
Influence Function Learning in Information Diffusion Networks

Nan Du, Yingyu Liang
Maria-Florina Balcan
Le Song

{DUNAN,YLIANG39}@GATECH.EDU
NINAMF@CC.GATECH.EDU
LSONG@CC.GATECH.EDU

College of Computing, Georgia Institute of Technology, 266 Ferst Drive, Atlanta, 30332 USA

Abstract

Can we learn the influence of a set of people in a social network from cascades of information diffusion? This question is often addressed by a two-stage approach: first learn a diffusion model, and then calculate the influence based on the learned model. Thus, the success of this approach relies heavily on the correctness of the diffusion model which is hard to verify for real world data. In this paper, we exploit the insight that the influence functions in many diffusion models are coverage functions, and propose a novel parameterization of such functions using a convex combination of random basis functions. Moreover, we propose an efficient maximum likelihood based algorithm to learn such functions directly from cascade data, and hence bypass the need to specify a particular diffusion model in advance. We provide both theoretical and empirical analysis for our approach, showing that the proposed approach can provably learn the influence function with low sample complexity, be robust to the unknown diffusion models, and significantly outperform existing approaches in both synthetic and real world data.

1. Introduction

Social networks are important in information diffusion, which has motivated the influence maximization problem: find a set of nodes whose initial adoptions of an idea can trigger the largest number of follow-ups. This problem has been studied extensively in literature from both modeling and algorithmic point of view (Kempe et al., 2003; Chen et al., 2010; Borgs et al., 2012; Rodriguez & Schölkopf, 2012; Du et al., 2013b). Essential to the influence maximization problem is the influence function of a set of nodes, which is an estimate of the expected number of triggered follow-ups from these nodes.

Proceedings of the 31st International Conference on Machine Learning, Beijing, China, 2014. JMLR: W&CP volume 32. Copyright 2014 by the author(s).

In practice, the influence function is not given to us, and we only observe the information diffusion traces, or cascades, originating from these nodes. In order to model the cascade data, many information diffusion models have been proposed in the literature, such as the discrete-time independent cascade model and linear threshold model (Kempe et al., 2003), and more recently the continuous-time independent cascade model (Gomez Rodriguez et al., 2011; Du et al., 2013b). To estimate the influence, we can employ a two-stage method: a particular diffusion model is first learned from cascade data, and then the influence function is evaluated or approximated from such learned model.

However, there still remain many challenges in these traditional two-stage approaches. First, real world information diffusion is complicated, and it is not easy to determine the most suitable diffusion model in practice. A chosen diffusion model may be misspecified compared to real world data and lead to large model bias. Second, the diffusion network structure can be also hidden from us, so we need to learn not only the parameters in the diffusion model, but also the diffusion network structure. This often leads to under-determined high dimensional estimation problem where specialized methods need to be designed (Du et al., 2012; 2013a). Third, calculating the influence based on learned diffusion models often leads to difficult graphical model inference problem where extra approximation algorithms need to be carefully designed (Du et al., 2013b).

If the sole purpose is to estimate the influence, can we avoid the challenging diffusion model learning and influence computation problem? In this paper, we provide a positive answer to the question and propose an approach which estimates the influence function directly from cascade data. Our approach will exploit the observation that the influence functions in many diffusion models are coverage functions. Instead of learning a particular diffusion model, we will aim to learn a coverage function instead, which will then naturally subsume many diffusion models as special cases. Furthermore, in the information diffusion context, we show that the coverage function can be represented as a sum of simpler functions, each of which is an expectation over random binary functions. Based on these structures of the problem, we propose a maximum-

likelihood based approach to learn the influence function directly from cascade data. More precisely,

Direct and robust approach. Our algorithm does not rely on the assumption of a particular diffusion model, and can be more robust to model misspecification than two-stage approaches. Furthermore, directly learning the coverage function also allows us to avoid the difficulty involved in diffusion model estimation and influence computation.

Novel Parameterization. We propose a parametrization of the coverage function using a convex combination of random basis function. Similar parameterization has been used in classification and kernel methods setting (Rahimi & Recht, 2008), but its usage in the information diffusion and coverage function estimation context is novel.

Approximation guarantee. We show that our parameterization using K random basis functions generates a rich enough family of functions which can approximate the true influence function within an error of $O(\frac{1}{\sqrt{K}})$. This allows us to work with a small number of parameters without creating too much bias at the same time.

Efficient algorithm. We propose a maximum likelihood based convex formulation to estimate the parameters, which allows us to leverage existing convex optimization techniques (Kivinen & Warmuth, 1997) to solve the problem efficiently. The time required to evaluate each gradient is $O(dmK)$, linear in the number of nodes d , the number of cascades m , and the number of basis functions K .

Sample complexity. We prove that to learn the influence function to an ϵ error, we only need $O(\frac{d^3}{\epsilon^3})$ cascades where d is the number of nodes in the diffusion networks. This is no obvious since the number possible source configurations can be exponential in the number of nodes in the network. Our approach is able to make use of the structure of the coverage function and be able to generalize only after seeing a polynomial number of cascades.

Superior performance. We evaluate our algorithms using large-scale datasets, and show that it achieves significantly better performance in both the synthetic cases where there is known model misspecification, and in real world data where the true model is completely unknown in advance.

2. Diffusion Models and Influence Function

Several commonly used models exist for information diffusion over networks. Interestingly, although these models are very different in nature, the derived influence functions belong to the same type of combinatorial functions — coverage functions. Such commonality allows us later to approach the problem of learning influence functions directly without assuming a particular diffusion model.

More specifically, a diffusion model is often associated

with a directed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, and a cascade from a model is just a set of influenced nodes according to the model given a set of source nodes $\mathcal{S} \subseteq \mathcal{V}$. In general, we have the following typical types of diffusion models :

Discrete-time independent cascade model (Kempe et al., 2003). Each edge is associated with a weight in $[0, 1]$. When a cascade is being generated from the source nodes \mathcal{S} , independently for each edge according to the edge weight, a binary random variable is sampled, indicating whether the edge is included in a “live edge graph” or not. The influenced nodes are those reachable from at least one of the source nodes in the resulting “live edge graph”.

Discrete-time linear threshold model (Kempe et al., 2003). Each edge is again associated with a weight in $[0, 1]$, but the sum of the incoming edge weights for each node is smaller or equal to 1. When a cascade is being generated from the source nodes \mathcal{S} , each node first independently sample one of its incoming edges with probability proportional to the edge weight. The chosen edges are then used to form the “live edge graph”. The influenced nodes are those reachable from at least one of the source nodes.

Continuous-time independent cascade model (Du et al., 2013b). Being different from the discrete-time models, this model associates each edge (j, i) with a transmission function, $f_{ji}(\tau_{ji})$, a density over time. The source nodes are assumed to be initially influenced at time zero. Then a diffusion time is sampled independently for each edge according to the transmission function and is viewed as the length of the edge. The influenced nodes are those within shortest distance T from at least one of the source nodes.

Being common to these diffusion models, the influence function, $\sigma(\mathcal{S}) : 2^{\mathcal{V}} \mapsto \mathbb{R}_+$, of a set of nodes \mathcal{S} , is defined as the expected number of influenced nodes with respect to the generative process of each model. This influence function is a combinatorial function which maps a subset \mathcal{S} of \mathcal{V} to a nonnegative number.

Although these diffusion models are very different in nature, their corresponding influence functions belong to the same type of functions — *coverage functions*, and share very interesting combinatorial structures (Kempe et al., 2003; Rodríguez & Schölkopf, 2012). This means that the influence function can be written as

$$\sigma(\mathcal{S}) = \sum_{u \in \bigcup_{s \in \mathcal{S}} \mathcal{A}_s} a_u, \quad (1)$$

with three sets of objects :

- (i) a ground set \mathcal{U} which may be different from the set \mathcal{V} of nodes in the diffusion network,
- (ii) a set of nonnegative weights $\{a_u\}_{u \in \mathcal{U}}$, each associated with an item in the ground set \mathcal{U} ,
- (iii) and a collection of subsets $\{\mathcal{A}_s : \mathcal{A}_s \subseteq \mathcal{U}\}_{s \in \mathcal{V}}$, one for each source node in diffusion network.

Essentially, each source node $s \in \mathcal{S}$ covers a set \mathcal{A}_s of items from \mathcal{U} , and the function value $\sigma(\mathcal{S})$ is the weighted sum over the union of items covered by all nodes in \mathcal{S} .

The combinatorial structures of coverage functions allow them to be potentially learned directly from cascades. However, the problem of learning coverage functions is very challenging for several reasons. First, there are an exponential number of different \mathcal{S} from the power set of \mathcal{V} , while one typically only observes a small number of cascades polynomial in the number of nodes, $d = |\mathcal{V}|$, in the network. Second, both the ground set \mathcal{U} , the weights $\{a_u\}$ and the subsets $\{\mathcal{A}_s\}$ are unknown, and one has to estimate a very large set of parameters if one wants to use the definition in (1) directly.

In fact, learning such combinatorial functions in general settings has attracted many recent research efforts (Balcan & Harvey, 2011; Badanidiyuru et al., 2012; Feldman & Kothari, 2013; Feldman & Vondrak, 2013), many of which show that coverage functions can be learned from just polynomial number of samples. However, existing algorithms are mostly of theoretical interest and impractical for real world problem yet. To tackle this challenge, we will exploit additional structure of the coverage function in the information diffusion context which allows us to derive compact parameterization of the function, and design a simple and efficient algorithm with provable guarantees.

3. Structure of the Influence Function

Besides being coverage functions, the influence functions, $\sigma(\mathcal{S})$, in the diffusion models discussed in Section 2 share additional structures. In all models, a random graph \mathcal{G} is first sampled from the distribution induced by a particular diffusion model; and then a function is defined for computing node reachability in the sampled graph; finally the influence is the expectation of this function with respect to the distribution of the random graphs.

3.1. Random reachability function

We represent each sampled random graph \mathcal{G} as a binary reachability matrix $\mathbf{R} \in \{0, 1\}^{d \times d}$ with (s, j) -th entry

$$\mathbf{R}_{sj} = \begin{cases} 1, & j \text{ is reachable from source } s, \\ 0, & \text{otherwise.} \end{cases} \quad (2)$$

Essentially, the s -th row of \mathbf{R} , denoted as $\mathbf{R}_{s,:}$, records the information that if s is the source, which node is reachable given sampled graph \mathcal{G} . Furthermore, the j -th column of \mathbf{R} , denoted as $\mathbf{R}_{:,j}$, records the information that whether j is reachable from each of the other nodes. Then given a set \mathcal{S} of sources, we can calculate whether a node j will be influenced or not in graph \mathcal{G} through a simple nonlinear function ϕ defined below.

First, we represent the set \mathcal{S} as an indicator vector $\chi_{\mathcal{S}} \in$

$\{0, 1\}^d$, with its i -th entry

$$\chi_{\mathcal{S}}(s) := \begin{cases} 1, & s \in \mathcal{S}, \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

Then the inner product $\chi_{\mathcal{S}}^{\top} \mathbf{R}_{:,j} \in \mathbb{Z}_+$ will give us an indication whether a target node j is reachable from any of the sources in \mathcal{S} . More specifically, $\chi_{\mathcal{S}}^{\top} \mathbf{R}_{:,j} \geq 1$ if the target node j is reachable, and 0 otherwise. Finally, using a concave function $\phi(u) = \min\{u, 1\} : \mathbb{Z}_+ \mapsto \{0, 1\}$, we can transform $\chi_{\mathcal{S}}^{\top} \mathbf{R}_{:,j}$ into a binary function of $\chi_{\mathcal{S}}$

$$\phi(\chi_{\mathcal{S}}^{\top} \mathbf{R}_{:,j}) : 2^{\mathcal{V}} \mapsto \{0, 1\}. \quad (4)$$

We note that $\phi(\chi_{\mathcal{S}}^{\top} \mathbf{R}_{:,j})$ itself is a coverage function where (i) the ground set \mathcal{U} contains a single item u_j , (ii) the weight on u_j is 1, (iii) and the collection of subset is either $\mathcal{A}_s = \{u_j\}$ if $\mathbf{R}_{sj} = 1$ and otherwise $\mathcal{A}_s = \emptyset$ if $\mathbf{R}_{sj} = 0$.

