

## Appendix 1. Likelihood and evidence lower bound

Before deriving the likelihood for the mixture of MHPs model, we start with the following Lemma about the relationship between the PDF and the conditional intensity function.

**Lemma 1.** Given the conditional intensity function of a Hawkes process,  $\lambda^*(t)$ , the probabilistic density function  $f^*$ , the survival function  $S^*$  and the cumulative distribution function  $F^*$  are given by:

$$\begin{aligned} f^*(t) &= \lambda^*(t) \exp\left(-\int_{t_n}^t \lambda^*(s) ds\right) \\ S^*(t) &= \exp\left(-\int_{t_n}^t \lambda^*(s) ds\right) \\ F^*(t) &= 1 - \exp\left(-\int_{t_n}^t \lambda^*(s) ds\right). \end{aligned}$$

*Proof.* By definition, we have:

$$\lambda^*(t) = \frac{f^*(t)}{1 - F^*(t)} = \frac{\frac{d}{dt} F^*(t)}{1 - F^*(t)} = -\frac{d}{dt} \log(1 - F^*(t))$$

The equation for  $F^*$  follows from the fundamental theorem of calculus, and the other two relationships can be derived using  $f^* = \frac{d}{dt} F^*$  and  $S^* = 1 - F^*$ .  $\square$

The likelihood for the complete data  $\{(t_n, Z_n, W_n)\}$  is given by:

$$L(t, Z, W) = \prod_{n=1}^N p(Z_n) p(W_n | Z_n) p^*(t_n | Z_n).$$

The first two terms are straightforward, i.e., the smoothing prior and the language model. Let  $t_0 = 0$  and  $t_{N+1} = T$ , the last term can be derived as follows:

$$\begin{aligned} L_{t|Z} &= \prod_{n=1}^N p^*(t_n | Z_n) \times S^*(T) \\ &= \prod_{n=1}^N \prod_{m=1}^M \lambda_{i_n, m}^*(t_n)^{Z_{nm}} \times \prod_{n=1}^N \prod_{m=1}^M \prod_{i=1}^I S_{i, m}^*(t_n) \times \prod_{m=1}^M \prod_{i=1}^I S_{i, m}^*(T) \\ &= \prod_{n=1}^N \prod_{m=1}^M \lambda_{i_n, m}^*(t_n)^{Z_{nm}} \times \prod_{n=1}^N \prod_{m=1}^M \prod_{i=1}^I \exp\left(-\int_{t_{n-1}}^{t_n} \lambda_{i, m}^*(s) ds\right) \times \prod_{m=1}^M \prod_{i=1}^I \exp\left(-\int_{t_N}^T \lambda_{i, m}^*(s) ds\right) \\ &= \prod_{n=1}^N \prod_{m=1}^M \lambda_{i_n, m}^*(t_n)^{Z_{mn}} \times \exp\left(-\sum_{i=1}^I \sum_{m=1}^M \int_0^T \lambda_{i, m}^*(s) ds\right). \end{aligned}$$

which gives the Eq(11) in the paper.

The normalization term in the evidence lower bound (ELBO) is derived as follows:

$$\begin{aligned}
& \sum_{i=1}^I \sum_{m=1}^M \int_0^T \mathbb{E}_q[\lambda_{i,m}(s)] ds \\
&= \sum_{i=1}^I \sum_{m=1}^M \int_0^T (\gamma_m \mu_i + \sum_{t_l < t} \phi_{tm} \alpha_{ii} \kappa(t - t_l)) dt \\
&= T \sum_{i=1}^I \sum_{m=1}^M \gamma_m \mu_i + \sum_{i=1}^I \int_0^T \sum_{t_l < t} \alpha_{ii} \kappa(t - t_l) dt \sum_{m=1}^M \phi_{tm} \\
&= T \sum_{i=1}^I \sum_{m=1}^M \gamma_m \mu_i + \sum_{i=1}^I \int_0^T \sum_{t_l < t} \alpha_{ii} \kappa(t - t_l) dt \\
&= T \sum_{i=1}^I \sum_{m=1}^M \gamma_m \mu_i + \sum_{i=1}^I \sum_{n=1}^{N+1} \sum_{l=1}^{n-1} \alpha_{ii} \int_{t_{n-1}}^{t_n} \kappa(t - t_l) dt \\
&= T \sum_{i=1}^I \sum_{m=1}^M \gamma_m \mu_i + \sum_{i=1}^I \sum_{l=1}^N \alpha_{ii} \sum_{n=l+1}^{N+1} (K(t_n - t_l) - K(t_{n-1} - t_l)) \\
&= T \sum_{i=1}^I \sum_{m=1}^M \gamma_m \mu_i + \sum_{i=1}^I \sum_{l=1}^N \alpha_{ii} K(T - t_l),
\end{aligned}$$

