

---

# Learning Time-Varying Coverage Functions

---

Nan Du<sup>†</sup>, Yingyu Liang<sup>‡</sup>, Maria-Florina Balcan<sup>◊</sup>, Le Song<sup>†</sup>

<sup>†</sup>College of Computing, Georgia Institute of Technology

<sup>‡</sup>Department of Computer Science, Princeton University

<sup>◊</sup>School of Computer Science, Carnegie Mellon University

dunan@gatech.edu, yingyul@cs.princeton.edu

ninamf@cs.cmu.edu, lsong@cc.gatech.edu

## Abstract

Coverage functions are an important class of discrete functions that capture the law of diminishing returns arising naturally from applications in social network analysis, machine learning, and algorithmic game theory. In this paper, we propose a new problem of learning time-varying coverage functions, and develop a novel parametrization of these functions using random features. Based on the connection between time-varying coverage functions and counting processes, we also propose an efficient parameter learning algorithm based on likelihood maximization, and provide a sample complexity analysis. We applied our algorithm to the influence function estimation problem in information diffusion in social networks, and show that with few assumptions about the diffusion processes, our algorithm is able to estimate influence significantly more accurately than existing approaches on both synthetic and real world data.

## 1 Introduction

Coverage functions are a special class of the more general submodular functions which play important role in combinatorial optimization with many interesting applications in social network analysis [1], machine learning [2], economics and algorithmic game theory [3], etc. A particularly important example of coverage functions in practice is the influence function of users in information diffusion modeling [1] — news spreads across social networks by word-of-mouth and a set of influential sources can collectively trigger a large number of follow-ups. Another example of coverage functions is the valuation functions of customers in economics and game theory [3] — customers are thought to have certain requirements and the items being bundled and offered fulfill certain subsets of these demands.

Theoretically, it is usually assumed that users' influence or customers' valuation are known in advance as an oracle. In practice, however, these functions must be learned. For example, given past traces of information spreading in social networks, a social platform host would like to estimate how many follow-ups a set of users can trigger. Or, given past data of customer reactions to different bundles, a retailer would like to estimate how likely customer would respond to new packages of goods. Learning such combinatorial functions has attracted many recent research efforts from both theoretical and practical sides (*e.g.*, [4, 5, 6, 7, 8]), many of which show that coverage functions can be learned from just polynomial number of samples.

However, the prior work has widely ignored an important dynamic aspect of the coverage functions. For instance, information spreading is a dynamic process in social networks, and the number of follow-ups of a fixed set of sources can increase as observation time increases. A bundle of items or features offered to customers may trigger a sequence of customer actions over time. These real world problems inspire and motivate us to consider a novel *time-varying coverage function*,  $f(\mathcal{S}, t)$ , which is a coverage function of the set  $\mathcal{S}$  when we fix a time  $t$ , and a continuous monotonic function of time  $t$  when we fix a set  $\mathcal{S}$ . While learning time-varying combinatorial structures has been ex-

explored in graphical model setting (*e.g.*, [9, 10]), as far as we are aware of, learning of time-varying coverage function has not been addressed in the literature. Furthermore, we are interested in estimating the entire function of  $t$ , rather than just treating the time  $t$  as a discrete index and learning the function value at a small number of discrete points. From this perspective, our formulation is the generalization of the most recent work [8] with even less assumptions about the data used to learn the model.

Generally, we assume that the historical data are provided in pairs of a set and a collection of timestamps when caused events by the set occur. Hence, such a collection of temporal events associated with a particular set  $\mathcal{S}_i$  can be modeled principally by a counting process  $N_i(t), t \geq 0$  which is a stochastic process with values that are positive, integer, and increasing along time [11]. For instance, in the information diffusion setting of online social networks, given a set of earlier adopters of some new product,  $N_i(t)$  models the time sequence of all triggered events of the followers, where each jump in the process records the timing  $t_{ij}$  of an action. In the economics and game theory setting, the counting process  $N_i(t)$  records the number of actions a customer has taken over time given a particular bundled offer. This essentially raises an interesting question of how to estimate the time-varying coverage function from the angle of counting processes. We thus propose a novel formulation which builds a connection between the two by modeling the cumulative intensity function of a counting process as a time-varying coverage function. The key idea is to parametrize the intensity function as a weighted combination of random kernel functions. We then develop an efficient learning algorithm TCOVERAGELARNER to estimate the parameters of the function using maximum likelihood approach. We show that our algorithm can provably learn the time-varying coverage function using only polynomial number of samples. Finally, we validate TCOVERAGELARNER on both influence estimation and maximization problems by using cascade data from information diffusion. We show that our method performs significantly better than alternatives with little prior knowledge about the dynamics of the actual underlying diffusion processes.

## 2 Time-Varying Coverage Function

We will first give a formal definition of the time-varying coverage function, and then explain its additional properties in details.

**Definition.** Let  $\mathcal{U}$  be a (potentially uncountable) domain. We endow  $\mathcal{U}$  with some  $\sigma$ -algebra  $\mathcal{A}$  and denote a probability distribution on  $\mathcal{U}$  by  $\mathbb{P}$ . A coverage function is a combinatorial function over a finite set  $\mathcal{V}$  of items, defined as

$$f(\mathcal{S}) := Z \cdot \mathbb{P} \left( \bigcup_{s \in \mathcal{S}} \mathcal{U}_s \right), \quad \text{for all } \mathcal{S} \in 2^{\mathcal{V}}, \quad (1)$$

where  $\mathcal{U}_s \subset \mathcal{U}$  is the subset of domain  $\mathcal{U}$  covered by item  $s \in \mathcal{V}$ , and  $Z$  is the additional normalization constant. For time-varying coverage functions, we let the size of the subset  $\mathcal{U}_s$  to grow monotonically over time, that is

$$\mathcal{U}_s(t) \subseteq \mathcal{U}_s(\tau), \quad \text{for all } t \leq \tau \text{ and } s \in \mathcal{V}, \quad (2)$$

which results in a combinatorial temporal function

$$f(\mathcal{S}, t) = Z \cdot \mathbb{P} \left( \bigcup_{s \in \mathcal{S}} \mathcal{U}_s(t) \right), \quad \text{for all } \mathcal{S} \in 2^{\mathcal{V}}. \quad (3)$$

In this paper, we assume that  $f(\mathcal{S}, t)$  is smooth and continuous, and its first order derivative with respect to time,  $f'(\mathcal{S}, t)$ , is also smooth and continuous.

**Representation.** We now show that a time-varying coverage function,  $f(\mathcal{S}, t)$ , can be represented as an expectation over random functions based on multidimensional step basis functions. Since  $\mathcal{U}_s(t)$  is varying over time, we can associate each  $u \in \mathcal{U}$  with a  $|\mathcal{V}|$ -dimensional vector  $\tau_u$  of change points. In particular, the  $s$ -th coordinate of  $\tau_u$  records the time that source node  $s$  covers  $u$ . Let  $\tau$  to be a random variable obtained by sampling  $u$  according to  $\mathbb{P}$  and setting  $\tau = \tau_u$ . Note that given all  $\tau_u$  we can compute  $f(\mathcal{S}, t)$ ; now we claim that the distribution of  $\tau$  is sufficient.

We first introduce some notations. Based on  $\tau_u$  we define a  $|\mathcal{V}|$ -dimensional step function  $\mathbf{r}_u(t) : \mathbb{R}_+ \mapsto \{0, 1\}^{|\mathcal{V}|}$ , where the  $s$ -th dimension of  $\mathbf{r}_u(t)$  is 1 if  $u$  is covered by the set  $\mathcal{U}_s(t)$  at time  $t$ , and 0 otherwise. To emphasize the dependence of the function  $\mathbf{r}_u(t)$  on  $\tau_u$ , we will also write  $\mathbf{r}_u(t)$  as  $\mathbf{r}_u(t|\tau_u)$ . We denote the indicator vector of a set  $\mathcal{S}$  by  $\chi_{\mathcal{S}} \in \{0, 1\}^{|\mathcal{V}|}$  where the  $s$ -th dimension of  $\chi_{\mathcal{S}}$  is 1 if  $s \in \mathcal{S}$ , and 0 otherwise. Then  $u \in \mathcal{U}$  is covered by  $\bigcup_{s \in \mathcal{S}} \mathcal{U}_s(t)$  at time  $t$  if  $\chi_{\mathcal{S}}^\top \mathbf{r}_u(t) \geq 1$ .

**Lemma 1.** *There exists a distribution  $\mathbb{Q}(\boldsymbol{\tau})$  over the vector of change points  $\boldsymbol{\tau}$ , such that the time-varying coverage function can be represented as*

$$f(\mathcal{S}, t) = Z \cdot \mathbb{E}_{\boldsymbol{\tau} \sim \mathbb{Q}(\boldsymbol{\tau})} [\phi(\boldsymbol{\chi}_{\mathcal{S}}^{\top} \mathbf{r}(t|\boldsymbol{\tau}))] \quad (4)$$

where  $\phi(x) := \min\{x, 1\}$ , and  $\mathbf{r}(t|\boldsymbol{\tau})$  is a multidimensional step function parameterized by  $\boldsymbol{\tau}$ .

*Proof.* Let  $\mathcal{U}_{\mathcal{S}} := \bigcup_{s \in \mathcal{S}} \mathcal{U}_s(t)$ . By definition (3), we have the following integral representation

$$f(\mathcal{S}, t) = Z \cdot \int_{\mathcal{U}} \mathbb{I}\{u \in \mathcal{U}_{\mathcal{S}}\} d\mathbb{P}(u) = Z \cdot \int_{\mathcal{U}} \phi(\boldsymbol{\chi}_{\mathcal{S}}^{\top} \mathbf{r}_u(t)) d\mathbb{P}(u) = Z \cdot \mathbb{E}_{u \sim \mathbb{P}(u)} [\phi(\boldsymbol{\chi}_{\mathcal{S}}^{\top} \mathbf{r}_u(t))].$$

We can define the set of  $u$  having the same  $\boldsymbol{\tau}$  as  $\mathcal{U}_{\boldsymbol{\tau}} := \{u \in \mathcal{U} \mid \boldsymbol{\tau}_u = \boldsymbol{\tau}\}$  and define a distribution over  $\boldsymbol{\tau}$  as  $d\mathbb{Q}(\boldsymbol{\tau}) := \int_{\mathcal{U}_{\boldsymbol{\tau}}} d\mathbb{P}(u)$ . Then the integral representation of  $f(\mathcal{S}, t)$  can be rewritten as

$$Z \cdot \mathbb{E}_{u \sim \mathbb{P}(u)} [\phi(\boldsymbol{\chi}_{\mathcal{S}}^{\top} \mathbf{r}_u(t))] = Z \cdot \mathbb{E}_{\boldsymbol{\tau} \sim \mathbb{Q}(\boldsymbol{\tau})} [\phi(\boldsymbol{\chi}_{\mathcal{S}}^{\top} \mathbf{r}(t|\boldsymbol{\tau}))],$$

which proves the lemma.  $\square$

### 3 Model for Observations

In general, we assume that the input data are provided in the form of pairs,  $(\mathcal{S}_i, N_i(t))$ , where  $\mathcal{S}_i$  is a set, and  $N_i(t)$  is a counting process in which each jump of  $N_i(t)$  records the timing of an event. We first give a brief overview of a counting process [11] and then motivate our model in details.

**Counting Process.** Formally, a counting process  $\{N(t), t \geq 0\}$  is any nonnegative, integer-valued stochastic process such that  $N(t') \leq N(t)$  whenever  $t' \leq t$  and  $N(0) = 0$ . The most common use of a counting process is to count the number of occurrences of temporal events happening along time, so the index set is usually taken to be the nonnegative real numbers  $\mathbb{R}_+$ . A counting process is a submartingale:  $\mathbb{E}[N(t) \mid \mathcal{H}_{t'}] \geq N(t')$  for all  $t > t'$  where  $\mathcal{H}_{t'}$  denotes the history up to time  $t'$ . By Doob-Meyer theorem [11],  $N(t)$  has the unique decomposition:

$$N(t) = \Lambda(t) + M(t) \quad (5)$$

where  $\Lambda(t)$  is a nondecreasing predictable process called the compensator (or cumulative intensity), and  $M(t)$  is a mean zero martingale. Since  $\mathbb{E}[dM(t) \mid \mathcal{H}_{t-}] = 0$ , where  $dM(t)$  is the increment of  $M(t)$  over a small time interval  $[t, t + dt)$ , and  $\mathcal{H}_{t-}$  is the history until just before time  $t$ ,

$$\mathbb{E}[dN(t) \mid \mathcal{H}_{t-}] = d\Lambda(t) := a(t) dt \quad (6)$$

where  $a(t)$  is called the intensity of a counting process.

**Model formulation.** We assume that the cumulative intensity of the counting process is modeled by a time-varying coverage function, *i.e.*, the observation pair  $(\mathcal{S}_i, N_i(t))$  is generated by

$$N_i(t) = f(\mathcal{S}_i, t) + M_i(t) \quad (7)$$

in the time window  $[0, T]$  for some  $T > 0$ , and  $df(\mathcal{S}, t) = a(\mathcal{S}, t)dt$ . In other words, the time-varying coverage function controls the propensity of occurring events over time. Specifically, for a fixed set  $\mathcal{S}_i$ , as time  $t$  increases, the cumulative number of events observed grows accordingly for that  $f(\mathcal{S}_i, t)$  is a continuous monotonic function over time; for a given time  $t$ , as the set  $\mathcal{S}_i$  changes to another set  $\mathcal{S}_j$ , the amount of coverage over domain  $\mathcal{U}$  may change and hence can result in a different cumulative intensity. This abstract model can be mapped to real world applications. In the information diffusion context, for a fixed set of sources  $\mathcal{S}_i$ , as time  $t$  increases, the number of influenced nodes in the social network tends to increase; for a given time  $t$ , if we change the sources to  $\mathcal{S}_j$ , the number of influenced nodes may be different depending on how influential the sources are. In the economics and game theory context, for a fixed bundle of offers  $\mathcal{S}_i$ , as time  $t$  increases, it is more likely that the merchant will observe the customers' actions in response to the offers; even at the same time  $t$ , different bundles of offers,  $\mathcal{S}_i$  and  $\mathcal{S}_j$ , may have very different ability to drive the customers' actions.

