiliary network structure in a static [12] or dynamic environment [13]. In addition to the traditional causal effect estimation, a new study of causal inference, including multiple treatments and a single outcome, has emerged, namely, Multiple Causal Inference. Researchers have shown that compared with traditional causal inference, it requires weaker assumptions [35]. ITE estimation would considerably benefit decision-making as it can provide potential outcomes with different treatment options. Our work introduces ITE estimation to societal event studies and exploits event-related causal information for event forecasting.

2.3 Knowledge Guided Machine Learning

Purely data-driven approaches might lead to unsatisfactory results when limited data are available to train well-performing and sufficiently generalized models. Such models may also break natural laws or other guidelines [36]. These problems have led to an increasing amount of research that focuses on incorporating additional prior knowledge into the learning process to improve machine learning models. For example, logic rules [37, 38] or algebraic equations [39, 40, 41] have been added as constraints to loss functions. Knowledge graphs are adopted to enhance neural networks with information about relations between instances [42]. The growth of this research suggests that the combination of data and knowledge-driven approaches is becoming relevant and showing benefits in a growing number of areas. Existing work has typically focused on pre-existing knowledge obtained by human experts. However, such approaches fail when prior knowledge is not available, e.g., for societal events. Some researchers explored causal knowledge-guided methods in health prediction [9] and image-to-video adaptation [10]. In this work, we study causal effects between societal events and use the learned causal information as prior knowledge for event prediction.

3 Problem Formulation

The objective of this study is two-fold: (1) given multiple pre-defined treatment events (e.g., appeal, investigation, etc.), estimate their causal effect on a target event (i.e., protest) individually; (2) predict the probability of the target event occurring in the future with the help of estimated causal information. In the following, we will introduce the observational event data, individual treatment effect learning, and event prediction.

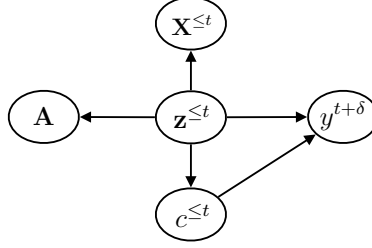


Figure 1: Causal graph defined for the ITE estimation on observational event data. The edges represent causal relations. For a location, we use $\mathbf{X}^{\leq t}$, $\mathbf{c}^{\leq t}$, $\mathbf{z}^{\leq t}$, $y^{t+\delta}$ to denote the covariates, assignment for a treatment event, representation of hidden confounders before time t , and outcome at time $t + \delta$, respectively. \mathbf{A} denotes the connectivity among locations.

3.1 Observational Event Data

In this work, we focus on modeling the occurrence of one type of societal event (i.e., “protest”) by exploring the possible effects it might receive from other types of events (e.g., “appeals” and “investigation”). A total of E categories of societal events are studied. These events happen at different locations and times. We use M, T to denote the sets of locations and timestamps of interest, respectively. The observational event data can be denoted as $\mathcal{D} = \{ \{ \mathbf{X}_i^{\leq t}, \mathbf{c}_i^{\leq t}, y_i^{t+\delta} \}_{i=1}^M, \mathbf{A} \}_{t \in T}$, where $\mathbf{X}_i^{\leq t}$, $\mathbf{c}_i^{\leq t}$, $y_i^{t+\delta}$ denote the pre-treatment covariates/features, observed treatments, and outcome, respectively. $\mathbf{A} \in \mathbb{R}^{M \times M}$ represents the connectivity of M locations, where each element can denote a fixed geographic distance or the degree of influence of events between locations. Important notations are presented in Table 1.

Covariates: We define the covariates $\mathbf{X}_i^{\leq t} = (\mathbf{x}_i^{t-\Delta+1}, \dots, \mathbf{x}_i^t) \in \mathbb{R}^{E \times \Delta}$ to be the historical events at location i with size Δ up to time t . $\mathbf{x}_i^t \in \mathbb{R}^E$ is a vector representing the frequencies of E types of events that occurred at location i at time t .

Treatments: The treatments $\mathbf{c}_i^{\leq t} \in \{0, 1\}^E$ can be represented by a binary vector with dimension E where each element indicates the occurrence states of a type of events (e.g., appeal). Specifically, the j -th element $c_{i(j)}^{\leq t} = 1$ indicates a notable (i.e., 50%) increase of the j -th event type at window $[t - \Delta + 1, t]$ from the previous window $[t - 2\Delta + 1, t - \Delta]$.¹ A value of 1 means getting treated and 0 means getting controlled. For convenience, we refer to each element in the treatment vector as a treatment event.²

Observed Outcome: The observed/factual outcome $y_i^{t+\delta} \in \{0, 1\}$ is a binary variable denoting if an event of interest (i.e., protest) occurs at location i in the future ($t + \delta$). $\delta \geq 1$ is the lead time indicating the number of timestamps in advance for a prediction.

3.2 Individual Treatment Effects Learning

We first define potential outcomes in observational event data following well-studied causal inference frameworks [45, 46]. We ignore the location subscript i for simplicity unless otherwise stated.

Potential Outcomes: In general, the potential outcome $Y(C)$ denotes what the outcome an instance would receive, if the instance was to take treatment C . A potential outcome $Y(C)$ is distinct from the observed/factual outcome Y in that not all potential outcomes are observed in the real world. In our problem, there are two potential outcomes for each treatment event. Given a location at time $t + \delta$ and the j -th treatment event, we denote by $y_{(j)}^{t+\delta}(1)$ the potential outcome (i.e., occurrence of protest) if the j -th treatment event is getting treated, i.e., $c_{(j)}^{\leq t} = 1$. Similarly, we denote by $y_{(j)}^{t+\delta}(0)$ the potential outcome we would observe if the treatment event is under control, i.e., $c_{(j)}^{\leq t} = 0$.

The factual outcome is $y^{t+\delta}$ when the location has already received the treatment assignment before time t . The counterfactual outcome is defined if the location obtains the opposite treatment assignment. In the observational study, only the factual outcomes are available, while the counterfactual outcomes can never be observed.

¹The comparison of two windows is motivated by studies showing that short-term historical data can lead to favorable performance in event prediction [5, 43]. A threshold of 50% is selected heuristically. We leave variant treatment settings for future work.

²Our setup differs from multiple causal inference [35, 44], which estimates the potential outcome of a combination of multiple treatments. We are more interested in studying the potential outcome of each element in the treatment vector.

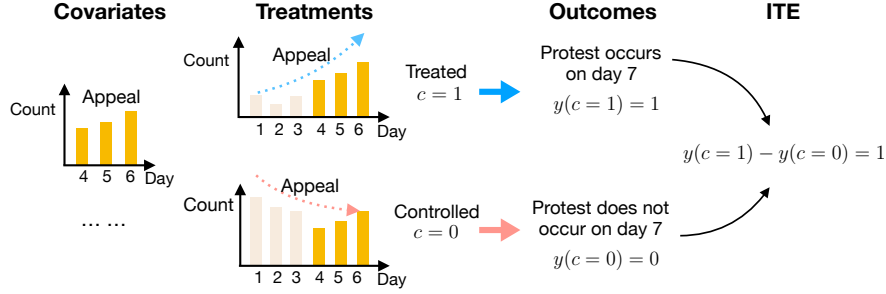


Figure 2: An example of ITE estimation on observational event data. We use the horizontal axis to indicate time and the bars to indicate the count of events in a location. Taking the Appeal event as the treatment event, the treated and controlled scenarios denote whether the amount of appeal events in the current window (on average) is greater or less than the previous window. The window size is 3. The light orange bars indicate the appeal events of the past window. The outcomes denote whether the protest will occur at a future time given each treatment assignment. ITE is estimated from the two potential outcomes.

The Individual Treatment Effect (ITE) is the difference between two potential outcomes of an instance, examining whether the treatment affects the outcome of the instance. In observational event data, for the j -th treatment event, we formulate the ITE for the location at time $t + \delta$ in the form of the Conditional Average Treatment Effect (CATE) [33, 13]:

$$\tau_{(j)}^{t+\delta} = \mathbb{E}[y_{(j)}^{t+\delta}(1) - y_{(j)}^{t+\delta}(0) | \mathbf{X}^{\leq t}, \mathbf{A}]. \quad (1)$$

We provide a toy example to illustrate ITE estimation on observational event data, as shown in Fig. 2.

