1  We thank all the reviewers for their efforts and constructive comments, which help improve the quality of our paper.

2  **To Reviewer #1**
3  **Q1. Lack of the evaluation of the uncertainty of the estimate.** We actually derived an error bound of the empirical
4  risk. Based on the analysis of the first and second moment of the estimator presented in Theorems 5.1 and 5.2, a
5  Chebyshev's type error bound can be easily obtained:

$$\mathbb{P}\Big(\big|\nabla\tilde{R}^{\mathcal{G}} - \mathbb{E}\nabla\tilde{R}^{\mathcal{G}}(\theta)\big| > \epsilon\Big) \leqslant \frac{\mathbb{V}\big[\nabla\tilde{R}^{\mathcal{G}}(\theta)\big]}{\epsilon^2} \leqslant \frac{\mathbb{V}\big[\nabla R_\ell(\theta)\big]}{\epsilon^2} = \frac{\frac{1-p}{p}\mathbb{E}\big[\nabla R(\theta)\big]^2 + \mathbb{V}\big[\nabla R(\theta)\big]}{\epsilon^2},$$

6  for any $\epsilon > 0$. As we reckon this error bound of the gradient estimation is a direct deduction of Theorems 5.1 and 5.2,
7  we didn't include this result due to page limit. We will present it as a corollary in the final version.
8  **Q2. Thinning bootstrap procedure.** Bootstrap is not suitable for point processes as it is a resampling method with
9  replacement, that is, an event could be sampled more than once. This violates the fundamental assumption of point
10  processes that events do not arrive simultaneously. A point cannot be sampled twice, otherwise the log-likelihood of the
11  point process will be infinite. Therefore, we propose thinning as a downsampling method, rather than bootstrap as a
12  resampling method with replacement. Alternatively, Jackknife resampling seems feasible for point processes, and we
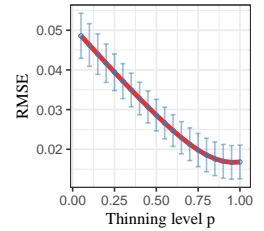13  plan to explore it in the future. We will add this discussion to the final version.
14  **Q3. Better figure representing MSE vs p.** Thanks for the suggestion and we will revise our paper accordingly. We
15  have conducted an experiment on this with the same setting in the paper and will add a new figure (shown below) in the
16  final version. It shows that RMSE decreases as $p$ increases.

17  **To Reviewer #2**
18  **Q1. The definition of R.** R is the likelihood/squared error loss, aka the residue. We will
19  make this clear in the final version.
20  **Q2. Make the transition to multi-variate processes a bit smoother.** This is indeed a
21  helpful suggestion, we will add a discussion to make it better.
22  **Q3. Discuss the gradient for stochastic intensities calculated on sub-samples potentially
23  being unbiased.** The gradient estimation is unbiased if and only if $\mathbb{E}[\mathcal{H}(t;\theta)\lambda(t;\theta)] =$



24  $\mathbb{E}\mathcal{H}(t;\theta)\mathbb{E}\lambda(t;\theta)$, as we shown in the proof of Theorem 5.1. However, this result is somehow inaccessible, therefore
25  we particularly indicate, in the Theorem 5.1, that the gradient estimation for NHPPs are unbiased, but not for all the
26  stochastic intensities. We also illustrated this result by the experiment shown in Fig. 3 and discussed in Line 266-269,
27  in order to provide an intuitive understanding. We will add a remark after the theorem to clarify this result.

28  **To Reviewer #3**
29  We want to highlight that the aim of this paper is NOT to propose a specific optimization/learning algorithm, but to
30  answer one of the most fundamental questions in point processes: what is the best sampling method for point processes?
31  We propose thinning as a solution, and derive important theoretical results of thinning for parameter estimation, gradient
32  estimation, and stochastic optimization, for point processes with decouplable intensities. In our empirical study, we
33  applied thinning to learn various state-of-the-art models (MMEL, Granger Causality for Hawkes, and Sparse Low-rank
34  Hawkes); to estimate the gradient in different types of point processes (NHPP and Hawkes) with different estimators
35  (MLE and LSE); and also to perform stochastic optimization under different algorithms (SGD and Adam). Through
36  this study, we show that thinning is a general downsampling solution for point processes with decouplable intensities,
37  and is not restricted to a specific learning/optimization algorithm or estimator.
38  **Q1. The applicability of the proposed model. Hawkes is a weak baseline.** The thinning method is applicable to
39  most, if not all, state-of-the-art models related to parametric point processes, including MMEL, GC, and sparse low-rank
40  Hawkes, as we shown in the synthetic experiment in Table 1. In the experiment with real datasets for stochastic
41  optimization, the aim is to test the applicability of thinning to typical and popular optimization methods such as SGD
42  and Adam. We picked Hawkes as the model to learn as it is the basis of many derivative models. Note that, **Hawkes is
43  not used as a baseline here**; our baselines are other sampling/optimization methods. The main purpose here is NOT
44  to show that Hawkes is good, but rather to show that thinning is effective when coupled with various optimization
45  algorithms (SGD and Adam).
46  **Q2. The focus on decouplable intensities is a bit limited.** As we mentioned in the paper, most state-of-the-art models
47  of parametric point processes are decouplable (we would be grateful if the reviewer can point out if we have missed any).
48  Besides, Netcodec (Long Tran, et al., 2015), parametric Hawkes (Liangda Li, et al., 2014), Hawkes with Stochastic
49  Excitations (Young Lee, et al., 2016) and SLANT (Abir De, et al., 2016) are also decouplable. This general category
50  covers all NHPPs, compound Poisson processes, renewal processes, marked point processes with independent markers,
51  and a certain part of doubly stochastic Poisson processes (Cox processes). We believe this is a general assumption with
52  many interesting properties. We choose this specific class of point processes in this first attempt of the problem as it
53  facilitates the mathematical preciseness and rigorous theoretical proofs. We agree that a broader assumption will serve
54  more scenarios, and we will investigate more general classes.