Compared to a regression model  $y_i = g(\mathcal{S}_i) + \epsilon_i$  with *i.i.d.* input data  $(\mathcal{S}_i, y_i)$ , our model outputs a special random function over time, that is, a counting process  $N_i(t)$  with the noise being a zero mean martingale  $M_i(t)$ . In contrast to functional regression models, our model exploits much more interesting structures of the problem. For instance, the random function representation in the last section can be used to parametrize the model. Such special structure of the counting process allows us to estimate the parameter of our model using maximum likelihood approach efficiently, and the martingale noise enables us to use exponential concentration inequality in analyzing our algorithm.

## 4 Parametrization

Based on the following two mild assumptions, we will show how to parametrize the intensity function as a weighted combination of random kernel functions, learn the parameters by maximum likelihood estimation, and eventually derive a sample complexity.

- (A1)  $a(\mathcal{S}, t)$  is smooth and bounded on  $[0, T]$ :  $0 < a_{\min} \leq a \leq a_{\max} < \infty$ , and  $\ddot{a} := d^2a/dt^2$  is absolutely continuous with  $\int \ddot{a}(t)dt < \infty$ .  
(A2) There is a known distribution  $\mathbb{Q}'(\boldsymbol{\tau})$  and a constant  $C$  with  $\mathbb{Q}'(\boldsymbol{\tau})/C \leq \mathbb{Q}(\boldsymbol{\tau}) \leq C\mathbb{Q}'(\boldsymbol{\tau})$ .

**Kernel Smoothing** To facilitate our finite dimensional parameterization, we first convolve the intensity function with  $K(t) = k(t/\sigma)/\sigma$  where  $\sigma$  is the bandwidth parameter and  $k$  is a kernel function (such as the Gaussian RBF kernel  $k(t) = e^{-t^2/2}/\sqrt{2\pi}$ ) with

$$0 \leq k(t) \leq \kappa_{\max}, \quad \int k(t) dt = 1, \quad \int t k(t) dt = 0, \quad \text{and} \quad \sigma_k^2 := \int t^2 k(t) dt < \infty. \quad (8)$$

The convolution results in a smoothed intensity  $a^K(\mathcal{S}, t) = K(t) \star (df(\mathcal{S}, t)/dt) = d(K(t) \star \Lambda(\mathcal{S}, t))/dt$ . By the property of convolution and exchanging derivative with integral, we have that

$$\begin{aligned} a^K(\mathcal{S}, t) &= d(Z \cdot \mathbb{E}_{\boldsymbol{\tau} \sim \mathbb{Q}(\boldsymbol{\tau})} [K(t) \star \phi(\boldsymbol{\chi}_{\mathcal{S}}^\top \boldsymbol{r}(t|\boldsymbol{\tau}))]) / dt && \text{by definition of } f(\cdot) \\ &= Z \cdot \mathbb{E}_{\boldsymbol{\tau} \sim \mathbb{Q}(\boldsymbol{\tau})} [d(K(t) \star \phi(\boldsymbol{\chi}_{\mathcal{S}}^\top \boldsymbol{r}(t|\boldsymbol{\tau}))) / dt] && \text{exchange derivative and integral} \\ &= Z \cdot \mathbb{E}_{\boldsymbol{\tau} \sim \mathbb{Q}(\boldsymbol{\tau})} [K(t) \star \delta(t - t(\mathcal{S}, \boldsymbol{\tau}))] && \text{by property of convolution and function } \phi(\cdot) \\ &= Z \cdot \mathbb{E}_{\boldsymbol{\tau} \sim \mathbb{Q}(\boldsymbol{\tau})} [K(t - t(\mathcal{S}, \boldsymbol{\tau}))] && \text{by definition of } \delta(\cdot) \end{aligned}$$

where  $t(\mathcal{S}, \boldsymbol{\tau})$  is the time when function  $\phi(\boldsymbol{\chi}_{\mathcal{S}}^\top \boldsymbol{r}(t|\boldsymbol{\tau}))$  jumps from 0 to 1. If we choose small enough kernel bandwidth,  $a^K$  only incurs a small bias from  $a$ . But the smoothed intensity still results in infinite number of parameters, due to the unknown distribution  $\mathbb{Q}(\boldsymbol{\tau})$ . To address this problem, we design the following random approximation with finite number of parameters.

**Random Function Approximation** The key idea is to sample a collection of  $W$  random change points  $\boldsymbol{\tau}$  from a known distribution  $\mathbb{Q}'(\boldsymbol{\tau})$  which can be different from  $\mathbb{Q}(\boldsymbol{\tau})$ . If  $\mathbb{Q}'(\boldsymbol{\tau})$  is not very far way from  $\mathbb{Q}(\boldsymbol{\tau})$ , the random approximation will be close to  $a^K$ , and thus close to  $a$ . More specifically, we will denote the space of weighted combination of  $W$  random kernel function by

$$\mathcal{A} = \left\{ a_{\boldsymbol{w}}^K(\mathcal{S}, t) = \sum_{i=1}^W w_i K(t - t(\mathcal{S}, \boldsymbol{\tau}_i)) : \boldsymbol{w} \geq 0, \frac{Z}{C} \leq \|\boldsymbol{w}\|_1 \leq ZC \right\}, \{\boldsymbol{\tau}_i\} \stackrel{i.i.d.}{\sim} \mathbb{Q}'(\boldsymbol{\tau}). \quad (9)$$

**Lemma 2.** *If  $W = \tilde{O}(Z^2/(\epsilon\sigma)^2)$ , then with probability  $\geq 1 - \delta$ , there exists an  $\tilde{a} \in \mathcal{A}$  such that  $\mathbb{E}_{\mathcal{S}} \mathbb{E}_t [(a(\mathcal{S}, t) - \tilde{a}(\mathcal{S}, t))^2] := \mathbb{E}_{\mathcal{S} \sim \mathbb{P}(\mathcal{S})} \int_0^T [(a(\mathcal{S}, t) - \tilde{a}(\mathcal{S}, t))^2] dt / T = O(\epsilon^2 + \sigma^4)$ .*

The lemma then suggests to set the kernel bandwidth  $\sigma = O(\sqrt{\epsilon})$  to get  $O(\epsilon^2)$  approximation error.

## 5 Learning Algorithm

We develop a learning algorithm, referred to as TCOVERAGELARNER, to estimate the parameters of  $a_{\boldsymbol{w}}^K(\mathcal{S}, t)$  by maximizing the joint likelihood of all observed events based on convex optimization techniques as follows.

**Maximum Likelihood Estimation** Instead of directly estimating the time-varying coverage function, which is the cumulative intensity function of the counting process, we turn to estimate the intensity function  $a(\mathcal{S}, t) = \partial \Lambda(\mathcal{S}, t) / \partial t$ . Given  $m$  i.i.d. counting processes,  $\mathcal{D}^m := \{(\mathcal{S}_1, N_1(t)), \dots, (\mathcal{S}_m, N_m(t))\}$  up to observation time  $T$ , the log-likelihood of the dataset is [11]

$$\ell(\mathcal{D}^m | a) = \sum_{i=1}^m \left\{ \int_0^T \{\log a(\mathcal{S}_i, t)\} dN_i(t) - \int_0^T a(\mathcal{S}_i, t) dt \right\}. \quad (10)$$

Maximizing the log-likelihood with respect to the intensity function  $a(\mathcal{S}, t)$  then gives us the estimation  $\hat{a}(\mathcal{S}, t)$ . The  $W$ -term random kernel function approximation reduces a function optimization problem to a finite dimensional optimization problem, while incurring only small bias in the estimated function.

---

**Algorithm 1** TCOVERAGELERNER

---

INPUT :  $\{(\mathcal{S}_i, N_i(t))\}, i = 1, \dots, m$ ;  
 Sample  $W$  random features  $\tau_1, \dots, \tau_W$  from  $\mathbb{Q}'(\tau)$ ;  
 Compute  $\{t(\mathcal{S}_i, \tau_w)\}, \{\mathbf{g}_i\}, \{\mathbf{k}(t_{ij})\}, i \in \{1, \dots, m\}, w = 1, \dots, W, t_{ij} < T$ ;  
 Initialize  $\mathbf{w}^0 \in \Omega = \{\mathbf{w} \geq 0, \|\mathbf{w}\|_1 \leq 1\}$ ;  
 Apply projected quasi-newton algorithm [12] to solve 11;  
 OUTPUT :  $a_w^K(\mathcal{S}, t) = \sum_{i=1}^W w_i K(t - t(\mathcal{S}, \tau_i))$

---

**Convex Optimization.** By plugging the parametrization  $a_w^K(\mathcal{S}, t)$  (9) into the log-likelihood (10), we formulate the optimization problem as :

$$\min_{\mathbf{w}} \sum_{i=1}^m \left\{ \mathbf{w}^\top \mathbf{g}_i - \sum_{t_{ij} < T} \log(\mathbf{w}^\top \mathbf{k}(t_{ij})) \right\} \quad \text{subject to} \quad \mathbf{w} \geq 0, \|\mathbf{w}\|_1 \leq 1, \quad (11)$$

where we define

$$\mathbf{g}_{ik} = \int_0^T K(t - t(\mathcal{S}_i, \tau_k)) dt \quad \text{and} \quad \mathbf{k}_l(t_{ij}) = K(t_{ij} - t(\mathcal{S}_i, \tau_l)), \quad (12)$$

$t_{ij}$  when the  $j$ -th event occurs in the  $i$ -th counting process. By treating the normalization constant  $Z$  as a free variable which will be tuned by cross validation later, we simply require that  $\|\mathbf{w}\|_1 \leq 1$ . By applying the Gaussian RBF kernel, we can derive a closed form of  $\mathbf{g}_{ik}$  and the gradient  $\nabla \ell$  as

$$\mathbf{g}_{ik} = \frac{1}{2} \left\{ \operatorname{erfc} \left( -\frac{t(\mathcal{S}_i, \tau_k)}{\sqrt{2}h} \right) - \operatorname{erfc} \left( \frac{T - t(\mathcal{S}_i, \tau_k)}{\sqrt{2}h} \right) \right\}, \quad \nabla \ell = \sum_{i=1}^m \left\{ \mathbf{g}_i - \sum_{t_{ij} < T} \frac{\mathbf{k}(t_{ij})}{\mathbf{w}^\top \mathbf{k}(t_{ij})} \right\}. \quad (13)$$

A pleasing feature of this formulation is that it is convex in the argument  $\mathbf{w}$ , allowing us to apply various convex optimization techniques to solve the problem efficiently. Specifically, we first draw  $W$  random features  $\tau_1, \dots, \tau_W$  from  $\mathbb{Q}'(\tau)$ . Then, we precompute the jumping time  $t(\mathcal{S}_i, \tau_w)$  for every source set  $\{\mathcal{S}_i\}_{i=1}^m$  on each random feature  $\{\tau_w\}_{w=1}^W$ . Because in general  $|\mathcal{S}_i| \ll n$ , this computation costs  $O(mW)$ . Based on the achieved  $m$ -by- $W$  jumping-time matrix, we preprocess the feature vectors  $\{\mathbf{g}_i\}_{i=1}^m$  and  $\mathbf{k}(t_{ij}), i \in \{1, \dots, m\}, t_{ij} < T$ , which costs  $O(mW)$  and  $O(mLW)$  where  $L$  is the maximum number of events caused by a particular source set before time  $T$ . Finally, we apply the projected quasi-newton algorithm [12] to find the weight  $\mathbf{w}$  that minimizes the negative log-likelihood of observing the given event data. Because the evaluation of the objective function and the gradient, which costs  $O(mLW)$ , is much more expensive than the projection onto the convex constraint set, and  $L \ll n$ , the worst case computation complexity is thus  $O(mnW)$ . Algorithm 1 summarizes the above steps in the end.

**Sample Strategy.** One important constitution of our parametrization is to sample  $W$  random change points  $\tau$  from a known distribution  $\mathbb{Q}'(\tau)$ . Because given a set  $\mathcal{S}_i$ , we can only observe the jumping time of the events in each counting process without knowing the identity of the covered items (which is a key difference from [8]), the best thing we can do is to sample from these historical data. Specifically, let the number of counting processes that a single item  $s \in \mathcal{V}$  is involved to induce be  $N_s$ , and the collection of all the jumping timestamps before time  $T$  be  $\mathcal{J}_s$ . Then, for the  $s$ -th entry of  $\tau$ , with probability  $|\mathcal{J}_s|/nN_s$ , we uniformly draw a sample from  $\mathcal{J}_s$ ; and with probability  $1 - |\mathcal{J}_s|/nN_s$ , we assign a time much greater than  $T$  to indicate that the item will never be covered until infinity. Given the very limited information, although this  $\mathbb{Q}'(\tau)$  might be quite different from  $\mathbb{Q}(\tau)$ , by drawing sufficiently large number of samples and adjusting the weights, we expect it still can lead to good results, as illustrated in our experiments later.

## 6 Sample Complexity

Suppose we use  $W$  random features and  $m$  training examples to compute an  $\epsilon_\ell$ -MLE solution  $\hat{a}$ , i.e.,

$$\ell(\mathcal{D}^m | \hat{a}) \geq \max_{a' \in \mathcal{A}} \ell(\mathcal{D}^m | a') - \epsilon_\ell.$$

The goal is to analyze how well the function  $\hat{f}$  induced by  $\hat{a}$  approximates the true function  $f$ . This sections describes the intuition and the complete proof is provided in the appendix.