In this study, we aim to estimate ITEs and then use them for event prediction. The challenge of ITE estimation lies in how to estimate the missing counterfactual outcome. Our estimation of ITE is built upon some essential assumptions. For simplicity and readability, we omit the subscripts for the location i and the treatment event j and use $c^{\leq t}$ to represent the j -th treatment event.

Assumption 1. No Interference. Assuming that one instance is defined as a location at a time in observational event data, the potential outcome on one instance should be unaffected by the particular assignment of treatments on other instances.

Assumption 2. Consistency. The potential outcome of treatment $c^{\leq t}$ equals to the observed outcome if the actual treatment received $c^{\leq t}$, i.e., $y^{t+\delta}(c^{\leq t}) = y^{t+\delta}$.

Assumption 3. Positivity. If the probability $\mathbb{P}(\mathbf{X}^{\leq t}, \mathbf{A}) \neq 0$, then the probability to receive treatment assignment 0 or 1 is positive, i.e., $0 < \mathbb{P}(c^{\leq t} = 1 | \mathbf{X}^{\leq t}, \mathbf{A}) < 1$, $c^{\leq t} \in \{0, 1\}$.

The Positivity assumption indicates that before time t , each treatment assignment has a non-zero probability of being given to a location. This assumption is testable in practice. In addition to these assumptions, most existing work [33, 11, 47] relies on the Ignorability assumption, which assumes that all confounding variables are observed and reliably measured by a set of features for each instance, i.e., hidden confounders do not exist.

Definition 1. Ignorability Assumption. Given pre-treatment covariates $\mathbf{X}^{\leq t}$, the outcome variables are independent of its treatment assignment, $y^{t+\delta}(0), y^{t+\delta}(1) \perp\!\!\!\perp c^{\leq t} | \mathbf{X}^{\leq t}$.

However, this assumption is untenable in societal event studies due to the complex environment in which societal events occur. We relax this assumption by introducing the existence of hidden confounders [12]. Note that hidden confounders are unobserved in observational event data but will be learned in our approach through a spatiotemporal model. We define a causal graph, as shown in Fig. 1. The hidden confounders $\mathbf{z}^{\leq t}$ causally affect the treatment and outcome.³ The potential outcomes are independent of the observed treatment, given the hidden confounders: $y^{t+\delta}(0), y^{t+\delta}(1) \perp\!\!\!\perp c^{\leq t} | \mathbf{z}^{\leq t}$. In addition, we assume the features $\mathbf{X}^{\leq t}$ and the connectivity of locations \mathbf{A} are proxy variables for hidden confounders $\mathbf{z}^{\leq t}$. Unobservable hidden confounders can be measured with $\mathbf{X}^{\leq t}$ and \mathbf{A} . Based on the temporal and spatial characteristics of our observational event data. We introduce the following assumption [13]:

Assumption 4. Spatiotemporal Dependencies in Hidden Confounders. In observational event data, hidden confounders capture spatial information among locations, reflected by \mathbf{A} , and show temporal dependencies of events across multiple historical steps (i.e., Δ).

³For the j -th treatment event, the hidden confounders can be written as $\mathbf{z}_{(j)}^{\leq t}$.

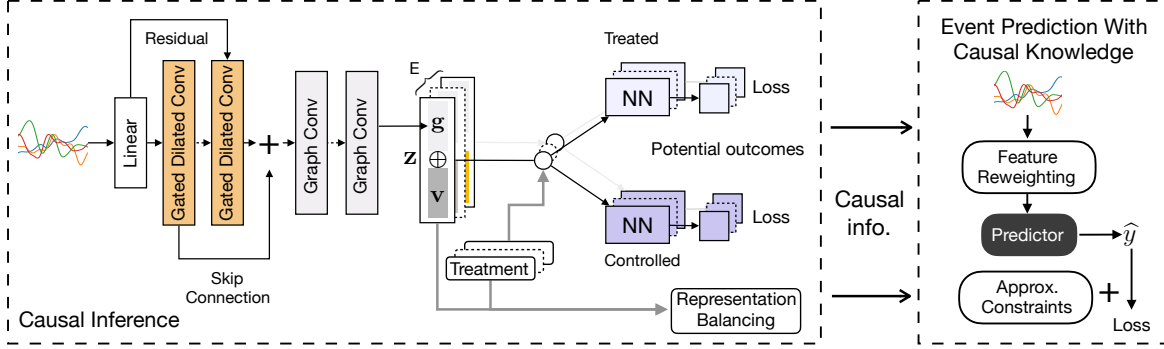


Figure 3: The overall architecture of the proposed framework. The proposed causal inference method learns to estimate ITEs in the presence of multiple treatment events of interest. The representation of hidden confounders is learned using a spatiotemporal model. Then, the vector representation of hidden confounders corresponding to each treatment event is fed into two neural networks (NNs) for individual treatment effect learning. Next, the framework takes the output of the causal inference model as prior knowledge to forecast events.

Note that this assumption does not contradict the *No Interference* assumption. We focus on the scenario in which spatiotemporal information can be exploited to control confounding bias.

3.3 Event Prediction

We present traditional event prediction and event prediction with causal knowledge proposed in this work.

Definition 2. Event Prediction. Learn a classifier that predicts the probability of the target event occurring at a location at time $t + \delta$ based on available data: $\mathbb{P}(y^{t+\delta} | \mathbf{X}^{\leq t}, \mathbf{A})$.

Instead of learning a mapping function from input features to event labels, we are interested in estimating treatment effects under different treatment events individually and exploiting such causal information to enhance event prediction.

Definition 3. Event Prediction with Causal Knowledge. Build an event forecaster using available data with causal information as prior knowledge: $\mathbb{P}(y^{t+\delta} | \mathbf{X}^{\leq t}, \mathbf{A}, \mathcal{C}(\mathbf{X}^{\leq t}, \mathbf{c}^{\leq t}, \mathbf{A}))$, where \mathcal{C} is the trained causal inference model that takes features, multiple treatments and the connectivity information of locations as input and outputs potential outcomes.

The multiple treatment setting (i.e., $\mathbf{c}^{\leq t}$) aims to produce informative causal knowledge to assist event prediction. We will discuss the proposed method of event prediction with causal knowledge in the following sections.

4 Methodology

We propose a novel framework **CAPE**, which incorporates **causal inference** into the **prediction** of future **event** occurrences in a spatiotemporal environment.⁴ In our framework, ITEs with different treatment events are jointly modeled in a spatiotemporal causal inference model. It will contribute to the final event prediction by feeding the causal output (e.g., potential outcomes) to a non-causal data-driven prediction model. The overall framework, as illustrated by Fig. 3, consists of two parts: (1) causal inference and (2) event prediction. The causal inference component is designed to estimate the ITE, including two essential modules: hidden confounder learning and potential outcome prediction. For each treatment event, it learns the representation of hidden confounders by capturing spatiotemporal dependencies and outputs the potential outcomes under different treatment assignments. The event prediction part comprises two robust learning modules, a feature reweighting module and an approximation constraint loss. They take the causal information learned from the causal inference model as prior knowledge to assist the training of a data-driven event prediction model. Next, we will elaborate on these components.

⁴Code is available at <https://github.com/amy-deng/cape>

4.1 Causal Inference

4.1.1 Hidden Confounder Learning

Hidden confounders are common in real-world observational data [48]. Assuming spatiotemporal dependencies exist in hidden confounders, we introduce a novel and effective network that models spatial and temporal information for each location at each time step. It consists of several temporal feature learning layers and spatial feature learning layers. Our network is based on the success of previous work [49, 50]. It is designed to be adaptable to a multi-task setting to learn hidden confounders of multiple treatments.