which proves Eq(15) in the paper.

## Appendix 2. Inference and learning algorithm

In this section, we derive the inference and learning algorithms.

**Inference for  $\phi$ .** The variational ELBO objective is decomposable for  $\phi$ s. Isolating the terms containing  $\phi_n$  we have:

$$\begin{aligned}
\max \mathcal{L}_{\phi_n} &= \sum_{m=1}^M \phi_{nm} \left( \log \pi_m + \sum_{v=1}^V w_{nv} \log \beta_{mv} - \sum_{l=1}^{n-1} \phi_{lm} \eta_{ln}^m \log(\eta_{ln}^m) - \log \phi_{nm} \right) \\
&\quad + \sum_{m=1}^M \phi_{nm} \left( \eta_{nn}^m \log(\gamma_m \mu_{i_n}) + \sum_{\ell=1}^{n-1} \phi_{lm} \eta_{ln}^m \log(\alpha_{ii} k(t_n - t_l)) + \sum_{\ell=n+1}^N \phi_{lm} \eta_{nl}^m \log(\alpha_{ii} k(t_n - t_l)) \right) \\
s.t. \quad &\sum_{m=1}^M \phi_{nm} = 1.
\end{aligned}$$

Eq(19) follows from the first-order optimality of the above optimization, i.e., by forming the Lagrangian and setting the derivative to zero.

**Inference for  $\eta$ .** The branching variables are optimal when the lower-bound in Eq(16) is tight, i.e., when the equality satisfies. This leads to:

$$\begin{aligned}\eta_{nn}^m &\propto \gamma_m \mu_{i_n}, \\ \eta_{ln}^m &\propto \alpha_{i_l i_n} k(t_n - t_l), \\ \eta_{nn}^m + \sum_{l=1}^{n-1} \phi_{lm} \eta_{ln}^m &= 1,\end{aligned}$$

which leads to Eq(20–21).

**Update  $\alpha$ ,  $\mu$  and  $\gamma$ .** The update formulas for  $\pi$  and  $\beta$  are trivial. Here we derive the update equations for  $\mu$ ,  $\gamma$  and  $\alpha$ , given that the variational parameters are optimal. We have, for  $\alpha$ , the ELBO:

$$\mathcal{L}_\alpha = \sum_{n=1}^N \sum_{m=1}^M \phi_{nm} \sum_{l=1}^{n-1} \eta_{ln}^m \phi_{lm} \log(\alpha_{i_l i_n} k(t_n - t_l)) - \sum_{i=1}^I \sum_{l=1}^N \alpha_{i_l} K(T - t_l).$$

Eq(25) follows by setting the derivative of  $\mathcal{L}_\alpha$  to zero. Eq(24) and Eq(26) can be derived similarly.

### Appendix 3. An alternative inference algorithm

In Eq(16), we lower-bound the expectation of the log-intensity by applying Jensen’s inequality with a set of branchings that satisfy Eq(17–18). Here, we present an alternative way to define these branchings and lower-bound the expected log-intensity, which leads to another version of the inference algorithm. This algorithm is slightly flawed in its mathematic form but in our experiments, we found that it converges faster than the one presented in the paper.

