

# Electronic Companion to “A Graphical Point Process Framework for Understanding Removal Effects in Multi-Touch Attribution”

In this Electronic Companion, we first list the important notations in Table EC.1. Then, we present the details of the proposed ADMM algorithm in Section EC.1, the assumptions for consistency in Section EC.2, and the proofs of Theorem 3, Theorem 4, and Theorem 5 in Sections EC.3-EC.5.

Notation	Type	Description
$p$	Integer	The number of event types (conversion and touchpoints) in the attribution problem.
$\mathcal{E}$	Set of labels	$\{1, \dots, p\}$ , the set of all event types, used as the node set for the Granger causality graph.
$e$	Label in $\mathcal{E}$	An event type, standing for a node on the Granger causality graph.
$i$	Integer	The index of an event on a path.
$t_i$	Number in $\mathbb{R}_{\geq 0}$	The relative timestamp of the $i$ -th event on a given path.
$e_i$	Label in $\mathcal{E}$	The event label of the $i$ -th event on a given path.
$(t_i, e_i)$	Element in $\mathbb{R}_{\geq 0} \times \mathcal{E}$	A detailed representation of an event.
$i^*$	Integer	The index of a conversion event on a path. $(t_{i^*}, e_{i^*})$ is the conversion event.
$m$	Integer	The number of events on a path.
$D$	Set of events	$\{(t_i, e_i)\}_{i=1}^m$ , a detailed representation of a path.
$T$	Number in $\mathbb{R}_{> 0}$	The time length of observation or the upper bound of $t_i$ .
$n$	Integer	The number of observed paths, i.e. sample size.
$j$	Integer	The index of a path among all the observations. When emphasized, the notations in the above rows are adapted correspondingly as in $D_j = \{(t_i^j, e_i^j)\}_{i=1}^{m_j}$ with length $T_j$ .
$N_e(t)$	Counting function $\mathbb{R}_{\geq 0} \rightarrow \mathbb{N}$	$\sum_{i: e_i=e} \mathbb{1}_{\{t_i \leq t\}}$ , the point process of event type $e$ .
$\mathbf{N}(t)$	Vector of counting functions	$(N_1(t), \dots, N_p(t))^\top$ , the multivariate point process of event type $1, \dots, p$ .
$\mathbf{N}_W(t)$	Vector of counting functions	$(N_e(t))_{e \in W}$ the subprocess of $\mathbf{N}(t)$ for event label subset $W \subseteq \mathcal{E}$ .
$\mathbf{N}_W^D(t)$	Vector of counting functions	$(N_e^D(t))_{e \in W}$ with $N_e^D(t) = \sum_{(t_i, e_i) \in D: e_i=e} \mathbb{1}_{\{t_i \leq t\}}$ , the observed subprocess for path $D$ .
$\lambda_e(t   \mathcal{H}_t)$	Non-negative function	The conditional intensity function for event type $e$ .
$(e' \rightarrow e)$	Element in $\mathcal{E} \times \mathcal{E}$	The directed edge from node $e'$ to node $e$ , representing Granger causality.
$\mathcal{A}$	Set of edges	The set of all edges where the Granger causality relations exist.
$q$	Integer	$1 \leq q < p$ , the number of customer-initiated event types.
$\mathcal{E}_c$	Set of labels	The set of customer-initiated event types, assumed to be $\{1, \dots, q\}$ .
$\mathcal{E}_f$	Set of labels	The set of firm-initiated event types, assumed to be $\{q+1, \dots, p\}$ .
$\mu_e$	Number in $\mathbb{R}_{\geq 0}$	The baseline intensity of event type $e$ .
$\psi_{e'e}(\cdot)$	Non-negative function	The kernel/transfer function describing the shape of the Granger causality impact ( $e' \rightarrow e$ ).
$\alpha_{e'e}$	Number in $\mathbb{R}_{\geq 0}$	The coefficient describing the scale of the Granger causality impact ( $e' \rightarrow e$ ).
$\boldsymbol{\alpha}_e$	Vector in $\mathbb{R}_{\geq 0}^p$	$(\alpha_{1e}, \dots, \alpha_{pe})^\top$ , the vector of Granger causality coefficients for node $e$ .
$A$	Matrix in $\mathbb{R}_{\geq 0}^{p \times q}$	$(\alpha_{e'e})_{e' \in \mathcal{E}, e \in \mathcal{E}_c}$ , the Granger causality coefficient matrix.
$R$	Set of observed events	A touchpoint event subset of a converted path $D$ .
$F_t^W(D)$	Set of observed events	$\{(t_i, e_i) \in D : t_i < t, e_i \in W\}$ , a subset of $D$ filtered by timestamp $t$ and event type set $W$ .
$\lambda_1(t_{i^*}   \mathcal{H}_{t_{i^*}}^D)$	Number in $\mathbb{R}_{\geq 0}$	The conditional intensity value of conversion ( $e=1$ ) at $t=t_{i^*}$ given path $D$ .
$\lambda_1(t_{i^*}   \mathcal{H}_{t_{i^*}}^{D \setminus R})$	Number in $\mathbb{R}_{\geq 0}$	The conditional intensity value of conversion at $t=t_{i^*}$ given path $D \setminus R$ .
$R^\circ$	Random set of observed events	A superset of $R$ obtained from thinning. The set of removed events in $R$ and additional events possibly to lose.
$\gamma_e$	Number in $\mathbb{R}_{> 0}$	The regularization parameters to control the sparsity of $\boldsymbol{\alpha}_e$ .
$\boldsymbol{\theta}$	Vector in $\mathbb{R}_{> 0}^{p+1}$	$(\mu_e, \alpha_{1e}, \dots, \alpha_{pe})^\top$ , the vector of parameters to learn for a given node $e$ .
$S$	Set of indices	$\{0\} \cup \{e' \in \mathcal{E} : \alpha_{e'e} > 0\}$ , the active set of signals corresponding to nonzero values of $\boldsymbol{\theta}$ .
$S^c$	Set of indices	$\{e' \in \mathcal{E} : \alpha_{e'e} = 0\}$ , the set of signals corresponding to zero values of $\boldsymbol{\theta}$ .
$s$	Integer	The cardinality of the active set $S$ .
$r$	Integer	The cardinality of the removal set $R$ .
$\hat{\boldsymbol{\theta}}$	Vector in $\mathbb{R}_{> 0}^{p+1}$	The regularized estimator of $\boldsymbol{\theta}$ .
$z$	Label in $\{1, \dots, Z\}$	The channel label for a touchpoint event type, where $Z$ is the number of channel types.
$\mathcal{C}_z$	Set of labels	The set of touchpoint event types belonging to channel $z$ .
$\mathbf{p}_z$	Number in $[0, 1]$	The proportion of channel-level conversion count (CCC) for channel $z$ .
$\hat{\mathbf{p}}_z$	Number in $[0, 1]$	The proportion of channel-level aggregated score (CAS) for channel $z$ .
$\mathbf{p}$	Vector in $[0, 1]^Z$	$(\mathbf{p}_1, \dots, \mathbf{p}_Z)^\top$ , the vector of proportions of CCC.
$\hat{\mathbf{p}}$	Vector in $[0, 1]^Z$	$(\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_Z)^\top$ , the vector of proportions of CAS.