Temporal Feature Learning Dilated casual convolution networks [51] handle long sequences in a non-recursive manner, which facilitates parallel computation and alleviates the gradient explosion problem. Gating mechanisms have shown benefits to control information flow through layers for convolution networks [49, 50, 52]. We employ the dilated causal convolution with a gating mechanism in temporal feature learning to capture a location’s temporal dependencies. For a location before time t , the multivariate time series of historical event occurrences is a matrix $\mathbf{X}^{\leq t} \in \mathbb{R}^{E \times \Delta}$, where each row indicates the frequency sequence of one type of events in the historical window with size Δ . We use a linear transformation to map the event frequency matrix into a latent space, i.e., $\mathbf{X}'^{\leq t} \in \mathbb{R}^{d_s \times \Delta}$. d_s indicates the feature dimension in the latent space. Then, we apply the dilated convolution to the sequence. For simplicity, we use $\mathbf{r} \in \mathbb{R}^{\Delta}$ to denote a row in the matrix $\mathbf{X}'^{\leq t}$. Formally, for the 1-D sequence input \mathbf{r} and a filter $\mathbf{\Gamma} \in \mathbb{R}^K$, the dilated causal convolution operation on element s of the sequence is defined as:

$$(\mathbf{r} *_d \mathbf{\Gamma})(s) = \sum_{k=0}^{K-1} \mathbf{\Gamma}(k) \mathbf{r}(s - d \cdot k), \quad (2)$$

where $*_d$ is a d -dilated convolution. K is the filter size. $(s), (k), (s - d \cdot k)$ are indices of vectors. $(\mathbf{r} *_d \mathbf{\Gamma})$ is the output vector.

We further incorporate a gated dilated convolutional layer which consists of two parallel dilated convolution layers:

$$\mathbf{h} = \tanh(\mathbf{r} *_d \mathbf{\Gamma}_1) \odot \sigma(\mathbf{r} *_d \mathbf{\Gamma}_2), \quad (3)$$

where $\mathbf{\Gamma}_1, \mathbf{\Gamma}_2$ are filters for dilated convolutional layers, and \odot is the Hadamard product. $\tanh(\cdot)$ is to regularize the features. $\sigma(\cdot)$ is the sigmoid function that determines the ratio of information passed to the next layer. Specifically, we stack multiple gated dilated convolutional layers (Eq. 3) with increasing dilation factors (e.g., $d = 1, 2, 4$). Residual and skip connections are applied to avoid the vanishing gradient problem [49, 52]. To this end, the temporal dependencies are captured, and we use $\mathbf{H}^t \in \mathbb{R}^{M \times d_s}$ to denote the learned temporal features for M locations at time t .

Spatial Feature Learning Graph convolution is a powerful operation to learn representations of nodes given the graph structure. To capture the spatial dependencies, we adopt the graph convolutional network (GCN) [53] to learn the spatial influence from locations by treating each location as a node in graph:

$$\mathbf{G}^t = \text{ReLU}(\mathbf{A}' \mathbf{H}^t \mathbf{W}^g), \quad (4)$$

where \mathbf{W}^g is the weight matrix for a GCN layer. \mathbf{G}^t denotes the spatiotemporal feature matrix referring to all locations, where each row \mathbf{g}^t captures the historical information of a specific location as well as the neighboring locations. \mathbf{A}' is a learnable adjacency matrix. The geographical adjacency matrix of locations usually cannot represent the connectivity of locations in the context of societal event forecasting. Therefore, we adopt the self-adaptive adjacency matrix [52], which does not require any prior knowledge and is learned through training. We randomly initialize two node embedding matrices with learnable parameters $\mathbf{E}_1, \mathbf{E}_2 \in \mathbb{R}^{M \times d_a}$. The self-adaptive adjacency matrix is defined as:

$$\mathbf{A}' = \text{Softmax}(\text{ReLU}(\mathbf{E}_1 \mathbf{E}_2^T)), \quad (5)$$

where the ReLU activation function eliminates weak connections and the Softmax applies normalization.

Hidden Confounder Learning To learn the representation of hidden confounders, we utilize the spatiotemporal feature and a learnable embedding specific to each treatment event (i.e., $\mathbf{v}_{(j)}, 1 \leq j \leq E$). It is worth pointing out that in the proposed framework, we include multiple treatment events and expect to estimate the ITE corresponding to each treatment event. The treatment-specific embedding aims to capture latent information of each treatment event and distinguish the hidden confounder representations learned for each treatment effect learning task. Similar ideas of task embeddings are studied in prior work [54]. Given a location and a time t , the representation of hidden confounders for the j -th treatment is:

$$\mathbf{z}_{(j)}^{\leq t} = \mathbf{g}^t \oplus \mathbf{v}_{(j)}, \quad (6)$$

where \oplus is concatenation.

4.1.2 Potential Outcome Prediction

Using the above components, we obtain the representation of hidden confounders $\mathbf{z}_{(j)}^{\leq t}$. Following the predefined causal graph in Fig. 1, the learned hidden confounders can be used to estimate potential outcomes. We use two networks that output two potential outcomes of the j -th treatment event, respectively:

$$\hat{y}_{(j)}^{t+\delta}(1) = \Phi_{(j)}^1(\mathbf{z}_{(j)}^{\leq t}), \quad \hat{y}_{(j)}^{t+\delta}(0) = \Phi_{(j)}^0(\mathbf{z}_{(j)}^{\leq t}), \quad (7)$$

where $\hat{y}_{(j)}^{t+\delta}(1/0)$ denotes the inferred potential outcomes when the j -th treatment event is getting treated or controlled. $\Phi_{(j)}^1(\cdot), \Phi_{(j)}^0(\cdot)$ are parameterized by deep neural networks with a sigmoid function at the last layer. The networks are trained end-to-end, and one can estimate the potential outcomes under multiple treatment events.

4.1.3 Objective Function

Potential Outcome Loss We use the binary cross-entropy loss as the objective factual loss for predicting potential outcomes. When only the j -th treatment event is considered (i.e., the general case for treatment effect learning [33, 12, 13]), the factual loss is:

$$\mathcal{L}_{(j)}^{\text{fact}} = - \sum_{t \in T} \sum_{i \in M} y_i^{t+\delta} \log \hat{y}_{i(j)}^{t+\delta} + (1 - y_i^{t+\delta}) \log (1 - \hat{y}_{i(j)}^{t+\delta}), \quad (8)$$

where $y_i^{t+\delta}$ is the observed outcome for location i at time $t + \delta$. $\hat{y}_{i(j)}^{t+\delta} = \hat{y}_{i(j)}^{t+\delta}(c_{i(j)}^{\leq t})$ is the predicted outcome given the observed treatment $c_{i(j)}^{\leq t} \in \{0, 1\}$. Since our model predicts potential outcomes for multiple treatment events, we express the total factual loss as follows:

$$\mathcal{L}^{\text{fact}} = \sum_{1 \leq j \leq E} \mathcal{L}_{(j)}^{\text{fact}} + \eta \cdot \Omega(\Theta), \quad (9)$$

where $\Omega(\Theta)$ stands for the ℓ_2 -norm regularization for all training parameters and η is the weight for scaling the regularization term.

Representation Balancing Studies have proved that balancing the representations of treated and control groups would help mitigate the confounding bias and minimize the upper bound of the outcome inference error [32, 33]. Therefore, we incorporate a representation balancing layer to force the distributions of hidden confounders of treated and controlled groups to be similar. Specifically, we adopt the integral probability metric (IPM) [33] to measure the difference between the distributions of the treated instances and the controlled instances in terms of their hidden confounder representations:

$$\mathcal{L}^{\text{disc}} = \alpha \cdot \text{IPM}(\mathcal{Z}_1, \mathcal{Z}_0), \quad (10)$$

where $\mathcal{Z}_1 = \{\mathbf{z}_{i(j)}^{\leq t}\}_{i,t,j:c_{i(j)}^{\leq t}=1}$, $\mathcal{Z}_0 = \{\mathbf{z}_{i(j)}^{\leq t}\}_{i,t,j:c_{i(j)}^{\leq t}=0}$ indicate the sets of hidden confounders for samples (in a batch) in the treated group and controlled group, respectively. The IPM can be Wasserstein or Maximum Mean Discrepancy (MMD) distances. α is a hyperparameter that indicates the imbalance penalty.

