
Counterfactual Neural Temporal Point Process for Misinformation Impact Estimation on Social Media

Anonymous Author(s)

Affiliation

Address

email

1 A Appendix

2 A.1 Data and Code

3 Together with this appendix, we also submit the synthetic data and necessary code. However, the
4 real-world data (from Twitter) is too huge for OpenReview platform. Thus, we can not submit it right
5 now. We will opensource a version of the data with anonymous processing to protect privacies after
6 this paper is accepted.

7 A.2 Proof of the theorems

8 **Theorem 1.** *For a user u , if the intensity function $\lambda(\mathbf{f}, t|Tr \cup X)$ is known, then we have:*

$$\mu(t, t+T, \lambda_1, Tr \cup X) = \int_{sup(\mathbf{f})} d\mathbf{f} \int_t^{t+T} \lambda(\mathbf{f}, t|Tr \cup X) dt \quad (1)$$

$$\phi(t, t+T, \lambda_1, Tr \cup X) = \int_{sup(\mathbf{f})} \mathbf{f} d\mathbf{f} \int_t^{t+T} \lambda(\mathbf{f}, t|Tr \cup X) dt \quad (2)$$

10 *Proof.* The first equation can be trivially proved by replacing the \mathbf{F} in the definition of λ with the
11 support set. The second one can be proved with the Campbell's Theorem [1]:

12 **Lemma 1.** *Campbell's Theorem (Campbell, 1909): For a point process S , and a measurable function*
13 *$f: \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$:*

$$\mathbb{E}[\sum_{x \in S} f(x)] = \int_{\mathbb{R}^d} f(x) \mathbb{E}(N(x, x+dx)) \quad (3)$$

14 Note that the above theorem is for all point process. For the specific scenario of temporal point
15 process, we also need to include the time range into consideration:

$$\mathbb{E}[\sum_{x \in S[t_1:t_2]} f(x)] = \int_{\mathbb{R}^d} f(x) \mathbb{E}(N(x, x+dx, t_1, t_2)) \quad (4)$$

16 From the definition of ϕ , we know that the $f(x) = x$. Recall the definition of intensity function
17 (Equation. 1 in the main content):

$$\mathbb{E}(N(\mathbf{F}, T_1, T_2)|S_h) = \int_{\mathbf{F}} d\mathbf{f} \int_{T_1}^{T_2} \lambda(\mathbf{f}, t|S_h) dt \quad (5)$$

18 We can convert this integral to a derivative formula:

$$\mathbb{E}(N(\mathbf{f}, \mathbf{f} + d\mathbf{f}, T_1, T_2)|S_h) = d\mathbf{f} \int_{T_1}^{T_2} \lambda(\mathbf{f}, t|S_h) dt \quad (6)$$

19 With the above formula and Campbell's theorem, we have:

$$\phi(t, t + T, \lambda_1, Tr \cup X) = \int_{\sup(\mathbf{f})} \mathbf{f} d\mathbf{f} \int_t^{t+T} \lambda(\mathbf{f}, t|Tr \cup X) dt \quad (7)$$

20

□

21 **Theorem 2.** *Given the following min-max game:*

$$\min_H \max_{\hat{p}} \mathbb{E}_{Tr, X \sim p(Tr, X)} \log \hat{p}(Tr|H(X)) \quad (8)$$

22 *the min gamer's Nash balanced solution H^* , ensures for any X_1, X_2 , the following equation holds:*

$$p(Tr|H^*(X_1)) = p(Tr|H^*(X_2)) \quad (9)$$

23 *where p denote the ground-truth conditional distribution of treatment given encoding.*

24 *Proof.* For the objective of the max gamer, we have:

$$\begin{aligned} \mathbb{E}_{Tr, X \sim p(Tr, X)} \log \hat{p}(Tr|H(X)) &= \mathbb{E}_{Tr, H(X) \sim p(Tr, H(X))} \log \hat{p}(Tr|H(X)) \\ &= \mathbb{E}_{H(X) \sim p(H(X))} \mathbb{E}_{Tr \sim p(Tr|H(X))} \log \hat{p}(Tr|H(X)) \\ &= \mathbb{E}_{H(X) \sim p(H(X))} - KL(p(Tr|H(X)), \hat{p}(Tr|H(X))) \\ &\quad - Q(p(Tr|H(X))) \end{aligned} \quad (10)$$