Again, we define a set of branching variables  $\{\eta^m\}$ , each of which fill a lower-triangular matrix, i.e.,  $\eta_n^m = [\eta_{1,n}^m, \dots, \eta_{n,n}^m]^T$ . But different from Eq(17–18), the branchings here define a multinomial distribution such that each  $\eta_{l,n}^m \geq 0$  and  $\sum_{l=1}^n \eta_{ln}^m = 1$ . Let  $\mathcal{E}[\eta_n^m] = -\sum_{l=1}^n \eta_{ln}^m \log(\eta_{ln}^m)$  be the entropy of  $\eta_n^m$ , we have:

$$\begin{aligned}&\mathbb{E}_q[\log(\lambda_{i_n, m}(t_n))] \\ &= \mathbb{E}_q[\log(\eta_{nn}^m \frac{\gamma_m \mu_{i_n}}{\eta_{nn}^m} + \sum_{l=1}^{n-1} \eta_{ln}^m \frac{Z_{lm} \alpha_{i_l, i_n} k(t_n - t_l)}{\eta_{ln}^m})] \\ &\geq \eta_{nn}^m \log(\gamma_m \mu_{i_n}) + \sum_{l=1}^{n-1} \eta_{ln}^m (\mathbb{E}_q[\log(Z_{lm})] + \log(\alpha_{i_l, i_n} k(t_n - t_l))) + \mathcal{E}[\eta_n^m]\end{aligned}$$

Unfortunately, in the Equation above, the term  $\mathbb{E}_q[\log(Z_{lm})]$ , i.e., the expected logarithm of a Bernoulli variable, is *not well-defined*. As a proxy, we substitute it with  $\log(\mathbb{E}_q[Z_{lm}]) = \log \phi_{lm}$ . Note that this is not mathematically sound.

This new lower bound leads to the following inference and learning algorithms.

**Variational inference:**

$$\begin{aligned}
\phi_{nm}^{\text{new}} &\propto \pi_m : \text{prior} \\
&\times \left( \prod_{v=1}^V \beta_{mv}^{w_{nv}} \right) : \text{content} \\
&\times (\gamma_m \mu_{i_n})^{\eta_{nn}^m} : \text{self triggering} \\
&\times \prod_{l=1}^{n-1} \left( \phi_{lm}^{\text{old}} \alpha_{i_l i_n} \kappa(t_n - t_l) \right)^{\eta_{ln}^m} : \text{influences from the past} \\
&\times \prod_{l=n+1}^N \exp\left(\eta_{nl}^m \frac{\phi_{lm}^{\text{old}}}{\phi_{nm}^{\text{old}}}\right) : \text{influences to the future}^1
\end{aligned}$$

where the branchings are updated via:

$$\begin{aligned}
\eta_{nn}^m &= \frac{\gamma_m \mu_{i_n}}{\gamma_m \mu_{i_n} + \sum_{l=1}^{n-1} \phi_{lm} \alpha_{i_l i_n} k(t_n - t_l)} \\
\eta_{ln}^m &= \frac{\phi_{lm} \alpha_{i_l i_n} k(t_n - t_l)}{\gamma_m \mu_{i_n} + \sum_{l=1}^{n-1} \phi_{lm} \alpha_{i_l i_n} k(t_n - t_l)}
\end{aligned}$$

**Learning:**

$$\begin{aligned}
\pi_m &\propto \sum_{n=1}^N \phi_{nm} \\
\beta_{mv} &\propto \sum_{n=1}^N w_{nv} \phi_{nm} \\
\mu_i &= \frac{\sum_{n=1}^N \delta(i_n = i) \sum_{m=1}^M \phi_{nm} \eta_{nn}^m}{\sum_{m=1}^M \gamma_m T} \\
\alpha_{ij} &= \frac{\sum_{m=1}^M \sum_{n=1}^N \sum_{l=1}^{n-1} \delta(i_l = i) \delta(i_n = j) \phi_{nm} \eta_{ln}^m}{\sum_{n=1}^N \delta(i_n = i) K(T - t_n)} \\
\gamma_m &= \frac{\sum_{n=1}^N \phi_{nm} \eta_{nn}^m}{\sum_{i=1}^I \mu_i T}
\end{aligned}$$

Note that the formulas for  $\pi$  and  $\beta$  are the same.