Formally, we present the loss function of the proposed causal inference model as:

$$\mathcal{L}^{\text{cau}} = \mathcal{L}^{\text{fact}} + \mathcal{L}^{\text{disc}}. \quad (11)$$

4.2 Event Prediction with Causal Knowledge

To improve the robustness of event predictions with imperfect real-world data, we incorporate causal information output by the causal inference model as priors to forecast future events. We introduce two robust learning modules into the training of event predictors: (1) feature reweighting, which involves causal information to weight the original input features to obtain causally enhanced features, and (2) approximation constraints, which use the predicted potential outcomes as value range constraints applied to event prediction scores. Next, we introduce these two modules in detail.

4.2.1 Feature Reweighting

Feature reweighting was introduced in object detection [55], where a reweighting vector is learned to indicate the importance of meta features for detecting objects. Here, we introduce a new feature reweighting method that leverages causal information. We use the ITE estimated from the causal inference model to reweight the event frequency features to predict future events.

Causal Feature Gates We define a feature gate based on ITE calculated using predicted potential outcomes from the causal inference model. For the j -th treatment event, the estimated ITE of a location at time $t + \delta$ is as follows:

$$\hat{\tau}_{(j)}^{t+\delta} = \hat{y}_{(j)}^{t+\delta}(1) - \hat{y}_{(j)}^{t+\delta}(0). \quad (12)$$

When considering multiple treatment events, we obtain the ITE vector $\hat{\boldsymbol{\tau}}^{t+\delta} \in \mathbb{R}^E$, where each element indicates a treatment event. A linear layer f_τ with sigmoid function is then applied to model the association between the effects of different treatment events:

$$\boldsymbol{\rho}^{t+\delta} = \sigma(f_\tau(\hat{\boldsymbol{\tau}}^{t+\delta})), \quad (13)$$

where $\boldsymbol{\rho}^{t+\delta} \in \mathbb{R}^E$ is the gating variables that will be applied to the original event frequency features. The sigmoid function σ converts the gating variable into a soft gated signal with a range of $(0, 1)$.

Reweighting Feature We reweight the event frequency features using the gating variables defined above. It is worth emphasizing that the event frequency vector \mathbf{x}^t has the same dimension as $\boldsymbol{\rho}^{t+\delta}$, and their corresponding elements represent the same event type. Nevertheless, we prefer not to apply the gating variables directly to the feature vector. ITE examines whether the binary treatment variable affects the outcome of an instance, while the event frequency vector refers to discrete numbers. To address this issue, we transform the event frequency feature into a latent vector using a position-wise feed-forward network (FFN) [56]. It maps the features into a continuous space, assuming that the gating variables can be aligned with the variables in this space. The formal procedures are defined as follows:

$$\text{FFN}(\mathbf{x}^t) = \text{ReLU}(\mathbf{x}^t \mathbf{W}^{r_0} + \mathbf{b}^{r_0}) \mathbf{W}^{r_1} + \mathbf{b}^{r_1}, \quad (14)$$

$$\tilde{\mathbf{x}}^t = \text{FFN}(\mathbf{x}^t) \odot \boldsymbol{\rho}^{t+\delta} + \mathbf{x}^t, \quad (15)$$

where $\{\mathbf{W}, \mathbf{b}\}^{\{r_0, r_1\}}$ are learnable parameters. A residual connection is added to ensure that the causally weighted elements still contain some original information. We denote the causality enhanced features across Δ historical steps as $\tilde{\mathbf{X}}^{\leq t} \in \mathbb{R}^{E \times \Delta}$. Such features are fed into a predictor to perform event prediction, denoted as $\mathcal{P}(\tilde{\mathbf{X}}^{\leq t}, \mathbf{A})$.

4.2.2 Approximation Constraints

The approximation constraints method was proposed to limit the reasonable range of the target variable during the model training process to generate a more robust model [41]. We follow this idea and propose a new method of integrating learned causal information into variable constraints. Given an event predictor \mathcal{P} , we denote the model's event prediction for a location at time $t + \delta$ as $\hat{y}^{t+\delta}$. Then, we assume that the causal range of the target variable, i.e., the event prediction, is $\hat{y}^{t+\delta} \in [l^{t+\delta}, u^{t+\delta}]$. The sample-wise boundaries are defined as:

$$l^{t+\delta} = \text{Min}(\hat{\mathbf{y}}^{t+\delta}), \quad u^{t+\delta} = \text{Max}(\hat{\mathbf{y}}^{t+\delta}), \quad (16)$$

where $\hat{\mathbf{y}}^{t+\delta} = \{\hat{y}_{(j)}^{t+\delta}(1), \hat{y}_{(j)}^{t+\delta}(0) | 1 \leq j \leq E\}$ is the set of potential outcomes for all treatment events. The minimum and maximum values are the lower and upper limits of the target variable for a given sample. Based on the range obtained from causal knowledge, we define a constraint loss term:

$$\mathcal{L}^{\text{cstr}} = \sum_{t \in T} \sum_{i \in M} \text{ReLU}(l_i^{t+\delta} - \hat{y}_i^{t+\delta}) + \text{ReLU}(\hat{y}_i^{t+\delta} - u_i^{t+\delta}). \quad (17)$$

The loss term can be involved during the training of the predictor \mathcal{P} . Given the proposed robust learning modules for event prediction, we train the predictor by minimizing the following loss function:

$$\mathcal{L}^{\text{evt}} = \mathcal{L}^{\text{pred}} + \mu \cdot \mathcal{L}^{\text{cstr}}, \quad (18)$$

where $\mathcal{L}^{\text{pred}}$ is the loss function defined by the predictor \mathcal{P} and μ is a hyperparameter. The training steps of the proposed method are shown in Algorithm 1.

5 Experimental Evaluation

The goal of the experimental evaluation is to answer the following research questions: **RQ1**: How well does **CAPE** estimate ITEs in observational event data? **RQ2**: Can **CAPE** improve the robustness of event prediction models? **RQ3**: What causal information can we learn from studies of causally related event prediction?

Next, we will describe the experimental setup and then show the experimental results to address the above questions.

Algorithm 1: CAPE

Input: Observational event data $\mathcal{D} = \{\{\mathbf{X}_i^{\leq t}, \mathbf{c}_i^{\leq t}, y_i^{t+\delta}\}_{i=1}^M, \mathbf{A}\}_{t \in \mathcal{T}}$, a predictor \mathcal{P} with randomly initialized parameters, and initialized model **CAPE**, including causal inference model \mathcal{C} and robust learning modules \mathcal{R} for event prediction.

▷ Train the causal inference model \mathcal{C} .

1 **while** \mathcal{C} has not converged **do**

2 Input \mathcal{D} to \mathcal{C} to obtain predicted potential outcomes. Calculate the loss \mathcal{L}^{cau} .

3 Update \mathcal{C} by optimizing Eq. 11.

4 Freeze the causal inference model \mathcal{C} .

 ▷ Train the predictor \mathcal{P} and robust modules \mathcal{R} with \mathcal{C} .

5 **while** \mathcal{P}, \mathcal{R} have not converged **do**

6 Input \mathcal{D} to \mathcal{C} to obtain predicted potential outcomes.

7 Calculate causally reweighed features $\tilde{\mathbf{X}}$ from robust learning modules \mathcal{R} using predicted potential outcomes.

8 Input $\{\tilde{\mathbf{X}}, \mathbf{A}\}$ to the predictor \mathcal{P} to obtain the event prediction. Calculate the prediction loss $\mathcal{L}^{\text{pred}}$.

9 Calculate the constraint loss $\mathcal{L}^{\text{cstr}}$ from \mathcal{R} using predicted potential outcomes.

10 Update the predictor \mathcal{P} and causal modules \mathcal{R} by optimizing Eq. 18.

Output: the updated model **CAPE**.

Table 2: Dataset Statistics. M is the number of locations in each dataset. **Positive** indicates the ratio of positive samples, i.e., the protest event has occurred. For India, we select top locations based on the total number of events. **Location** represents the geographical level of events.