A natural choice for connecting the error between  $f$  and  $\hat{f}$  with the log-likelihood cost used in MLE is the Hellinger distance [22]. So it suffices to prove an upper bound on the Hellinger distance  $h(a, \hat{a})$  between  $\hat{a}$  and the true intensity  $a$ , for which we need to show a high probability bound on the (total) empirical Hellinger distance  $\hat{H}^2(a, a')$  between the two. Here,  $h$  and  $\hat{H}$  are defined as

$$h^2(a, a') := \frac{1}{2} \mathbb{E}_{\mathcal{S}} \mathbb{E}_t \left[ \sqrt{a(\mathcal{S}, t)} - \sqrt{a'(\mathcal{S}, t)} \right]^2,$$

$$\hat{H}^2(a, a') := \frac{1}{2} \sum_{i=1}^m \int_0^T \left[ \sqrt{a(\mathcal{S}_i, t)} - \sqrt{a'(\mathcal{S}_i, t)} \right]^2 dt.$$

The key for the analysis is to show that the empirical Hellinger distance can be bounded by a martingale plus some other additive error terms, which we then bound respectively. This martingale is defined based on our hypotheses and the martingales  $M_i$  associated with the counting process  $N_i$ :

$$M(t|g) := \int_0^t g(t) d \left( \sum_i M_i(t) \right) = \sum_{i=1}^m \int_0^t g(t) dM_i(t)$$

where  $g \in \mathcal{G} = \left\{ g_{a'} = \frac{1}{2} \log \frac{a+a'}{2a} : a' \in \mathcal{A} \right\}$ . More precisely, we have the following lemma.

**Lemma 3.** *Suppose  $\hat{a}$  is an  $\epsilon_\ell$ -MLE. Then*

$$\hat{H}^2(\hat{a}, a) \leq 16M(T; g_{\hat{a}}) + 4 \left[ \ell(\mathcal{D}^m | a) - \max_{a' \in \mathcal{A}} \ell(\mathcal{D}^m | a') \right] + 4\epsilon_\ell.$$

The right hand side has three terms: the martingale (estimation error), the likelihood gap between the truth and the best one in our hypothesis class (approximation error), and the optimization error. We then focus on bounding the martingale and the likelihood gap.

To bound the martingale, we first introduce a notion called  $(d, d')$ -covering dimension measuring the complexity of the hypothesis class, generalizing that in [25]. Based on this notion, we prove a uniform convergence inequality, combining the ideas in classic works on MLE [25] and counting process [13]. Compared to the classic uniform inequality, our result is more general, and the complexity notion has more clear geometric interpretation and are thus easier to verify. For the likelihood gap, recall that by Lemma 2, there exists an good approximation  $\tilde{a} \in \mathcal{A}$ . The likelihood gap is then bounded by that between  $a$  and  $\tilde{a}$ , which is small since  $a$  and  $\tilde{a}$  are close.

Combining the two leads to a bound on the Hellinger distance based on bounded dimension of the hypothesis class. We then show that the dimension of our specific hypothesis class is at most the number of random features  $W$ , and convert  $\hat{H}^2(\hat{a}, a)$  to the desired  $\ell_2$  error bound on  $f$  and  $\hat{f}$ .

**Theorem 4.** *Suppose  $W = \tilde{O} \left( Z^2 \left[ \left( \frac{ZT}{\epsilon} \right)^{5/2} + \left( \frac{ZT}{\epsilon a_{\min}} \right)^{5/4} \right] \right)$  and  $m = \tilde{O} \left( \frac{ZT}{\epsilon} [W + \epsilon_\ell] \right)$ . Then with probability  $\geq 1 - \delta$  over the random sample of  $\{\tau_i\}_{i=1}^W$ , we have that for any  $0 \leq t \leq T$ ,*

$$\mathbb{E}_{\mathcal{S}} \left[ \hat{f}(\mathcal{S}, t) - f(\mathcal{S}, t) \right]^2 \leq \epsilon.$$

The theorem shows that the number of random functions needed to achieve  $\epsilon$  error is roughly  $O(\epsilon^{-5/2})$ , and the sample size is  $O(\epsilon^{-7/2})$ . They also depend on  $a_{\min}$ , which means with more random functions and data, we can deal with intensities with more extreme values. Finally, they increase with the time  $T$ , i.e., it is more difficult to learn the function values at later time points.

## 7 Experiments

We evaluate TCOVERAGELARNER on both synthetic and real world information diffusion data. We show that our method can be more robust to model misspecification than other state-of-the-art alternatives by learning a temporal coverage function all at once.

### 7.1 Competitors

Because our input data only include pairs of a source set and the temporal information of its triggered events  $\{(\mathcal{S}_i, N_i(t))\}_{i=1}^m$  with unknown identity, we first choose the general kernel ridge regression model as the major baseline, which directly estimates the influence value of a source set

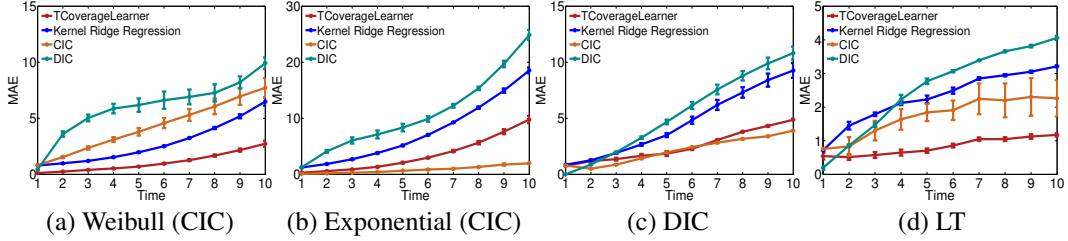


Figure 1: MAE of the estimated influence on test data along time with the true diffusion model being continuous-time independent cascade with pairwise Weibull (a) and Exponential (b) transmission functions, (c) discrete-time independent cascade model and (d) linear-threshold cascade model.

$\chi_S$  by  $f(\chi_S) = \mathbf{k}(\chi_S)(\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$  where  $\mathbf{k}(\chi_S) = K(\chi_{S_i}, \chi_S)$ , and  $\mathbf{K}$  is the kernel matrix. We discretize the time into several steps and fit a separate model to each of them. Between two consecutive time steps, the predictions are simply interpolated. In addition, to further demonstrate the robustness of TCOVERAGELEARNER, we compare it to the two-stage methods which must know the *identity* of the nodes involved in an information diffusion process to first learn a specific diffusion model based on which they can then estimate the influence. We give them such an advantage and study three well-known diffusion models: (I) Continuous-time Independent Cascade model(CIC)[14, 15]; (II) Discrete-time Independent Cascade model(DIC)[1]; and (III) Linear-Threshold cascade model(LT)[1].

## 7.2 Influence Estimation on Synthetic Data

We generate Kronecker synthetic networks ( $[0.9\ 0.5; 0.5\ 0.3]$ ) which mimic real world information diffusion patterns [16]. For CIC, we use both Weibull distribution (Wbl) and Exponential distribution (Exp) for the pairwise transmission function associated with each edge, and randomly set their parameters to capture the heterogeneous temporal dynamics. Then, we use NETRATE [14] to learn the model by assuming an exponential pairwise transmission function. For DIC, we choose the pairwise infection probability uniformly from 0 to 1 and fit the model by [17]. For LT, we assign the edge weight  $w_{uv}$  between  $u$  and  $v$  as  $1/d_v$ , where  $d_v$  is the degree of node  $v$  following [1]. Finally, 1,024 source sets are sampled with power-law distributed cardinality (with exponent 2.5), each of which induces eight independent cascades(or counting processes), and the test data contains another 128 independently sampled source sets with the ground truth influence estimated from 10,000 simulated cascades up to time  $T = 10$ . Figure 1 shows the MAE(Mean Absolute Error) between the estimated influence value and the true value up to the observation window  $T = 10$ . The average influence is 16.02, 36.93, 9.7 and 8.3. We use 8,192 random features and two-fold cross validation on the train data to tune the normalization  $Z$ , which has the best value 1130, 1160, 1020, and 1090, respectively. We choose the RBF kernel bandwidth  $h = 1/\sqrt{2\pi}$  so that the magnitude of the smoothed approximate function still equals to 1 (or it can be tuned by cross-validation as well), which matches the original indicator function. For the kernel ridge regression, the RBF kernel bandwidth and the regularization  $\lambda$  are all chosen by the same two-fold cross validation. For CIC and DIC, we learn the respective model up to time  $T$  for once.

Figure 1 verifies that even though the underlying diffusion models can be dramatically different, the prediction performance of TCOVERAGELEARNER is robust to the model changes and consistently outperforms the nontrivial baseline significantly. In addition, even if CIC and DIC are provided with extra information, in Figure 1(a), because the ground-truth is continuous-time diffusion model with Weibull functions, they do not have good performance. CIC assumes the right model but the wrong family of transmission functions. In Figure 1(b), we expect CIC should have the best performance for that it assumes the correct diffusion model and transmission functions. Yet, TCOVERAGELEARNER still has comparable performance with even less information. In Figure 1(c), although DIC has assumed the correct model, it is hard to determine the correct step size to discretize the time line, and since we only learn the model once up to time  $T$  (instead of at each time point), it is harder to fit the whole process. In Figure 1(d), both CIC and DIC have the wrong model, so we have similar trend as Figure synthetic(a). Moreover, for kernel ridge regression, we have to first partition the timeline with arbitrary step size, fit the model to each of time, and interpolate the value between neighboring time legs. Not only will the errors from each stage be accumulated to the error of the final prediction, but also we cannot rely on this method to predict the influence of a source set beyond the observation window  $T$ .

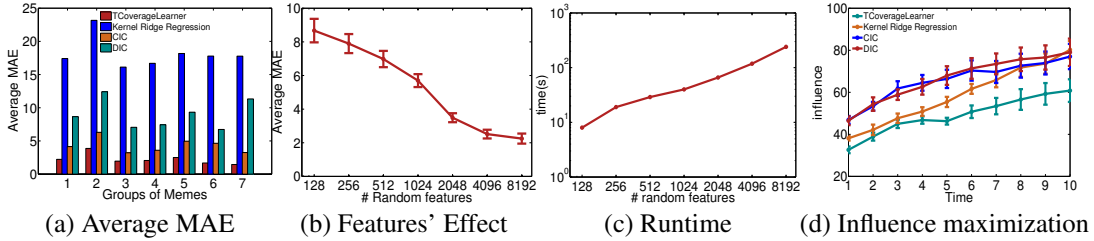


Figure 2: (a) Average MAE from time 1 to 10 on seven groups of real cascade data; (b) Improved estimation with increasing number of random features; (c) Runtime in log-log scale; (d) Maximized influence of selected sources on the held-out testing data along time.

Overall, compared to the kernel ridge regression, `TCOVERAGELEARNER` only needs to be trained once given all the event data up to time  $T$  in a compact and principle way, and then can be used to infer the influence of any given source set at any particular time much more efficiently and accurately. In contrast to the two-stage methods, `TCOVERAGELEARNER` is able to address the more general setting with much less assumption and information but still can produce consistently competitive performance.

### 7.3 Influence Estimation on Real Data

MemeTracker is a real-world dataset [18] to study information diffusion. The temporal flow of information was traced using quotes which are short textual phrases spreading through the websites. We have selected seven typical groups of cascades with the representative keywords like ‘apple and jobs’, ‘tsunami earthquake’, etc., among the top active 1,000 sites. Each set of cascades is split into 60%-train and 40%-test. Because we often can observe cascades only from single seed node, we rarely have cascades produced from multiple sources simultaneously. However, because our model can capture the correlation among multiple sources, we challenge `TCOVERAGELEARNER` with sets of randomly chosen multiple source nodes on the independent hold-out data. Although the generation of sets of multiple source nodes is simulated, the respective influence is calculated from the real test data as follows : Given a source set  $\mathcal{S}$ , for each node  $u \in \mathcal{S}$ , let  $\mathcal{C}(u)$  denote the set of cascades generated from  $u$  on the testing data. We uniformly sample cascades from  $\mathcal{C}(u)$ . The average length of all sampled cascades is treated as the true influence of  $\mathcal{S}$ . We draw 128 source sets and report the average MAE along time in Figure 2(a). Again, we can observe that `TCOVERAGELEARNER` has consistent and robust estimation performance across all testing groups. Figure 2(b) verifies that the prediction can be improved as more random features are exploited, because the representational power of `TCOVERAGELEARNER` increases to better approximate the unknown true coverage function. Figure 2(c) indicates that the runtime of `TCOVERAGELEARNER` is able to scale linearly with large number of random features. Finally, Figure 2(d) shows the application of the learned coverage function to the influence maximization problem along time, which seeks to find a set of source nodes that maximize the expected number of infected nodes by time  $T$ . The classic greedy algorithm[19] is applied to solve the problem, and the influence is calculated and averaged over the seven held-out test data. It shows that `TCOVERAGELEARNER` is very competitive to the two-stage methods with much less assumption. Because the greedy algorithm mainly depends on the relative rank of the selected sources, although the estimated influence value can be different, the selected set of sources could be similar, so the performance gap is not large.

## 8 Conclusions

We propose a new problem of learning temporal coverage functions with a novel parametrization connected with counting processes and develop an efficient algorithm which is guaranteed to learn such a combinatorial function from only polynomial number of training samples. Empirical study also verifies our method outperforms existing methods consistently and significantly.

**Acknowledgments** This work was supported in part by NSF grants CCF-0953192, CCF-1451177, CCF-1101283, and CCF-1422910, ONR grant N00014-09-1-0751, AFOSR grant FA9550-09-1-0538, Raytheon Faculty Fellowship, NSF IIS1116886, NSF/NIH BIGDATA 1R01GM108341, NSF CAREER IIS1350983 and Facebook Graduate Fellowship 2014-2015.