**Table EC.1** Summary of important notations.

## EC.1. Details of the Proposed ADMM Algorithm

For each node  $e \in \mathcal{E}_c$ , learning its parent nodes yields

$$\min_{\mu_e \geq 0, \boldsymbol{\alpha}_e \in \mathbb{R}_{\geq 0}^p} \frac{1}{n} \sum_{j=1}^n \phi_e(D_j; \mu_e, \boldsymbol{\alpha}_e) + \gamma_e \|\boldsymbol{\alpha}_e\|_1, \quad (\text{EC.1})$$

where

$$\phi_e(D_j; \mu_e, \boldsymbol{\alpha}_e) = \frac{1}{2} \int_0^{T_j} [\lambda_e(t | \mathcal{H}_t^{D_j})]^2 dt - \int_0^{T_j} \lambda_e(t | \mathcal{H}_t^{D_j}) dN_e^{D_j}(t).$$

In the following context, we write  $\boldsymbol{\theta} = (\theta_0, \theta_1, \dots, \theta_p)^\top = (\mu_e, \boldsymbol{\alpha}_e^\top)^\top \in \mathbb{R}_{\geq 0}^{p+1}$ . Let  $\mathbf{X}_j(t) = (X_{j,0}(t), X_{j,1}(t), \dots, X_{j,p}(t))^\top$  with

$$X_{j,0}(t) = 1, \quad X_{j,k}(t) = \int_0^t \psi_{ke}(t-u) dN_k^{D_j}(u), \quad k = 1, \dots, p.$$

Now the conditional intensity can be written as  $\lambda_e(t | \mathcal{H}_t^{D_j}) = \boldsymbol{\theta}^\top \mathbf{X}_j(t)$ . Let  $V = (V_{kk'})_{k,k'=0}^p \in \mathbb{R}^{(p+1) \times (p+1)}$  and  $\mathbf{b} = (b_0, \dots, b_p)^\top \in \mathbb{R}^{p+1}$ , where for  $k = 0, \dots, p$ ,

$$V_{kk'} = \frac{1}{n} \sum_{j=1}^n \int_0^{T_j} X_{j,k}(t) X_{j,k'}(t) dt, \quad b_k = \frac{1}{n} \sum_{j=1}^n \int_0^{T_j} X_{j,k}(t) dN_e^{D_j}(t). \quad (\text{EC.2})$$

Then the regularized solution satisfies

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}_{\geq 0}^{p+1}} \frac{1}{2} \boldsymbol{\theta}^\top V \boldsymbol{\theta} - \mathbf{b}^\top \boldsymbol{\theta} + \gamma_e \|\boldsymbol{\alpha}_e\|_1. \quad (\text{EC.3})$$

The above problem is equivalent to a linearly constrained one

$$\min_{\boldsymbol{\theta} \in \mathbb{R}_{\geq 0}^{p+1}, \boldsymbol{\alpha}_e = \boldsymbol{\alpha}'_e} \frac{1}{2} \boldsymbol{\theta}^\top V \boldsymbol{\theta} - \mathbf{b}^\top \boldsymbol{\theta} + \gamma_e \|\boldsymbol{\alpha}'_e\|_1.$$

For any  $p$ -dimensional vector  $\mathbf{x} = (x_1, \dots, x_p)^\top \in \mathbb{R}^p$ , let  $\|\mathbf{x}\|_2 := \sqrt{\sum_{k=1}^p x_k^2}$  denote its  $L_2$ -norm.

We design an alternating direction method of multipliers (ADMM). The corresponding augmented Lagrangian function is

$$\mathcal{L}_\eta(\boldsymbol{\theta}, \boldsymbol{\alpha}'_e, \boldsymbol{\omega}) = \frac{1}{2} \boldsymbol{\theta}^\top V \boldsymbol{\theta} - \mathbf{b}^\top \boldsymbol{\theta} + \gamma_e \|\boldsymbol{\alpha}'_e\|_1 + \boldsymbol{\omega}^\top (\boldsymbol{\alpha}_e - \boldsymbol{\alpha}'_e) + \frac{1}{2} \eta \|\boldsymbol{\alpha}_e - \boldsymbol{\alpha}'_e\|_2^2,$$

where  $\|\boldsymbol{\alpha}_e - \boldsymbol{\alpha}'_e\|_2^2 = (\boldsymbol{\alpha}_e - \boldsymbol{\alpha}'_e)^\top (\boldsymbol{\alpha}_e - \boldsymbol{\alpha}'_e)$  represents the squared  $L_2$ -distance between  $\boldsymbol{\alpha}_e$  and  $\boldsymbol{\alpha}'_e$ .

The learning algorithm is shown in Algorithm 3. It consists of two steps, pre-computation and optimization, corresponding to (EC.2) and (EC.3) respectively.

**Algorithm 3** Graphical point process learning by ADMM

---

Input: Paths  $D_1, \dots, D_n$  and the regularization parameter  $\gamma_e$ .
Pre-compute:  $V, \mathbf{b}$ .Initialize: Set  $\eta > 0$  and proper initial values of  $\mu_e, \boldsymbol{\alpha}_e, \boldsymbol{\alpha}'_e$ , and  $\boldsymbol{\omega}$ .**while** not converge **do**

$$\begin{pmatrix} \mu_e \\ \boldsymbol{\alpha}_e \end{pmatrix} = \left[ \left( V + \begin{pmatrix} 0 \\ \eta I_p \end{pmatrix} \right)^{-1} \left( \mathbf{b} + \begin{pmatrix} 0 \\ \eta \boldsymbol{\alpha}'_e - \boldsymbol{\omega} \end{pmatrix} \right) \right]_+.$$

$$\boldsymbol{\alpha}'_e = (\boldsymbol{\alpha}_e + \eta^{-1} \boldsymbol{\omega} - \eta^{-1} \gamma_e \mathbb{1}_p)_+.$$

$$\boldsymbol{\omega} = \boldsymbol{\omega} + \eta(\boldsymbol{\alpha}_e - \boldsymbol{\alpha}'_e).$$

**end while**Return:  $\boldsymbol{\theta} = (\mu_e, \boldsymbol{\alpha}_e^\top)^\top$ .

For typical attribution use cases with sufficient observations, the main computational cost of Algorithm 3 is the pre-computation step, which requires traversing all the paths to collect the first-order and second-order statistics. Recall the assumption that the number of events  $m_j$  on each path  $D_j$  satisfies  $m_j < \bar{m}$ ,  $j = 1, \dots, n$ . The time cost for this step is  $O(\bar{m}^2 n)$  for each node.