Dataset	M	Positive	Location	Time	Time Unit	Source
India	14	30.1%	State	2000-2017	3 days	ICEWS
Nigeria	6	65.7%	Geopolitical zone	2015-2020	1 day	GDELT
Australia	8	44.4%	State	2015-2020	1 day	GDELT
Canada	13	26.8%	State	2015-2020	1 day	GDELT

5.1 Datasets

Experimental evaluation is conducted on two data sources: Integrated Conflict Early Warning System (ICEWS) [57], and Global Database of Events, Language, and Tone (GDELT) [58]. These two data sources include daily events encoded from news reports.⁵ We construct event datasets for four countries, i.e., **India**, **Nigeria**, **Australia** and **Canada**, based on their large volume of events. Events are categorized into 20 main categories (e.g., appeal, demand, protest, etc.) according to CAMEO methodology [59]. Each event is encoded with geolocation, time (day, month, year), category, etc. In this work, we focus on predicting one category of events: *protest*, as the target variable, and using event historical data of all event types as feature variables. Data statistics are shown in Table 2. Positive in the table refers to the proportion of positive samples.

5.2 Evaluation Metrics

For the ITE estimation, since there is no ground truth counterfactual outcomes, we report the ATT error $\epsilon_{\text{ATT}} = |\text{ATT} - \mathbb{E}_{\mathcal{S}}[\hat{y}^{t+\delta}(1) - \hat{y}^{t+\delta}(0)|c^{\leq t} = 1]|$ [33], where ATT is the true average treatment effect on the treated, i.e., $\text{ATT} = \mathbb{E}_{\mathcal{S}}[y^{t+\delta}|c^{\leq t} = 1] - \mathbb{E}_{\mathcal{S}}[y^{t+\delta}|c^{\leq t} = 0]$. \mathcal{S} denotes the subset of samples simulating a randomized controlled trial. Specifically, given the treatment event, we employ a 1-nearest neighbor algorithm [60] to find a matching control instance (without replacement) for each treated instance. Euclidean distance is adopted to measure feature vectors. The matching process is performed for each location.

We quantify the predictive performance of event prediction based on Balanced Accuracy (BACC), i.e., $\text{BACC} = (\text{TPR} + \text{TNR})/2$. TPR and TNR are the true positive rate and true negative rate, respectively. BACC is a good metric when the classes are imbalanced.

⁵For event data from GDELT, we only select root events identified in news reports.

5.3 Comparative Methods

For the ITE estimation, we compare our causal inference model, notated as CAPE_C , with two groups of baselines: (i) Tree based methods: Bayesian Additive Regression Trees (**BART**) [31] and Causal Forest (**CF**) [47]; (ii) Representation learning based methods: Counterfactual regression with MMD (**CFR-MMD**) [33] and Wasserstein metric (**CFR-WASS**) [33], Causal Effect Variational Autoencoder (**CEVAE**) [11], Network Deconfounder (**Net-Deconf**) [12], and Similarity Preserved Individual Treatment Effect (**SITE**) [61].

We study three variants of our model to examine the impact of different components in our model: (i) CAPE_{C-G} which removes the spatial feature learning. (ii) CAPE_{C-T} which replaces the temporal feature learning with a simple linear transformation. (iii) CAPE_{C-B} removes the loss term $\mathcal{L}^{\text{disc}}$.

To evaluate the effectiveness of proposed robust learning modules in event prediction, we adopt two spatiotemporal models as the predictor \mathcal{P} : (i) **Cola-GNN** [62]: A graph-based framework for long-term Influenza-like illness prediction; (ii) **GWNet** [52]: A state-of-the-art spatiotemporal graph model for traffic prediction. Given the spatiotemporal characteristics of societal event data, these models can be well applied to our problem. Note that we do not adopt protest event prediction models [4, 21] because they model on more complex data, such as text and knowledge graphs. We leave the causal exploration of such complex data to future work.

6 Implementation Details

For the causal inference model, we use three gated temporal convolutional layers with dilation factors $d = 1, 2, 4$, and two graph convolutional layers. The dimension d_a is set to 10. The feature dimensions of all other hidden layers including d_s are set to be equal and searched from $\{16, 32, 64\}$.

The number of treatment events E is 20, where each treatment event corresponds to an event type, such as appeal and protest. Following previous work [4, 5], we set the historical window size Δ to 7 and the lead time δ to 1. The hyperparameter η used for parameter regularization is fixed to $1e-5$. We use the squared linear MMD for representation balancing [33]. The imbalance penalty α is searched from $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$. The scaling term μ in Eq. 18 is searched from $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$. All parameters are initialized with Glorot initialization [63] and trained using the Adam [64] optimizer with learning rate 10^{-3} and dropout rate 0.5. The batch size is set to 64. We use the objective value on the validation set for early stopping.

For causal inference baselines, **CF**⁶ and **BART**⁷ are implemented using R packages. We implement the causal inference models **CFR-MMD**, **CFR-WASS**, **SITE** by ourselves and use the source code of **CEVAE**⁸ and **Net-Deconf**⁹. We apply parameter searching on all baseline models. For representation learning based approach, the dimension of hidden layers are searched from $\{32, 64, 128\}$ and the number of hidden layers are searched from $\{1, 2\}$. For models that introduce balancing representation learning, we search the hyperparameter from $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$. The model **Net-Deconf** involves an auxiliary network and we use the geographic adjacency matrix for locations.

For the experiments on event forecasting, we run the source code of **Cola-GNN**¹⁰ and **GWNet**¹¹. For event prediction models, we fixed the dimension of hidden layers to 32. **Cola-GNN** takes the geographic adjacency matrix as input and **GWNet** learns the adaptive adjacency matrix.

We report the average of 5 randomized trials for all experiments. At each training, we randomly split the data into training, validation, and test sets at a ratio of 70%-15%-15% with a fixed seed value. All python code is implemented using Python 3.7.7 and Pytorch 1.5.0 with CUDA 9.2.

7 Experimental Results

7.1 Results of ITE Estimation (RQ1)

To evaluate the effectiveness of our proposed causal inference framework, we limit the number of treatment events to be one and compare our model with other baselines. We focus on two treatment events: *appeal* and *reject*. The motivation is

⁶https://rdrr.io/cran/grf/man/causal_forest.html

⁷<https://rdrr.io/cran/BART/>

⁸https://github.com/rik-helweg/CEVAE_pytorch

⁹<https://github.com/rguo12/network-deconfounder-wsdm20>

¹⁰<https://github.com/amy-deng/colagnn>

¹¹<https://github.com/nnzhan/Graph-WaveNet>

Table 3: ITE estimation results showing the mean and standard deviation of ATT errors on all datasets with treatment event being *Appeal*. Lower is better.

Treatment event: Appeal				
	India	Nigeria	Australia	Canada
BART	.114 ± .012	.024 ± .015	.034 ± .018	.035 ± .017
CF	.112 ± .011	.034 ± .015	.038 ± .021	.038 ± .016
CFR-MMD	.014 ± .012	.020 ± .013	.016 ± .011	.017 ± .010
CFR-WASS	.021 ± .010	.025 ± .016	.019 ± .013	.012 ± .009
CEVAE	.012 ± .007	.018 ± .014	.018 ± .011	.017 ± .016
SITE	.013 ± .009	.019 ± .015	.012 ± .008	.016 ± .013
Net-Deconf	.034 ± .027	.030 ± .017	.022 ± .023	.034 ± .020
CAPE_C	.011 ± .008	.010 ± .007	.010 ± .007	.008 ± .006
CAPE_{C-G}	.019 ± .009	.019 ± .011	.019 ± .008	.016 ± .011
CAPE_{C-T}	.011 ± .009	.027 ± .021	.013 ± .010	.011 ± .007
CAPE_{C-B}	.018 ± .012	.026 ± .013	.020 ± .013	.019 ± .012

Table 4: ITE estimation results showing the mean and standard deviation of ATT errors on all datasets with treatment event being *Reject*. Lower is better.