## References

- [1] David Kempe, Jon Kleinberg, and Éva Tardos. Maximizing the spread of influence through a social network. In *SIGKDD 2003*, pages 137–146. ACM, 2003.
- [2] C. Guestrin, A. Krause, and A. Singh. Near-optimal sensor placements in gaussian processes. In *International Conference on Machine Learning ICML'05*, 2005.
- [3] Benny Lehmann, Daniel Lehmann, and Noam Nisan. Combinatorial auctions with decreasing marginal utilities. In *EC '01*, pages 18–28, 2001.
- [4] Maria-Florina Balcan and Nicholas JA Harvey. Learning submodular functions. In *Proceedings of the 43rd annual ACM symposium on Theory of computing*, pages 793–802. ACM, 2011.
- [5] A. Badanidiyuru, S. Dobzinski, H. Fu, R. D. Kleinberg, N. Nisan, and T. Roughgarden. Sketching valuation functions. In *Annual ACM-SIAM Symposium on Discrete Algorithms*, 2012.
- [6] Vitaly Feldman and Pravesh Kothari. Learning coverage functions. *arXiv preprint arXiv:1304.2079*, 2013.
- [7] Vitaly Feldman and Jan Vondrak. Optimal bounds on approximation of submodular and xos functions by juntas. In *FOCS*, 2013.
- [8] Nan Du, Yingyu Liang, Nina Balcan, and Le Song. Influence function learning in information diffusion networks. In *ICML 2014*, 2014.
- [9] L. Song, M. Kolar, and E. P. Xing. Time-varying dynamic bayesian networks. In *Neural Information Processing Systems*, pages 1732–1740, 2009.
- [10] M. Kolar, L. Song, A. Ahmed, and E. P. Xing. Estimating time-varying networks. *Ann. Appl. Statist.*, 4(1):94–123, 2010.
- [11] Odd Aalen, Oernulf Borgan, and Håkon K Gjessing. *Survival and event history analysis: a process point of view*. Springer, 2008.
- [12] M. P. Friedlander K. Murphy M. Schmidt, E. van den Berg. Optimizing costly functions with simple constraints: A limited-memory projected quasi-newton algorithm. In *AISTATS 2009*.
- [13] Sara van de Geer. Exponential inequalities for martingales, with application to maximum likelihood estimation for counting processes. *The Annals of Statistics*, pages 1779–1801, 1995.
- [14] Manuel Gomez Rodriguez, David Balduzzi, and Bernhard Schölkopf. Uncovering the temporal dynamics of diffusion networks. *arXiv preprint arXiv:1105.0697*, 2011.
- [15] Nan Du, Le Song, Hongyuan Zha, and Manuel Gomez Rodriguez. Scalable influence estimation in continuous time diffusion networks. In *NIPS 2013*, 2013.
- [16] Jure Leskovec, Deepayan Chakrabarti, Jon Kleinberg, Christos Faloutsos, and Zoubin Ghahramani. Kronecker graphs: An approach to modeling networks. 11(Feb):985–1042, 2010.
- [17] Praneeth Netrapalli and Sujay Sanghavi. Learning the graph of epidemic cascades. In *SIGMETRICS/PERFORMANCE*, pages 211–222. ACM, 2012.
- [18] Jure Leskovec, Lars Backstrom, and Jon Kleinberg. Meme-tracking and the dynamics of the news cycle. In *SIGKDD2009*, pages 497–506. ACM, 2009.
- [19] G. Nemhauser, L. Wolsey, and M. Fisher. An analysis of the approximations for maximizing submodular set functions. *Mathematical Programming*, 14:265–294, 1978.
- [20] L. Wasserman. *All of Nonparametric Statistics*. Springer, 2006.
- [21] Ali Rahimi and Benjamin Recht. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. In *Neural Information Processing Systems*, 2009.
- [22] Sara van de Geer. Hellinger-consistency of certain nonparametric maximum likelihood estimators. *The Annals of Statistics*, pages 14–44, 1993.
- [23] G.R. Shorack and J.A. Wellner. *Empirical Processes with Applications to Statistics*. Wiley, New York, 1986.
- [24] Wing Hung Wong and Xiaotong Shen. Probability inequalities for likelihood ratios and convergence rates of sieve mles. *The Annals of Statistics*, pages 339–362, 1995.
- [25] L. Birgé and P. Massart. Minimum Contrast Estimators on Sieves: Exponential Bounds and Rates of Convergence. *Bernoulli*, 4(3), 1998.
- [26] Kenneth S Alexander. Rates of growth and sample moduli for weighted empirical processes indexed by sets. *Probability Theory and Related Fields*, 75(3):379–423, 1987.

## A Approximation Error

Recall that we view our data as a marked counting process

$$N_i(t) = f(\mathcal{S}_i, t) + M_i(t).$$

where  $t \in [0, T]$  and  $T$  is the time window,  $\mathcal{S}_i \subseteq \mathcal{V}$  is the marker, and  $M_i(t)$  is a zero mean local martingale.

We make the following assumptions for our analysis of the parametrization and estimation.

- (A1)  $f(\mathcal{S}, t)$  has derivative  $a(\mathcal{S}, t)$  with respect to  $t$ . For any  $\mathcal{S}$ ,  $a(\mathcal{S}, t) = df(\mathcal{S}, t)/dt$  is smooth and bounded on  $[0, T]$ ;  $a(\mathcal{S}, t)$  is smooth and bounded on  $[0, T]$ :  $0 < a_{\min} \leq a \leq a_{\max} < \infty$ , and  $\ddot{a} := d^2a/dt^2$  is absolutely continuous with  $\int \ddot{a}(t)dt < \infty$ .
- (A2) There is a known distribution  $\mathbb{Q}'(\boldsymbol{\tau})$  and a constant  $C$  with  $\mathbb{Q}'(\boldsymbol{\tau})/C \leq \mathbb{Q}(\boldsymbol{\tau}) \leq C\mathbb{Q}'(\boldsymbol{\tau})$ .

Let  $a^K$  denote the convolution of  $a$  with a kernel smoothing function  $K$  with bandwidth  $\sigma$ . More precisely,  $K(t) = \frac{1}{\sigma}k(\frac{t}{\sigma})$  and  $k$  is a kernel with

$$0 \leq k(t) \leq \kappa_{\max}, \quad \int k(t) dt = 1, \quad \int t k(t) dt = 0, \quad \text{and} \quad \sigma_k^2 := \int t^2 k(t) dt < \infty.$$

Let

$$\mathcal{A} = \left\{ a_{\mathbf{w}}^K = \sum_{i=1}^W w_i K(t - t(\mathcal{S}_i, \boldsymbol{\tau}_i)) : \mathbf{w} \geq 0, \frac{Z}{C} \leq \|\mathbf{w}\|_1 \leq ZC \right\}$$

denote our hypothesis class. In the following, we show that there exists  $\tilde{a} \in \mathcal{A}$  that is close to  $a$  when the number of features  $W$  is sufficiently large. We first show that  $a$  is close to  $a^K$  and then show that there exists  $\tilde{a} \in \mathcal{A}$  close to  $a^K$ . The first step follows directly from a classic result in kernel density estimation.

**Lemma 5** (e.g., Theorem 6.28 in [20]). *For any  $\mathcal{S}$  and  $t$ ,  $a^K(\mathcal{S}, t) - a(\mathcal{S}, t) = O(\sigma^4)$ .*

For the second step, we have the following lemma based on the quantitative  $C$  measuring the difference between the true distribution  $\mathbb{Q}$  of the features and the sample distribution  $\mathbb{Q}'$ .

**Lemma 6.** *Let  $\mathbb{P}(\mathcal{S})$  be any distribution of  $\mathcal{S}$ . Suppose  $\boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_W$  are drawn i.i.d. from  $\mathbb{Q}'(\boldsymbol{\tau})$ , and  $W = O\left(\left(\frac{CZ\kappa_{\max}}{\epsilon\sigma}\right)^2 \log \frac{1}{\delta\delta_1}\right)$ . Then with probability at least  $1 - \delta$  over  $\boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_W$ , there exists  $\tilde{a} \in \mathcal{A}$  such that,*

$$\Pr_{\mathcal{S} \sim \mathbb{P}(\mathcal{S})} \left\{ \mathbb{E}_t [\tilde{a}(\mathcal{S}, t) - a^K(\mathcal{S}, t)]^2 \geq \epsilon^2 \right\} \leq \delta_1.$$

*Proof.* Let  $a_i(\mathcal{S}, t) = Z \frac{\mathbb{Q}(\boldsymbol{\tau}_i)}{\mathbb{Q}'(\boldsymbol{\tau}_i)} K(t - t(\mathcal{S}, \boldsymbol{\tau}_i))$  for  $i = 1, \dots, W$ . Then  $\mathbb{E}_{\boldsymbol{\tau}_i \sim \mathbb{Q}'(\boldsymbol{\tau}_i)} [a_i] = a^K$ . Let  $\tilde{a}(\mathcal{S}, t) = \frac{Z}{W} \sum_{i=1}^W \frac{\mathbb{Q}(\boldsymbol{\tau}_i)}{\mathbb{Q}'(\boldsymbol{\tau}_i)} K(t - t(\mathcal{S}, \boldsymbol{\tau}_i))$  be the sample average of these functions. Then  $\tilde{a} \in \mathcal{A}$  since  $\frac{Z}{CW} \leq \frac{Z}{W} \frac{\mathbb{Q}(\boldsymbol{\tau}_i)}{\mathbb{Q}'(\boldsymbol{\tau}_i)} \leq \frac{ZC}{W}$ .

Fix  $\mathcal{S}$ , and consider the Hilbert space with the inner product

$$\langle f, g \rangle = \mathbb{E}_t [f(\mathcal{S}, t)g(\mathcal{S}, t)] = \frac{1}{T} \int_0^T f(\mathcal{S}, t)g(\mathcal{S}, t)dt.$$

We now apply the following lemma, which states that the average of bounded vectors in a Hilbert space concentrates towards its expectation in the Hilbert norm exponentially fast.

**Claim 1** (Lemma 4 in [21]). *Let  $\mathbf{X} = \{x_1, \dots, x_W\}$  be iid random variables in a ball  $\mathcal{A}$  of radius  $M$  centered around the origin in a Hilbert space. Denote their average by  $\bar{\mathbf{X}} = \frac{1}{W} \sum_{i=1}^W x_i$ . Then for any  $\delta > 0$ , with probability  $1 - \delta$ ,*

$$\|\bar{\mathbf{X}} - \mathbb{E}\bar{\mathbf{X}}\| \leq \frac{M}{\sqrt{W}} \left( 1 + \sqrt{2 \log \frac{1}{\delta}} \right).$$

Since  $\|\mathbf{w}\|_1 \leq CZ$  and  $|K| \leq \frac{\kappa_{\max}}{\sigma}$ , the norm  $\|a_i\| \leq \frac{CZ\kappa_{\max}}{\sigma}$ . Then when  $W = O\left(\left(\frac{CZ\kappa_{\max}}{\epsilon\sigma}\right)^2 \log \frac{1}{\delta\delta_1}\right)$ , for any fixed  $\mathcal{S}$  we have

$$\Pr_{\boldsymbol{\tau}} \left[ \mathbb{E}_t [\tilde{a}(\mathcal{S}, t) - a^K(\mathcal{S}, t)]^2 \geq \epsilon^2 \right] \leq \delta\delta_1$$

where  $\Pr_{\boldsymbol{\tau}}$  is over the random sample of  $\boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_W$ . This leads to

$$\Pr_{\mathcal{S} \sim \mathbb{P}(\mathcal{S})} \Pr_{\boldsymbol{\tau}} \left[ \mathbb{E}_t [\tilde{a}(\mathcal{S}, t) - a^K(\mathcal{S}, t)]^2 \geq \epsilon^2 \right] \leq \delta\delta_1$$

Exchanging  $\Pr_{\mathcal{S} \sim \mathbb{P}(\mathcal{S})}$  and  $\Pr_{\boldsymbol{\tau}}$  by Fubini's theorem, and then by Markov's inequality, we have

$$\Pr_{\boldsymbol{\tau}} \left\{ \Pr_{\mathcal{S} \sim \mathbb{P}(\mathcal{S})} \left[ \mathbb{E}_t [\tilde{a}(\mathcal{S}, t) - a^K(\mathcal{S}, t)]^2 \geq \epsilon^2 \right] \geq \delta_1 \right\} \leq \delta$$

This means with probability at least  $1 - \delta$  over the random features, on at least  $1 - \delta_1$  probability mass of the distribution of  $\mathcal{S}$ ,  $\mathbb{E}_t [\tilde{a}(\mathcal{S}, t) - a^K(\mathcal{S}, t)]^2 \leq \epsilon^2$ .  $\square$

Combining the two, we have the following approximation error bound.

