INTRODUCTION

Georgia Tech

Self- and mutually-exciting point processes are popular models for dependent discrete event data. To date, most existing models assume stationary kernels (including the classical Hawkes processes) and simple parametric models. Modern applications with complex event data require more general models that can incorporate contextual information of the events, called marks, besides the temporal and location information. Moreover, such applications often require non-stationarity to capture more complex spatio-temporal dependence. In this paper, we introduce a novel and general neural network-based non-stationary influence kernel with high expressiveness for handling complex discrete events data while providing theoretical performance guarantees. We demonstrate the superior performance of our proposed method compared with the state-of-the-art on synthetic and real data.



Figure: An example of non-stationary influence kernel k(t', t) of event time t'and future time t > t'.

MARKED TEMPORAL POINT PROCESS

A data point in the discrete event sequence:

$$x = (t, m), \quad t \in [0, T), \quad m \in \mathcal{M},$$

Conditional intensity function:

$$\lambda(x)dx = \mathbb{E}\left(d\mathbb{N}(x)|\mathcal{H}_{t(x)}\right).$$

Hawkes process with kernel $k \in \mathcal{K} \subset C^0(\mathcal{X} \times \mathcal{X})$:

$$\lambda[k](x) = \mu + \int_{x' \in \mathcal{X}_{t(x)}} k(x', x) d\mathbb{N}(x'), \qquad (1)$$

Log-likelihood given M sequences $\{x_{i,j}\}, i = 1, ..., N_j, j =$ $1,\ldots,M$

$$\ell[k] \coloneqq \frac{1}{M} \sum_{j=1}^{M} \left(\int_{\mathcal{X}} \log \lambda_j[k](x) d\mathbb{N}_j(x) - \int_{\mathcal{X}} \lambda_j[k](x) dx \right),$$
(2)

CONTRIBUTION

- \checkmark The kernel function is represented by a **spectral decom**position of the influence kernel.
- \checkmark The spectral decomposition of asymmetric influence kernel consists of a sum of the product of feature maps, which can be **parameterized by neural networks**.
- \checkmark We establish theoretical guarantees of the MLE for the true kernel function based on functional variational analysis and finite-dimensional asymptotic analysis, which shed light on theoretical understanding of neural network-based kernel functions.

NEURAL SPECTRAL MARKED POINT PROCESSES

Shixiang Zhu[†], Haoyun Wang[†], Zheng Dong[†], Xiuyuan Cheng^{*}, Yao Xie[†]

[†] H. Milton Stewart School of Industrial and Systems Engineering, Georgia Institute of Technology * Department of Mathematics, Duke University

NEURAL SPECTRAL KERNEL

Influence kernel k is represented using finite-rank decomposition:

$$k(x', x) = \sum_{r=1}^{R} \nu_r \psi_r(x') \phi_r(x), \quad \nu_r \ge 0, \quad (3) \quad c_1,$$

where

$$u_r: \mathcal{X} \to \mathbb{R}, \quad \phi_r: \mathcal{X} \to \mathbb{R}, \quad r = 1, \cdots, R,$$

are two sets of *feature functions* in some smooth functional space $\mathcal{F} \subset C^0(\mathcal{X})$, and ν_r is the corresponding weight or *spectrum*.



OPTIMAL FIT USING FEATURE FUNCTION BASIS REPRESENTATION

• Assume that the feature functions of the kernel can be wellapproximated by a linear combination of basis functions: $b_i(x)$: $\mathcal{X} \to \mathbb{R}, i = 1, \dots, S$:

$$\psi_r(x) = \sum_{i=1}^S \alpha_{ri} b_i(x), \ \phi_r(x) = \sum_{i=1}^S \beta_{ri} b_i(x), \ r = 1, \cdots, R.$$

The kernel function in (3) can be written as $k_A(x', x) =$ $b(x')^T A b(x)$, where $A_{pq} = \sum_{r=1}^R \nu_r \alpha_{rp} \beta_{rq}$.