Then the influence of \mathcal{S} in graph \mathcal{G} is the number of target nodes reachable from the source set \mathcal{S}

$$\#(\mathcal{S}|\mathbf{R}) := \sum_{j=1}^d \phi(\chi_{\mathcal{S}}^{\top} \mathbf{R}_{:,j}). \quad (5)$$

$\#(\mathcal{S}|\mathbf{R})$ is also a coverage function where (i) the ground set \mathcal{U} contains d items u_1, \dots, u_d , (ii) the weight on each u_j is 1, (iii) and $\mathcal{A}_s = \{u_j | \mathbf{R}_{sj} = 1\}$. Since the graph \mathcal{G} and the associated \mathbf{R} are random quantities, the Φ function is a random function.

3.2. Expectation of random functions

Each diffusion model will induce a distribution over random graph \mathcal{G} and hence a distribution $p_{\mathbf{R}}$ over the random binary matrix \mathbf{R} . Then the overall influence of a source set \mathcal{S} in a diffusion model is the expected value of $\#(\mathcal{S}|\mathbf{R})$, i.e.,

$$\sigma(\mathcal{S}) := \mathbb{E}_{\mathbf{R} \sim p_{\mathbf{R}}} [\#(\mathcal{S}|\mathbf{R})], \quad (6)$$

which is also a coverage function, since non-negative combinations of coverage functions are still coverage functions (See Appendix A).

Next we will manipulate expression (6) to expose its structure as a sum over a set of conditional probabilities

$$\begin{aligned} & \mathbb{E}_{\mathbf{R} \sim p_{\mathbf{R}}} [\#(\mathcal{S}|\mathbf{R})] \\ &= \mathbb{E}_{\mathbf{R} \sim p_{\mathbf{R}}} \left[\sum_{j=1}^d \phi(\chi_{\mathcal{S}}^{\top} \mathbf{R}_{:,j}) \right] \quad (\text{by definition (5)}) \\ &= \sum_{j=1}^d \mathbb{E}_{\mathbf{R} \sim p_{\mathbf{R}}} [\phi(\chi_{\mathcal{S}}^{\top} \mathbf{R}_{:,j})] \quad (\text{sum} \Leftrightarrow \text{expectation}) \\ &= \sum_{j=1}^d \underbrace{\Pr \{ \phi(\chi_{\mathcal{S}}^{\top} \mathbf{R}_{:,j}) = 1 | \chi_{\mathcal{S}} \}}_{:= f_j(\chi_{\mathcal{S}})} \quad (\phi(\cdot) \text{ is binary}), \end{aligned}$$

where $f_j(\chi_{\mathcal{S}})$ is the conditional probability of $\phi(\chi_{\mathcal{S}}^{\top} \mathbf{R}_{:,j})$ being 1 given the set indicator $\chi_{\mathcal{S}}$.

Strategy for learning: The form of the influence function as a sum over conditional probabilities suggests a simple strategy for learning the influence function:

1. we learn each $f_j(\chi_S)$ separately,
2. and then sum them together,

which we will elaborate in subsequent sections.

4. Random Basis Function Approximation

In this section, we will provide a novel parameterization of function $f_j(\chi_S)$ using random basis functions. Recall from the derivation in (7) that

$$f_j(\chi_S) = \mathbb{E}_{r \sim p_j(r)} [\phi(\chi_S^\top r)] \quad (8)$$

where $r := \mathbf{R}_{:,j}$ and $p_j(r)$ is the marginal distribution of column j of \mathbf{R} induced by $p_{\mathbf{R}}$. Since $f_j(\chi_S)$ is an expectation w.r.t. a distribution $p_j(r)$ over the binary vectors $\{0, 1\}^d$, we will use a convex combination of random basis functions to parameterize $f_j(\chi_S)$. A similar idea, called random kitchen sinks (Rahimi & Recht, 2008), has appeared in the classification and kernel methods context. Our use of such parameterization is novel in the information diffusion and coverage function learning context, and our analysis is also different.

Specifically, consider drawing a set of K random binary vectors (random features) $\{r_1, r_2, \dots, r_K\}$ from some distribution $q(r)$ on $\{0, 1\}^n$, and build functions of the form

$$f^w(\chi_S) = \sum_{k=1}^K w_k \phi(\chi_S^\top r_k) = w^\top \phi(\chi_S), \quad (9)$$

$$\text{subject to } \sum_{k=1}^K w_k = 1, w_k \geq 0 \quad (10)$$

where $w := (w_1, \dots, w_K)^\top$ are parameters to be learned, r_k is the sampled random feature, and $\phi(\chi_S) := (\phi(\chi_S^\top r_1), \dots, \phi(\chi_S^\top r_K))^\top$. Since each random basis function $\phi(\chi_S^\top r_k)$ takes value either 0 or 1, the above combination of such functions will qualify as a probability in $[0, 1]$. We will denote the class of functions defined by equations (9) and (10) as $\hat{\mathcal{F}}^w$.

How well can the random basis function $f^w(\chi_S)$ approximate the original function $f_j(\chi_S)$? We can show that there exists some w such that $f^w(\chi_S)$ approximates $f_j(\chi_S)$ well when K is sufficiently large. More specifically, let C be the minimum value such that

$$p_j(r) \leq C q_j(r), \quad \forall j \in [d], \quad \forall r \in \{0, 1\}^n.$$

Intuitively, C measures how far away the sampling distribution $q_j(r)$ is from the true distribution $p_j(r)$.

Lemma 1. *Let $p_\chi(\chi_S)$ be a distribution of χ_S . If $K = O(\frac{C^2}{\epsilon^2} \log \frac{C}{\epsilon\delta})$ and r_1, \dots, r_K are drawn i.i.d. from $q_j(r)$, then with probability at least $1 - \delta$, there exists an $f^w \in \hat{\mathcal{F}}^w$ such that $\mathbb{E}_{\chi_S \sim p_\chi} [(f_j(\chi_S) - f^w(\chi_S))^2] \leq \epsilon^2$.*

Alternatively, the lemma can also be interpreted as the approximation error ϵ scales as $O(\frac{C}{\sqrt{K}})$. Note that we require that w lie in a simplex, i.e., $w_k \geq 0$ and $\sum_{k=1}^K w_k = 1$, and it is slightly different from that in Rahimi & Recht (2008).

5. Efficient Learning Algorithm

After generating the random features, we can learn the weights $w = (w_1, \dots, w_K)$ by fitting $f^w(\chi_S)$ to training data. Since the target function $f_j(\chi_S)$ is a conditional probability, l_2 or l_1 error metric may not be suitable loss functions to optimize. A natural approach is maximum conditional likelihood estimation. We use an efficient exponentiated gradient algorithm for performing the estimation for the weights in f^w . Here, we describe the algorithm, and then present the sample complexity analysis in the next section.

Suppose we observe a dataset of m i.i.d. cascades

$$\mathcal{D}^m := \{(\mathcal{S}_1, \mathcal{I}_1), \dots, (\mathcal{S}_m, \mathcal{I}_m)\}, \quad (11)$$

where each cascade is a pair of observation of the source set \mathcal{S}_i and the corresponding set \mathcal{I}_i of influenced nodes. Each cascade $(\mathcal{S}_i, \mathcal{I}_i)$ in the dataset is obtained by first sampling a source set \mathcal{S}_i from a distribution $p_\chi(\chi_S)$ (e.g., power law), then sampling a random reachability matrix \mathbf{R} from $p_{\mathbf{R}}$, and finally calculating $\mathcal{I}_i := \{j : \phi(\chi_S^\top \mathbf{R}_{:,j}) = 1\}$. We note that \mathbf{R} is an intermediate quantity which is not observed in the dataset. In our setting, we let $\mathcal{S}_i \subseteq \mathcal{I}_i$ which means the nodes in the source set are also considered as influenced nodes.

For a particular cascade $(\mathcal{S}_i, \mathcal{I}_i)$ and a particular target node j , we can define a binary variable indicating whether the target node j is influenced in this cascade, $y_{ij} := \mathbb{I}\{j \in \mathcal{I}_i\}$. Then the conditional likelihood of the status of node j (influenced or not) can be expressed using $f_j(\chi_S)$

$$f_j(\chi_{\mathcal{S}_i})^{y_{ij}} (1 - f_j(\chi_{\mathcal{S}_i}))^{1-y_{ij}}. \quad (12)$$

So in the following, we will focus on learning individual function f^w which is an approximation of $f_j(\chi_S)$.

5.1. Maximum conditional likelihood estimation

In a way very similar to logistic regression and conditional random fields by Lafferty et al. (2001), we will maximize the conditional log-likelihood the y_{ij} given the $\chi_{\mathcal{S}_i}$. In contrast to logistic regression and conditional random fields where the models usually take the exponential family form, we will use a form of a convex combination of random basis function (f^w). The additional challenge for this parameterization is that the conditional probability may be zero for some \mathcal{S} . To address this challenge, we will use a truncated or Winsorized version of the function f^w

$$f^{w,\lambda}(\chi_S) = (1 - 2\lambda)f^w(\chi_S) + \lambda \quad (13)$$

which squashes the function output to the range of $[\lambda, 1 - \lambda]$. We will denote this new class of functions as $\hat{\mathcal{F}}^{w,\lambda}$. Although this transformation introduces additional bias to the function class, we show in later analysis that it is fine if we choose λ to be about the same level as the approximation error. In practice, λ is selected via cross-validation.

Then the log-likelihood of the data \mathcal{D}^m can be written as

$$\ell_j(w) := \sum_{i=1}^m y_{ij} \log f^{w,\lambda}(\chi_{S_i}) + (1 - y_{ij}) \log(1 - f^{w,\lambda}(\chi_{S_i})), \quad (14)$$

and we can find w by maximizing the log-likelihood

$$\hat{w} := \underset{w}{\operatorname{argmax}} \ell_j(w) \quad (15)$$

subject to $\sum_{k=1}^K w_k = 1, w_k \geq 0.$

One can easily show that the optimization problem in (15) is a convex optimization problem over a probability simplex. Hence we can leverage existing techniques from convex optimization by [Kivinen & Warmuth \(1997\)](#) and [Schmidt et al. \(2009\)](#) to find \hat{w} efficiently.

5.2. Exponentiated gradient algorithm

We describe a simple exponentiated gradient (EG) algorithm, originally introduced by [Kivinen & Warmuth \(1997\)](#) in the online learning context. The EG updates involve the following simple multiplicative modification

$$w_k^{t+1} = \frac{1}{Z^t} w_k^t \exp(-\eta \nabla_k(w^t)) \quad (16)$$

where $Z^t = \sum_{k=1}^K w_k^t \exp(-\eta \nabla_k(w^t))$ is the normalization constant, the parameter $\eta > 0$ is the learning rate, and the gradient $\nabla(w^t)$ is given by

$$\nabla(w) = (1 - 2\lambda) \sum_{i=1}^m \left(\frac{1 - y_{ij}}{1 - \lambda - (1 - 2\lambda)w^\top \phi(\chi_{S_i})} - \frac{y_{ij}}{\lambda + (1 - 2\lambda)w^\top \phi(\chi_{S_i})} \right) \phi(\chi_{S_i}) \quad (17)$$

Algorithm 1 summarizes algorithm for learning the influence function. We first generate K random features $\{r_1, \dots, r_K\}$ from the given distribution $q_j(r)$. Then, we precompute m feature vectors $\phi(\chi_{S_i}) = (\phi(\chi_{S_i}^\top r_1), \dots, \phi(\chi_{S_i}^\top r_K))^\top$. Because χ_{S_i} is usually very sparse, this preprocessing costs $O(K \sum_{i=1}^m |\mathcal{S}_i|)$, where $|\mathcal{S}_i|$ is the cardinality of the set \mathcal{S}_i . Then we use the exponentiated gradient algorithm to find the weight w that maximizes the log-likelihood of the training data. According to [Kivinen & Warmuth \(1997\)](#), to get within ϵ of the optimum, we need $O(\frac{1}{\epsilon\eta})$ iterations, where the main work of each iteration is evaluating the gradient with complexity $O(dmK)$. The final estimate $\hat{\sigma}(\mathcal{S})$ is the sum of all the functions learned for each node. The learning task for each node is independent of those for the other nodes (except that we use the same set of training data), so the algorithm can be easily parallelized. We refer to our algorithm as INFLUARNER.