**Lemma 2** *Let  $\mathbb{P}(\mathcal{S})$  be any distribution of  $\mathcal{S}$ . Suppose  $\boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_W$  are drawn i.i.d. from  $\mathbb{Q}'(\boldsymbol{\tau})$ , and  $W = O\left(\left(\frac{CZ\kappa_{\max}}{\epsilon\sigma}\right)^2 \log \frac{1}{\delta\delta_1}\right)$ . Then with probability at least  $1 - \delta$  over  $\boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_W$ , there exists  $\tilde{a} \in \mathcal{A}$  such that with probability at least  $1 - \delta_1$  over  $\mathcal{S}$ ,*

$$\mathbb{E}_t [\tilde{a}(\mathcal{S}, t) - a(\mathcal{S}, t)]^2 \leq \epsilon^2 + O(\sigma^4).$$

*Consequently, if  $W = O\left(\left(\frac{CZ\kappa_{\max}}{\epsilon\sigma}\right)^2 \log \frac{a_{\max} + CZ\kappa_{\max}}{\delta\epsilon}\right)$ , with probability at least  $1 - \delta$  over  $\boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_W$ , there exists  $\tilde{a} \in \mathcal{A}$  such that*

$$\mathbb{E}_{\mathcal{S}} \mathbb{E}_t [\tilde{a}(\mathcal{S}, t) - a(\mathcal{S}, t)]^2 = O(\epsilon^2 + \sigma^4).$$

*Proof.* The first statement follows from Lemma 5 and 6. Since  $[\tilde{a}(\mathcal{S}, t) - a(\mathcal{S}, t)]^2 \leq C_1 := (a_{\max} + CZ\kappa_{\max})^2$ , we can set  $\delta_1 = \epsilon^2/C_1$ . Then

$$\mathbb{E}_{\mathcal{S}} \mathbb{E}_t [\tilde{a}(\mathcal{S}, t) - a(\mathcal{S}, t)]^2 \leq (1 - \delta_1)(\epsilon^2 + O(\sigma^4)) + \delta_1 C_1 = O(\epsilon^2 + \sigma^4)$$

which completes the proof.  $\square$

For convenience, let  $\epsilon_a^2 := O(\epsilon^2 + \sigma^4)$  denote the  $\ell_2$  approximation error.

## B Sample Complexity

**Setup** Recall that the true intensity  $a$  is bounded on  $[0, T]$ :

$$0 < a_{\min} \leq a \leq a_{\max} < \infty.$$

The kernel  $K$  is also bounded on  $[0, T]$ :

$$0 < \kappa_{\min} \leq K(t) \leq \kappa_{\max}, \forall t \in [0, T]$$

where  $\kappa_{\min} := \min_{t \in [0, T]} K(t) > 0$  is satisfied for typical kernels, e.g., the Gaussian kernel. Our hypothesis class is

$$\mathcal{A} = \left\{ a_{\mathbf{w}}^K = \sum_{i=1}^W w_i K(t - t(\mathcal{S}_i, \boldsymbol{\tau}_i)) : \mathbf{w} \geq 0, \frac{Z}{C} \leq \|\mathbf{w}\|_1 \leq ZC \right\}$$

and thus  $a_{\mathbf{w}}^K$  is also bounded:  $\forall a' \in \mathcal{A}$ ,

$$0 < a_{\min}^w := \frac{Z\kappa_{\min}}{C} \leq a'(\mathcal{S}, t) \leq a_{\max}^w := CZ\kappa_{\max}, \forall \mathcal{S}, t \in [0, T].$$

With the exception of  $\kappa_{\min}$  and  $a_{\min}^w$  that depend on  $\sigma$ , all other parameters are treated as constants.

We observe  $\mathcal{D}^m = \{(\mathcal{S}_i, N_i(t))\}_{i=1}^m$ , and we want to fit  $a(\mathcal{S}, t)$  by  $a_w^K(\mathcal{S}, t)$  by using maximum likelihood estimation (MLE). The log-likelihood is defined as

$$\ell(\mathcal{D}^m|a') := \sum_{i=1}^m \int_0^T [\log a'(\mathcal{S}_i, t)] dN_i(t) - \sum_{i=1}^m \int_0^T a'(\mathcal{S}_i, t) dt$$

and we optimize the log-likelihood to find an approximate solution.

**Definition 7.** We say that  $\hat{a} \in \mathcal{A}$  is an  $\epsilon_\ell$ -MLE if

$$\ell(\mathcal{D}^m|\hat{a}) \geq \max_{a' \in \mathcal{A}} \ell(\mathcal{D}^m|a') - \epsilon_\ell.$$

**Analysis Roadmap** Our final goal is to bound the  $\ell_2$  error between the truth  $f(t)$  and the function  $\hat{f}(t) = \int_0^t \hat{a}(s) ds$  induced by the MLE output  $\hat{a}$ . A natural choice for connecting  $\ell_2$  error with the log-likelihood cost used in MLE is the Hellinger distance. So it suffices to prove an upper bound on the hellinger distance between the MLE output  $\hat{a}$  and the truth  $a$ , for which we need to show a high probability bound on the empirical Hellinger distance between the two. The key for the analysis is to show that the empirical Hellinger distance can be bounded by a martingale plus some additive error terms. This martingale is defined based on the martingales  $M_i$  associated with the counting process  $N_i$ . The additive error terms are the optimization error and the likelihood gap between the truth and the best one in our hypothesis class. Therefore, our analysis focuses on two parts: a high probability bound for the martingale, and a high probability bound on the likelihood gap.

To bound the martingale, we need to show a uniform convergence inequality. We first introduce a dimension notion measuring the complexity of the hypothesis class, and then prove the uniform convergence based on this notion. Compared to classic uniform inequality for (unmarked) counting process [13], our uniform inequality is for marked counting processes, and the complexity notion and the related conditions have more clear geometric interpretation and are thus easier to verify.

To bound the likelihood gap, we decompose it into three terms, related respectively to the martingale part of the counting processes, the compensate part of the counting processes, and the cumulative difference between the two intensities  $a$  and  $\hat{a}$ . The first term can be bounded by bounding its variance and applying a classic martingale inequality. The second term reduces to the KL-divergence, which can be bounded by the  $\ell_2$  approximation error between the truth and the hypotheses. Similarly, the cumulative difference between the two intensities can be bounded by the  $\ell_2$  approximation error.

We then combine the two to get a bound on the Hellinger distance between the MLE output and the truth based on the dimension of the hypothesis class. This bound is for general hypothesis class, so we bound the dimension of our specific hypothesis class. Finally, we convert the Hellinger distance between the MLE output and the truth to the desired  $\ell_2$  error bound on  $f$  and  $\hat{f}$ .

The rest of the section is organized as follows. We first show the construction of the key martingale upper bound for the Hellinger distance in Section B.1, and then show how to bound the martingale and the likelihood gap in Section B.2 and Section B.3 respectively. In Section B.4 we provide the general bound for the Hellinger distance based on the dimension of the hypothesis class. Finally, in Section B.5 we bound the dimension of our hypothesis class and convert the Hellinger distance to  $\ell_2$  error, achieving the final bound for learning time varying coverage functions.

## B.1 Upper Bound the Hellinger Distance

More precisely, the Hellinger distance is defined as

$$h^2(a, a') = \frac{1}{2} \mathbb{E}_S \mathbb{E}_t \left[ \sqrt{a(\mathcal{S}, t)} - \sqrt{a'(\mathcal{S}, t)} \right]^2$$

where  $\mathbb{E}_S$  is with respect to the random drawing of  $\mathcal{S}$ , and  $\mathbb{E}_t [g(t)]$  denotes  $\frac{1}{T} \int_0^T g(t) dt$ . Define the (total) empirical Hellinger distance as

$$\hat{H}^2(a, a') = \frac{1}{2} \sum_{i=1}^m \int_0^T \left[ \sqrt{a(\mathcal{S}_i, t)} - \sqrt{a'(\mathcal{S}_i, t)} \right]^2 dt$$

and note that  $\mathbb{E}_S \mathbb{E}_t \left[ \hat{H}^2(a, a') \right] = mT h^2(a, a')$ .

Define a martingale

$$M(t|g) := \int_0^t g(t)d \left( \sum_i M_i(t) \right) = \sum_{i=1}^m \int_0^t g(t)dM_i(t) \quad (14)$$

where  $M_i(t)$  is the martingale in the marked counting process  $(\mathcal{S}_i, N_i(t))$ , and  $g \in \mathcal{G}$  where  $\mathcal{G}$  is a set of functions defined as

$$\mathcal{G} = \left\{ g_{a'} = \frac{1}{2} \log \frac{a+a'}{2a} : a' \in \mathcal{A} \right\}.$$

Let  $V_n(t|g)$  denote the  $n$ -th order variation process corresponding to  $M(t|g)$ .

Define two distances on  $\mathcal{G}$ :

$$d_{2,t}^2(g, g') = \frac{1}{2} \sum_{i=1}^m \int_0^t [\exp(g) - \exp(g')]^2 d\Lambda_i(t)$$

where  $\Lambda_i(t) = f(\mathcal{S}_i, t)$  is the compensator of  $N_i(t)$  and

$$d_{\infty,t}(g, g') = \max_{\tau \in [0,t], \mathcal{S}} |\exp(g(\mathcal{S}, \tau)) - \exp(g'(\mathcal{S}, \tau))|.$$

Now we show that  $\widehat{H}(\cdot, \cdot)$  can be bounded by a martingale plus some additive error terms.

**Lemma 3** *Suppose  $\widehat{a}$  is an  $\epsilon_\ell$ -MLE. Then*

$$\begin{aligned} \widehat{H}^2 \left( \frac{\widehat{a} + a}{2}, a \right) &\leq M(T|g_{\widehat{a}}) + \frac{1}{4} \left[ \ell(\mathcal{D}^m|a) - \max_{a' \in \mathcal{A}} \ell(\mathcal{D}^m|a') \right] + \frac{1}{4} \epsilon_\ell, \\ \widehat{H}^2(\widehat{a}, a) &\leq 16M(T|g_{\widehat{a}}) + 4 \left[ \ell(\mathcal{D}^m|a) - \max_{a' \in \mathcal{A}} \ell(\mathcal{D}^m|a') \right] + 4\epsilon_\ell. \end{aligned}$$

*Proof.* This is a generalization of Lemma 4.1 in [13] and the proof largely follows their arguments.

**Claim 2.** *For any  $b \geq 0$ ,*

$$\frac{1}{2} [\ell(\mathcal{D}^m|b) - \ell(\mathcal{D}^m|a)] \leq M \left( T \left| \frac{1}{2} \log \frac{b}{a} \right. \right) - \widehat{H}^2(b, a).$$

*Proof.* Let  $h_b := \frac{1}{2} \log \frac{b}{a}$ .

$$\frac{1}{2} [\ell(\mathcal{D}^m|b) - \ell(\mathcal{D}^m|a)] = M(T|h_b) + \sum_{i=1}^m \int_0^T h_b d\Lambda_i(t) - \frac{1}{2} \sum_{i=1}^m \int_0^T (b-a) dt$$

and

$$\begin{aligned} \sum_{i=1}^m \int_0^T h_b d\Lambda_i(t) - \frac{1}{2} \sum_{i=1}^m \int_0^T (b-a) dt &= \sum_{i=1}^m \int_0^T \log \sqrt{\frac{b}{a}} d\Lambda_i(t) - \frac{1}{2} \sum_{i=1}^m \int_0^T (b-a) dt \\ &\leq \sum_{i=1}^m \int_0^T \left( \sqrt{\frac{b}{a}} - 1 \right) d\Lambda_i(t) - \frac{1}{2} \sum_{i=1}^m \int_0^T (b-a) dt \\ &= \sum_{i=1}^m \int_0^T (\sqrt{ba} - a) dt - \frac{1}{2} \sum_{i=1}^m \int_0^T (b-a) dt \\ &= -\widehat{H}^2(b, a). \end{aligned}$$

This completes the proof.  $\square$

**Claim 3.** *For any  $\widehat{a} \geq 0$ ,*

$$\ell \left( \mathcal{D}^m \left| \frac{\widehat{a} + a}{2} \right. \right) - \ell(\mathcal{D}^m|a) \geq \frac{1}{2} [\ell(\mathcal{D}^m|\widehat{a}) - \ell(\mathcal{D}^m|a)].$$

*Proof.* By the concavity of the log function,

$$\begin{aligned}
\ell\left(\mathcal{D}^m \mid \frac{\hat{a}+a}{2}\right) - \ell(\mathcal{D}^m \mid a) &= \sum_{i=1}^m \int_0^T \log\left(\frac{\hat{a}+a}{2a}\right) dN_i(t) - \sum_{i=1}^m \int_0^T \left(\frac{\hat{a}+a}{2} - a\right) dt \\
&\geq \frac{1}{2} \sum_{i=1}^m \int_0^T \log\left(\frac{\hat{a}}{a}\right) dN_i(t) - \frac{1}{2} \sum_{i=1}^m \int_0^T (\hat{a} - a) dt \\
&= \frac{1}{2} [\ell(\mathcal{D}^m \mid \hat{a}) - \ell(\mathcal{D}^m \mid a)].
\end{aligned} \tag{15}$$

□

We let  $b = \frac{\hat{a}+a}{2}$  in Claim 2 and combine with Claim 3, which then leads to

$$\frac{1}{2} [\ell(\mathcal{D}^m \mid \hat{a}) - \ell(\mathcal{D}^m \mid a)] \leq M \left( T \mid \frac{1}{2} \log \frac{\hat{a}+a}{2a} \right) - \hat{H}^2(b, a).$$

Note that  $\frac{1}{2} \log \frac{\hat{a}+a}{2a}$  is just  $g_{\hat{a}}$ . This, together with the definition of  $\epsilon_\ell$ -MLE, completes the proof for the first statement.

For the second statement, we use the following claim.

**Claim 4** ([22]).  $2\hat{H}^2(\frac{a+b}{2}, a) \leq \hat{H}^2(b, a) \leq 16\hat{H}^2(\frac{a+b}{2}, a)$ .

The second statement then follows from the first statement. □

## B.2 Bounding the Martingale

We begin with some basics about martingales. Here, for a martingale  $M(t)$ , let  $V_n(t)$  denote its  $n$ -th order variation process for  $n \geq 2$ , and let  $V(t) := V_2(t)$ . In particular,

$$V(t) := \lim_{j \rightarrow \infty} \sum_{k=1}^n \text{Var}(\Delta M_k \mid \mathcal{H}_{(k-1)t/j}) \tag{16}$$

where the time interval  $[0, t]$  is partitioned into  $j$  subintervals each of length  $t/j$ , and  $\Delta M_k := M(kt/j) - M((k-1)t/j)$  is the increment of the martingale over the  $k$ th of these intervals. The higher order moments are defined similarly.

