

2 We would like to thank the reviewers for their detailed feedback and insightful comments, we will incorporate the  
3 suggested clarifications in the paper.

4 **Remarks for Reviewer 1.** We will add the reference of the SVI approach of Linderman and Adams in the introduction.

5 **Remarks for Reviewer 2.**

6 • *On the normalization coefficient and prior  $p_\alpha$  in (4):* It is known that the regularized-MLE objective is equivalent to  
7 MAP objective up to the constant normalization coefficient. The regularization term  $R(\mu, W)/\alpha$  in (4) can be seen  
8 as the negative of the logarithm of the unnormalized prior. So to derive the prior  $p_\alpha(\mu, W)$  we only need to compute  
9 the normalization term, which is the integral of  $\exp(-R(\mu, W)/\alpha)$  over  $\mu$  and  $W$ , and which therefore cannot be a  
10 function of the integration variables  $\mu$  and  $W$ .

11 • *On Optimizing  $\alpha$  in (4) and (8) directly:* In line 142 we actually mean that the MAP estimator cannot be optimized  
12 over  $\alpha$ . Indeed, as demonstrated in the example of Appendix B, the MAP objective is an unbounded function (from  
13 above) of  $\alpha$ .

14 • *On line 147:* Indeed it should be “maximum likelihood”, we apologize for the confusion caused by the typo.

15 • *On performance if only 1 dimension has few observations:* In our setting, data scarcity comes from the short length  
16 of the observation window and not from missing data. So, if one dimension  $i$  has much fewer timestamps than  
17 others, it means that overall, it has smaller intensity (i.e. small  $\mu_i$  and incoming  $W_{ij}$ ). Therefore the likelihood  
18 function naturally enforces small values in these parameters to explain the observed intervals with no timestamps.  
19 The setting where the data scarcity comes from both short observation window and missing data requires extending  
20 our probabilistic model and is an interesting direction for future work.

21 **Remarks for Reviewer 3.**

22 • *Computational complexity of the algorithm:* Our gradient-based method is computationally efficient and scales well  
23 to large data regimes. For small data-regimes, the state-of-the-art methods empirically seem to converge faster as  
24 they need fewer iterations (even if there is no proof of convergence rate in the papers). However, when the number of  
25 nodes gets large and the number of observations increases, the per-iteration cost of the state-of-the-art methods grows  
26 faster than our gradient-based approach, which we expect to converge faster for such settings. Indeed, the complexity  
27 of MLE-ADM4 is  $O(N_{iter}n^3d^2)$  (see Table 1, Achab et al. 2016, [1]), whereas the complexity of our approach can  
28 be reduced to  $O(N_{iter}nd^2 + d^2n^2)$  where the reduced cost comes from efficiently pre-computing some constant  
29 terms in the log-likelihood function (at the cost of memory), which is a one shot cost of  $O(d^2n^2)$ .

30 • *To evaluate the scalability of our approach, we ran additional simulations on increasingly large-dimensional problems.*  
31 As shown in Figure 1, the per-iteration running time of our approach VI-EXP (implemented in python) scales  
32 better than the one of MLE-ADM4 (implemented in C++). In addition, even if our gradient descent algorithm  
33 requires more iterations to converge, we show in Figure 2 that VI-EXP reaches the same F1-score as MLE-ADM4  
34 faster. Empirically, we expect similar results for the non-parametric setting. We performed simulations only for the  
35 parametric setting due to the time constraint of the rebuttal.

36 • *Optimizing the decay parameters:* It is possible to optimize the decay parameters but we chose to use this particular  
37 form of exponential kernel as an example designed to match the efficient C++ implementation of MLE-ADM4,  
38 which takes advantage of the convexity of the problem. In addition, considering a fixed decay enables the use of  
39 caching as discussed above.

40 • *Visualizing the estimated causal-networks:* In high-dimensions, the resulting networks are not easy to visualize. We  
41 tried to draw the learned networks on top of a map of the Ebola dataset, but the figure needs to be rendered too large  
42 to be clear. Given space limitation, we did not plot any.

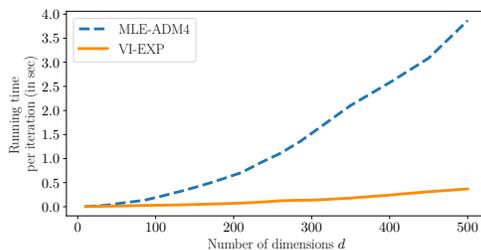


Figure 1: Comparison of per-iteration running time.

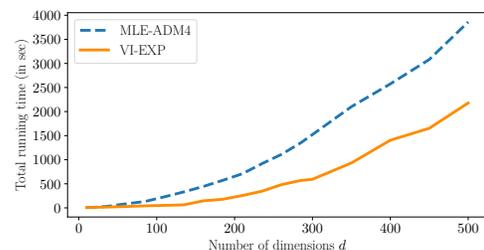


Figure 2: Running time required for our approach VI-EXP to reach the same F1-Score as MLE-ADM4.