5.3. How to choose random basis function

By our analysis in Lemma 1, the number of random features needed for node j depends on the sampling distribu-

Algorithm 1 INFLUARNER

input training data $\{(\mathcal{S}_i, \mathcal{I}_i)\}_{i=1}^m$, $\lambda \in (0, \frac{1}{4})$
for each node $j \in [d]$ **do**
 sample K random features $\{r_1, \dots, r_K\}$ from $q_j(r)$;
 compute $\phi(\chi_{S_i}) = (\phi(\chi_{S_i}^\top r_1), \dots, \phi(\chi_{S_i}^\top r_K))$, $\forall i$;
 initialize w^1 to a interior point of a K -simplex;
for $t = 1, \dots, T$ **do**
 calculate $\nabla(w^t)$ using (17)
 update $w^{t+1} \propto w^t \exp(-\eta \nabla(w^t))$ using (16)
end for
 $\hat{f}_j^{w,\lambda}(\chi_S) = \lambda + (1 - 2\lambda)(w^T)^\top \phi(\chi_S)$
end for
output $\hat{\sigma}(\mathcal{S}) = \sum_{j=1}^d \hat{f}_j^{w,\lambda}(\chi_S).$

tion $q_j(r)$. More precisely, it has quadratic dependence on C where $p_j(r) \leq Cq_j(r)$ for all r . If we know $p_j(r)$, then by sampling random features from $p_j(r)$, we have $C = 1$ so that much fewer features are needed. However, in practice, $p_j(r)$ is often unknown, so we consider estimating $p_j(r)$ by $q_j(r)$ using the following simple approach.

Inspired by the empirical success of Naïve Bayes algorithm in classification by [Bishop \(2006\)](#) and the mean field approximation in graphical model inference ([Wainwright & Jordan, 2003](#)), we assume that $q_j(r)$ is fully factorized, *i.e.*,

$$q_j(r) = \prod_{s=1}^d q_j(r(s)).$$

where $q_j(r(s))$ means the marginal distribution of the i -th dimension of r . Given a training dataset \mathcal{D}^m as in equation (11), we estimate each $q_j(r(s))$ using the frequency of node j being influenced by source node i , *i.e.*, $q_j(r(s)) = \frac{1}{|\mathcal{D}_s^m|} \sum_{i \in \mathcal{D}_s^m} y_{ij}$ where $\mathcal{D}_s^m := \{i : s \in \mathcal{S}_i\}$. Although this $q_j(r)$ may be quite different from $p_j(r)$, by the additional steps of drawing random features and adjusting the corresponding weights, it leads to very good results, as illustrated in our experiments.

A more intelligent approach for choosing $q_j(r)$ may be first learning a diffusion model outlined in Section 2 and then using samples from the diffusion model to generate the random basis functions. This approach requires more computation and is left for future study.

6. Sample Complexity of MLE

Here we analyze Algorithm 1 and provide sample complexity bounds for the number of random basis functions and the size of the training data needed to get a solution close to the truth. We describe our results here and provide the proof in the appendix.

We note that existing analysis for random kitchen sink ([Rahimi & Recht, 2008](#)) does not apply to the maximum likelihood estimation. Therefore, we use a general framework by [Birgé & Massart \(1998\)](#) for maximum like-

likelihood estimation. Loosely speaking, the error of the maximum likelihood estimator $\hat{f}_j^{w,\lambda}(\chi_S) \in \hat{\mathcal{F}}^{w,\lambda}$ is bounded by the best possible in the hypothesis class plus a term scale roughly as $\tilde{O}(D/m)$, where D is the dimension of the set of candidate models based on a covering approach. Hence, to get sample complexity bounds for our problem, we need to bound the dimension of $\hat{\mathcal{F}}^{w,\lambda}$. We consider the mapping from the weight w to the corresponding hypothesis $f \in \hat{\mathcal{F}}^{w,\lambda}$, and show that the distance between two functions f and f' are approximately the distance between their corresponding weights w and w' . Then a covering on the space of w induces a covering on the function space $\hat{\mathcal{F}}^{w,\lambda}$, and thus the dimensions of the two spaces are approximately the same, which is $O(K)$. Combined the dimension bound with Lemma 1, we arrive at the following:

Lemma 2. *Assume the statement in Lemma 1 is true. If $m = \tilde{O}(\frac{K}{\epsilon})$, then the maximum likelihood estimator $\hat{f}_j^{w,\lambda} \in \hat{\mathcal{F}}^{w,\lambda}$ satisfies*

$$\mathbb{E}_{\mathcal{D}^m} \mathbb{E}_{p_\chi} \left[(\hat{f}_j^{w,\lambda}(\chi_S) - f_j(\chi_S))^2 \right] \leq \tilde{O} \left(\frac{\epsilon^2 + \lambda^2}{\lambda} \right).$$

This means to get ϵ accuracy, it suffices to choose $\lambda = \epsilon$ and choose K large enough to make sure that the l_2 error between the true function and the set of candidate functions in $\hat{\mathcal{F}}^{w,\lambda}$ is at most ϵ^2 . The bound then follows by applying the above argument on each node with accuracy $O(\epsilon/d)$.

Theorem 3. *Suppose in Algorithm 1, we set $\lambda = \tilde{O}(\frac{\epsilon}{d})$, $K = \tilde{O}(\frac{C^2 d^2}{\epsilon^2})$, and $m = \tilde{O}(\frac{C^2 d^3}{\epsilon^3})$. Then with probability at least $1 - \delta$ over the drawing of the random features, the output of Algorithm 1 satisfies*

$$\mathbb{E}_{\mathcal{D}^m} \mathbb{E}_{p_\chi} \left[\left(\sum_{j=1}^d \hat{f}_j^{w,\lambda}(\chi_S) - \sigma(\mathcal{S}) \right)^2 \right] \leq \epsilon.$$

Intuitively, the l_2 error of the function $\sum_{j=1}^d \hat{f}_j^{w,\lambda}$ learned is small if the number K of random features and the size m of the training data are sufficiently large. Both quantities have a quadratic dependence on C , since if C is large, then the difference between p_j and q_j could be large, and thus we need more random features to approximate f_j and also more training data to learn the weights. K and m also depend on the number d of nodes in the network, for the reason that we need to estimate each f_j up to accuracy ϵ/d so that their sum is estimated to accuracy ϵ . This is far too pessimistic, as we observe in our experiment that much smaller K or m is needed.

7. Experiments

We evaluate INFLUARNER in synthetic and real world data. We compare it to the state-of-the-art two-stage approaches, as well as methods based on linear regression and logistic regression, and show that INFLUARNER is more robust to model misspecification than these alternatives.

7.1. Competitors

Two-stage methods. Two-stage learning methods depend on the diffusion model assumptions, families of pairwise temporal dynamics, and whether network structures are given or not. We design the following four representative competitors :

1. Continuous-time Independent Cascade model with exponential pairwise transmission function (CIC).
2. Continuous-time Independent Cascade model with exponential pairwise transmission function and given network Structure (CIC-S).
3. Discrete-time Independent Cascade model (DIC).
4. Discrete-time Independent Cascade model with given network Structure (DIC-S).

For the methods CIC and CIC-S, we use NETRATE (Gomez Rodriguez et al., 2011) to learn the structure and parameters of the pairwise transmission functions. For DIC and DIC-S, we learn the pairwise infection probability based on the method of (Netrapalli & Sanghavi, 2012).

Approach based on logistic regression. Instead of using random features, we represent $f_j(\chi_S)$ using a modified logistic regression

$$f_j(\chi_S) = \frac{2 \exp(w^\top \chi_S)}{1 + \exp(w^\top \chi_S)} - 1, \text{ where } w \geq 0. \quad (18)$$

Since the sigmoid function is concave in \mathbb{R}_+ and $w^\top \chi_S$ is a linear function of χ_S , the representation in (18) is also a submodular function of the set \mathcal{S} . We learn w by maximizing the log-likelihood subject to the nonnegative constraint. We also experimented with the original logistic regression model which does not lead to a submodular function, and thus does not perform as well as the representation in (18) (and hence not reported).

Approach based on linear regression. We use the linear regression model, $w^\top \chi_S + b$, to directly regress from χ_S to the cascade size $|\mathcal{Z}|$. This approach does use the knowledge that the influence function is a coverage function.

7.2. Synthetic Data

We generate Kronecker type of synthetic networks with the parameter matrix $[0.9 \ 0.5; 0.5 \ 0.3]$, which mimics the information diffusion traces in real world networks (Leskovec et al., 2010). The generated networks consist of 1,024 nodes and 2,048 edges. Given a generated network structure, we apply the continuous-time independent cascade, the discrete-time independent cascades and the linear-threshold model to generate the cascades, respectively.

For the continuous-time diffusion model, we used both Weibull distribution (Wbl) and exponential distribution (Exp) for the pairwise transmission function, and set their parameters at random to capture the heterogeneous temporal dynamics. For the Weibull distribution, $f(t; \alpha, \beta) =$

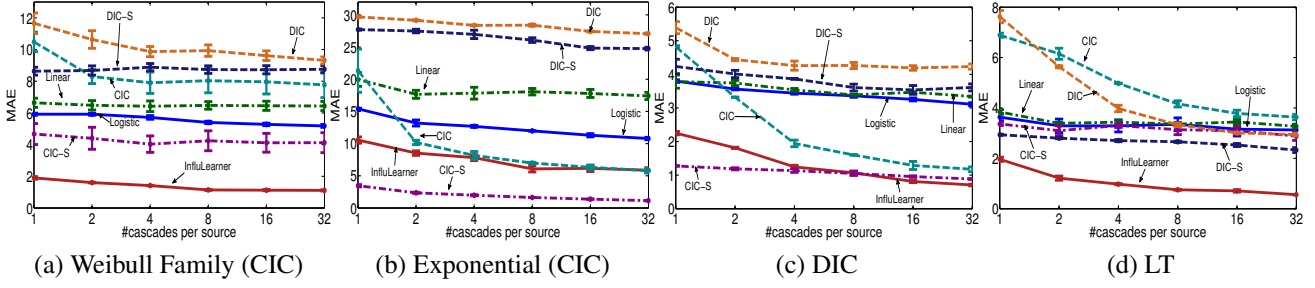


Figure 1. Over the generated synthetic networks with 1,024 nodes and 2,048 edges, we present the mean absolute error of the estimated influence on the testing data by increasing the number of training data when the true diffusion model is (a) continuous-time independent cascade with pairwise Weibull transmission functions, (b) continuous-time independent cascade with pairwise exponential transmission functions, (c) discrete-time independent cascade model and (d) linear-threshold cascade model.

$\frac{\beta}{\alpha} \left(\frac{t}{\alpha}\right)^{\beta-1} e^{-(t/\alpha)^\beta}$, $t \geq 0$, where $\alpha > 0$ is a scale parameter and $\beta > 0$ is a shape parameter. We choose α and β from 1 to 10 uniformly at random for each edge in order to have heterogeneous temporal dynamics. The true influence value range is from 1 to 235, and the average value is 15.78 with the time window $T = 10$. For the exponential distribution, the average influence is 37.81.

For the discrete-time independent cascade model, the pairwise infection probability is chosen uniformly from 0 to 1. For the discrete-time linear-threshold model, we followed Kempe et al. (2003) where the edge weight w_{uv} between u and v is $1/d_v$, and d_v is the degree of node v . We run these generative models for 10 time steps. The average influence values are 9.2 and 8.9 respectively.