25 where $Q(p(Tr|H(X)))$ is the entropy of $p(Tr|H(X))$. Note that for a fixed H , $Q(p(Tr|H(X)))$
26 is a constant. Thus, the max objective is the same as minimizing the KL-divergence between
27 $p(Tr|H(X))$ and $\hat{p}(Tr|H(X))$ for every X :

$$\max_{\hat{p}} \mathbb{E}_{Tr, X \sim p(Tr, X)} \log \hat{p}(Tr|H(X)) \iff \min_{\hat{p}} \mathbb{E}_{H(X) \sim p(H(X))} KL(p(Tr|H(X)), \hat{p}(Tr|H(X))) \quad (11)$$

28 It is known that the KL-divergence is minimized if and only if $\hat{p}(Tr|H(X))$ perfectly fit $p(Tr|H(X))$
29 at each point. Thus, the balance solution of the max gamer is to learn $\hat{p}(Tr|H(X)) = p(Tr|H(X))$.
30 Meanwhile, if there are two X_1 and X_2 such that $p(Tr|H(X_1)) \neq p(Tr|H(X_2))$, then the min
31 gamer can always update the H by swapping $H(X_1)$ and $H(X_2)$ while maintaining other points the
32 same. This is doable when H is a universal function approximator. As a result, the KL-divergence
33 will increase, which violates the definition of Nash balance. Therefore, at the balance point for any
34 X_1 and X_2 , $p(Tr|H(X_1)) = p(Tr|H(X_2))$ □

35 A.3 Synthetic Dataset

36 In this section, we introduce the how to generate the synthetic data in details. We present the basic
37 ideas in Figure 1. To generate the data, we design a system to simulate the social media activities and
38 events in real world. The basic elements in this system are **users** and **news**.

- 39 • **User:** Each user i is represented with a **hidden** vector u_i , which represent the hidden status
40 of a social media user, such as interests, ideas and political trends in reality. If a user engage
41 with a piece of news, then this user's hidden vector will change accordingly. We randomly
42 assign each user with an randomly initialized vector to simulate the fact that users may hold
43 different ideas and interests before starting to use social media.
- 44 • **News:** Each piece of news n is characterized by two randomly generated **hidden** feature
45 vectors: a topic vector $v_{topic}(n)$ and an inherent influence vector $v_{in}(n)$. These two vectors
46 decide which users will be attracted by the news and how they will be influenced by the
47 news.

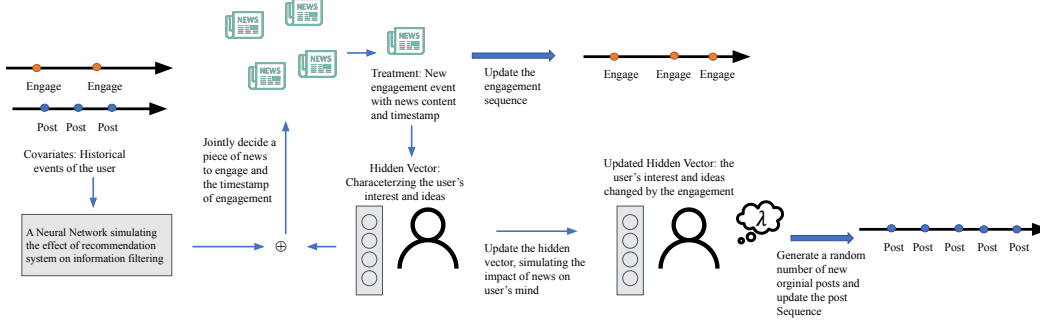


Figure 1: The general idea of generating the synthetic dataset

In this system, a user has two kinds of activities: (1) engaging with a piece of news and (2) posting a post with original contents. These activities are characterized by two temporal point process respectively. For the process of engaging with a piece of news, the intensity function of a user's (i) engagement with a specific piece of news n is defined as:

$$\lambda_e(n, t | u_i, S_h^{(e)}, S_h^{(g)}) = v_{topic}(n) \cdot u_i + \sum_{(\mathbf{f}_j, t_j) \in S_h^{(e)} \cup S_h^{(g)}} \exp(t_j - t) NN_{rec}(\mathbf{f}_j, v_{topic}(n)) \quad (12)$$

where $S_h^{(e)}$ and $S_h^{(g)}$ denotes the historical sequences of a user's engagement and posting activities respectively. In this equation, the first term models the consistency of news topic and user interest and the $NN_{rec}(\cdot, \cdot)$ in the second term is a random parameterized neural network simulating the recommendation system that recommend contents based on past history. $\exp(t_j - t)$ make sure that the influence of a history event decrease as time goes by, simulating the property that recommendation system usually give higher attention to recent events.