Informally, the increment  $dV(t)$  of the predictable variation process can be written as  $dV(t) = \text{Var}(dM(t) \mid \mathcal{H}_{t-}) = \text{Var}(dN(t) \mid \mathcal{H}_{t-})$ , since  $a(t)$  is predictable given  $\mathcal{H}_{t-}$ . Finally,  $dN(t)$  may only take the value 0 or 1, and it follows that  $dV(t) = a(t)dt(1 - a(t)dt) \approx a(t)dt = d\Lambda(t)$ . This motivates the following claim, which will be useful in our later analysis.

**Claim 5** ([11]).  $V(t) = \int_0^t a(s) ds = \Lambda(t)$ .

The following two classic martingale inequalities will also be useful.

**Lemma 8** ([23]). *Suppose that  $|dM(t)| \leq C_M$  for all  $t \geq 0$  and some  $0 \leq C_M < \infty$ , and let  $V(t)$  denote its variation process. Then for each  $x > 0, y > 0$ ,*

$$\Pr [M(t) \geq x \text{ and } V(t) \leq y^2 \text{ for some } t] < \exp \left[ -\frac{x^2}{2(xC_M + y^2)} \right].$$

**Lemma 9** ([13]). *Suppose for all  $t \geq 0$  and some constant  $0 < C_M < \infty$ ,*

$$V_n(t) \leq \frac{n!}{2} C_M^{n-2} R(t), \quad \forall n \geq 2,$$

where  $R(t)$  is a predictable process. Then for each  $x > 0, y > 0$ ,

$$\Pr [M(t) \geq x \text{ and } R(t) \leq y^2 \text{ for some } t] < \exp \left[ -\frac{x^2}{2(xC_M + y^2)} \right].$$

**Uniform Inequality for Marked Counting Processes** Now, we will prove a uniform inequality for the martingale  $M(t|g)$  defined in (14), which is based on the marked counting process and the function  $g \in \mathcal{G}$ . Consider the following complexity notion for  $\mathcal{G}$  based on a covering argument.

**Definition 10.** Suppose  $d$  and  $d'$  are two families of metrics on  $\mathcal{G}$  which are indexed by  $t$ , that is, for any  $t \geq 0$ ,  $d_t$  and  $d'_t$  are two metrics on  $\mathcal{G}$ . The  $(d, d')$ -covering dimension of  $\mathcal{G}$  is the minimum  $D \geq 1$  such that there exist  $c_1 \geq 1$  and  $c_2 \geq 1$  satisfying the following. For each  $\epsilon > 0$  and each ball  $\mathcal{B} \subseteq \mathcal{G}$  with radius  $R \geq \epsilon$ , one can find  $\mathcal{C} \subseteq \mathcal{G}$  with

$$|\mathcal{C}| \leq (c_1 R / \epsilon)^D$$

that is an  $\epsilon$ -covering of  $\mathcal{B}$  for the  $d_t$  metric and a  $(c_2 \epsilon)$ -covering for the  $d'_t$  metric for each  $t \geq 0$ .

Based on this notion we have the following uniform inequality.

**Theorem 11.** Let  $D$  be the  $(d, d')$ -covering dimension of  $\mathcal{G}$ . Suppose for any  $g, g' \in \mathcal{G}$ , any  $n \geq 2$ ,

$$V_n(t|g - g') \leq C_1 \frac{n!}{2} C_2^{n-2} d_t^2(g, g'),$$

and

$$V_n(t|g - g') \leq C_3 \frac{n!}{2} [C_4 d'_t(g, g')]^{n-2} d_t^2(g, g')$$

where  $C_1, C_2, C_3, C_4 > 0$  are some constants. Then there exists a constant  $C_0 > 0$ , such that for any  $g^* \in \mathcal{G}$ , any  $y, z > 0$  and  $x \geq C_0(y+1)(z+D)$ ,

$$\Pr [M(t|g - g^*) \geq x \text{ and } d_t(g^*, g) \leq y \text{ for some } t \text{ and some } g \in \mathcal{G}] \leq \exp[-z].$$

**Corollary 12.** Let  $D, V_n$  as specified in Theorem 11. Then there exists a constant  $C_0 > 0$ , such that for any  $g^* \in \mathcal{G}$ , any  $y > 0$  and  $0 < \delta < 1$ , we have that with probability  $\geq 1 - \delta$ ,

$$M(t|g - g^*) \leq C_0(y+1) \left( D + \log \frac{1}{\delta} \right)$$

for all  $g$  and  $t$  satisfying  $d_t(g^*, g) \leq y$ .

*Proof.* Let  $M(\cdot)$  denote  $M(t|\cdot)$  for short. For each  $k = 0, 1, 2, \dots$ , for the ball  $\mathcal{B}(g^*, y)$  and  $\delta_k = O(2^{-k})y$ , there exists a subset  $\mathcal{C}_k$  of size  $\exp\{O(kD)\}$  that is both a  $\delta_k$ -covering with respect to  $d_t$  and a  $(r\delta_k)$ -covering with respect to  $d'_t$  for some constant  $r > 0$ . Let  $g_k$  denote the one in  $\mathcal{C}_k$  closest to  $g$ . Since  $g = g_0 + \sum_{k=0}^{\infty} (g_{k+1} - g_k)$ , we have

$$\begin{aligned} & \Pr [M(g - g^*) \geq x \text{ and } d_t(g^*, g) \leq y \text{ for some } t \text{ and some } g \in \mathcal{G}] \\ & \leq \sum_{g_0 \in \mathcal{C}_0} \Pr [M(g_0 - g^*) > \eta \text{ and } d_t(g_0, g^*) \leq 2y \text{ for some } t] \\ & \quad + \sum_{k=0}^{\infty} \sum_{g_k, g_{k+1}} \Pr [M(g_k - g_{k+1}) > \eta_k \text{ and } d_t(g_k, g_{k+1}) \leq 2\delta_k \text{ for some } t] \end{aligned}$$

as long as  $\eta + \sum_{k=0}^{\infty} \eta_k \leq x$ .

We have by Lemma 9 that

$$\Pr [M(g_0 - g^*) > \eta \text{ and } d_t(g_0, g^*) \leq 2y \text{ for some } t] \leq \exp \left[ -O \left( \frac{\eta^2}{\eta + y^2} \right) \right].$$

Also, for  $g_k, g_{k+1}$  we have

$$\begin{aligned} V_n(t|g_k - g_{k+1}) & \leq C_3 \frac{n!}{2} d_t^2(g_k, g_{k+1}) [C_4 d'_t(g_k, g_{k+1})]^{n-2} \\ & \leq C_3 \frac{n!}{2} d_t^2(g_k, g_{k+1}) [C_4 d'_t(g_k, g) + C_4 d'_t(g, g_{k+1})]^{n-2} \\ & \leq C_3 \frac{n!}{2} d_t^2(g_k, g_{k+1}) [C_4 r \delta_k + C_4 r \delta_{k+1}]^{n-2} \\ & \leq C_3 \frac{n!}{2} d_t^2(g_k, g_{k+1}) [2C_4 r \delta_k]^{n-2}. \end{aligned}$$

Then by Lemma 9,

$$\Pr [M(g_k - g_{k+1}) > \eta_k \text{ and } d_t(g_k, g_{k+1}) \leq 2\delta_k \text{ for some } t] \leq \exp \left[ -O \left( \frac{\eta_k^2}{\eta_k \delta_k + \delta_k^2} \right) \right].$$

Note that  $\eta = O(y\sqrt{z} + cz)$  ensures  $\frac{\eta^2}{c\eta + y^2} \geq z$ . So we can choose  $\eta = O(y\sqrt{z + D} + z + D)$  and  $\eta_k = O(\delta_k(z + kD))$  so that the final statement holds.

We still need to verify  $\eta + \sum_{k=0}^{\infty} \eta_k \leq x$ . Since  $\eta = O(y\sqrt{z + D} + z + D)$  and  $\eta_k = O(\delta_k(z + kD)) = O(2^{-k}y(z + kD))$ , it suffices to have  $x = O((y + 1)(z + D))$ .  $\square$

### B.3 Bounding the Likelihood Gap

**Lemma 13.** *Suppose there exists an  $\tilde{a} \in \mathcal{A}$  such that with probability at least  $1 - \delta_1$  over  $\mathcal{S}$ ,  $\mathbb{E}_t [a'(\mathcal{S}, t) - a(\mathcal{S}, t)]^2 \leq \epsilon_a^2$ . With probability  $\geq 1 - m\delta_1$  over  $\{\mathcal{S}_i\}_{i=1}^m$ , we have that with probability  $\geq 1 - \delta_2$  over  $\{M_i\}_{i=1}^m$ ,*

$$\ell(\mathcal{D}^m|a) - \ell(\mathcal{D}^m|\tilde{a}) \leq B(\delta_2) := O\left(\sqrt{c_\ell^2 Q \log \frac{1}{\delta_2}} + \log \frac{1}{\delta_2} + Q \log \frac{a_{\max}^w}{a_{\min}^w} + mT\epsilon_a\right)$$

where

$$Q = \frac{mT\epsilon_a^2}{a_{\min} + a_{\min}^w}, \text{ and } c_\ell^2 = \frac{4\left(\sqrt{\frac{a_{\max}^w}{a_{\min}^w}} - 1 - \frac{1}{2} \log\left(\frac{a_{\min}^w}{a_{\max}^w}\right)\right)}{\left(\sqrt{\frac{a_{\min}^w}{a_{\max}^w}} - 1\right)^2}.$$

**Corollary 14.** *Under the condition of Lemma 13,  $\ell(\mathcal{D}^m|a) - \max_{a' \in \mathcal{A}} \ell(\mathcal{D}^m|a') \leq B(\delta_2)$ .*

*Proof.* With probability  $\geq 1 - m\delta_1$ ,  $\mathbb{E}_t [a(\mathcal{S}_i, t) - a(\tilde{\mathcal{S}}_i, t)]^2 \leq \epsilon_a^2$  for all  $\mathcal{S}_i$ . Assume this is true.

$$\begin{aligned} & \ell(\mathcal{D}^m|a) - \ell(\mathcal{D}^m|\tilde{a}) \\ &= \sum_{i=1}^m \left[ \int_0^T (\log a - \log \tilde{a}) dN_i(t) - \int_0^T (a - \tilde{a}) dt \right] \\ &= \sum_{i=1}^m \left[ \underbrace{\int_0^T \log\left(\frac{a}{\tilde{a}}\right) dM_i(t)}_{T_{i1}} + \underbrace{\int_0^T \log\left(\frac{a}{\tilde{a}}\right) d\Lambda_i(t)}_{T_{i2}} - \underbrace{\int_0^T (a - \tilde{a}) dt}_{T_{i3}} \right] \end{aligned}$$

where  $\Lambda_i(t) := f(\mathcal{S}_i, t)$  is the compensator of  $N_i(t)$ . There are three terms under the sum, each of which is bounded in the following.

**Bounding  $T_{i1}$**  The first term  $T_{i1}$  has zero expectation, and its variance is  $\text{Var}(T_{i1}) = \mathbb{E}_M [T_{i1}^2]$ . Then

$$\mathbb{E}_M [T_{i1}^2] = \int_0^T \log^2\left(\frac{\tilde{a}}{a}\right) dV_i(t) = \int_0^T \log^2\left(\frac{\tilde{a}}{a}\right) d\Lambda_i(t) = 4 \int_0^T \left[\frac{1}{2} \log\left(\frac{\tilde{a}}{a}\right)\right]^2 d\Lambda_i(t).$$

We now apply the following claim:

**Claim 6** ([24]). *If  $g \geq -L$  for some constant  $L > 0$ , then*

$$|g|^n \leq \frac{n!}{2} C_L^2 \frac{1}{2} [\exp(g) - 1]^2, \text{ for any } n \geq 2,$$

where  $C_L^2 = \frac{4(e^L - 1 - L)}{(e^{-L} - 1)^2}$ .

Since  $\frac{1}{2} \log\left(\frac{\tilde{a}}{a}\right) \geq \frac{1}{2} \log\left(\frac{a_{\min}^w}{a_{\max}^w}\right)$ , by the above claim we have

$$\left[\frac{1}{2} \log\left(\frac{\tilde{a}}{a}\right)\right]^2 \leq O(c_\ell^2) \left(\sqrt{\frac{\tilde{a}}{a}} - 1\right)^2$$



and thus

$$\begin{aligned}
\mathbb{E}_M [T_{i1}^2] &\leq O(c_\ell^2) \int_0^T a \left( \sqrt{\frac{\tilde{a}}{a}} - 1 \right)^2 dt = O(c_\ell^2) \int_0^T \left( \sqrt{\tilde{a}} - \sqrt{a} \right)^2 dt \\
&\leq O(c_\ell^2 T) \mathbb{E}_t \left( \frac{\tilde{a} - a}{\sqrt{\tilde{a}} + \sqrt{a}} \right)^2 \\
&\leq B_1 := O \left( \frac{c_\ell^2 T \epsilon_a^2}{a_{\min} + a_{\min}^w} \right).
\end{aligned}$$

We have that the variance of  $\sum_i T_{i1}$  is bounded by  $mB_1$  and that  $|\sum_i dM(t|\mathcal{S}_i)| \leq 1$  almost surely. By martingale inequality in Lemma 8,

$$\Pr_M \left[ \sum_i T_{i1} \geq C_1 \left( \sqrt{mB_1 \log \frac{1}{\delta_2}} + \log \frac{1}{\delta_2} \right) \right] \leq \frac{\delta_2}{2}$$

for sufficiently large  $C_1$ .