The source locations are sampled uniformly without replacement from \mathcal{V} , and the source set sizes conform to a power law distribution with parameter 2.5. For the training set, we independently sample 1,024 source sets, and independently generate 8 to 128 cascades for each source set. The test set contains 128 independently sampled source sets with the ground truth influence estimated from 10,000 simulated cascades.

7.3. Robustness to model misspecification

The cascades used in Figure 1(a) are generated from the continuous-time independent cascade model with pairwise Weibull transmission functions. We expect that the four two-stage methods are not doing well due to model misspecification of one form or the other. Figure 1(a) shows the MAE (Mean Absolute Error) between the estimated value and the true value. Both CIC-S and CIC used the correct continuous-time diffusion model but the wrong family of pairwise transmission functions, so their performance lies in the middle. However, CIC-S has the prior knowledge about the true network structure, so it is reduced to a much simpler learning problem and is thus better than CIC. DIC-S and DIC used the wrong diffusion model with unit time step (which is hard to determine in practice), so they have the lowest performance overall. In contrast, INFLUERNER does not explicitly make assumptions about

diffusion models or transmission functions but only learns the influence function directly from the data. Thus, it is much more robust and better than the two-stage methods. Since INFLUERNER has better representational power than the logistic regression based approach, it is able to better approximate the true influence function and thus can achieve the best performance overall.

The cascades used in Figure 1(b) are generated from the continuous-time independent cascade model with pairwise exponential transmission functions. Note that in this case we expect CIC-S and CIC to do well, since they have the correct assumptions about both the diffusion model and the family of transmission functions. Particularly, with the prior knowledge of the true network structure, CIC-S simply fits the model parameters for each edge, and thus the estimates converge to the true influence function quickly. Still, we see that the performance of INFLUERNER is close to that of CIC-S and CIC. Figure 1(b) again shows that INFLUERNER is robust to diffusion model changes.

In Figure 1(c, d), we generate cascades according to discrete-time independent cascade model and linear threshold model respectively. In Figure 1(c), DIC-S and DIC assumes the correct model, so their performance improves a lot. However, in Figure 1(d), because CIC-S, CIC, DIC-S, and DIC all assume the wrong diffusion model, we observe a similar trend as in Figure 1(a): INFLUERNER is robust and obtain the best results. Note that in this case, the gap between different methods is not as big since the average influence value is small.

7.4. Scalability

Figure 2(a) reports the parallel runtime of INFLUERNER as we increase the number of training cascades per source set. We arbitrarily divide the 1,024 independent learning problems into 32 individual jobs running on a cluster of 32 cores (AMD Opteron(tm) Processor, 2.5GHz). It shows that the runtime grows almost linearly as the number of cascades increases.

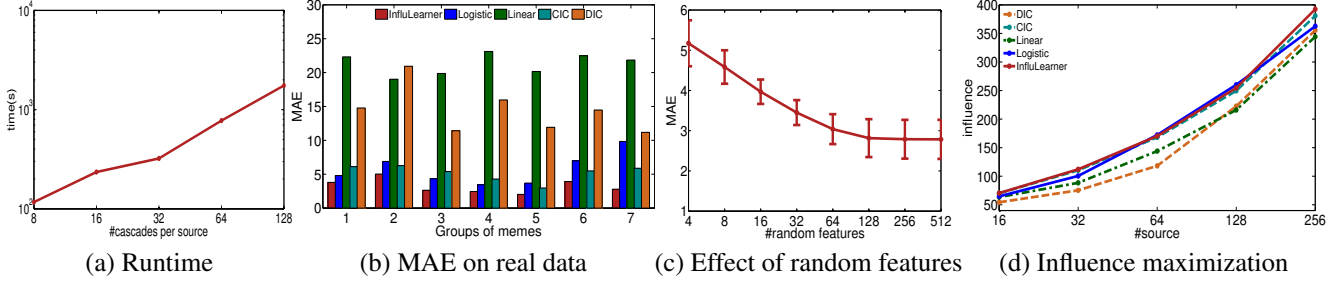


Figure 2. (a) Runtime in log-log scale; (b) MAE on seven sets of real cascade data; (c) The performance gain of using different number of random features; (d) Maximized expected influence of different selected sources on the real hold-out testing data.

7.5. Influence estimation on real data

We further evaluate the performance of our proposed method on the MemeTracker dataset which includes 300 million blog posts and articles collected from 5,000 active media sites between March 2011 and February 2012 (Leskovec et al., 2009). The flow of information was traced using quotes which are short textual phrases spreading through the websites. Because all published documents containing a particular quote are time-stamped, a cascade induced by the same quote like ‘apple and jobs’ is a collection of times when the media site first mentioned it. We have selected seven groups of cascades with the typical keywords like ‘apple and jobs’, ‘tsunami earthquake’, ‘william kate marriage’, ‘occupy wall-street’, ‘airstrikes’, ‘egypt’ and ‘elections’. We split each set of cascades into 60%-train and 40%-test. Because we do not have any prior knowledge about either the diffusion structure or the underlying diffusion mechanism on the real cascades data, we only compare INFLUARNER with the Logistic regression, Linear regression, CIC and DIC.

We evaluate the performance on the held-out testing cascades as follows : we randomly select 10 sources from the testing cascades, which represents one particular source set S . For each node $u \in S$, let $\mathcal{C}(u)$ denote the set of cascades generated from u on the testing data. For each $u \in S$, we uniformly sample one cascade from $\mathcal{C}(u)$. Thus, the union of all sampled cascades is the set of nodes infected by source set S . We repeat the process for 1,000 times and take the average of the number of infected nodes as the true influence of source set S . Finally, we have generated 100 source sets and report the MAE of each method in Figure 2(b). We can see that the performance of INFLUARNER is robust and consistent across all groups of testing cascades, and is significantly better than the other competitors.

Moreover, Figure 2(c) demonstrates the effect of the number of random features on the performance of INFLUARNER by showing the average MAE over the seven sets of cascade data as the number of random features increases. As the number of random features grows, INFLUARNER approximates the true influence better, and

thus the MAE decreases. It seems that 128 to 256 random features are sufficient to achieve good performance overall.

7.6. Influence maximization on real data

Finally, we use the learned influence function (from INFLUARNER, Logistic, Linear, CIC and DIC) for solving the influence maximization problem Kempe et al. (2003); Du et al. (2013b). Here we want to find a set \mathcal{S}^* of C source nodes which maximizes the influence, i.e., $\mathcal{S}^* = \operatorname{argmax}_{|\mathcal{S}| \leq C} \sigma(\mathcal{S})$. We will use a greedy algorithm framework of Nemhauser et al. (1978) to solve the problem. We use the held-out test cascade to estimate the influence achieved by selected source nodes. The observation time window used is $T = 14$.

Figure 2(d) shows the influence achieved in Meme group 1 (the rest of the testing groups has similar results as in the Appendix). INFLUARNER, Logistic and CIC perform consistently better than DIC and linear regression. The source nodes selected by INFLUARNER, Logistic and CIC are very similar, though the estimated influence value can be different. As a result, the influence value of INFLUARNER, Logistic and CIC are very close.

8. Conclusion

Based on the observation that the influence function in many diffusion models are coverage functions, we propose to directly learn the influence from cascade data. In this paper, we provide a novel parameterization of the influence function as a convex combination of random basis functions, and an efficient maximum likelihood based algorithm for learning the weighting of the random basis functions. Theoretically, we show that the algorithm can learn the influence with low sample complexity, and our empirical study also shows our method outperforms traditional two-stage approaches.

Acknowledgement

The research was supported in part by NSF/NIH BIGDATA 1R01GM108341, NSF IIS-1116886, NSF CAREER IIS-1350983 to L.S.; NSF CCF-1101283, NSF CAREER CCF-0953192, AFOSR FA9550-09-1-0538, ONR N00014-09-1-0751, and a Microsoft Research Faculty Fellowship to M.B.; a Raytheon Faculty Fellowship to M.B. and L.S.; and a Facebook Fellowship to N.D.

References

- Badanidiyuru, A., Dobzinski, S., Fu, H., Kleinberg, R. D., Nisan, N., and Roughgarden, T. Sketching valuation functions. In *Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms*, 2012.
- Balcan, Maria-Florina and Harvey, Nicholas JA. Learning submodular functions. In *Proceedings of the 43rd annual ACM symposium on Theory of computing*, pp. 793–802. ACM, 2011.
- Birgé, L. and Massart, P. Minimum Contrast Estimators on Sieves: Exponential Bounds and Rates of Convergence. *Bernoulli*, 4(3), 1998.
- Bishop, Christopher. *Pattern Recognition and Machine Learning*. Springer, 2006.
- Borgs, Christian, Brautbar, Michael, Chayes, Jennifer, and Lucier, Brendan. Influence maximization in social networks: Towards an optimal algorithmic solution. *arXiv preprint arXiv:1212.0884*, 2012.
- Chen, Wei, Wang, Chi, and Wang, Yajun. Scalable influence maximization for prevalent viral marketing in large-scale social networks. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 1029–1038. ACM, 2010.
- Du, N., Song, L., Smola, A., and Yuan, M. Learning networks of heterogeneous influence. In *Advances in Neural Information Processing Systems 25*, pp. 2789–2797, 2012.
- Du, N., Song, L., Woo, H., and Zha, H. Uncover topic-sensitive information diffusion networks. In *Artificial Intelligence and Statistics (AISTATS)*, 2013a.
- Du, Nan, Song, Le, Rodriguez, Manuel Gomez, and Zha, Hongyuan. Scalable influence estimation in continuous-time diffusion networks. In *Advances in Neural Information Processing Systems 26*, 2013b.
- Feldman, Vitaly and Kothari, Pravesh. Learning coverage functions. *arXiv preprint arXiv:1304.2079*, 2013.
- Feldman, Vitaly and Vondrak, Jan. Optimal bounds on approximation of submodular and xos functions by juntas. In *FOCS*, 2013.
- Gomez Rodriguez, Manuel, Balduzzi, David, and Schölkopf, Bernhard. Uncovering the temporal dynamics of diffusion networks. *arXiv preprint arXiv:1105.0697*, 2011.
- Kempe, David, Kleinberg, Jon, and Tardos, Éva. Maximizing the spread of influence through a social network. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 137–146. ACM, 2003.
- Kivinen, J. and Warmuth, M. K. Exponentiated gradient versus gradient descent for linear predictors. *Information and Computation*, 132(1):1–64, January 1997.
- Lafferty, J. D., McCallum, A., and Pereira, F. Conditional random fields: Probabilistic modeling for segmenting and labeling sequence data. In *Proceedings of International Conference on Machine Learning*, volume 18, pp. 282–289, San Francisco, CA, 2001. Morgan Kaufmann.
- Leskovec, Jure, Backstrom, Lars, and Kleinberg, Jon. Meme-tracking and the dynamics of the news cycle. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 497–506. ACM, 2009.
- Leskovec, Jure, Chakrabarti, Deepayan, Kleinberg, Jon, Faloutsos, Christos, and Ghahramani, Zoubin. Kronecker graphs: An approach to modeling networks. *Journal of Machine Learning Research*, 11(Feb):985–1042, 2010.
- Nemhauser, G., Wolsey, L., and Fisher, M. An analysis of the approximations for maximizing submodular set functions. *Mathematical Programming*, 14:265–294, 1978.
- Netrapalli, Praneeth and Sanghavi, Sujay. Learning the graph of epidemic cascades. In *SIGMETRICS/PERFORMANCE*, pp. 211–222. ACM, 2012. ISBN 978-1-4503-1097-0.
- Rahimi, Ali and Recht, Benjamin. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. In *Advances in neural information processing systems*, pp. 1313–1320, 2008.
- Rodriguez, M.G. and Schölkopf, B. Influence maximization in continuous time diffusion networks. In *Proceedings of the International Conference on Machine Learning*, 2012.
- Schmidt, M., van den Berg, E., Friedlander, M. P., and Murphy, K. Optimizing costly functions with simple constraints: A limited-memory projected quasi-newton algorithm. In van Dyk, D. and Welling, M. (eds.), *Proceedings of The Twelfth International Conference on Artificial Intelligence and Statistics (AISTATS) 2009*, volume 5, pp. 456–463, Clearwater Beach, Florida, April 2009.
- Wainwright, M. J. and Jordan, M. I. Graphical models, exponential families, and variational inference. Technical Report 649, UC Berkeley, Department of Statistics, September 2003.