Treatment event: Reject				
	India	Nigeria	Australia	Canada
BART	.184 ± .012	.072 ± .021	.045 ± .030	.070 ± .021
CF	.181 ± .013	.062 ± .021	.039 ± .031	.067 ± .022
CFR-MMD	.022 ± .011	.021 ± .020	.024 ± .013	.018 ± .009
CFR-WASS	.016 ± .012	.018 ± .009	.026 ± .021	.016 ± .011
CEVAE	.015 ± .010	.021 ± .018	.020 ± .012	.019 ± .015
SITE	.010 ± .008	.019 ± .010	.022 ± .019	.024 ± .012
Net-Deconf	.026 ± .021	.020 ± .018	.024 ± .019	.023 ± .014
CAPE_C	.016 ± .011	.016 ± .005	.012 ± .009	.015 ± .012
CAPE_{C-G}	.010 ± .013	.020 ± .012	.015 ± .010	.011 ± .011
CAPE_{C-T}	.015 ± .014	.014 ± .013	.007 ± .004	.016 ± .011
CAPE_{C-B}	.017 ± .013	.028 ± .013	.018 ± .014	.018 ± .009

that *appeal* events might be a potential cause of protest events, as they express a serious or urgent request, typically to the public. *Reject* events represent verbal conflicts [59], which contain dissatisfaction with the current state and may lead to a future occurrence of protest. Table 3 and Table 4 report the ATT errors of all causal inference models on four datasets when treatment variable to be *appeal* and *reject*, respectively. The results show that the tree-based model performs worse than the representational learning-based model. The findings reflect the limitations of the tree-based models and highlight the benefits of representation learning for estimating ITE for observational event data. **CFR-MMD** and **CFR-WASS** learn a balanced representation such that the induced treated and control distributions look similar. Both models achieved good results in most cases, demonstrating the importance of controlling for representation distributions to predict potential outcomes. **CEVAE** learns latent variables based on variational autoencoders and **SITE** focuses on capturing local similarities to estimate ITE. These two models present the most stable and relatively small ATT errors in all settings. This suggests that learning the latent variables and considering similarity information is useful for estimating ITE for observational event data. The model **Net-Deconf** learns hidden confounders by leveraging network/spatial information. However, it does not outperform representation-based baselines. This may be because the model was designed for semi-synthetic datasets and the spatial characteristics of observational event data are different from the network used in the original paper. Our proposed causal inference framework learns hidden confounders while capturing spatial and temporal information and achieves the best performance. For our model variants, we observe that removing the representation balancing makes the results worse. Ignoring the temporal or spatial feature learning can also deteriorate the results. This reflects the possible spatiotemporal dependencies underlying the hidden confounders. It also demonstrates the capability of the proposed model in capturing the spatiotemporal information of the observational event data.

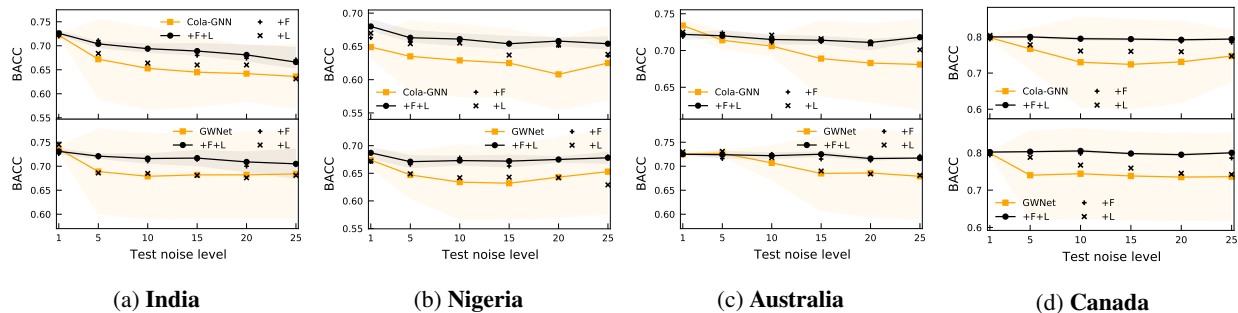


Figure 4: BACC value of event prediction when varying the noise level in validation and test sets. Higher is better.

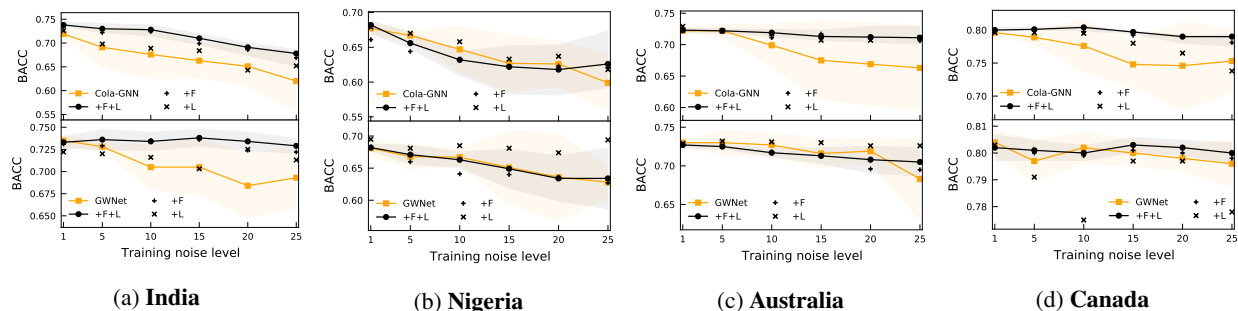


Figure 5: BACC value of event prediction when varying the noise level in the training set. Higher is better.

7.2 Robustness Tests in Event Prediction (RQ2)

In this subsection, we perform two robustness tests on event prediction for all datasets and conduct a case study on the proposed feature reweighting module.

7.2.1 Robustness to Test Noise

A model is considered to be robust if its output variable is consistently accurate when one or more input variables drastically change due to unforeseen circumstances. In this setting, we add Poisson noise into the validation and test sets while keeping the training data noise-free. We aim to verify whether our method guarantees good prediction performance when the test input features are biased. We vary the rate parameter (aka expectation) of the Poisson distribution from 1 to 25 and provide the comparison results for different noise levels in Fig 4. We notice that training with the proposed robust learning module leads to higher average BACC results and lower variance over multiple runs. In most cases, the feature reweighting module (+F) contributes more in improving the prediction performance. Incorporating these two modules (+F+L) can lead to better overall results. The results suggest that forecasting events with learned causal information is beneficial to improve the robustness of the prediction.

7.2.2 Robustness to Training Noise

Human errors or machine failures in real-world data collection usually reduce data accuracy. With this motivation, we assume that only the training data are biased and test whether our method can achieve decent event prediction results on unbiased test data. As shown in Fig. 5, applying robust learning modules can help the prediction model achieve better performance in BACC when the noise level increases. Adding the approximation constraint loss (+L) can lead to a higher BACC than adding the two modules (Fig. 5b and Fig. 5c). The results also illustrate that even with biased data (with corrupted features), the trained causal inference model learns valuable information that contributes to event prediction.

7.2.3 Case Study of Feature Reweighting

To illustrate the functionality of the proposed feature reweighting on robust event prediction, we provide several examples in the India dataset, as shown in Fig. 6. We use **Cola-GNN** for analysis, given the more apparent improvements when it is applied with the feature reweighting module. Specifically, we first train an event prediction model on the India

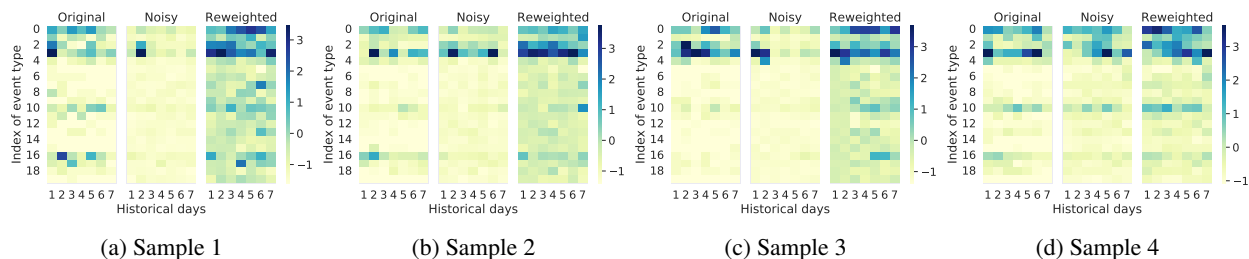


Figure 6: Distributions of original, noisy, and reweighted features of corrupted samples on the India dataset.