According to Hong and Luo (2017), ADMM converges linearly in this problem. During each iteration, the computation cost depends on  $(p+1)$ -dimensional matrix operations. Then the complexity of the optimization step is  $O(p^3)$  for each node.

DRE requires the neighborhood of conversion only with  $p+1$  parameters and thus the learning complexity is  $O(\bar{m}^2 n + p^3)$ . As a comparison, TRE needs the information of the entire graph. There are  $q$  unique nodes to be learned, with a total of  $(p+1)q$  parameters. Therefore, the overall graph learning complexity is  $O(\bar{m}^2 n q + p^3 q)$ . It is worth pointing out that these  $q$  nodes can be learned parallelly.

**EC.2. Assumptions**

In this section, we provide the assumptions for establishing the rates of convergence of the proposed estimators and consistency results.

ASSUMPTION EC.1. (*Independence*) *The process of the customer-initiated event types  $\mathbf{N}_{\mathcal{E}_c}^{D_j}(t)$ ,  $j = 1, \dots, n$  are independent and follow model (1).*

This assumption basically says the prospective customers behave independently of the ads, which is generally true due to interactions between them being sparse. The convergence properties of the estimator rely on this independent distribution of the observations. We do not need to assume the external (firm-initiated) events are *I.I.D.* across paths.

ASSUMPTION EC.2. (*Boundedness*) *There exist constants  $\bar{T}$  and  $\bar{X}$  such that  $T_j < \bar{T}$  and  $\sup_{t \in (0, T_j)} \|\mathbf{X}_j(t)\|_\infty < \bar{X}$  a.s. for  $j = 1, \dots, n$ .*

The first condition in this assumption requires a limit to the length of observation. Assuming  $\psi_{ke}(\cdot)$ ,  $k = 1, \dots, p$  are universally bounded, the second condition reduces to the overall frequency of each type of event is bounded.

For a matrix  $M = (M_{ij})_{i,j}$ , the matrix  $L_1$ -norm is defined as  $\|M\|_{1,\infty} := \max_i(\sum_j |M_{ij}|)$ . Let  $G = (G_{kk'})_{k,k'=0}^p = \mathbb{E}V$  be the population version of  $V$ . Assume the sub-matrix  $G_{SS} = (G_{kk'})_{k,k' \in S}$  is non-singular and define  $\kappa := \|G_{SS}^{-1}\|_{1,\infty}$ .

ASSUMPTION EC.3. (*Irrepresentability*) *There exists a constant  $\xi \in (0, 1)$  such that*

$$\|G_{S^c S} G_{SS}^{-1}\|_{1,\infty} < 1 - \xi.$$

This condition is similar to Condition 3 for sparse survival models in Lin and Lv (2013). It is a generalization of the *irrepresentable* condition (Zhao and Yu 2006) that is almost necessary and sufficient for the LASSO to achieve the model selection consistency. For example, Wainwright (2009) used their Condition (15) to analyze the recovery of sparsity pattern using the LASSO.

Below we repeat Theorem 3 with constants given in detail.

THEOREM 3. *Under Assumptions EC.1-EC.3, there exist  $C_3 > 0$  and  $C_4 > 0$  such that the regularized estimator  $\hat{\boldsymbol{\theta}}$  satisfies the following properties:*

(i) *If  $p$  is fixed, then for any constant  $0 < \nu < 1$ , when*

$$\begin{cases} n > \max(C_3^{-1}, 8\xi^{-2}\kappa^2\bar{X}^4\bar{T}^2 s^2)^{\frac{1}{1-\nu}} \\ \gamma_e = 4\xi^{-1}C_3^{-\frac{1}{2}}C_4 n^{-\frac{1-\nu}{2}} \end{cases},$$

*with probability at least  $1 - 2(p+1)(p+2)\exp(-n^\nu)$ ,*

- (*Edge selection*)  $\hat{\boldsymbol{\theta}}_{S^c} = 0$
- ( *$L_\infty$ -error*)  $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_\infty \leq 10\xi^{-1}\kappa C_3^{-\frac{1}{2}}C_4 n^{-\frac{1-\nu}{2}};$

(ii) *If  $p$  diverges, then for any constant  $\zeta > 2$ , when*

$$\begin{cases} n > \max(C_3^{-1}, 8\xi^{-2}\kappa^2\bar{X}^4\bar{T}^2 s^2) \cdot \zeta \log(p+1) \\ \gamma_e = 4\xi^{-1}C_4 \sqrt{\frac{\zeta \log(p+1)}{C_3 n}} \end{cases},$$

*with probability at least  $1 - 3/(p+1)^{\zeta-2}$ ,*

- (*Edge selection*)  $\hat{\boldsymbol{\theta}}_{S^c} = 0$
- ( *$L_\infty$ -error*)  $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_\infty \leq 10\xi^{-1}\kappa C_3^{-\frac{1}{2}}C_4 \sqrt{\frac{\zeta \log(p+1)}{n}}.$

### EC.3. Proof of Theorem 3

In this proof, let  $\boldsymbol{\theta}^*$  denote the ground truth of  $\boldsymbol{\theta}$  to distinguish from the variable  $\boldsymbol{\theta}$  in the optimization problem. In the sequel, we first prove part (i) and then prove part (ii) of Theorem 3.