A. Proofs for Structure of the Influence Function

To prove that the influence function $\sigma(\mathcal{S})$ is a coverage function, the key is that non-negative combinations of coverage functions are still coverage functions. We state and prove the property for the case of combining two coverage functions, while for the general case we can simply repeat the argument.

Lemma 4. Suppose $c^{(1)}$ and $c^{(2)}$ are two coverage functions mapping from $2^{\mathcal{V}}$ to \mathbb{R}_+ . If $\alpha^{(1)} \geq 0$ and $\alpha^{(2)} \geq 0$, then $\sum_{\ell=1}^2 \alpha^{(\ell)} c^{(\ell)}$ is also a coverage function mapping from $2^{\mathcal{V}}$ to \mathbb{R}_+ .

Proof. By definition, for $\ell = 1, 2$, there exists a universe $\mathcal{U}^{(\ell)}$, a set of weights $\{a_u^{(\ell)}\}_{u \in \mathcal{U}^{(\ell)}}$, and a family of subsets $\{\mathcal{A}_v^{(\ell)} : \mathcal{A}_v^{(\ell)} \subseteq \mathcal{U}^{(\ell)}\}_{v \in \mathcal{V}}$ such that for any $\mathcal{S} \subseteq \mathcal{V}$,

$$c^{(\ell)}(\mathcal{S}) = \sum_{u \in \bigcup_{v \in \mathcal{S}} \mathcal{A}_v^{(\ell)}} a_u^{(\ell)}.$$

Define a new universe $\mathcal{U} = \bigcup_{\ell=1}^2 \mathcal{U}^{(\ell)}$, where elements in $\mathcal{U}^{(\ell)}$ ($\ell = 1, 2$) are treated as different elements. Define the corresponding weights a_u for $u \in \mathcal{U}$ as follows: if $u \in \mathcal{U}^{(\ell)}$, then $a_u = \alpha^{(\ell)} a_u^{(\ell)}$. Define a family of subsets $\{\mathcal{A}_v : \mathcal{A}_v \subseteq \mathcal{U}\}_{v \in \mathcal{V}}$ where $\mathcal{A}_v = \bigcup_{\ell=1}^2 \mathcal{A}_v^{(\ell)}$. Then the corresponding coverage function is

$$\begin{aligned} c(\mathcal{S}) &= \sum_{u \in \bigcup_{v \in \mathcal{S}} \mathcal{A}_v} a_u = \sum_{u \in \bigcup_{\ell=1}^2 \bigcup_{v \in \mathcal{S}} \mathcal{A}_v^{(\ell)}} a_u = \sum_{\ell=1}^2 \sum_{u \in \bigcup_{v \in \mathcal{S}} \mathcal{A}_v^{(\ell)}} a_u = \sum_{\ell=1}^2 \sum_{u \in \bigcup_{v \in \mathcal{S}} \mathcal{A}_v^{(\ell)}} \alpha^{(\ell)} a_u^{(\ell)} \\ &= \sum_{\ell=1}^2 \alpha^{(\ell)} \sum_{u \in \bigcup_{v \in \mathcal{S}} \mathcal{A}_v^{(\ell)}} a_u^{(\ell)} = \sum_{\ell=1}^2 \alpha^{(\ell)} c^{(\ell)}(\mathcal{S}). \end{aligned}$$

Therefore, $c(\mathcal{S}) = \sum_{\ell=1}^2 \alpha^{(\ell)} c^{(\ell)}$ is a coverage function. \square

Since $\Phi(\mathcal{S}|\mathbf{R})$ is a coverage function for any fixed \mathbf{R} , and the influence

$$\sigma(\mathcal{S}) = \mathbb{E}_{\mathbf{R} \sim p_{\mathbf{R}}} [\Phi(\mathcal{S}|\mathbf{R})]$$

is a convex combination of $\Phi(\mathcal{S}|\mathbf{R})$, we have the following corollary.

Corollary 5. The influence function $\sigma(\mathcal{S})$ is a coverage function.

Note If we naively construct the universe for the influence function as in the proof of Lemma 4, this will lead to a universe of size 2^d , which is exponential in d . It seems to imply that the function is difficult to learn. However, as shown in (Badanidiyuru et al., 2012), there exists a coverage function that is a $(1 + \epsilon)$ multiplicative approximation to σ , and is defined on a universe of size $O\left(\frac{d^2}{\epsilon^2}\right)$. This suggests that there are structures in a coverage function that make learning tractable, even if it is defined on an exponentially large universe. On the other hand, the proof in (Badanidiyuru et al., 2012) does not immediately lead to an efficient learning algorithm, since the construction explicitly makes use of the weights of the elements in the universe defining σ .

B. Proofs for Random Basis Function Approximation

In this section, we fix a node j and $f_j(\chi_{\mathcal{S}}) = \mathbb{E}_{r \sim p_j(r)} [\phi(\chi_{\mathcal{S}}^\top r)]$. Suppose a set of K random features $\{r_{j1}, \dots, r_{jK}\}$ is drawn from the distribution $q_j(r)$ over $\{0, 1\}^n$. We show that given sufficiently many random features, there exists a convex combination of the random basis functions that approximates the truth f_j .

The number of random features needed depends on how close the sample distribution q_j is to the true distribution p_j . The “distance” between the two is formalized in the following definition.

Definition 6. Let C be the minimum value such that

$$p_j(r) \leq C q_j(r) \text{ for all } j \in [d], r \in \{0, 1\}^n.$$

We first introduce an intermediate class $\tilde{\mathcal{F}}^w$ that depends on C , and show that there exists a function in $\tilde{\mathcal{F}}^w$ that is close to f_j . We then utilize the structure of our problem to show that the same is true for a class $\hat{\mathcal{F}}^w$ that does not depend on C . In

particular, define

$$\tilde{\mathcal{F}}^w := \left\{ f^w(\chi_S) = \sum_{k=1}^K w_k \phi(\chi_S^\top r_{jk}) \mid 0 \leq w_k \leq \frac{C}{K} \right\}, \quad (19)$$

$$\hat{\mathcal{F}}^w := \left\{ f^w(\chi_S) = \sum_{k=1}^K w_k \phi(\chi_S^\top r_{jk}) \mid w_k \geq 0, \sum_{k=1}^K w_k \leq 1 \right\}. \quad (20)$$

Lemma 7. Let p_χ be any distribution of χ_S . If r_{j1}, \dots, r_{jK} are drawn i.i.d. from $q_j(r)$, then with probability at least $1 - \delta$ over r_{j1}, \dots, r_{jK} , there exists $\tilde{f} \in \tilde{\mathcal{F}}^w$ such that

$$\Pr_{\chi_S \sim p_\chi} \left[\left| \tilde{f}(\chi_S) - f_j(\chi_S) \right| \geq \epsilon \right] \leq \epsilon^2 / C$$

when $K = O(\frac{C^2}{\epsilon^2} \log \frac{C}{\epsilon \delta})$. Consequently,

$$\mathbb{E}_{\chi_S \sim p_\chi} \left[(f_j(\chi_S) - \tilde{f}(\chi_S))^2 \right] \leq 3\epsilon^2.$$

Proof. Here we prove the first statement, which is stronger and implies the second one. Let $f^k(\chi_S) = \frac{p(r_{jk})}{q(r_{jk})} \phi(\chi_S^\top r_{jk})$ for $k = 1, \dots, K$. Then $\mathbb{E}_{r_{jk} \sim q_j(r)}[f^k] = f_j$. Let $\tilde{f}(\chi_S) = \frac{1}{K} \sum_{i=1}^K \frac{p(r_{jk})}{q(r_{jk})} \phi(\chi_S^\top r_{jk})$ be the sample average of these functions. Then $\tilde{f} \in \tilde{\mathcal{F}}^w$ since $0 \leq \frac{1}{K} \frac{p(r_{jk})}{q(r_{jk})} \leq \frac{C}{K}$.

By Hoeffding's inequality, when $K = O(\frac{C^2}{\epsilon^2} \log \frac{C}{\epsilon \delta})$, for any fixed S we have

$$\Pr_r \left[\left| \tilde{f}(\chi_S) - f_j(\chi_S) \right| \geq \epsilon \right] \leq \delta \epsilon^2 / C$$

where \Pr_r is over the random sample of r_{j1}, \dots, r_{jK} . This leads to

$$\Pr_{\chi_S \sim p_\chi} \Pr_r \left[\left| \tilde{f}(\chi_S) - f_j(\chi_S) \right| \geq \epsilon \right] \leq \delta \epsilon^2 / C.$$

Exchanging $\Pr_{\chi_S \sim p_\chi}$ and \Pr_r by Fubini's theorem, and then by Markov's inequality, we have

$$\Pr_r \left\{ \Pr_{\chi_S \sim p_\chi} \left[\left| \tilde{f}(\chi_S) - f_j(\chi_S) \right| \geq \epsilon \right] \geq \epsilon^2 / C \right\} \leq \delta.$$

This means with probability at least $1 - \delta$ over the random sample of r_{j1}, \dots, r_{jK} , on at least $1 - \epsilon^2 / C$ probability mass of the distribution of S , $[\tilde{f}(\chi_S) - f_j(\chi_S)]^2 \leq \epsilon^2$. Since $|\tilde{f}(\chi_S)| \leq C$ and $|f_j(\chi_S)| \leq 1$,

$$\mathbb{E}_{\chi_S \sim p_\chi} \left[(f_j(\chi_S) - \tilde{f}(\chi_S))^2 \right] \leq \epsilon^2(1 - \epsilon^2) + (C + 1)\epsilon^2 / C < 3\epsilon^2.$$

□

Note that for learning over $\tilde{\mathcal{F}}^w$, the parameter C needs to be determined. However, there are additional structures in our problem that can be utilized to further restrict $\tilde{\mathcal{F}}^w$ and get rid of the dependence on C .

Lemma 1. Let p_χ be any distribution of χ_S . If $K = O(\frac{C^2}{\epsilon^2} \log \frac{C}{\epsilon \delta})$ and r_{j1}, \dots, r_{jK} are drawn iid from $q_j(r)$, then with probability at least $1 - \delta$ over r_{j1}, \dots, r_{jK} , there exists $\hat{f} \in \hat{\mathcal{F}}^w$ such that

$$\mathbb{E}_{\chi_S \sim p_\chi} \left[(f_j(\chi_S) - \hat{f}(\chi_S))^2 \right] \leq \epsilon^2.$$

Proof. Construct a distribution Δ_1 that assigns probability 1 to $\chi_S = \mathbf{1}$ and probability 0 to all other source sets. Note that the definition of \tilde{f} is independent of the distribution of χ_S , so that we can apply Lemma 7 for \tilde{f} on both Δ_1 and p_χ .