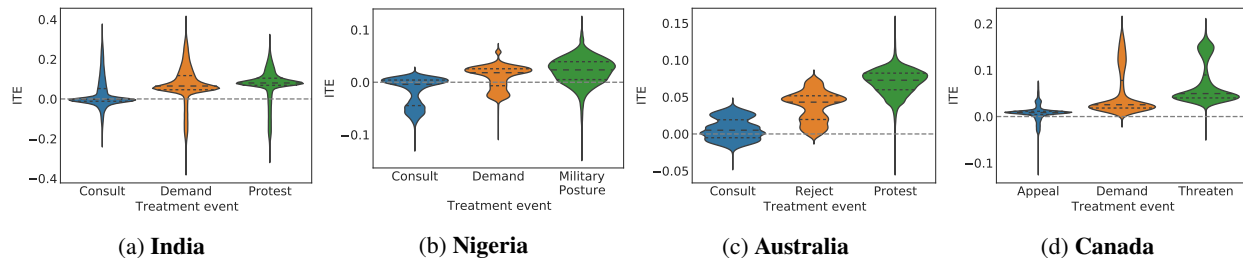


Figure 7: ITE distribution of different treatment events on the outcome *protest* in violin plot. The loosely dashed line represents the median.

dataset using **Cola-GNN** with the feature reweighting module. We select four corrupted test samples with random noise added to their input features (noise level of 5). We visualize the original features, the noisy features, and the ones obtained from the feature reweighting module. We can observe that the reweighted features can encode similar patterns of original features. It highlights the advantages of the ITE used in the feature reweighting module and demonstrates its ability to capture crucial information underlying the data distribution.

7.3 Causal Effect in Societal Events (RQ3)

In our study, whether there is a significant increase in certain types of events (e.g., appeal) over the past window is defined as the treatment of a location. The outcome is the future occurrence of a target event, i.e., protest. In this case, the ITE measures the difference in the outcome of the protest occurring between the two scenarios of the treatment event (i.e., increased or not). Thus, when the necessary assumptions hold, it implies a causal effect of the treatment event on the protest. A higher ITE suggests that an increase in a treatment event will be more influential in the occurrence of future protests, compared to a decrease or no change. To better illustrate the effect of treatment events on future protests, we visualize the predicted ITEs based on Eq. 12. Violin plots for the four datasets are shown in Fig. 7. We select three treatment events for each dataset. They have relatively low, moderate, and high ITE on average, respectively. The results vary from datasets due to different social environments. In India and Australia, massive historical protests may lead to future protests. In Nigeria and Canada, events related to military posture and threats, respectively, are likely to be more dominant factors in future protests. Nevertheless, we hardly conclude that protests will occur when the treatment event increases substantially because both types of events can be affected by hidden variables (i.e., unknown social factors). These results can provide supporting evidence for conjectures on protest triggers and generate hypotheses for future experiments.

8 Conclusion and Future Work

Learning causal effects of societal events is beneficial to decision-making and helps practitioners understand the underlying dynamics of events. In this paper, we introduce a deep learning framework that can estimate the causal effects of societal events and predict societal events simultaneously. We design a novel spatiotemporal causal inference model for estimating ITEs and propose two robust learning modules that use the learned causal information as prior knowledge for societal event prediction. We conducted extensive experiments on several real-world event datasets and showed that our approach achieves the best results in ITE estimation and robust event prediction. One future direction is to examine other potential causes of event occurrence, such as events with specific themes and potentially biased media coverage.

9 Broader Impacts

This work aims to advance computational social science by investigating causal effects among societal events from observational data. The causal effects among different types of societal events have not been extensively studied. In this work, we provide preliminary results on estimating the individual causal effects of one type of event on another and incorporate this causal information to improve the predictive power of event prediction models. We hope to provide a way to understand human behavior from the societal and causal inference aspects and broaden the possibilities for future work on societal event studies.

References

- [1] Liang Zhao, Feng Chen, Chang-Tien Lu, and Naren Ramakrishnan. Spatiotemporal event forecasting in social media. In *SIAM*, pages 963–971. SIAM, 2015.
- [2] Liang Zhao, Jieping Ye, Feng Chen, Chang-Tien Lu, and Naren Ramakrishnan. Hierarchical incomplete multi-source feature learning for spatiotemporal event forecasting. In *KDD*, pages 2085–2094. ACM, 2016.
- [3] Yue Ning, Sathappan Muthiah, Huzefa Rangwala, and Naren Ramakrishnan. Modeling precursors for event forecasting via nested multi-instance learning. In *KDD*, pages 1095–1104. ACM, 2016.
- [4] Songgaojun Deng, Huzefa Rangwala, and Yue Ning. Learning dynamic context graphs for predicting social events. In *KDD*, pages 1007–1016. ACM, 2019.
- [5] Songgaojun Deng, Huzefa Rangwala, and Yue Ning. Dynamic knowledge graph based multi-event forecasting. *KDD*, page 1585–1595, New York, NY, USA, 2020. Association for Computing Machinery.
- [6] Kira Radinsky, Sagie Davidovich, and Shaul Markovitch. Learning causality for news events prediction. In *WWW*, pages 909–918, 2012.
- [7] Kira Radinsky and Eric Horvitz. Mining the web to predict future events. In *WSDM*, pages 255–264, 2013.
- [8] Stephen Bonner and Flavian Vasile. Causal embeddings for recommendation. In *RecSys*, pages 104–112, 2018.
- [9] Jia Li, Xiaowei Jia, Haoyu Yang, Vipin Kumar, Michael Steinbach, and Gyorgy Simon. Teaching deep learning causal effects improves predictive performance. *arXiv preprint arXiv:2011.05466*, 2020.
- [10] Jin Chen, Xinxiao Wu, Yao Hu, and Jiebo Luo. Spatial-temporal causal inference for partial image-to-video adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1027–1035, 2021.
- [11] Christos Louizos, Uri Shalit, Joris M Mooij, David Sontag, Richard Zemel, and Max Welling. Causal effect inference with deep latent-variable models. In *NIPS*, pages 6446–6456, 2017.
- [12] Ruocheng Guo, Jundong Li, and Huan Liu. Learning individual treatment effects from networked observational data. *arXiv preprint arXiv:1906.03485*, 2019.
- [13] Jing Ma, Ruocheng Guo, Chen Chen, Aidong Zhang, and Jundong Li. Deconfounding with networked observational data in a dynamic environment. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, *WSDM '21*, page 166–174, New York, NY, USA, 2021. Association for Computing Machinery.
- [14] Andranik Tumasjan, Timm Oliver Sprenger, Philipp G Sandner, and Isabell M Welp. Predicting elections with twitter: What 140 characters reveal about political sentiment. *ICWSM*, 10(1):178–185, 2010.
- [15] Brendan O’Connor, Ramnath Balasubramanian, Bryan R Routledge, Noah A Smith, et al. From tweets to polls: Linking text sentiment to public opinion time series., 2010.
- [16] Johan Bollen, Huina Mao, and Xiaojun Zeng. Twitter mood predicts the stock market. *Journal of computational science*, 2(1):1–8, 2011.
- [17] Alessio Signorini, Alberto Maria Segre, and Philip M Polgreen. The use of twitter to track levels of disease activity and public concern in the us during the influenza a h1n1 pandemic. *PloS one*, 6(5):e19467, 2011.
- [18] Harshvardhan Achrekar, Avinash Gandhe, Ross Lazarus, Ssu-Hsin Yu, and Benyuan Liu. Predicting flu trends using twitter data. In *IEEE Conference on Computer Communications Workshops*, pages 702–707. IEEE, 2011.
- [19] Xiaofeng Wang, Matthew S Gerber, and Donald E Brown. Automatic crime prediction using events extracted from twitter posts. In *International conference on social computing, behavioral-cultural modeling, and prediction*, pages 231–238. Springer, 2012.
- [20] Nathan Kallus. Predicting crowd behavior with big public data. In *Proceedings of the 23rd International Conference on World Wide Web*, pages 625–630, 2014.