**EC.3.1. Proof of (i) of Theorem 3.**

*Proof of Part (i) of Theorem 3.* Let  $U(\boldsymbol{\theta})$  be the negative derivative of the loss function, that is

$$U(\boldsymbol{\theta}) = \mathbf{b} - V\boldsymbol{\theta} = \frac{1}{n} \sum_{j=1}^n \int_0^{T_j} \mathbf{X}_j(t) \left( dN_e^{D_j}(t) - \boldsymbol{\theta}^\top \mathbf{X}_j(t) dt \right).$$

For any  $K$ -dimensional vector  $\mathbf{x} = (x_1, \dots, x_K)^\top \in \mathbb{R}^K$ , where  $K \geq 1$  is an arbitrary integer, let  $\text{sign}(\mathbf{x})$  denote the vector  $(\text{sign}(x_1), \dots, \text{sign}(x_K))^\top$ . Then the vector  $\hat{\boldsymbol{\theta}} = (\hat{\mu}_e, \hat{\boldsymbol{\alpha}}_e) \in \mathbb{R}_{\geq 0}^{p+1}$  is a strict local minimizer of problem (EC.3), if the following conditions hold:

$$\begin{cases} U_{\hat{S}}(\hat{\boldsymbol{\theta}}) - \gamma_e \boldsymbol{\delta}_{\hat{S}} \circ \text{sign}(\hat{\boldsymbol{\theta}}_{\hat{S}}) = 0 \\ \|U_{\hat{S}^c}(\hat{\boldsymbol{\theta}})\|_\infty < \gamma_e \end{cases}, \quad (\text{EC.4})$$

where  $\boldsymbol{\delta} = (0, \mathbf{1}_p)^\top$  and  $\circ$  denotes the Hadamard product.

We study the oracle solution:

$$\tilde{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}_{\geq 0}^{p+1}, \boldsymbol{\theta}_{S^c} = 0} \frac{1}{2} \boldsymbol{\theta}^\top V \boldsymbol{\theta} - \mathbf{b}^\top \boldsymbol{\theta} + \gamma_e \|\boldsymbol{\theta}_{S \setminus \{0\}}\|_1. \quad (\text{EC.5})$$

Note that  $\tilde{\boldsymbol{\theta}}$  will be a solution to the original regularized problem (EC.3) as long as  $\|U_{S^c}(\tilde{\boldsymbol{\theta}})\|_\infty < \gamma_e$  is satisfied.

For  $\tilde{\boldsymbol{\theta}}$  and the true  $\boldsymbol{\theta}^*$ , we have  $U_S(\tilde{\boldsymbol{\theta}}) = \mathbf{b}_S - V_{SS} \tilde{\boldsymbol{\theta}}_S$  and  $U_S(\boldsymbol{\theta}^*) = \mathbf{b}_S - V_{SS} \boldsymbol{\theta}_S^*$ , and thus

$$U_S(\tilde{\boldsymbol{\theta}}) - U_S(\boldsymbol{\theta}^*) = -V_{SS}(\tilde{\boldsymbol{\theta}}_S - \boldsymbol{\theta}_S^*).$$

Then by (EC.4),

$$\tilde{\boldsymbol{\theta}}_S - \boldsymbol{\theta}_S^* = V_{SS}^{-1} [U_S(\boldsymbol{\theta}^*) - \gamma_e \boldsymbol{\delta}_S \circ \text{sign}(\tilde{\boldsymbol{\theta}}_S)]. \quad (\text{EC.6})$$

Need to guarantee  $V_{SS}$  is strictly positive definite and thus invertible, which requires that the least eigenvalue of  $V_{SS}$ , denoted by  $\Lambda_{\min}(V_{SS})$ , is positive. Since  $0 \leq \int_0^{T_j} X_{j,k}(t) X_{j,k'}(t) dt \leq \overline{X^2 T}$  is a bounded random variable, by Hoeffding's inequality, for any  $\epsilon_1 > 0$ ,

$$\mathbb{P}(|V_{kk'} - G_{kk'}| > \epsilon_1) \leq 2 \exp\left(-\frac{2n\epsilon_1^2}{X^4 T^2}\right).$$

Now, we use the union bound to obtain that

$$\mathbb{P}(\|V - G\|_\infty > \epsilon_1) \leq 2(p+1)^2 \exp\left(-\frac{2n\epsilon_1^2}{X^4 T^2}\right).$$

Let  $\epsilon_1 = \frac{1}{\sqrt{2}} \overline{X^2 T} n^{-\frac{1-\nu}{2}}$ , then with probability at least  $1 - 2(p+1)^2 \exp(-n^\nu)$ , where  $0 < \nu < 1$  is a constant,

$$\|V - G\|_\infty \leq \frac{1}{\sqrt{2}} \overline{X^2 T} n^{-\frac{1-\nu}{2}}.$$

We let  $\epsilon_1 < \frac{\xi}{4s\kappa}$ , which suggests the choice of  $n$ :

$$n > (8\xi^{-2}\kappa^2\bar{X}^4\bar{T}^2s^2)^{\frac{1}{1-\nu}}. \quad (\text{EC.7})$$

Let  $\|\cdot\|_{\text{op}}$  denote matrix operator norm. Given the fact that

$$\begin{aligned} \Lambda_{\min}(G_{SS})^{-1} &= \|G_{SS}^{-1}\|_{\text{op}} \leq \|G_{SS}^{-1}\|_{1,\infty} = \kappa, \\ \|V_{SS} - G_{SS}\|_{\text{op}} &\leq \|V_{SS} - G_{SS}\|_{1,\infty} \leq s\epsilon_1 < \frac{\xi}{4\kappa}, \end{aligned}$$

we have

$$\Lambda_{\min}(V_{SS}) \geq \Lambda_{\min}(G_{SS}) - \|V_{SS} - G_{SS}\|_{\text{op}} \geq (1 - \frac{\xi}{4})\kappa^{-1} > 0.$$

Then we give a bound of  $\|U(\boldsymbol{\theta}^*)\|_{\infty}$ , the gradient of the loss function at  $\boldsymbol{\theta}^*$ . By the definition of  $U(\boldsymbol{\theta})$ ,

$$U(\boldsymbol{\theta}^*) = \frac{1}{n} \sum_{j=1}^n \int_0^{T_j} \mathbf{X}_j(t) \left( dN_e^{D_j}(t) - (\boldsymbol{\theta}^*)^\top \mathbf{X}_j(t) dt \right).$$

Note that  $N_e^{D_j}(t) - \int_0^t (\boldsymbol{\theta}^*)^\top \mathbf{X}_j(u) du$  is a counting process martingale with respect to  $\mathcal{H}_t$ . On the other hand, the component

$$u_{j,k}(\boldsymbol{\theta}^*) := \int_0^{T_j} X_{j,k}(t) \left( dN_e^{D_j}(t) - (\boldsymbol{\theta}^*)^\top \mathbf{X}_j(t) dt \right)$$

for  $k = 0, \dots, p$  is a martingale as well. By Theorem 3 in Hansen et al. (2015), the random variable  $u_{j,k}(\boldsymbol{\theta}^*)$  is sub-exponential. For any  $\epsilon_2 > 0$ ,

$$\mathbb{P}(|u_{j,k}(\boldsymbol{\theta}^*)| \geq \sqrt{2\|\boldsymbol{\theta}^*\|_1 \bar{X}^2 \bar{T} \epsilon_2} + \frac{1}{3}\epsilon_2) \leq 2 \exp\left(-\frac{\epsilon_2}{\bar{X}}\right).$$