**Bounding  $T_{i2}$**  Since  $T_{i2}$  is just the KL-divergence between  $a(\mathcal{S}_i, \cdot)$  and  $\tilde{a}(\mathcal{S}_i, \cdot)$ , we can apply the following claim.

**Claim 7** (Eqn (7.6) in Lemm 5 in [25]). *The KL-divergence between  $g(\cdot)$  and  $\tilde{g}(\cdot)$  is at most  $4 + 2 \log \left[ \max_t \left| \frac{g(t)}{\tilde{g}(t)} \right| \right]$  times their Hellinger distance  $\frac{1}{2} \int_0^T (\sqrt{g(t)} - \sqrt{\tilde{g}(t)}) dt$ .*

By this claim, we have

$$\begin{aligned}
\int_0^T \log \left( \frac{a}{\tilde{a}} \right) d\Lambda_i(t) &\leq \left( 4 + 2 \log \left[ \max_t \left| \sqrt{\frac{a(\mathcal{S}_i, t)}{\tilde{a}(\mathcal{S}_i, t)}} \right| \right] \right) \int_0^T \left( \sqrt{a(\mathcal{S}_i, t)} - \sqrt{\tilde{a}(\mathcal{S}_i, t)} \right)^2 dt \\
&\leq \left( 4 + 2 \log \frac{a_{\max}}{a_{\min}^w} \right) \int_0^T \left( \frac{a(\mathcal{S}_i, t) - \tilde{a}(\mathcal{S}_i, t)}{\sqrt{a(\mathcal{S}_i, t)} + \sqrt{\tilde{a}(\mathcal{S}_i, t)}} \right)^2 dt \\
&\leq B_2 := \left( 4 + 2 \log \frac{a_{\max}}{a_{\min}^w} \right) \frac{T \epsilon_a^2}{a_{\min} + a_{\min}^w}.
\end{aligned}$$

**Bounding  $T_{i3}$**  For  $T_{i3}$ , we have

$$|T_{i3}| \leq \int_0^T |a(\mathcal{S}_i, t) - \tilde{a}(\mathcal{S}_i, t)| dt \leq T \sqrt{\mathbb{E}_t |a(\mathcal{S}_i, t) - \tilde{a}(\mathcal{S}_i, t)|^2} = T \epsilon_a =: B_3.$$

Combining the bounds together, we have that  $\ell(\mathcal{D}^m|a) - \max_{a' \in \mathcal{A}} \ell(\mathcal{D}^m|a')$  is bounded by  $O \left( \sqrt{B_1 \log \frac{1}{\delta_2}} + \log \frac{1}{\delta_2} + m(B_2 + B_3) \right)$ .  $\square$

#### B.4 MLE for Marked Counting Processes

Here we apply Theorem 11 to bound the empirical Hellinger distance between an approximate MLE and the truth.

**Theorem 15.** *Let  $D$  be the  $(d_2, d_\infty)$ -covering dimension of  $\mathcal{G}$ , and  $\hat{a}$  be an  $\epsilon_\ell$ -MLE. There exist constants  $C_1, C_2 > 1$  such that for any  $\{\mathcal{S}_i\}_{i=1}^m$ , if  $z \geq C_1 [D + \Delta + \epsilon_\ell]$ , then we have*

$$\Pr_M \left[ \hat{H}^2(\hat{a}, a) \geq z \right] \leq \exp[-z/C_2] + \Pr_M \left[ \ell(\mathcal{D}^m|a) - \max_{a' \in \mathcal{A}} \ell(\mathcal{D}^m|a') > \Delta \right]$$

where  $\Pr_M$  is with respect to the randomness in  $\{M_i\}_{i=1}^m$ .

*Proof.* We first verify the conditions of Theorem 11 is satisfied and then apply it to prove the claim.

Since  $g, g' \in \mathcal{G}$  are lower bounded by  $\frac{1}{2} \log \frac{1}{2}$ , Claim 6 leads to

$$|g - g'|^n \leq C'_1 \frac{n!}{2} \frac{1}{2} [\exp(g) - \exp(g')]^2, \text{ for any } n \geq 2$$

for some constant  $C'_1 > 0$ . Since  $(\mathcal{S}_i, N_i(t))$  are independent, and the counting process  $|dM_i(t)| \leq C'_2 = 1$  for all  $t$  and  $\mathcal{S}$ , then

$$\begin{aligned} V_n(t|g - g') &= \sum_i \int_0^t |g - g'|^n dV_{i,n} \\ &\leq (C'_2)^{n-2} \sum_i \int_0^t |g - g'|^n dV_{i,2} = (C'_2)^{n-2} \sum_i \int_0^t |g - g'|^n d\Lambda_i \\ &\leq C'_1 (C'_2)^{n-2} \frac{n!}{2} d_t^2(g, g') \end{aligned}$$

where  $V_{i,n}$  are the  $n$ -th order variation processes for  $M_i$ , and  $\Lambda_i$  is the compensator of  $M_i$ . This verifies the first condition. For the second condition, by Claim 6 we have

$$|g(\mathcal{S}, t) - g'(\mathcal{S}, t)|^2 \leq C'_1 \frac{2!}{2} \frac{1}{2} [\exp(g(\mathcal{S}, t)) - \exp(g'(\mathcal{S}, t))]^2 \leq C'_4 d_{\infty,t}^2(g, g')$$

where  $(C'_4)^2 = C'_1 \frac{2!}{2} \frac{1}{2}$ . Then

$$|g(\mathcal{S}, t) - g'(\mathcal{S}, t)|^{n-2} = (|g(\mathcal{S}, t) - g'(\mathcal{S}, t)|^2)^{(n-2)/2} \leq [C'_4 d_{\infty,t}(g, g')]^{n-2}$$

and

$$\begin{aligned} V_n(t|g - g') &= \sum_i \int_0^t |g - g'|^n dV_{i,n} \leq (C'_2)^{n-2} \sum_i \int_0^t |g - g'|^2 |g - g'|^{n-2} dV_{i,2} \\ &\leq [C'_2 C'_4 d_{\infty,t}(g, g')]^{n-2} \sum_i \int_0^t |g - g'|^2 dV_{i,2} \\ &\leq [C'_2 C'_4 d_{\infty,t}(g, g')]^{n-2} \sum_i \int_0^t |g - g'|^2 d\Lambda_i \\ &= 2d_{2,t}^2(g, g') [C'_2 C'_4 d_{\infty,t}(g, g')]^{n-2} \leq 2 \frac{n!}{2} d_{2,t}^2(g, g') [C'_2 C'_4 d_{\infty,t}(g, g')]^{n-2}. \end{aligned}$$

We are now ready to apply Theorem 11. The argument is classic, see for example, in [26]. By Lemma 3 and Lemma 4, it suffices to prove

$$\Pr_M \left[ M(T|g_{\hat{a}}) \geq \hat{H}^2 \left( \frac{\hat{a} + a}{2}, a \right) - \Delta \text{ and } \hat{H} \left( \frac{\hat{a} + a}{2}, a \right) > \frac{z}{4} \right] \leq \exp[-O(z)].$$

Let  $\bar{b} := \frac{a+b}{2}$  for  $b \in \mathcal{A}$ . The left hand side of the above inequality is bounded by

$$\begin{aligned} &\Pr_M \left[ M(T|g_b) \geq \hat{H}^2(\bar{b}, a) - \Delta \text{ and } \hat{H}(\bar{b}, a) > \frac{z}{4} \text{ for some } b \right] \\ &\leq \sum_{j=1}^{\infty} \Pr_M \left[ M(T|g_b) \geq \left( 2^{j-1} \frac{z}{4} \right)^2 - \Delta \text{ and } \hat{H}(\bar{b}, a) > 2^j \frac{z}{4} \text{ for some } b \right]. \end{aligned}$$

Denote the  $j$ -th term on the right hand side as  $P_j$ . Note that  $g_a = 0$  and  $M(T|g_b) = M(T|g_b - g_a)$ , and  $\hat{H}(\bar{b}, a) = d_{2,T}^2(g_b, g_a)$ . So we can apply Theorem 11 on  $P_j$ . By setting  $z = \Omega(\max\{D, \Delta\})$  and  $z_j = O(2^j z)$ , we have  $P_j \leq \exp[-z_j]$  and thus  $\sum_{j=1}^{\infty} P_j \leq \exp[-O(z)]$ .  $\square$

## B.5 Sample Complexity of MLE for Learning Time Varying Coverage Functions

To apply Theorem 15 in our case, we need: 1) to bound the dimension of our hypothesis class; 2) to transfer the Hellinger distance to  $\ell_2$  error to get the final bound.

**Lemma 16.** *The  $(d_2, d_{\infty})$ -covering dimension of  $\mathcal{G}$  is at most the number of random features  $W$ .*

*Proof.* Note that  $d_{2,t}$  and  $d_{\infty,t}$  are both nondecreasing with respect to  $t$ . So it suffices to show the existence of a covering of size exponential in  $W$  with respect to both  $d_{2,T}$  and  $d_{\infty,T}$ . In the following, we only consider the time  $T$ , and write  $d_{2,T}$  ( $d_{\infty,T}$  respectively) as  $d_2$  ( $d_\infty$  respectively). Note that

$$d_2^2(g_{a'}, g_{a''}) = \widehat{H}^2 \left( \frac{a' + a}{2}, \frac{a'' + a}{2} \right) \quad \text{and} \quad d_\infty(g_{a'}, g_{a''}) = \max_{t,S} \left| \frac{a' + a}{2a} - \frac{a'' + a}{2a} \right| = \left| \frac{a' - a''}{2a} \right|.$$

Then, the covering dimension of  $\mathcal{G}$  is just the  $(d_2, d_\infty)$ -covering dimension of  $\mathcal{A}$  on which the distances are (overloading notations):

$$d_2^2(a', a'') := d_2^2(g_{a'}, g_{a''}), \quad d_\infty(a', a'') := d_\infty(g_{a'}, g_{a''}).$$

Then we can apply the same argument as Lemma 15 in [8] to show the dimension is at most  $W$ . That is, define a mapping from  $\mathbf{w}$  to  $a_{\mathbf{w}}^K$ , and show that the  $\ell_\infty$  distance of the former is approximately the  $d_2$  distance of the latter, and the  $d_\infty$  distance is bounded by the  $d_2$  distance (up to constant factors).

We will need to introduce the following definition and then prove a claim showing that the  $\ell_\infty$  distance on  $\mathbf{w}$  is approximately the  $d_2$  distance on  $a_{\mathbf{w}}^K$ .

**Definition 17.** Define  $\xi = \min_{\mathbf{w} \neq 0} \frac{\mathbf{w}^\top \mathbf{A} \mathbf{w}}{\mathbf{w}^\top \mathbf{w}}$ , where  $\mathbf{A} = \frac{1}{2T} \sum_{\mathcal{S}} \mathbb{P}(\mathcal{S}) \Phi \Phi^\top$  and

$$\Phi = \int_0^T \phi dt, \quad \text{and} \quad \phi = [K(t - t(\mathcal{S}, \tau_1)), \dots, K(t - t(\mathcal{S}, \tau_W))]^\top.$$

**Claim 8.** For an  $\mathbf{w}, \mathbf{w}'$ ,

$$\sqrt{\frac{\xi}{2T a_{\max}^w}} \|\mathbf{w} - \mathbf{w}'\|_\infty \leq d_2(a_{\mathbf{w}}, a_{\mathbf{w}'}) \leq \frac{W \kappa_{\max}}{2\sqrt{a_{\min}^w}} \|\mathbf{w} - \mathbf{w}'\|_\infty.$$

*Proof.* (1) By definition, we have

$$\begin{aligned} d_2^2(a_{\mathbf{w}}, a_{\mathbf{w}'}) &= \frac{1}{2} \mathbb{E}_{\mathcal{S}} \mathbb{E}_t \left[ \sqrt{\mathbf{w}^\top \phi} - \sqrt{\mathbf{w}'^\top \phi} \right]^2 = \frac{1}{2} \mathbb{E}_{\mathcal{S}} \mathbb{E}_t \left[ \frac{\mathbf{w}^\top \phi - \mathbf{w}'^\top \phi}{\sqrt{\mathbf{w}^\top \phi} + \sqrt{\mathbf{w}'^\top \phi}} \right]^2 \\ &\geq \frac{1}{2a_{\max}^w} \mathbb{E}_{\mathcal{S}} \mathbb{E}_t \left[ \mathbf{w}^\top \phi - \mathbf{w}'^\top \phi \right]^2 \\ &= \frac{1}{2a_{\max}^w T} (\mathbf{w} - \mathbf{w}')^\top \mathbf{A} (\mathbf{w} - \mathbf{w}') \\ &\geq \frac{\xi}{2a_{\max}^w T} \|\mathbf{w} - \mathbf{w}'\|_2^2 \geq \frac{\xi}{2a_{\max}^w T} \|\mathbf{w} - \mathbf{w}'\|_\infty^2. \end{aligned}$$

(2) By definition we have

$$\begin{aligned} d_2^2(a_{\mathbf{w}}, a_{\mathbf{w}'}) &= \frac{1}{2} \mathbb{E}_{\mathcal{S}} \mathbb{E}_t \left[ \sqrt{\mathbf{w}^\top \phi} - \sqrt{\mathbf{w}'^\top \phi} \right]^2 = \frac{1}{2} \mathbb{E}_{\mathcal{S}} \mathbb{E}_t \left[ \frac{\mathbf{w}^\top \phi - \mathbf{w}'^\top \phi}{\sqrt{\mathbf{w}^\top \phi} + \sqrt{\mathbf{w}'^\top \phi}} \right]^2 \\ &\leq \frac{1}{4a_{\min}^w} \mathbb{E}_{\mathcal{S}} \mathbb{E}_t \left[ \mathbf{w}^\top \phi - \mathbf{w}'^\top \phi \right]^2 \\ &\leq \frac{1}{4a_{\min}^w} W^2 \kappa_{\max}^2 \|\mathbf{w} - \mathbf{w}'\|_\infty^2. \end{aligned}$$

□

To bound the dimension, the key is to construct coverings of small sizes. By the above claim, the  $d_2$  metric on  $\mathcal{A}$  approximately corresponds to the  $\ell_\infty$  metric on the set of weights. So based on coverings for the weights with respect to the  $\ell_\infty$  metric, we can construct coverings for  $\mathcal{A}$  with respect to the  $d_2$  metric. We then show that they are also coverings with respect to the  $d_\infty$  metric. The bound on the dimension then follows from the sizes of these coverings.