Similarly, the intensity function of a user's (i) process generating original posts is defined as:

$$\lambda_g(t | u_i, S_h^{(e)}, S_h^{(g)}) = u_i \cdot u_i + \sum_{(\mathbf{f}_j, t_j) \in S_h^{(e)} \cup S_h^{(g)}} \exp(t_j - t) NN_{post}(\mathbf{f}_j, u_i) \quad (13)$$

where $u_i \cdot u_i$ models the user's reliance to social media and the second term reflect the influence of history events to future event.

In real world, the news topic, influence on user mind and user's mind status are not observable. Instead we can only observe the tweet content text features, such as the representation extracted with a pre-trained language model. To simulate this fact, we use two neural networks NN_{news} and NN_{user} with random parameters to project the vectors to observable features:

$$\mathbf{f}_{news}(n) = NN_{news}(v_{topic}(n), v_{in}(n)), \quad \mathbf{f}_{user}(i, t) = NN_{user}(u_i, t) \quad (14)$$

where $\mathbf{f}_{news}(n)$ is the observable content feature of news n and $\mathbf{f}_{user}(i, t)$ is the observable content feature of the post generated by user i at time t . $\mathbf{f}_{news}(n)$ does not change over time because once a piece of news is reported, its content does not change significantly. Note that a normal user's activity is usually uncertain. To simulate such uncertainty, in each hidden layer of NN_{user} , we add a random noise. Since the two networks are with random parameters, the projection will be highly non-linear and thus brings enough complexity on the learning of causal impact.

The overall algorithm is shown in Algorithm 1:

A.4 Synthetic Data Experiments

We remove the users with only one event and then acquire around 13000 accounts. Then we split the training set, validation set and test set according to 8000/2000/3000 to train the proposed model. Specifically, we adopt the Adam optimizer with learning rate $1e - 3$ and the max epoch is set to be

Algorithm 1 Synthetic Data Generation

```
1: Set the ending time  $T$ . {The simulation will end at time  $T$ }
2: Randomly generated user hidden representations and news features and a posting time of each
   news. {Each piece of news will only be engaged with after it is posted.}
3: for news  $n$  do
4:    $\mathbf{f}_{news}(n) = NN_{news}(v_{topic}(n), v_{in}(n))$ 
5: end for
6: for user  $i$  do
7:    $t_{end} = 0$ 
8:    $S_h^{(e)}, S_h^{(g)} = \emptyset, \emptyset$ 
9:   repeat
10:    Draw a posting event timestamp  $t_g$  based on intensity function  $\lambda_g(t|u_i, S_h^{(e)}, S_h^{(g)})$ 
11:     $t_{cur} = t_g$ 
12:     $\mathbf{f}_{cur} = \mathbf{f}_{user}(i, t_g)$ 
13:     $flag = 'generating'$ 
14:    for news  $n$  do
15:      if  $t_p(n) \geq t_{cur}$  then
16:        Continue.
17:      end if
18:      Draw an event timestamp  $t_e$  based on intensity function  $\lambda_e(n, t|u_i, S_h^{(e)}, S_h^{(g)})$ 
19:      if  $t_e < t_{cur}$  then
20:         $t_{cur} = t_e$ 
21:         $\mathbf{f}_{cur} = \mathbf{f}_{news}(n)$ 
22:         $flag = 'engagement'$ 
23:         $n_{cur} = n$ 
24:      end if
25:    end for
26:     $t_{end} = t_{cur}$ 
27:    if  $flag$  is  $'engagement'$  then
28:       $S_h^{(e)} = S_h^{(e)} \cup (\mathbf{f}_{cur}, t_{cur})$ 
29:       $u_i = u_i + NN_{scale}(v_{topic}(n_{cur}), u_i)v_{in}(n_{cur})$  {Update the user status.  $NN_{scale}$  is a
        neural network with random parameters to adjust news impact scale.}
30:    else
31:       $S_h^{(g)} = S_h^{(g)} \cup (\mathbf{f}_{cur}, t_{cur})$ 
32:    end if
33:  until  $t_{end} \geq T$ 
34: end for
```

76 500 with batch size 128. In addition, we add the dropout with 0.1 rate. We train the models and
77 baselines and evaluate their inference time on 1 NVIDIA GeForce RTX 2080 TI GPUs for the fair
competition. The inference time with an error bar is presented in Table 1

| Method | Inference Time |
|-----------------------------------------|-------------------|
| FullyNN (CNTPP w/o Adversarial Balance) | 7.23±0.2 ms |
| CRN-VAE (Approximation) | 4.05±0.1ms |
| CRN-VAE (Sampling) | 29.34± 2.3ms |
| CNTPP(Ours) | 7.12± 0.2ms |