Without loss of generality, assume $r_{jk} \neq \mathbf{0}$ for any k , since otherwise we can remove r_{jk} without changing \tilde{f} . Then $\mathbf{1}^\top r_{jk} > 0$ and thus $\tilde{f}(\mathbf{1}) = \sum_{k=1}^K w_k$. By Lemma 7 on Δ_1 , with probability $1 - \delta/2$ we have

$$\sqrt{\mathbb{E}_{\chi_S \sim \Delta_1} \left[(f_j(\chi_S) - \tilde{f}(\chi_S))^2 \right]} = |f_j(\mathbf{1}) - \tilde{f}(\mathbf{1})| = \left| f_j(\mathbf{1}) - \sum_{k=1}^K w_k \right| \leq \frac{\epsilon}{2}.$$

Then $\sum_{k=1}^K w_k \leq f_j(\mathbf{1}) + \frac{\epsilon}{2} \leq 1 + \frac{\epsilon}{2}$. Define $\hat{f} = \tilde{f}/(1 + \epsilon/2)$. Then $\hat{f} \in \hat{\mathcal{F}}^w$ and

$$|\hat{f}(\chi_S) - \tilde{f}(\chi_S)| = \frac{\epsilon/2}{1 + \epsilon/2} \tilde{f}(\chi_S) \leq \frac{\epsilon/2}{1 + \epsilon/2} \sum_{k=1}^K w_k \leq \frac{\epsilon}{2}.$$

By Lemma 7 on p_χ , with probability $1 - \delta/2$ we have

$$\mathbb{E}_{\chi_S \sim p_\chi} \left[(f_j(\chi_S) - \tilde{f}(\chi_S))^2 \right] \leq \frac{\epsilon^2}{4}.$$

Then we have

$$\mathbb{E}_{\chi_S \sim p_\chi} \left[(f_j(\chi_S) - \hat{f}(\chi_S))^2 \right] \leq 2\mathbb{E}_{\chi_S \sim p_\chi} \left[(f_j(\chi_S) - \tilde{f}(\chi_S))^2 + (\tilde{f}(\chi_S) - \hat{f}(\chi_S))^2 \right] \leq 2 \left(\frac{\epsilon}{2} \right)^2 + \frac{2\epsilon^2}{4} = \epsilon^2$$

which completes the proof. \square

C. Proofs for Sample Complexity

In this section, we provide the complete proof for the sample complexity of learning the weights of the random basis functions by maximum likelihood estimation (MLE). We are not aware of any previous work providing the analysis of MLE for the hypothesis class in our problem (the weighted sum of the random basis functions). Therefore, we adopt the general framework in (Birgé & Massart, 1998), which analyzes the sample complexity based on a particular dimension notion for the hypothesis class. Then we bound the dimension of our hypothesis class, which then leads to our sample bound. The techniques used in bounding the dimension can be extended to other hypothesis classes, and thus may be of independent interest.

In the following, we first review the framework and paraphrase their result for distributions over a discrete domain, since this suffices for our purpose. We then apply the result to learning the conditional probability f_j for an individual node j , and finally prove the bound for the entire influence function.

C.1. Review of MLE for probability estimation

The MLE estimator is defined as follows. Suppose we observe m data points Z_1, \dots, Z_m independent identically distributed according to the true probability function p^* over a discrete domain \mathcal{Z} . The hypothesis class \mathcal{H} is a set of functions, each of which is the square root¹ of a probability function. That is, for each $h \in \mathcal{H}$, $h = \sqrt{p_h}$ where p_h is a probability over \mathcal{Z} . The MLE estimator is $\hat{h} = \operatorname{argmax}_{h \in \mathcal{H}} \sum_{i=1}^m \log [h(Z_i)]$. More generally, an approximate MLE estimator is \hat{h} such that

$$\sum_{i=1}^m \log [\hat{h}(Z_i)] + 1 \geq \sup_{h \in \mathcal{H}} \sum_{i=1}^m \log [h(Z_i)]. \quad (21)$$

The goal is to analyze how the difference between \hat{h} and the truth $h^* = \sqrt{p^*}$ decreases with the sample size m .

Complexity of the hypothesis class To analyze the sample complexity, we need to introduce some metric over the hypotheses and some notion bounding the complexity of the hypothesis class based on the metric. Given h, \tilde{h} that are the square roots of two probabilities, the ℓ_2 distance is

$$d(h, \tilde{h}) := \|h - \tilde{h}\| = \sqrt{\sum_{Z \in \mathcal{Z}} [h(Z) - \tilde{h}(Z)]^2}. \quad (22)$$

Note that $d(h, \tilde{h})/\sqrt{2}$ is just the Hellinger distance. Similar to the ℓ_2 distance, we can define ℓ_∞ distance:

$$d_\infty(h, \tilde{h}) := \|h - \tilde{h}\|_\infty = \max_{Z \in \mathcal{Z}} |h(Z) - \tilde{h}(Z)|. \quad (23)$$

Both the ℓ_2 and ℓ_∞ distances are bounded over all square roots of probabilities, so a hypothesis class with such metrics is always a bounded metric space. To measure the complexity of such a metric space, a common notion is the following:

Definition 8. Given a set \mathcal{B} equipped with metric d , and a real number $\epsilon > 0$, $\mathcal{T} \subseteq \mathcal{B}$ is an ϵ -covering of \mathcal{B} if the following

¹We will always talk about the square root of the probabilities. This is because the ℓ_2 distance over such hypotheses correspond to the Hellinger distance, which plays a key role in the analysis of MLE and appears in the final bound.

holds: for every $h \in \mathcal{B}$ there exists $\tilde{h} \in \mathcal{T}$ such that $d(h, \tilde{h}) < \epsilon$.

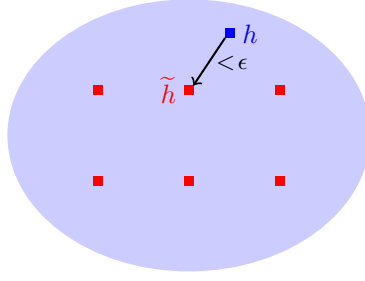


Figure 3. Illustration of ϵ -covering.

Intuitively, if we construct balls around points in \mathcal{T} with radius ϵ , then these balls can cover all points in \mathcal{B} . Note that the dimension depends on the metric d . The result in (Birgé & Massart, 1998) actually depends on both the ℓ_2 and ℓ_∞ metrics on \mathcal{H} . More precisely, we introduce the following $\ell_{2,\infty}$ dimension².

Definition 9 ((Birgé & Massart, 1998)). The $\ell_{2,\infty}$ dimension of \mathcal{H} is the minimum $D \geq 1$ such that there exist constants $c_0 \geq 1$ and $c_1 \geq 1$ satisfying the following. For each $\epsilon > 0$ and each ball $\mathcal{B} \subseteq \mathcal{H}$ with radius $R \geq 5\epsilon$, one can find \mathcal{T} with

$$|\mathcal{T}| \leq (c_0 R / \epsilon)^D$$

that is an ϵ -covering of \mathcal{B} for the ℓ_2 metric and a $c_1 \epsilon$ -covering for the ℓ_∞ metric.

The condition says that for any given distance threshold ϵ and any sufficiently large ball in \mathcal{H} , we can find a finite $O(\epsilon)$ -covering \mathcal{T} that is simultaneously with respect to both the ℓ_2 metric and the ℓ_∞ metric, and the size of the covering depends exponentially on the dimension D .

Sample complexity based on $\ell_{2,\infty}$ dimension The following result bounds the expected squared ℓ_2 distance between the MLE estimator and the truth, by a constant times the best Kullback-Leibler divergence from the truth to any hypothesis, plus a penalty term roughly $\tilde{O}(D/m)$ where D is the dimension of \mathcal{H} and m is the number of data points. The Kullback-Leibler divergence between $h, \tilde{h} \in \mathcal{H}$ is defined as

$$\text{KL}(h, \tilde{h}) := \mathbb{E}_{Z \sim p_h(Z)} \left[\log \frac{h^2(Z)}{\tilde{h}^2(Z)} \right]. \quad (24)$$

Theorem 10 (Theorem 3 in (Birgé & Massart, 1998)). Assume \mathcal{H} has $\ell_{2,\infty}$ dimension $D \in [1, m]$. Let \hat{h} be an approximate MLE estimator, i.e., it satisfies (21). Then there is a constant $c > 0$ such that

$$\mathbb{E}_{\mathcal{D}_m} [d^2(h^*, \hat{h})] \leq c \inf_{h \in \mathcal{H}} \text{KL}(h^*, h) + \frac{cD}{m} (1 + \log[c_0(1 + c_1)])$$

where $\mathbb{E}_{\mathcal{D}_m}$ is with respect to the randomness in the data Z_1, \dots, Z_m generated from the true distribution $(h^*)^2$.

On the right hand side of the bound is the Kullback-Leibler divergence, instead of the squared distance as on the left. The following lemma is useful for connecting the two.

Lemma 11 (Eqn. (7.5) and (7.6) in Lemma 5 in (Birgé & Massart, 1998)). If h and \tilde{h} are the square roots of two probabilities and $\|h/\tilde{h}\|_\infty < +\infty$, then

$$d^2(h, \tilde{h}) \leq \text{KL}(h, \tilde{h}) \leq 2[1 + \log \|h/\tilde{h}\|_\infty] d^2(h, \tilde{h}).$$

²The result (Birgé & Massart, 1998) actually depends on a covering property, which basically says that the $\ell_{2,\infty}$ dimension of \mathcal{H} is bounded by D . For our purpose, it is more convenient to introduce a definition of the dimension. Also note that in (Birgé & Massart, 1998), the covering property actually requires that \mathcal{T} is simultaneously an ϵ -net of \mathcal{B} for the ℓ_2 metric and a $c_1 \epsilon$ -net for the ℓ_∞ metric. But in fact, this requirement can be relaxed to that \mathcal{T} is a covering (instead of a net) as in our definition. See Assumption $\mathcal{M}'_{2,\infty}$ and Theorem 10 in their subsequent work (?).

C.2. Estimation for individual node

Here we consider learning f_j for a fixed node j . Assume that the event stated in Lemma 1 happens, and fix the set of random features r_{j1}, \dots, r_{jK} . We first formalize our hypothesis class for learning f_j , and then analyze the sample complexity.

Hypothesis class Recall that for learning f_j , we get training data in the form $Z_i = (\chi_{S_i}, y_{ij})$, where $\chi_{S_i} \in \{0, 1\}^d$ is the indicator vector of S_i and $y_{ij} \in \{0, 1\}$ indicates whether node j gets influenced by S_i . Let p^* denote the true distribution

$$p^*(\chi_S, y) = p_\chi(\chi_S)p(y|\chi = \chi_S)$$

where p_χ is the distribution of χ_S , and $p(y|\chi = \chi_S)$ is the conditional probability

$$p(y|\chi = \chi_S) = [f_j(\chi_S)]^y [1 - f_j(\chi_S)]^{1-y}.$$

Similarly, given a function f , define the distribution induced as

$$p(\chi_S, y|f) = p_\chi(\chi_S)p(y|\chi = \chi_S, f) \quad \text{where} \quad p(y|\chi = \chi_S, f) = [f(\chi_S)]^y [1 - f(\chi_S)]^{1-y}.$$

We could define our hypothesis class as the square roots of the probability distributions induced by functions in $\hat{\mathcal{F}}^w$. Unfortunately, there is some subtle technical difficulty: $p(\chi_S, y|f)$ can be arbitrarily close to 0, in which case our technique for bounding the dimension of our hypothesis class fails (in particular, we cannot construct coverings for our hypotheses based on coverings for the weights; see the proof of Lemma 15). Therefore, we add a small offset to functions in $\hat{\mathcal{F}}^w$ and ensure that they are bounded away from 0. More precisely, define

$$\hat{\mathcal{F}}^{w,\lambda} := \left\{ f^{w,\lambda} \mid f^{w,\lambda} = f^w + \lambda, f^w \in (1 - 2\lambda)\hat{\mathcal{F}}^w \right\} \quad (25)$$

where $\lambda \in (0, 1)$ is a constant whose value will be determined later. For any $f^{w,\lambda} \in \hat{\mathcal{F}}^{w,\lambda}$, we have $\lambda \leq f^{w,\lambda}(\chi_S) \leq 1 - \lambda$ for any χ_S . Then the probability $p(\chi_S, y|f^{w,\lambda})$ introduced by $f^{w,\lambda}$ satisfies that $p(\chi_S, y|f^{w,\lambda}) \geq \lambda$ for any χ_S and y , which will allow us to use our technique.

Still, for $\hat{\mathcal{F}}^{w,\lambda}$ to be meaningful, we need to show there exists a function in $\hat{\mathcal{F}}^{w,\lambda}$ close to f_j . The following lemma shows that this is true as long as λ is small.