By  $U_k(\boldsymbol{\theta}^*) = \frac{1}{n} \sum_{j=1}^n u_{j,k}(\boldsymbol{\theta}^*)$  and Bernstein's inequality (Vershynin 2018), there exist constant  $C_3 > 0$  and  $C_4 > 0$  depending on  $\bar{X}$  and  $\bar{T}$  such that

$$\mathbb{P}(|U_k(\boldsymbol{\theta}^*)| \geq \epsilon_2) \leq 2 \exp\left(-C_3 n \min\left(\frac{\epsilon_2^2}{C_4^2}, \frac{\epsilon_2}{C_4}\right)\right).$$

And thus

$$\mathbb{P}(\|U(\boldsymbol{\theta}^*)\|_{\infty} \geq \epsilon_2) \leq 2(p+1) \exp\left(-C_3 n \min\left(\frac{\epsilon_2^2}{C_4^2}, \frac{\epsilon_2}{C_4}\right)\right).$$

When  $\epsilon_2 < C_4$ , with probability at least  $1 - 2(p+1) \exp(-n^\nu)$ ,

$$\|U(\boldsymbol{\theta}^*)\|_{\infty} \leq C_3^{-\frac{1}{2}} C_4 n^{-\frac{1-\nu}{2}}. \quad (\text{EC.8})$$

Now we check the strict dual feasibility condition (EC.4). According to (EC.6),

$$\begin{aligned} U_{S^c}(\tilde{\boldsymbol{\theta}}) &= U_{S^c}(\boldsymbol{\theta}^*) - V_{S^c S}(\tilde{\boldsymbol{\theta}}_S - \boldsymbol{\theta}_S^*) \\ &= U_{S^c}(\boldsymbol{\theta}^*) - V_{S^c S} V_{SS}^{-1} [U_S(\boldsymbol{\theta}^*) - \gamma_e \boldsymbol{\delta}_S \circ \text{sign}(\tilde{\boldsymbol{\theta}}_S)]. \end{aligned}$$

If  $\|U(\boldsymbol{\theta}^*)\|_\infty \leq \epsilon_2$ , then

$$\begin{aligned} \|U_{S^c}(\tilde{\boldsymbol{\theta}})\|_\infty &\leq \|U_{S^c}(\boldsymbol{\theta}^*)\|_\infty + \|V_{S^c S} V_{SS}^{-1}\|_{1,\infty} \|U_S(\boldsymbol{\theta}^*) - \gamma_e \boldsymbol{\delta}_S \circ \text{sign}(\tilde{\boldsymbol{\theta}}_S)\|_\infty \\ &\leq \|U_{S^c}(\boldsymbol{\theta}^*)\|_\infty + \|V_{S^c S} V_{SS}^{-1}\|_{1,\infty} (\|U_S(\boldsymbol{\theta}^*)\|_\infty + \gamma_e) \\ &\leq \epsilon_2 + \|V_{S^c S} V_{SS}^{-1}\|_{1,\infty} (\epsilon_2 + \gamma_e). \end{aligned}$$

Let  $\epsilon_2 \leq \xi \gamma_e / 4 < C_4$ , which suggests the choice of  $n$  and  $\gamma_e$  by (EC.8):

$$\begin{cases} n > C_3^{-\frac{1}{1-\nu}} \\ \gamma_e = 4\xi^{-1} C_3^{-\frac{1}{2}} C_4 n^{-\frac{1-\nu}{2}} \end{cases} \quad (\text{EC.9})$$

Recall the choice of  $n$  in (EC.7), which implies  $\kappa s \epsilon_1 < \xi / 4 < 1/2$ . Then

$$\begin{aligned} \|V_{SS}^{-1} - G_{SS}^{-1}\|_{1,\infty} &= \|G_{SS}^{-1}(G_{SS} - V_{SS})V_{SS}^{-1}\|_{1,\infty} \\ &\leq \|G_{SS}^{-1}\|_{1,\infty} \|G_{SS} - V_{SS}\|_{1,\infty} \|V_{SS}^{-1}\|_{1,\infty} \\ &\leq \|G_{SS}^{-1}\|_{1,\infty} \|G_{SS} - V_{SS}\|_{1,\infty} (\|G_{SS}^{-1}\|_{1,\infty} + \|V_{SS}^{-1} - G_{SS}^{-1}\|_{1,\infty}) \\ &\leq \frac{\|G_{SS}^{-1}\|_{1,\infty}^2 \|G_{SS} - V_{SS}\|_{1,\infty}}{1 - \|G_{SS}^{-1}\|_{1,\infty} \|G_{SS} - V_{SS}\|_{1,\infty}} \\ &\leq \frac{\kappa^2 s \epsilon_1}{1 - \kappa s \epsilon_1} \\ &< \kappa. \end{aligned}$$

Therefore,

$$\|V_{SS}^{-1}\|_{1,\infty} \leq \|G_{SS}^{-1}\|_{1,\infty} + \|V_{SS}^{-1} - G_{SS}^{-1}\|_{1,\infty} < 2\kappa. \quad (\text{EC.10})$$

As a result,

$$\begin{aligned} &\|V_{S^c S} V_{SS}^{-1} - G_{S^c S} G_{SS}^{-1}\|_{1,\infty} \\ &\leq \|(V_{S^c S} - G_{S^c S})V_{SS}^{-1}\|_{1,\infty} + \|G_{S^c S} G_{SS}^{-1}(V_{SS} - G_{SS})V_{SS}^{-1}\|_{1,\infty} \\ &\leq (\|V_{S^c S} - G_{S^c S}\|_{1,\infty} + \|G_{S^c S} G_{SS}^{-1}\|_{1,\infty} \|V_{SS} - G_{SS}\|_{1,\infty}) \|V_{SS}^{-1}\|_{1,\infty} \\ &< 2s \epsilon_1 \kappa. \end{aligned}$$

Use the fact that  $\kappa s \epsilon_1 < \xi / 4$  again, we have

$$\|V_{S^c S} V_{SS}^{-1}\|_{1,\infty} \leq \|V_{S^c S} V_{SS}^{-1} - G_{S^c S} G_{SS}^{-1}\|_{1,\infty} + \|G_{S^c S} G_{SS}^{-1}\|_{1,\infty} < 1 - \xi / 2.$$

Therefore

$$\begin{aligned} \|U_{S^c}(\tilde{\boldsymbol{\theta}})\|_\infty &\leq \epsilon_2 + \|V_{S^c S} V_{SS}^{-1}\|_{1,\infty} (\epsilon_2 + \gamma_e) \\ &< \xi \gamma_e / 4 + (1 - \xi / 2) (\xi \gamma_e / 4 + \gamma_e) \\ &< \gamma_e. \end{aligned}$$

Thus  $\tilde{\boldsymbol{\theta}}$  satisfies the strict dual feasibility condition (EC.4). Therefore, the oracle solution  $\tilde{\boldsymbol{\theta}}$  is a finite solution to the original problem and any solution has a support  $\hat{S}$  as a subset of  $S$ . Thus  $\hat{\boldsymbol{\theta}}$  is also a solution to the oracle problem (EC.5) and we have  $\tilde{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}$  since the objective function of the oracle problem (EC.5) is strictly convex.