Table 1: Inference Time with Error Bar

78

79 A.5 Real World Dataset

80 To collect the Twitter dataset related to COVID-19 vaccines, we use the tracked keywords including
81 vaccine, Pfizer, BioNTech, Moderna, Janssen, AstraZeneca, Sinopharm to filtering all tweets via the
82 streaming Twitter API which returns a 1% sample of all tweets. We collect the tweets just before
83 Pfizer-BioNTech and Moderna were approved by the FDA for Emergency Use Authorization (EUA)

| Methods | ASD↓ | ASD _{in} ↓ | ASD _{mis} ↓ |
|---------------|--------------|---------------------|----------------------|
| Event Feature | 0.123 | 0.123 | 0.122 |
| FullyNN | 0.073 | 0.069 | 0.072 |
| CNTPP (ours) | 0.045 | 0.042 | 0.044 |

Table 2: Comparison of (normalized) Average Sum of Distances with different methods on the real-world dataset

from Dec 9, 2020 - April 24, 2021. After that, we only use the tweets with labels of 'MisInformation' and 'Information', including a total of 169,008 tweets from 24,192 users.

A.6 Real World Data Experiments

For the impact on sentiments and subjectivity experiments, we use first use TextBlob sentiment analyse tool [3] to generate the continuous sentiment and subjectivity score of the content user posts, where the former reflects the emotional charge of a statement and the latter reflects the degree of objectivity. After that, we use the tweets' timestamp to calculate the time interval, and use the sentiment and subjectivity score including whether the news is positive (1) or not (0) as the covariates of our causal structure. We split the training set, validation set and test set according to the ratio of 80%/10%/10% to train the proposed model. Specifically, we adopt the Adam optimizer with learning rate $1e-5$ and the epoch is set to be 500 with batch size 128. In addition, we add the dropout with 0.1 rate.

For the impact on text representation, we use the pre-trained BERT model [2] to extract the representation vectors of each tweet's content with 768 dimension. After that, we use PCA tools to reduce into 2 dimensions and use the new vectors as the covariates in a new model. The training setting is same with the impact on sentiments and subjectivity experiments.

Note that all the experiments are finished on the NVIDIA GeForce RTX 2080 TI GPUs for the fair competition.

Besides, we also evaluate FullyNN (a baseline) on the real-world dataset. Its results are visualized in Figure 2(a) and 2(b). As we can see, it can not acquire recognizable differences between information and misinformation. We further visualize the raw event features (BERT embedding with dimension reduction of PCA, not finetuned on fake news detection) in Figure 2(c). As we can see, it is also hard for a pre-trained language model to distinguish information and misinformation in an unsupervised manner.

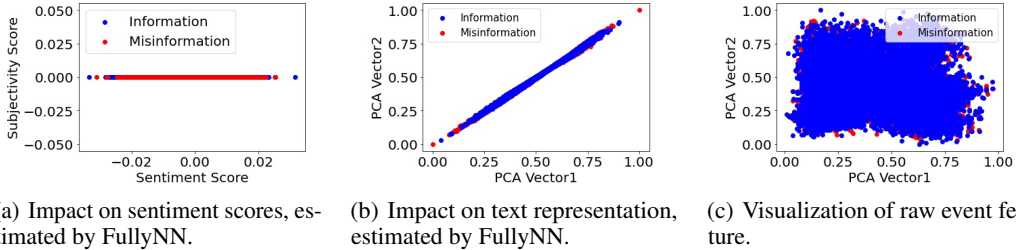


Figure 2: Visualization of the results from baselines on real world

We also calculate Normalized Averaged Sum of Distances (NASD) for information cluster and misinformation cluster. More specifically, we first normalize the scale range of each embedding dimension (ITE vector or raw input) to be same for different methods to ensure that the scales embedding spaces of different methods are comparable. Then, we calculate their Averaged Sum of Distances (ASD), which is defined as $ASD = \sum_{i=0}^k \sum_{p \in C_i} (p - m_i) / n$, where C_i is the set of ITE results (or event feature) and $i = 1$ means the information results, $i = 0$ means the misinformation results, p stands for an instance in C_i , m_i is the center of all instance representations in C_i , n is the

total number of instance. The lower these metric are, the better that information and misinformation are distinguished. We also separately report the ASD on two clusters $ASD_{in} = \sum_{p \in C_0} (p - m_0)/n_0$ and $ASD_{mis} = \sum_{p \in C_0} (p - m_1)/n_1$. The comparison of our model against FullyNN and event features on this metric is shown in Table 2.

References

- [1] N. Campbell. The study of discontinuous phenomena. *Proc Cambr. Phil. Soc*, vol. 15:pp. 117–136, 1909.
- [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [3] Steven Loria et al. textblob documentation. *Release 0.15*, 2:269, 2018.