• Let $A \in \mathcal{A}$ be the one which maximizes the expected loglikelihood function. i.e.

$$\widetilde{A} = \underset{A \in \mathcal{A}}{\operatorname{arg\,max}} \mathbb{E}\left(\ell_A\right). \tag{4}$$

Under Assumption, let the ℓ_2 -norm of a kernel be

$$||k||_{2}^{2} = \int_{\mathcal{X}} \int_{\mathcal{X}_{t(x)}} k(x', x)^{2} dx' dx.$$
 (5)

Then we have

$$\|k^* - k_{\widetilde{A}}\|_2^2 \le \frac{c_2^5 |\mathcal{M}| T + c_2^4}{c_1^4} \exp(2(c_2 - c_1) |\mathcal{M}| T) D(k^*, \mathcal{K}_{\text{finite}})^2,$$

where $D(k^*, \mathcal{K}_{\text{finite}})$ is the ℓ_2 -distance between the true kernel and the set $\mathcal{K}_{\text{finite}}$,

 $D(k^*, \mathcal{K}_{\text{finite}}) = \min_{k \in \mathcal{K}_{\text{finite}}} \|k^* - k\|_2.$

KERNEL IDENTIFIABILITY

Assumption: (A1) The kernel function family $\overline{\mathcal{K}} \subset C^0(\mathcal{X} \times \mathcal{X})$, Id kernel functions in $\overline{\mathcal{K}}$ is uniformly bounded; (A2) There exist , c_2 positive constants, such that for any $k \in \overline{\mathcal{K}}, c_1 \leq \lambda_i[k](x) \leq k$ $c_2, \forall x \in \mathcal{X} \text{ and } \forall j.$

Kernel identifiability using maximum likelihood: Under Assumption, the true kernel function k^* is locally identifiable in that k^* is a local minimum solution of maximum likelihood (2) in expectation.

Asymptotic normality of low-rank MLE

Let the singular value decomposition of \widetilde{A} be $\widetilde{A} = U\Lambda V^T$ and I be the expected Hessian matrix of the log-likelihood at A. Let J be the covariance matrix of a single trajectory's score function at \widetilde{A} , and $\widetilde{G} \in \mathbb{R}^{S \times S}$ be the expected score at \tilde{A} , Let $F = (\mathbb{I}_S \otimes U, V \otimes \mathbb{I}_S) \in \mathbb{R}^{S^2 \times 2SR}$ where \otimes is the Kronecker product, \mathbb{I}_S is the identity matrix of size S, $\widetilde{C} = (\widetilde{A}^{\dagger} \otimes \widetilde{G})Q_{S,S} + ((\widetilde{A}^{\dagger} \otimes \widetilde{G})Q_{S,S})^T$, where \dagger represents pseudo-inverse and $Q_{a,b} \in \mathbb{R}^{ab \times ab}$ is the permutation matrix such that $\operatorname{vec}(P^T) = Q_{a,b}\operatorname{vec}(P)$ for any *a*-by-*b* matrix *P*.

If $F^T(\widetilde{I} + \widetilde{C})F$ shares the same null-space with F, then the low-rank estimator \widehat{A}_{MLE} , solved from the constrained maximum likelihood problem, when $M \to \infty$, satisfies

 $\sqrt{M}(\operatorname{vec}(\widehat{A}_{\mathrm{MLE}}) - \operatorname{vec}(\widetilde{A})) \rightarrow$ $\mathcal{N}(0, F(F^T(\widetilde{I} + \widetilde{C})F)^{\dagger}F^T\widetilde{J}F(F^T(\widetilde{I} + \widetilde{C})F)^{\dagger}F),$

Real Data Experiments

Table: Predictive log-likelihood on real data.

ℓ	Earthquake $(2D)$	Robbery $(1D)$	#Parameters	Training/Testing time ²
MPP	-56.50	-74.47	171,555	0.766 / 0.84
TPP	-218.39	-132.55	$274,\!168$	0.245 / 7.29
ural Hawkes	-189.39	-96.10	282,755	0.204 / 6.09
wkes	NA	-197.84	2	0.021 / < 0.01







Figure: Predicted conditional intensity using our method and other baseline approaches for synthetic data sets.

Zhu, Shixiang, Haoyun Wang, Zheng Dong, Xiuyuan Cheng, and Yao Xie. "Neural Spectral Marked Point Processes." arXiv:2106.10773 (2021).





NUMERICAL EXPERIMENTS

REFERENCE