Finally, by (EC.6), (EC.7), (EC.8), (EC.9), and (EC.10), if

$$\begin{cases} n > \max(C_3^{-1}, 8\xi^{-2}\kappa^2\bar{X}^4\bar{T}^2 s^2)^{\frac{1}{1-\nu}} \\ \gamma_e = 4\xi^{-1}C_3^{-\frac{1}{2}}C_4 n^{-\frac{1-\nu}{2}} \end{cases},$$

then with probability at least  $1 - 2(p+1)(p+2)\exp(-n^\nu)$ ,

$$\begin{aligned} \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_\infty &= \|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_\infty \\ &= \|\tilde{\boldsymbol{\theta}}_S - \boldsymbol{\theta}^*_S\|_\infty \\ &= \|V_{SS}^{-1}[U_S(\boldsymbol{\theta}^*) - \gamma_e \boldsymbol{\delta}_S \circ \text{sign}(\tilde{\boldsymbol{\theta}}_S)]\|_\infty \\ &\leq \|V_{SS}^{-1}\|_{1,\infty} (\|U_S(\boldsymbol{\theta}^*)\|_\infty + \gamma_e) \\ &\leq 2\kappa(\epsilon_2 + \gamma_e) \\ &\leq 10\xi^{-1}\kappa C_3^{-\frac{1}{2}}C_4 n^{-\frac{1-\nu}{2}}. \end{aligned}$$

Now, we complete the proof of part (i). □

### EC.3.2. Proof of (ii) of Theorem 3.

*Proof of Part (ii) of Theorem 3.* Compared with the fixed  $p$  setting, the case where  $p$  diverges does not have a significant difference except for the tail probability.

Recall the  $L_\infty$ -error of the Hessian matrix obtained by the union bound, which is

$$\mathbb{P}(\|V - G\|_\infty > \epsilon_1) \leq 2(p+1)^2 \exp\left(-\frac{2n\epsilon_1^2}{\bar{X}^4\bar{T}^2}\right).$$

Let  $\epsilon_1 = \bar{X}^2\bar{T} \sqrt{\frac{\zeta \log(p+1)}{2n}}$ , then with probability at least  $1 - 2/(p+1)^{\zeta-2}$ , where  $\zeta > 2$  is a constant,

$$\|V - G\|_\infty \leq \bar{X}^2\bar{T} \sqrt{\frac{\zeta \log(p+1)}{2n}}.$$

We let  $\epsilon_1 < \frac{\xi}{4s\kappa}$ , which requires

$$n > 8\xi^{-2}\zeta\kappa^2\bar{X}^4\bar{T}^2 s^2 \log(p+1). \quad (\text{EC.11})$$

The martingale concentration gives

$$\mathbb{P}(\|U(\boldsymbol{\theta}^*)\|_\infty \geq \epsilon_2) \leq 2(p+1) \exp\left(-C_3 n \min\left(\frac{\epsilon_2^2}{C_4^2}, \frac{\epsilon_2}{C_4}\right)\right).$$

When  $\epsilon_2 < C_4$ , with probability at least  $1 - 2/(p+1)^{\zeta-1}$ , where  $\zeta > 1$  is a constant,

$$\|U(\boldsymbol{\theta}^*)\|_\infty \leq C_4 \sqrt{\frac{\zeta \log(p+1)}{C_3 n}}. \quad (\text{EC.12})$$

For the strict dual feasibility condition (EC.4), we let  $\epsilon_2 \leq \xi \gamma_e / 4 < C_4$ , which suggests the choice of  $n$  and  $\gamma_e$  by (EC.12):

$$\begin{cases} n > C_3^{-1} \zeta \log(p+1) \\ \gamma_e = 4\xi^{-1} C_4 \sqrt{\frac{\zeta \log(p+1)}{C_3 n}} \end{cases}. \quad (\text{EC.13})$$

Finally, combining (EC.6), (EC.10), (EC.11), (EC.12), (EC.13), if

$$\begin{cases} n > \max(C_3^{-1}, 8\xi^{-2} \kappa^2 \bar{X}^4 \bar{T}^2 s^2) \cdot \zeta \log(p+1) \\ \gamma_e = 4\xi^{-1} C_4 \sqrt{\frac{\zeta \log(p+1)}{C_3 n}} \end{cases},$$

then with probability at least  $1 - 3/(p+1)^{\zeta-2}$ ,

$$\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_\infty \leq \|V_{SS}^{-1}\|_{1,\infty} (\|U_S(\boldsymbol{\theta}^*)\|_\infty + \gamma_e) \leq 10\xi^{-1} \kappa C_4 \sqrt{\frac{\zeta \log(p+1)}{C_3 n}}.$$

Now, we have completed the proof of part (ii). □

Therefore, the proof of Theorem 3 is complete.

#### EC.4. Proof of Theorem 4

*Proof of Theorem 4.* By the triangle inequality,

$$\begin{aligned} & \left| \widehat{\text{att}}_{t_i^*}^{(\text{direct})}(R | D) - \text{att}_{t_i^*}^{(\text{direct})}(R | D) \right| \\ &= \left| \sum_{(t_i, e_i) \in R} \widehat{\text{att}}_{t_i^*}^{(\text{direct})}(\{(t_i, e_i)\} | D) - \sum_{(t_i, e_i) \in R} \text{att}_{t_i^*}^{(\text{direct})}(\{(t_i, e_i)\} | D) \right| \\ &\leq \sum_{(t_i, e_i) \in R} \left| \widehat{\text{att}}_{t_i^*}^{(\text{direct})}(\{(t_i, e_i)\} | D) - \text{att}_{t_i^*}^{(\text{direct})}(\{(t_i, e_i)\} | D) \right|. \end{aligned} \quad (\text{EC.14})$$