More precisely, given  $\epsilon > 0$  and a ball  $\mathcal{B} \subseteq \mathcal{A}$  with radius  $R \geq \epsilon$ , we construct an  $\epsilon$ -covering  $\mathcal{C}$  as follows. Define a mapping  $\pi : \mathbf{w} \mapsto a_{\mathbf{w}}$ , and define  $\mathcal{B}^w = \pi^{-1}(\mathcal{B})$ . By Claim 8, the radius of  $\mathcal{B}^w$

is at most  $R^w = \sqrt{\frac{2T a_{\max}^w}{\xi}} R$  (with respect to the  $\ell_\infty$  metric). Now consider finding an  $\epsilon^w$ -covering for  $\mathcal{B}^w$  with respect to the  $\ell_\infty$  metric, where  $\epsilon^w = \left(\frac{W \kappa_{\max}}{2\sqrt{a_{\min}^w}}\right)^{-1} \epsilon$ . Since  $\mathcal{B}^w \subseteq \mathbb{R}^W$ , by taking the grid with length  $\epsilon^w/2$  on each dimension, we can get such a covering  $\mathcal{C}^w$  with

$$|\mathcal{C}^w| \leq \left(\frac{4R^w}{\epsilon^w}\right)^W \leq \left(4\sqrt{\frac{2T a_{\max}^w}{\xi}} \frac{W \kappa_{\max}}{2\sqrt{a_{\min}^w}} \frac{R}{\epsilon}\right)^W.$$

Let  $\mathcal{C} = \pi(\mathcal{C}^w)$ , and for any  $b \in \mathcal{B}$  find  $\tilde{b}$  as follows. Suppose  $\mathbf{w}_b \in \mathcal{B}^w$  satisfies  $\pi(\mathbf{w}_b) = b$  and  $\mathbf{w}_{\tilde{b}}$  is the nearest neighbor of  $\mathbf{w}_b$  in  $\mathcal{C}^w$ , then we set  $\tilde{b} = \pi(\mathbf{w}_{\tilde{b}})$ .

First, we argue that  $\mathcal{C}$  is an  $\epsilon$ -covering w.r.t. the  $d_2$  metric, i.e.,  $d(b, \tilde{b}) < \epsilon$  for any  $b \in \mathcal{B}$ . It follows from Claim 8:

$$d_2(b, \tilde{b}) \leq \frac{W \kappa_{\max}}{2\sqrt{a_{\min}^w}} \|\mathbf{w}_b - \mathbf{w}_{\tilde{b}}\|_\infty < \frac{W \kappa_{\max}}{2\sqrt{a_{\min}^w}} \epsilon^w = \epsilon.$$

Second, we argue that  $\mathcal{C}$  is also an  $O(\epsilon)$ -covering w.r.t. the  $d_\infty$  metric, i.e.,  $d_\infty(b, \tilde{b}) = O(\epsilon)$  for any  $b \in \mathcal{B}$ .

$$\begin{aligned} d_\infty(\pi(\mathbf{w}_b), \pi(\mathbf{w}_{\tilde{b}})) &= \max_{t, \mathcal{S}} \left| \sqrt{\frac{b+a}{2a}} - \sqrt{\frac{\tilde{b}+a}{2a}} \right| \\ &= \max_{t, \mathcal{S}} \left| \frac{|b - \tilde{b}|}{\sqrt{2a} (\sqrt{b+a} + \sqrt{\tilde{b}+a})} \right| \\ &\leq \frac{\max_{t, \mathcal{S}} |(\mathbf{w}_{\tilde{b}} - \mathbf{w}_b)^\top \phi|}{2\sqrt{2a_{\min}} (a_{\min}^w + a_{\min})} \\ &\leq \frac{W \kappa_{\max}}{2\sqrt{2a_{\min}} (a_{\min}^w + a_{\min})} \|\mathbf{w}_b - \mathbf{w}_{\tilde{b}}\|_\infty. \end{aligned}$$

So the conditions in the definition of the dimension are satisfied with  $D = W$ ,  $c_1 = 4\sqrt{\frac{2T a_{\max}^w}{\xi}} \frac{W \kappa_{\max}}{2\sqrt{a_{\min}^w}}$  and  $c_2 = \frac{W \kappa_{\max}}{2\sqrt{2a_{\min}} (a_{\min}^w + a_{\min})}$ , and thus the dimension of  $\mathcal{A}$  is at most  $W$ .  $\square$

Now, we can plug the lemma into Theorem 15, and convert the Hellinger distance to the  $\ell_2$  distance between  $f$  and our output function  $\hat{f}$  defined by  $\hat{a}$ .

**Theorem 18.** *Suppose  $\hat{a}$  is an  $\epsilon_\ell$ -MLE, and  $\hat{f}$  is the corresponding function.*

(i) *Suppose there exists an  $\tilde{a} \in \mathcal{A}$  such that with probability at least  $1 - \delta_1$  over  $\mathcal{S}$ ,  $\mathbb{E}_t (a' - a)^2 \leq \epsilon_a^2$ . Then for any  $0 \leq t \leq T$ , and  $\nu > 0$ ,*

$$\begin{aligned} &\mathbb{E}_{\mathcal{S}} \left[ \hat{f}(\mathcal{S}, t) - f(\mathcal{S}, t) \right]^2 \\ &\leq O \left( t^2 \left\{ \nu A_{\max}^2 + \frac{A_{\max}}{mT} \left[ W + \log \frac{1}{\nu} + \epsilon_\ell \right] + A_{\max} \left[ \epsilon_a + \frac{\epsilon_a^2}{A_{\min}} \log \left( \frac{a_{\max}}{a_{\min}^w} \right) + \sqrt{\frac{c_\ell^2 \epsilon_a^2}{A_{\min} mT}} \log \frac{1}{\nu} \right] \right\} \right) \end{aligned}$$

where  $A_{\max} = a_{\max} + a_{\max}^w$ ,  $A_{\min} = a_{\min} + a_{\min}^w$ , and  $c_\ell^2$  is defined in Lemma 13.

(ii) *Consequently, if*

$$W = O \left( (CZ \kappa_{\max})^2 \left[ \left( \frac{A_{\max} T}{\epsilon} \right)^{5/2} + \left( \frac{A_{\max} T \log \frac{a_{\max}}{a_{\min}^w}}{\epsilon A_{\min}} \right)^{5/4} \right] \log \frac{m A_{\max} T}{\epsilon \delta} \right)$$

and

$$m = O \left( \frac{A_{\max} T}{\epsilon} \left[ W + \log \frac{A_{\max} T}{\epsilon} + \epsilon_\ell \right] + \frac{1}{A_{\min} \sqrt{a_{\min}^w} T} \log \frac{A_{\max} T}{\epsilon} \right).$$

then with probability  $\geq 1 - \delta$  over  $\{\tau_i\}_{i=1}^W$ , for any  $0 \leq t \leq T$ ,

$$\mathbb{E}_{\mathcal{S}} \left[ \widehat{f}(\mathcal{S}, t) - f(\mathcal{S}, t) \right]^2 \leq \epsilon.$$

*Proof.* (i) By Theorem 15 and Lemma 13, there exists  $\Omega_{\mathcal{S}}$  of probability at least  $1 - m\delta_1$  so that for any outcome of  $\{\mathcal{S}_i\}_{i=1}^m$  in it, we have that with probability  $\geq 1 - 2\delta_2$ ,

$$\widehat{H}^2(\widehat{a}, a) \leq z = O\left(D + B(\delta_2) + \epsilon_{\ell} + \log \frac{1}{\delta_2}\right)$$

where  $D \leq W$  by Lemma 16.

Since  $\widehat{H}^2(\widehat{a}, a) \leq mT(a_{\max} + a_{\max}^w)$  and  $\mathbb{E}_{\mathcal{D}^m} [\widehat{H}^2(\widehat{a}, a)] = mTh^2(\widehat{a}, a)$ , we have

$$h^2(\widehat{a}, a) \leq \epsilon^2(\delta_1, \delta_2) := (1 - m\delta_1)(1 - 2\delta_2) \frac{z}{mT} + (m\delta_1 + 2\delta_2)(a_{\max} + a_{\max}^w).$$

Now we convert the Hellinger distance between the intensities to the  $\ell_2$  distance between the function  $f$  and the output  $\widehat{f}$  defined by  $\widehat{a}$ . For any  $0 \leq \tau \leq T$ ,

$$\begin{aligned} \mathbb{E}_{\mathcal{S}} \left[ \widehat{f}(\mathcal{S}, \tau) - f(\mathcal{S}, \tau) \right]^2 &\leq \mathbb{E}_{\mathcal{S}} \left[ \int_0^{\tau} |\widehat{a}(\mathcal{S}, t) - a(\mathcal{S}, t)| dt \right]^2 \\ &\leq \tau \mathbb{E}_{\mathcal{S}} \int_0^{\tau} [\widehat{a}(\mathcal{S}, t) - a(\mathcal{S}, t)]^2 dt \\ &\leq \tau \mathbb{E}_{\mathcal{S}} \int_0^{\tau} \left[ \left( \sqrt{\widehat{a}(\mathcal{S}, t)} - \sqrt{a(\mathcal{S}, t)} \right) \left( \sqrt{\widehat{a}(\mathcal{S}, t)} + \sqrt{a(\mathcal{S}, t)} \right) \right]^2 dt \\ &\leq 2(a_{\max} + a_{\max}^w) \tau \mathbb{E}_{\mathcal{S}} \int_0^{\tau} \left[ \sqrt{\widehat{a}(\mathcal{S}, t)} - \sqrt{a(\mathcal{S}, t)} \right]^2 dt \\ &\leq 4(a_{\max} + a_{\max}^w) \tau^2 h^2(\widehat{a}, a) \leq 4(a_{\max} + a_{\max}^w) \tau^2 \epsilon^2(\delta_1, \delta_2). \end{aligned}$$

The first statement then follows from choosing  $\delta_1 = \nu/m$  and  $\delta_2 = \nu$ .

(ii) The second statement follows from the first statement and Lemma 2. More precisely, we check each error term and set the parameters as follows.

- To ensure  $t^2 A_{\max}^2 \nu = O(\epsilon)$ , let  $\nu = O\left(\frac{\epsilon}{A_{\max}^2 T^2}\right)$ .

- To ensure  $\frac{t^2 A_{\max}}{mT} [W + \log \frac{1}{\nu} + \epsilon_{\ell}] = O(\epsilon)$ , let

$$m = O\left(\frac{A_{\max} T}{\epsilon} \left[ W + \log \frac{A_{\max} T}{\epsilon} + \epsilon_{\ell} \right]\right). \quad (17)$$

- To ensure that  $\epsilon_a^2 = O(\epsilon_0^2)$ , let  $\sigma = \sqrt{\epsilon_0}$ , and

$$W = O\left(\left(\frac{CZ\kappa_{\max}}{\epsilon_0 \sigma}\right)^2 \log \frac{1}{\delta_1 \delta}\right) = O\left(\frac{(CZ\kappa_{\max})^2}{\epsilon_0^{5/2}} \log \frac{mA_{\max} T}{\epsilon \delta}\right).$$

- To ensure  $t^2 A_{\max} \epsilon_a = O(\epsilon)$ , we need  $\epsilon_0^2 = O\left(\left(\frac{\epsilon}{A_{\max} T}\right)^2\right)$ . To ensure  $\frac{t^2 A_{\max} \epsilon_a^2}{A_{\min}} \log\left(\frac{a_{\max}}{a_{\min}^w}\right) = O(\epsilon)$ , we need  $\epsilon_0^2 = O\left([A_{\min} \epsilon] / [A_{\max} T^2 \log\left(\frac{a_{\max}}{a_{\min}^w}\right)]\right)$ . Then we need

$$W = O\left((CZ\kappa_{\max})^2 \left[ \left(\frac{A_{\max} T}{\epsilon}\right)^{5/2} + \left(\frac{A_{\max} T \log \frac{a_{\max}}{a_{\min}^w}}{\epsilon A_{\min}}\right)^{5/4} \right] \log \frac{mA_{\max} T}{\epsilon \delta}\right). \quad (18)$$

- To ensure  $t^2 A_{\max} \sqrt{\frac{c_\ell^2 \epsilon_a^2}{A_{\min} m T}} \log \frac{1}{\nu} = O(\epsilon)$ , we need

$$m = O\left(\frac{c_\ell^2}{A_{\min} T} \log \frac{A_{\max} T}{\epsilon}\right) = O\left(\frac{1}{A_{\min} \sqrt{a_{\min}^w} T} \log \frac{A_{\max} T}{\epsilon}\right). \quad (19)$$

The bound for  $W$  and  $m$  then follows from (18) and (17) (19) respectively. The kernel bandwidth  $\sigma$  is chosen such that  $\sigma = \sqrt{\epsilon_0} = O\left(\min\left\{\left(\frac{\epsilon}{A_{\max} T}\right)^{1/2}, [A_{\min} \epsilon]^{1/4} / \left[A_{\max} T^2 \log\left(\frac{a_{\max}}{a_{\min}^w}\right)\right]^{1/4}\right\}\right)$ .  $\square$