Lemma 12. *Assume that the statement in Lemma 1 happens. Then there exists $\hat{f}^{w,\lambda} \in \hat{\mathcal{F}}^{w,\lambda}$ such that*

$$\mathbb{E}_{\chi_S \sim p_\chi} \left[(f_j(\chi_S) - \hat{f}^{w,\lambda}(\chi_S))^2 \right] \leq 2\epsilon^2 + 2\lambda^2.$$

Proof. Let $\hat{f}^w \in \hat{\mathcal{F}}^w$ be such that $\mathbb{E}_{\chi_S \sim p_\chi} \left[(f_j(\chi_S) - \hat{f}^w(\chi_S))^2 \right] \leq \epsilon^2$. Define $\hat{f}^{w,\lambda} = (1 - 2\lambda)\hat{f}^w + \lambda$. Then $|\hat{f}^w(\chi_S) - \hat{f}^{w,\lambda}(\chi_S)| = |\lambda - 2\lambda\hat{f}^w(\chi_S)| \leq \lambda$. The lemma then follows from

$$\mathbb{E}_{\chi_S \sim p_\chi} \left[(f_j(\chi_S) - \hat{f}^{w,\lambda}(\chi_S))^2 \right] \leq 2\mathbb{E}_{\chi_S \sim p_\chi} \left[(f_j(\chi_S) - \hat{f}^w(\chi_S))^2 \right] + 2\mathbb{E}_{\chi_S \sim p_\chi} \left[(\hat{f}^{w,\lambda}(\chi_S) - \hat{f}^w(\chi_S))^2 \right].$$

□

Therefore, our hypothesis class is defined as

$$\mathcal{H}_K := \left\{ \sqrt{p(\chi_S, y|f^{w,\lambda})} \mid f^{w,\lambda} \in \hat{\mathcal{F}}^{w,\lambda} \right\}. \quad (26)$$

In other words, \mathcal{H}_K is the square roots of the probabilities induced by $\hat{\mathcal{F}}^{w,\lambda}$. Let $h^* = \sqrt{p^*(\chi_S, y)}$ denote the element corresponding to the true distribution. Note that we do not assume h^* is in \mathcal{H}_K .

Sample complexity To bound the dimension of \mathcal{H}_K and apply Theorem 10, the key is to construct coverings for \mathcal{H}_K based on those for the weights, since the feasible set of weights is a subset of \mathbb{R}^K which has nice structure. We first relate the topology of \mathcal{H}_K to that of the weights w in Lemma 14, which makes the construction possible. We then bound on the dimension in Lemma 15, and subsequently bound the sample complexity in Lemma 2.

To begin with, let $\Delta := \{w \mid w \geq \mathbf{0}, \|w\|_1 \leq 1 - 2\lambda\}$ denote the feasible set of the weights w of the functions in $\hat{\mathcal{F}}^{w,\lambda}$, and consider a mapping $\pi : \Delta \rightarrow \mathcal{H}_K$ as follows:

$$\pi(w) := \sqrt{p(\cdot|f^{w,\lambda})}, \quad \text{where} \quad f^{w,\lambda}(\chi_S) = \sum_{k=1}^K w_k \phi(\chi_S^\top r_k) + \lambda.$$

Lemma 14 shows that the ℓ_2 distance between $\pi(w)$ and $\pi(w')$ is approximately the ℓ_∞ distance between w and w' , relating the topology of \mathcal{H}_K to that of the weights w . The following quantity is useful in the process:

Definition 13. Let $A^j = \Sigma \Phi^j$ where Σ is a $2^n \times 2^n$ diagonal matrix with entries $\Sigma_{\chi_S, \chi_S} = \sqrt{p_\chi(\chi_S)}$, and Φ^j is a $2^n \times K$ matrix with entries $\Phi_{\chi_S, k}^j = \phi(\chi_S^\top r_{jk})$. Define

$$\Lambda^j := \min_{w \neq 0} \frac{\|A^j w\|}{\|w\|}, \quad \Lambda = \min_{j \in [d]} \Lambda^j.$$

Intuitively, Λ reflects how the change in w affects $A^j w$, which subsequently affects the corresponding hypothesis in \mathcal{H}_K . This quantity thus goes into the relation between the distance on the set of w and the distance on \mathcal{H}_K , as shown in Lemma 14.

Lemma 14. For an $w, w' \in \Delta$,

$$\frac{\Lambda}{2} \|w - w'\|_\infty \leq \|\pi(w) - \pi(w')\| \leq \frac{K}{\sqrt{2\Lambda}} \|w - w'\|_\infty.$$

Proof. For simplicity, let f be a shorthand of $f^{w, \lambda}(\chi_S)$ and f' be a shorthand of $f^{w', \lambda}(\chi_S)$ in the proof.

(1) By definition of the norm in (22), we have

$$\begin{aligned} \|\pi(w) - \pi(w')\|^2 &= \sum_{(\chi_S, y)} (\sqrt{p(\chi_S, y|f)} - \sqrt{p(\chi_S, y|f')})^2 \\ &= \sum_{\chi_S} p_\chi(\chi_S) \sum_y (\sqrt{p(y|\chi = \chi_S, f)} - \sqrt{p(y|\chi = \chi_S, f')})^2 \\ &= \sum_{\chi_S} p_\chi(\chi_S) \left[(\sqrt{f} - \sqrt{f'})^2 + (\sqrt{1-f} - \sqrt{1-f'})^2 \right] \\ &= \mathbb{E}_{p_\chi} \left[(\sqrt{f} - \sqrt{f'})^2 + (\sqrt{1-f} - \sqrt{1-f'})^2 \right]. \end{aligned}$$

This leads to

$$\|\pi(w) - \pi(w')\|^2 \geq \mathbb{E}_{p_\chi} \left[(\sqrt{f} - \sqrt{f'})^2 \right] \geq \frac{1}{4} \mathbb{E}_{p_\chi} [(f - f')^2] = \frac{1}{4} \sum_{\chi_S} p_\chi(\chi_S) (f - f')^2$$

where the second inequality follows from Lemma 16. The right hand side expands to

$$\frac{1}{4} \sum_{\chi_S} p_\chi(\chi_S) (f - f')^2 = \frac{1}{4} \sum_{\chi_S} p_\chi(\chi_S) \left[\sum_{k=1}^K \phi(\chi_S^\top r_{jk})(w_k - w'_k) \right]^2 = \frac{1}{4} \|A^j w - A^j w'\|^2.$$

where the last step follows from the definition of A^j . So

$$\|\pi(w) - \pi(w')\|^2 \geq \frac{1}{4} \|A^j w - A^j w'\|^2 \geq \frac{\Lambda^2}{4} \|w - w'\|^2 \geq \frac{\Lambda^2}{4} \|w - w'\|_\infty^2$$

where the second inequality follows from the definition of Λ .

(2) By definition we have

$$|f(\chi_S) - f'(\chi_S)| \leq \|w - w'\|_1 \leq K \|w - w'\|_\infty$$

for any χ_S . Then

$$\begin{aligned} \|\pi(w) - \pi(w')\|^2 &= \mathbb{E}_{p_\chi} \left[(\sqrt{f} - \sqrt{f'})^2 + (\sqrt{1-f} - \sqrt{1-f'})^2 \right] \\ &\leq \mathbb{E}_{p_\chi} \left[\frac{(f - f')^2}{4\lambda} + \frac{((1-f) - (1-f'))^2}{4\lambda} \right] \leq \frac{K^2}{2\lambda} \|w - w'\|_\infty^2 \end{aligned}$$

where the first inequality follows from Lemma 16.(2) and the fact that $\lambda \leq f \leq 1 - \lambda$ and $\lambda \leq f' \leq 1 - \lambda$. \square

Lemma 15. The $\ell_{2, \infty}$ dimension of \mathcal{H}_K is at most K .

Proof. To bound the dimension, the key is to construct coverings of small sizes. By Lemma 14, the ℓ_2 metric on \mathcal{H}_K approximately corresponds to the ℓ_∞ metric on the set of weights. So based on coverings for the weights with respect to

the ℓ_∞ metric, we can construct coverings for \mathcal{H}_K with respect to the ℓ_2 metric. We then show that they are also coverings with respect to the ℓ_∞ metric. The bound on the dimension then follows from the sizes of these coverings.

More precisely, given $\epsilon > 0$ and a ball $\mathcal{B} \subseteq \mathcal{H}_K$ with radius $R > 5\epsilon$, we construct an ϵ -covering \mathcal{T} as follows. Define $\mathcal{B}^w = \pi^{-1}(\mathcal{B})$. By Lemma 14, the radius of \mathcal{B}^w is at most $R^w = \frac{2}{\lambda}R$ (with respect to the ℓ_∞ metric). Now consider finding an ϵ^w -covering for \mathcal{B}^w with respect to the ℓ_∞ metric, where $\epsilon^w = (\frac{K}{\sqrt{2\lambda}})^{-1}\epsilon$. Since $\mathcal{B}^w \subseteq \mathbb{R}^K$, by taking the grid with length $\epsilon^w/2$ on each dimension, we can get such a covering \mathcal{T}^w with

$$|\mathcal{T}^w| \leq \left(\frac{4R^w}{\epsilon^w} \right)^K \leq \left(\frac{8K}{\sqrt{2\lambda}\Lambda} \frac{R}{\epsilon} \right)^K.$$

Let $\mathcal{T} = \pi(\mathcal{T}^w)$, and for any $h \in \mathcal{B}$ find \tilde{h} as follows. Suppose $w_h \in \mathcal{B}^w$ satisfies $\pi(w_h) = h$ and $w_{\tilde{h}}$ is the nearest neighbor of w_h in \mathcal{T}^w , then we set $\tilde{h} = \pi(w_{\tilde{h}})$. See Figure 4 for an illustration.

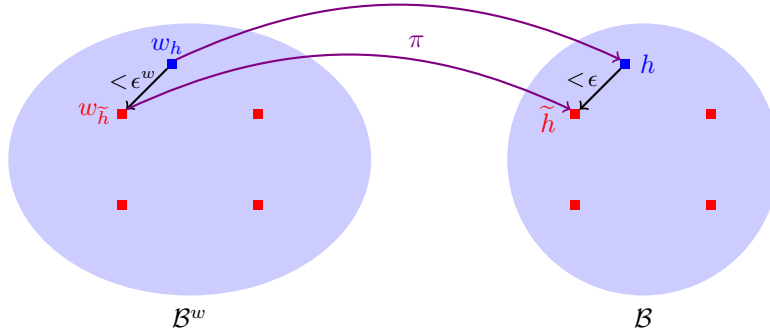


Figure 4. Illustration of the mapping.

First, we argue that \mathcal{T} is an ϵ -covering w.r.t. the ℓ_2 metric, i.e., $d(h, \tilde{h}) < \epsilon$ for any $h \in \mathcal{B}$. It follows from Lemma 14:

$$d(h, \tilde{h}) \leq \frac{K}{\sqrt{2\lambda}} \|w_h - w_{\tilde{h}}\|_\infty < \frac{K}{\sqrt{2\lambda}} \epsilon^w = \epsilon.$$

Second, we argue that \mathcal{T} is also an $O(\epsilon)$ -covering w.r.t. the ℓ_∞ metric, i.e., $d_\infty(h, \tilde{h}) = \|h - \tilde{h}\|_\infty = O(\epsilon)$ for any $h \in \mathcal{B}$. We have $\|h - \tilde{h}\| < \epsilon$, then $\|w_h - w_{\tilde{h}}\|_\infty < \frac{2}{\lambda}\epsilon$ by Lemma 14. Let $f_h := f^{w_h, \lambda}$ and $f_{\tilde{h}} := f^{w_{\tilde{h}}, \lambda}$. Then

$$|f_h(\chi_S) - f_{\tilde{h}}(\chi_S)| \leq \|w_h - w_{\tilde{h}}\|_1 \leq K \|w_h - w_{\tilde{h}}\|_\infty < \frac{2K}{\lambda}\epsilon$$

for any χ_S , and thus

$$\begin{aligned} \|\pi(w_h) - \pi(w_{\tilde{h}})\|_\infty &= \max_{\chi_S} \max \left\{ |\sqrt{f_h} - \sqrt{f_{\tilde{h}}}|, |\sqrt{1-f_h} - \sqrt{1-f_{\tilde{h}}}| \right\} \\ &\leq \max_{\chi_S} |\sqrt{f_h} - \sqrt{f_{\tilde{h}}}| \leq \max_{\chi_S} \frac{|f_h - f_{\tilde{h}}|}{2\sqrt{\lambda}} < \frac{K}{\Lambda\sqrt{\lambda}}\epsilon \end{aligned}$$

where the second inequality follows from Lemma 16.(2).