So it is sufficient to show the theorem is true for each touchpoint. Let  $\hat{\lambda}_1(t | \mathcal{H}_t^D)$  denote the estimated conditional intensity of conversion:

$$\hat{\lambda}_1(t | \mathcal{H}_t^D) = \hat{\mu}_1 + \sum_{i: t_i < t} \hat{\alpha}_{e_i 1} \psi_{e_i 1}(t - t_i).$$

Without loss of generality, assume  $\bar{\psi}_1 \geq 1$ . Then for the conversion intensity,

$$\begin{aligned} & \left| \hat{\lambda}_1(t_i^* | \mathcal{H}_{t_i^*}^D) - \lambda_1(t_i^* | \mathcal{H}_{t_i^*}^D) \right| \\ &= \left| \left( \hat{\mu}_1 + \sum_{i < i^*} \hat{\alpha}_{e_i 1} \psi_{e_i 1}(t_i^* - t_i) \right) - \left( \mu_1 + \sum_{i < i^*} \alpha_{e_i 1} \psi_{e_i 1}(t_i^* - t_i) \right) \right| \\ &\leq |\hat{\mu}_1 - \mu_1| + \sum_{i < i^*} |\hat{\alpha}_{e_i 1} - \alpha_{e_i 1}| |\psi_{e_i 1}(t_i^* - t_i)| \\ &\leq m \bar{\psi}_1 \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_\infty, \end{aligned}$$

where  $\boldsymbol{\theta} = (\mu_1, \boldsymbol{\alpha}_1^\top)^\top$ . By Theorem 3, we can take  $n$  to be sufficiently large such that  $m\bar{\psi}_1 \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_\infty < \lambda_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)/2$ . As a result, we have  $\hat{\lambda}_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D) > \lambda_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)/2$ . Then the estimated direct removal effect of  $\{(t_i, e_i)\}$  given  $D$  satisfies

$$\begin{aligned}
& \left| \widehat{\text{att}}_{t_{i^*}}^{(\text{direct})}(\{(t_i, e_i)\} | D) - \text{att}_{t_{i^*}}^{(\text{direct})}(\{(t_i, e_i)\} | D) \right| \\
&= \left| \frac{\hat{\alpha}_{e_i 1} \psi_{e_i 1}(t_{i^*} - t_i)}{\hat{\lambda}_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)} - \frac{\alpha_{e_i 1} \psi_{e_i 1}(t_{i^*} - t_i)}{\lambda_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)} \right| \\
&\leq \bar{\psi}_1 \left| \frac{\hat{\alpha}_{e_i 1}}{\hat{\lambda}_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)} - \frac{\alpha_{e_i 1}}{\lambda_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)} \right| \\
&\leq \bar{\psi}_1 \left| \frac{\hat{\alpha}_{e_i 1}}{\hat{\lambda}_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)} - \frac{\alpha_{e_i 1}}{\hat{\lambda}_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)} \right| + \bar{\psi}_1 \left| \frac{\alpha_{e_i 1}}{\hat{\lambda}_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)} - \frac{\alpha_{e_i 1}}{\lambda_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)} \right| \\
&= \bar{\psi}_1 \frac{|\hat{\alpha}_{e_i 1} - \alpha_{e_i 1}|}{\hat{\lambda}_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)} + \bar{\psi}_1 |\alpha_{e_i 1}| \frac{|\hat{\lambda}_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D) - \lambda_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)|}{\hat{\lambda}_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D) \lambda_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)} \\
&\leq 2\bar{\psi}_1 \frac{\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_\infty}{\lambda_1(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)} + 2\bar{\psi}_1 \|\boldsymbol{\theta}\|_\infty \frac{m\bar{\psi}_1 \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_\infty}{\lambda_1^2(t_{i^*} | \mathcal{H}_{t_{i^*}}^D)}. \tag{EC.15}
\end{aligned}$$

Combining (EC.14), (EC.15), and Theorem 3, we complete the proof of Theorem 4.  $\square$

## EC.5. Proof of Theorem 5

*Proof of Theorem 5.* By the definitions and the triangle inequality,

$$\begin{aligned}
& \left| \widehat{\text{att}}_{t_{i^*}}^{(\text{total})}(R | D) - \text{att}_{t_{i^*}}^{(\text{total})}(R | D) \right| \\
&= \left| \mathbb{E}[\widehat{\text{att}}_{t_{i^*}}^{(\text{direct})}(\widehat{R}^\diamond | D) | D] - \mathbb{E}[\text{att}_{t_{i^*}}^{(\text{direct})}(R^\diamond | D) | D] \right| \\
&= \left| \sum_{R \subseteq R' \subseteq \Omega} \widehat{\text{att}}_{t_{i^*}}^{(\text{direct})}(R' | D) \mathbb{P}(\widehat{R}^\diamond = R' | D) - \sum_{R \subseteq R' \subseteq \Omega} \text{att}_{t_{i^*}}^{(\text{direct})}(R' | D) \mathbb{P}(R^\diamond = R' | D) \right| \\
&\leq \sum_{R \subseteq R' \subseteq \Omega} \left| \widehat{\text{att}}_{t_{i^*}}^{(\text{direct})}(R' | D) \mathbb{P}(\widehat{R}^\diamond = R' | D) - \text{att}_{t_{i^*}}^{(\text{direct})}(R' | D) \mathbb{P}(R^\diamond = R' | D) \right| \\
&\leq \sum_{R \subseteq R' \subseteq \Omega} \left| \widehat{\text{att}}_{t_{i^*}}^{(\text{direct})}(R' | D) - \text{att}_{t_{i^*}}^{(\text{direct})}(R' | D) \right| \\
&\quad + \sum_{R \subseteq R' \subseteq \Omega} \left| \mathbb{P}(\widehat{R}^\diamond = R' | D) - \mathbb{P}(R^\diamond = R' | D) \right|, \tag{EC.16}
\end{aligned}$$

where the last inequality is due to the the attribution scores and the probabilities involved being bounded by 1. Each component of the first term is bounded in Theorem 4. For the second term, the estimated version of thinning yields

$$\mathbb{P}(\widehat{R}^\diamond = R' | D) = \prod_{(t_i, e_i) \in \Omega \setminus R} \mathbb{P}(\text{Bernoulli}(1 - \frac{\hat{\lambda}_{e_i}(t_i | \mathcal{H}_{t_i}^{D \setminus R'})}{\hat{\lambda}_{e_i}(t_i | \mathcal{H}_{t_i}^D)}) = \mathbb{1}_{\{(t_i, e_i) \in R'\}}.$$