- [21] Songgaojun Deng, Huzefa Rangwala, and Yue Ning. Understanding event predictions via contextualized multilevel feature learning. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pages 342–351, 2021.
- [22] Liang Zhao, Qian Sun, Jieping Ye, Feng Chen, Chang-Tien Lu, and Naren Ramakrishnan. Multi-task learning for spatio-temporal event forecasting. In *KDD*, pages 1503–1512. ACM, 2015.
- [23] Yuyang Gao, Liang Zhao, Lingfei Wu, Yanfang Ye, Hui Xiong, and Chaowei Yang. Incomplete label multi-task deep learning for spatio-temporal event subtype forecasting. In *AAAI*, volume 33, pages 3638–3646, 2019.
- [24] Matthew S Gerber. Predicting crime using twitter and kernel density estimation. *Decision Support Systems*, 61:115–125, 2014.
- [25] Xiaofeng Wang, Donald E Brown, and Matthew S Gerber. Spatio-temporal modeling of criminal incidents using geographic, demographic, and twitter-derived information. In *ISI*, pages 36–41. IEEE, 2012.
- [26] Liang Zhao, Feng Chen, Jing Dai, Ting Hua, Chang-Tien Lu, and Naren Ramakrishnan. Unsupervised spatial event detection in targeted domains with applications to civil unrest modeling. *PLoS one*, 9(10), 2014.
- [27] Andrew Anglemyer, Hacsı T Horvath, and Lisa Bero. Healthcare outcomes assessed with observational study designs compared with those assessed in randomized trials. *Cochrane Database of Systematic Reviews*, (4), 2014.
- [28] Jan-Eric Gustafsson. Causal inference in educational effectiveness research: A comparison of three methods to investigate effects of homework on student achievement. *School Effectiveness and School Improvement*, 24(3):275–295, 2013.
- [29] Wei Sun, Pengyuan Wang, Dawei Yin, Jian Yang, and Yi Chang. Causal inference via sparse additive models with application to online advertising. In *AAAI*, 2015.
- [30] Jennifer L Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20(1):217–240, 2011.
- [31] Hugh A Chipman, Edward I George, Robert E McCulloch, et al. Bart: Bayesian additive regression trees. *The Annals of Applied Statistics*, 4(1):266–298, 2010.
- [32] Fredrik Johansson, Uri Shalit, and David Sontag. Learning representations for counterfactual inference. In *ICML*, pages 3020–3029, 2016.
- [33] Uri Shalit, Fredrik D Johansson, and David Sontag. Estimating individual treatment effect: generalization bounds and algorithms. In *ICML*, pages 3076–3085. JMLR. org, 2017.
- [34] Paul R Rosenbaum and Donald B Rubin. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983.
- [35] Yixin Wang and David M Blei. The blessings of multiple causes. *Journal of the American Statistical Association*, (just-accepted):1–71, 2019.
- [36] Laura von Rueden, Sebastian Mayer, Katharina Beckh, Bogdan Georgiev, Sven Giesselbach, Raoul Heese, Birgit Kirsch, Julius Pfrommer, Annika Pick, Rajkumar Ramamurthy, et al. Informed machine learning—a taxonomy and survey of integrating knowledge into learning systems. *arXiv preprint arXiv:1903.12394*, 2019.
- [37] Michelangelo Diligenti, Soumali Roychowdhury, and Marco Gori. Integrating prior knowledge into deep learning. In *2017 16th IEEE international conference on machine learning and applications (ICMLA)*, pages 920–923. IEEE, 2017.
- [38] Jingyi Xu, Zilu Zhang, Tal Friedman, Yitao Liang, and Guy Broeck. A semantic loss function for deep learning with symbolic knowledge. In *International conference on machine learning*, pages 5502–5511. PMLR, 2018.
- [39] Anuj Karpatne, William Watkins, Jordan Read, and Vipin Kumar. Physics-guided neural networks (pgnn): An application in lake temperature modeling. *arXiv preprint arXiv:1710.11431*, 2017.
- [40] Russell Stewart and Stefano Ermon. Label-free supervision of neural networks with physics and domain knowledge. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [41] Nikhil Muralidhar, Mohammad Raihanul Islam, Manish Marwah, Anuj Karpatne, and Naren Ramakrishnan. Incorporating prior domain knowledge into deep neural networks. In *2018 IEEE International Conference on Big Data (Big Data)*, pages 36–45. IEEE, 2018.
- [42] Peter W Battaglia, Razvan Pascanu, Matthew Lai, Danilo Rezende, and Koray Kavukcuoglu. Interaction networks for learning about objects, relations and physics. *arXiv preprint arXiv:1612.00222*, 2016.
- [43] Woojeong Jin, Meng Qu, Xisen Jin, and Xiang Ren. Recurrent event network: Autoregressive structure inference over temporal knowledge graphs. In *EMNLP*, 2020.

- [44] Ioana Bica, Ahmed Alaa, and Mihaela Van Der Schaar. Time series deconfounder: Estimating treatment effects over time in the presence of hidden confounders. In *International Conference on Machine Learning*, pages 884–895. PMLR, 2020.
- [45] Donald B Rubin. Bayesian inference for causal effects: The role of randomization. *The Annals of statistics*, pages 34–58, 1978.
- [46] Donald B Rubin. Causal inference using potential outcomes: Design, modeling, decisions. *Journal of the American Statistical Association*, 100(469):322–331, 2005.
- [47] Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018.
- [48] Judea Pearl et al. Causal inference in statistics: An overview. *Statistics surveys*, 3:96–146, 2009.
- [49] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- [50] Yann N Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with gated convolutional networks. In *International conference on machine learning*, pages 933–941. PMLR, 2017.
- [51] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.
- [52] Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, and Chengqi Zhang. Graph wavenet for deep spatial-temporal graph modeling. *arXiv preprint arXiv:1906.00121*, 2019.
- [53] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017.
- [54] Risto Vuorio, Shao-Hua Sun, Hexiang Hu, and Joseph J Lim. Multimodal model-agnostic meta-learning via task-aware modulation. *arXiv preprint arXiv:1910.13616*, 2019.
- [55] Bingyi Kang, Zhuang Liu, Xin Wang, Fisher Yu, Jiashi Feng, and Trevor Darrell. Few-shot object detection via feature reweighting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8420–8429, 2019.
- [56] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, pages 5998–6008, 2017.
- [57] Elizabeth Boschee, Jennifer Lautenschlager, Sean O’Brien, Steve Shellman, James Starz, and Michael Ward. Icews coded event data, 2015.
- [58] Kalev Leetaru and Philip A Schrod. Gdelt: Global data on events, location, and tone, 1979–2012. In *ISA annual convention*, volume 2, pages 1–49. Citeseer, 2013.
- [59] Elizabeth Boschee, Jennifer Lautenschlager, Sean O’Brien, Steve Shellman, James Starz, and Michael Ward. *Cameo.cdb.09b5.pdf*. In *ICEWS Coded Event Data*. Harvard Dataverse, 2015.
- [60] Liu Yang and Rong Jin. Distance metric learning: A comprehensive survey. *Michigan State University*, 2(2):4, 2006.
- [61] Liuyi Yao, Sheng Li, Yaliang Li, Mengdi Huai, Jing Gao, and Aidong Zhang. Representation learning for treatment effect estimation from observational data. *Advances in Neural Information Processing Systems*, 31, 2018.
- [62] Songgaojun Deng, Shusen Wang, Huzefa Rangwala, Lijing Wang, and Yue Ning. Cola-gnn: Cross-location attention based graph neural networks for long-term ili prediction. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pages 245–254, 2020.
- [63] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, 2010.
- [64] D Kinga and J Ba Adam. A method for stochastic optimization. In *ICLR*, volume 5, 2015.