So the conditions in the definition of the dimension are satisfied with $D = K$ and $c_0 = c_1 = O\left(\frac{K}{\Lambda\sqrt{\lambda}}\right)$, and thus the dimension of \mathcal{H}_K is at most K . \square

Lemma 2. Assume the statement in Lemma 1 happens. Let \hat{h} be an approximate MLE estimator, i.e., it satisfies (21). Let $\hat{f}_j^{w, \lambda}$ be the corresponding function in $\hat{\mathcal{F}}^{w, \lambda}$. Then when $m = O\left(\frac{K}{\epsilon} \log \frac{K}{\lambda\Lambda}\right)$,

$$\mathbb{E}_{\mathcal{D}^m} \left[\mathbb{E}_{p_\chi} [(\hat{f}(\chi_S) - f_j(\chi_S))^2] \right] \leq 8c \left(\epsilon + \frac{\epsilon^2 + \lambda^2}{\lambda} \left[1 + \log \frac{1}{\lambda} \right] \right)$$

where c is the constant in Theorem 10.

Proof. The lemma follows from Theorem 10 and Lemma 15. On the left hand side of that bound in Theorem 10, we have

$$\begin{aligned} d^2(h^*, \hat{h}) &= \mathbb{E}_{p_{\mathcal{X}}} \left[\left(\sqrt{\hat{f}_j^{w,\lambda}} - \sqrt{f_j} \right)^2 + \left(\sqrt{1 - \hat{f}_j^{w,\lambda}} - \sqrt{1 - f_j} \right)^2 \right] \\ &\geq \mathbb{E}_{p_{\mathcal{X}}} \left[\left(\sqrt{\hat{f}_j^{w,\lambda}} - \sqrt{f_j} \right)^2 \right] \geq \frac{1}{4} \mathbb{E}_{p_{\mathcal{X}}} \left[\left(\hat{f}_j^{w,\lambda} - f_j \right)^2 \right] \end{aligned}$$

where the last inequality follows from Lemma 16.(1).

On the right hand side of the bound, the number of points m is sufficiently large so that the penalty term is at most ϵ . So it suffices to show that $\inf_{h \in \mathcal{H}_K} \text{KL}(h^*, h) \leq 2[1 + \log \frac{1}{\lambda}] \frac{\epsilon^2 + \lambda^2}{\lambda}$. By Lemma 12, there exists $\hat{f}^{w,\lambda} \in \hat{\mathcal{F}}^{w,\lambda}$ such that $E_{p_{\mathcal{X}}}[(f_j - \hat{f}^{w,\lambda})^2] \leq 2\epsilon^2 + 2\lambda^2$. Let $\hat{h}^{w,\lambda} = \sqrt{p(\cdot | \hat{f}^{w,\lambda})}$ denote the element in \mathcal{H}_K corresponding to $\hat{f}^{w,\lambda}$. Then

$$\text{KL}(h^*, \hat{h}^{w,\lambda}) \leq 2 \left[1 + \log \left\| \frac{h^*}{\hat{h}^{w,\lambda}} \right\|_{\infty} \right] d^2(h^*, \hat{h}^{w,\lambda}) \leq 2 \left[1 + \log \frac{1}{\lambda} \right] d^2(h^*, \hat{h}^{w,\lambda})$$

where the first inequality follows from Lemma 11, and the second inequality follows from the definition of h^* and $\hat{h}^{w,\lambda}$, and the fact that $p(\chi_{\mathcal{S}}, y | \hat{f}^{w,\lambda}) \geq \lambda$ for any $\chi_{\mathcal{S}}$ and y . The proof is completed by noting

$$\begin{aligned} d^2(h^*, \hat{h}^{w,\lambda}) &= \mathbb{E}_{p_{\mathcal{X}}} \left[\left(\sqrt{f_j} - \sqrt{\hat{f}^{w,\lambda}} \right)^2 + \left(\sqrt{1 - f_j} - \sqrt{1 - \hat{f}^{w,\lambda}} \right)^2 \right] \\ &\leq \frac{\mathbb{E}_{p_{\mathcal{X}}}[(f_j - \hat{f}^{w,\lambda})^2]}{4\lambda} + \frac{\mathbb{E}_{p_{\mathcal{X}}}[(1 - f_j - (1 - \hat{f}^{w,\lambda}))^2]}{4\lambda} \leq \frac{\epsilon^2 + \lambda^2}{\lambda} \end{aligned}$$

where the first inequality follows from Lemma 16.(2), and the last inequality follows from the choice of $\hat{f}^{w,\lambda}$ as in Lemma 12. \square

Below are some technical facts that are used in the analysis.

Lemma 16. (1) If $f_1, f_2 \in [0, 1]$, then $4(\sqrt{f_1} - \sqrt{f_2})^2 \geq (f_1 - f_2)^2$.

(2) If $f_1 \geq \lambda > 0$ and $f_2 \geq \lambda$, then $|\sqrt{f_1} - \sqrt{f_2}| \leq \frac{|f_1 - f_2|}{2\sqrt{\lambda}}$.

Proof. Both claims follow from the fact that $f_1 - f_2 = (\sqrt{f_1} - \sqrt{f_2})(\sqrt{f_1} + \sqrt{f_2})$. \square

C.3. Estimation of the entire influence function

We now combine the bounds for individual nodes to get the sample complexity for learning the entire influence function.

Theorem 3. Let $\epsilon \in (0, 1/4)$ and $\lambda = \frac{\epsilon}{c'd \log \frac{d}{\epsilon}}$ where $c' > 0$ is a sufficiently large constant. If $K = O(\frac{C^2 d^2}{\epsilon^2} \log^2 \frac{d}{\epsilon} [\log \frac{Cd}{\delta} + \log \frac{d}{\epsilon}])$,

$$m = O \left(\frac{C^2 d^3}{\epsilon^3} \log^3 \frac{d}{\epsilon} \left[\log \frac{1}{\Lambda} + \log \frac{Cd}{\epsilon} + \log \frac{d}{\delta} \right] \right)$$

then with probability $1 - \delta$ over the drawing of the random features,

$$\mathbb{E}_{\mathcal{D}_m} \left[\mathbb{E}_{p_{\mathcal{X}}} \left[\left(\sum_{j=1}^d \hat{f}_j^{w,\lambda}(\chi_{\mathcal{S}}) - \sigma(\mathcal{S}) \right)^2 \right] \right] \leq \epsilon$$

where $\mathbb{E}_{\mathcal{D}_m}$ is with respect to the randomness of $\{(\chi_{\mathcal{S}_i}, \mathbf{y}_i)\}_{i=1}^m$. The running time of the algorithm is $O(dmK)$.

Proof. Let $\lambda = \epsilon_0$ and $\epsilon_0 = \frac{\epsilon}{c'd \log \frac{d}{\epsilon}}$ where $c' > 0$ is a sufficiently large constant, so that $8c \left(\epsilon_0 + \frac{\epsilon_0^2 + \lambda^2}{\lambda} [1 + \log \frac{1}{\lambda}] \right) \leq \frac{\epsilon}{2d}$ where c is the constant in Lemma 2.

Apply Lemma 1 with error rate ϵ_0 and confidence parameter δ/d . Then when $K = O(\frac{C^2}{\epsilon_0^2} \log \frac{Cd}{\delta})$, with probability at least

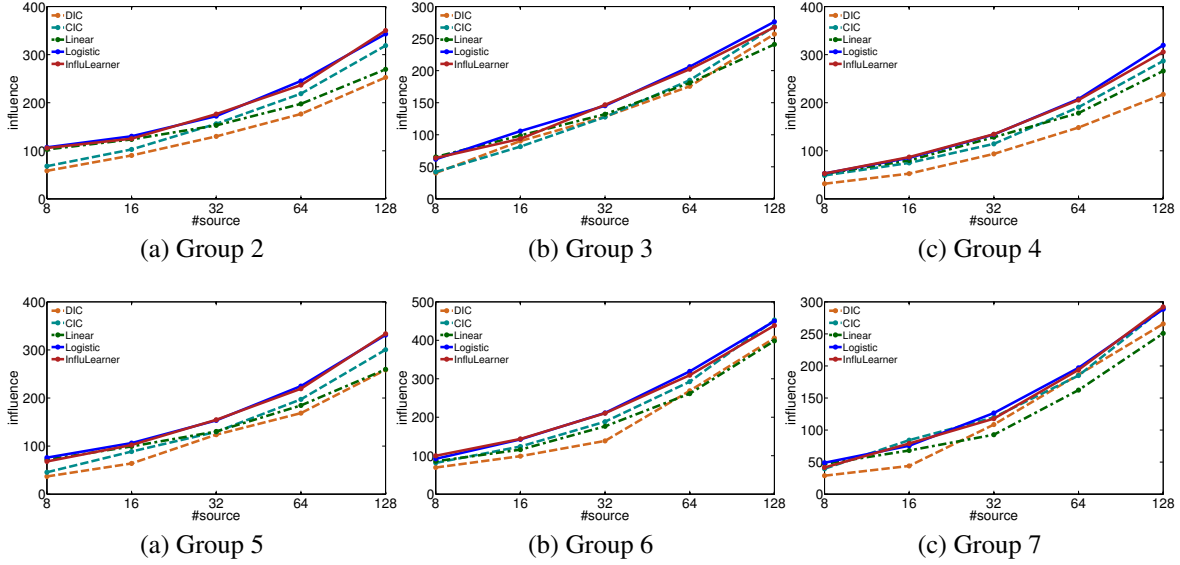


Figure 5. Expected influence vs. #sources on the real hold-out testing data.

$1 - \delta$, for each node $j \in [d]$ there exists $\hat{f}'_j \in \hat{\mathcal{F}}^w$ satisfying

$$\mathbb{E}_{p_{\chi}} \left[(\hat{f}'_j(\chi_S) - f_j(\chi_S))^2 \right] \leq \epsilon_0^2.$$

Then by Lemma 2, when $m = O\left(\frac{K}{\epsilon_0} \log \frac{K}{\lambda \lambda}\right)$, for each node $j \in [d]$ we find $\hat{f}_j^{w, \lambda}$ satisfying

$$\mathbb{E}_{\mathcal{D}_m} \left[\mathbb{E}_{p_{\chi}} \left[(\hat{f}_j^{w, \lambda}(\chi_S) - f_j(\chi_S))^2 \right] \right] \leq \frac{\epsilon}{2d}.$$

The theorem follows from the fact that $\mathbb{E}_{p_{\chi}} [(\sum_{j=1}^d \hat{f}_j^{w, \lambda}(\chi_S) - \sigma(\mathcal{S}))^2] \leq 2 \sum_{j=1}^d \mathbb{E}_{p_{\chi}} [(\hat{f}_j^{w, \lambda}(\chi_S) - f_j(\chi_S))^2]$.

Runtime. The maximum likelihood estimation only needs to be solved approximately. In particular, it suffices to get \hat{h} such that

$$\sum_{i=1}^m \log[\hat{h}(Z_i)] + 1 \geq \sup_{h \in \mathcal{H}_K} \sum_{i=1}^m \log[h(Z_i)].$$

By the convergence rate of EG (see Section 4.4 in (Kivinen & Warmuth, 1997)), we only need $O(1/\eta)$ iterations, where the learning rate η can be viewed as a constant. Each iteration takes time $O(mK)$, and we need to use EG for each of the d nodes. Hence, the total time is $O(dmK)$. \square

C.4. Additional experimental results

We report the additional experimental evaluations on the application of the learnt influence functions to the continuous-time influence maximization problem on the rest six groups of hold-out real testing cascades datasets. Compared to DIC and Linear regression, Figure 5 verifies that the performance of INFLUeR, Modified Logistic and CIC are better and more consistent with each other.