Therefore,

$$\begin{aligned}
& \left| \mathbb{P}(\widehat{R}^\diamond = R' \mid D) - \mathbb{P}(R^\diamond = R' \mid D) \right| \\
&= \left| \prod_{(t_i, e_i) \in \Omega \setminus R} \left( 1 - \frac{\widehat{\lambda}_{e_i}(t_i \mid \mathcal{H}_{t_i}^{D \setminus R'})}{\widehat{\lambda}_{e_i}(t_i \mid \mathcal{H}_{t_i}^D)} \right)^{\mathbb{1}_{\{(t_i, e_i) \in R'\}}} \left( \frac{\widehat{\lambda}_{e_i}(t_i \mid \mathcal{H}_{t_i}^{D \setminus R'})}{\widehat{\lambda}_{e_i}(t_i \mid \mathcal{H}_{t_i}^D)} \right)^{1 - \mathbb{1}_{\{(t_i, e_i) \in R'\}}} \right. \\
&\quad \left. - \prod_{(t_i, e_i) \in \Omega \setminus R} \left( 1 - \frac{\lambda_{e_i}(t_i \mid \mathcal{H}_{t_i}^{D \setminus R'})}{\lambda_{e_i}(t_i \mid \mathcal{H}_{t_i}^D)} \right)^{\mathbb{1}_{\{(t_i, e_i) \in R'\}}} \left( \frac{\lambda_{e_i}(t_i \mid \mathcal{H}_{t_i}^{D \setminus R'})}{\lambda_{e_i}(t_i \mid \mathcal{H}_{t_i}^D)} \right)^{1 - \mathbb{1}_{\{(t_i, e_i) \in R'\}}} \right| \\
&\leq \sum_{(t_i, e_i) \in \Omega \setminus R} \left| \frac{\widehat{\lambda}_{e_i}(t_i \mid \mathcal{H}_{t_i}^{D \setminus R'})}{\widehat{\lambda}_{e_i}(t_i \mid \mathcal{H}_{t_i}^D)} - \frac{\lambda_{e_i}(t_i \mid \mathcal{H}_{t_i}^{D \setminus R'})}{\lambda_{e_i}(t_i \mid \mathcal{H}_{t_i}^D)} \right| \\
&\leq (m-r) \max_{(t_i, e_i) \in \Omega \setminus R} \left| \frac{\widehat{\lambda}_{e_i}(t_i \mid \mathcal{H}_{t_i}^{D \setminus R'})}{\widehat{\lambda}_{e_i}(t_i \mid \mathcal{H}_{t_i}^D)} - \frac{\lambda_{e_i}(t_i \mid \mathcal{H}_{t_i}^{D \setminus R'})}{\lambda_{e_i}(t_i \mid \mathcal{H}_{t_i}^D)} \right|. \tag{EC.17}
\end{aligned}$$

Using similar tricks in Theorem 4, we can show that each component in (EC.17) has the same rate of convergence up to some constants. Combining (EC.16), (EC.17), and Theorem 3, we can derive the desired bound. It suffices to make sure Theorem 3 holds for each  $e = 1, \dots, q$ . The corresponding tail probabilities can be obtained by the union bound. Let  $\boldsymbol{\theta}^{(e)} = (\mu_e, \boldsymbol{\alpha}_e^\top)^\top$  denote the true parameter and  $\widehat{\boldsymbol{\theta}}^{(e)}$  be its regularized estimate for  $e = 1, \dots, q$  in the following context. Then there exists  $C_2 > 0$  such that

- (i) If  $p$  is fixed and  $n$  is sufficiently large, then for any constant  $0 < \nu < 1$ ,

$$\begin{aligned}
& \mathbb{P} \left( \left| \widehat{\text{att}}_{t_{i^*}}^{(\text{total})}(R \mid D) - \text{att}_{t_{i^*}}^{(\text{total})}(R \mid D) \right| > C_2 n^{-\frac{1-\nu}{2}} \right) \\
&\leq \sum_{e=1}^q \mathbb{P} \left( \|\widehat{\boldsymbol{\theta}}^{(e)} - \boldsymbol{\theta}^{(e)}\|_\infty > 10(\xi^{(e)})^{-1} \kappa^{(e)} (C_3^{(e)})^{-\frac{1}{2}} C_4^{(e)} n^{-\frac{1-\nu}{2}} \right) \\
&\leq q \cdot 2(p+1)(p+2) \exp(-n^\nu) \\
&\leq 2p(p+1)(p+2) \exp(-n^\nu);
\end{aligned}$$

- (ii) If  $p$  diverges and  $n$  is sufficiently large, then for any constant  $\zeta > 3$ ,

$$\begin{aligned}
& \mathbb{P} \left( \left| \widehat{\text{att}}_{t_{i^*}}^{(\text{total})}(R \mid D) - \text{att}_{t_{i^*}}^{(\text{total})}(R \mid D) \right| > C_2 \sqrt{\frac{\zeta \log(p+1)}{n}} \right) \\
&\leq \sum_{e=1}^q \mathbb{P} \left( \|\widehat{\boldsymbol{\theta}}^{(e)} - \boldsymbol{\theta}^{(e)}\|_\infty > 10(\xi^{(e)})^{-1} \kappa^{(e)} (C_3^{(e)})^{-\frac{1}{2}} C_4^{(e)} \sqrt{\frac{\zeta \log(p+1)}{n}} \right) \\
&\leq q \cdot 3/(p+1)^{\zeta-2} \\
&< 3/(p+1)^{\zeta-3}.
\end{aligned}$$

Therefore, the proof of Theorem 5 is complete.  $\square$

## References

- Hansen, N. R., Reynaud-Bouret, P. and Rivoirard, V. (2015), ‘Lasso and probabilistic inequalities for multivariate point processes’, *Bernoulli* **21**(1), 83–143.
- Hong, M. and Luo, Z.-Q. (2017), ‘On the linear convergence of the alternating direction method of multipliers’, *Mathematical Programming* **162**(1-2), 165–199.
- Lin, W. and Lv, J. (2013), ‘High-dimensional sparse additive hazards regression’, *Journal of the American Statistical Association* **108**(501), 247–264.
- Vershynin, R. (2018), *High-Dimensional Probability: An Introduction with Applications in Data Science*, Vol. 47, Cambridge University Press.
- Wainwright, M. J. (2009), ‘Sharp thresholds for high-dimensional and noisy sparsity recovery using  $\ell_1$ -constrained quadratic programming (lasso)’, *IEEE Transactions on Information Theory* **55**(5), 2183–2202.
- Zhao, P. and Yu, B. (2006), ‘On model selection consistency of lasso’, *The Journal of Machine Learning Research* **7**, 2541–2563.