## Appendix 4. Sparsity in MHPs and MMHP

In this appendix, we show that simply adding Lasso (i.e.,  $\ell_1$ ) or ElasticNet (i.e., mixture of  $\ell_1$  and  $\ell_2$ ) penalties to the ELBO objective cannot lead to topological sparsity. To illustrate this, let us consider the plain MHP model (i.e., without latent mixture variables). Suppose  $\lambda^*(t) = [\lambda_1^*(t), \dots, \lambda_I^*(t)]^\top$  be the conditional intensity function of an  $I$ -dimensional Hawkes process, where

$$\lambda_i^*(t) = \mu_i + \sum_{t_l < t} \alpha_{i i_l} \kappa(t - t_l)$$

Given a sequence of events  $\{(t_n, i_n) | n = 1, \dots, N\}$ . The log-likelihood is given by:

$$\mathcal{L} = \sum_{n=1}^N \log(\mu_{i_n} + \sum_{\ell=1}^{n-1} \alpha_{i_\ell, i_n} \kappa(t_n - t_\ell)) - T \sum_{i=1}^I \mu_i - \sum_{n=1}^N \sum_{i=1}^I \alpha_{i_n, i} K(T - t_n)$$

Introducing the branchings, we have the following lower-bound:

$$\mathcal{L} \geq \mathcal{J} = \sum_{n=1}^N (\eta_{nn} \log(\mu_{i_n}) + \sum_{\ell=1}^{n-1} \eta_{\ell n} \log(\alpha_{i_\ell, i_n} \kappa(t_n - t_\ell))) - T \sum_{i=1}^I \mu_i - \sum_{n=1}^N \sum_{i=1}^I \alpha_{i_n, i} K(T - t_n)$$

Optimizing this lower-bound given the branchings yields the following MLE for  $\alpha$ :

$$\alpha_{ij} = \frac{\sum_{n=1}^N \sum_{\ell=1}^{n-1} \delta_{i_\ell, i} \delta_{i_n, j} \eta_{\ell n}}{\sum_{n=1}^N \delta_{i_n, i} K(T - t_n)},$$

the inferred infectivity is not sparse.

We want to enforce sparsity of  $\alpha$ . Let  $\mathcal{J}(\alpha)$  be the objective involving  $\alpha$ , we have:

$$\mathcal{J}(\alpha) = \sum_{n=1}^N \sum_{\ell=1}^{n-1} \eta_{\ell n} \log(\alpha_{i_\ell, i_n} \kappa(t_n - t_\ell)) - \sum_{n=1}^N \sum_{i=1}^I \alpha_{i_n, i} K(T - t_n)$$

Note that this objective already includes a  $\ell_1$  regularization on  $\alpha$ . Simply adding another  $\ell_1$  term cannot give us a sparse solution, i.e.:

$$\min \mathcal{J}(\alpha) - C \sum_{i,j=1}^I |\alpha_{ij}| = \sum_{n=1}^N \sum_{\ell=1}^{n-1} \eta_{\ell n} \log(\alpha_{i_\ell, i_n} \kappa(t_n - t_\ell)) - \left( \sum_{n=1}^N K(T - t_n) - C \right) \sum_{i=1}^I \alpha_{i_n, i}$$

yields

$$\alpha_{ij} = \left( \frac{\sum_{n=1}^N \sum_{\ell=1}^{n-1} \delta_{i_\ell, i} \delta_{i_n, j} \eta_{\ell n}}{\sum_{n=1}^N \delta_{i_n, i} K(T - t_n) - C} \right)^+$$

Note that in the above, although increasing the regularization strength  $C$  could lead to zero  $\alpha$ s, but the regularization path is not well-behaved (see Figure 1).

Alternatively adding an  $\ell_2$  term (to make a elastic net type regularization) also doesn't work:

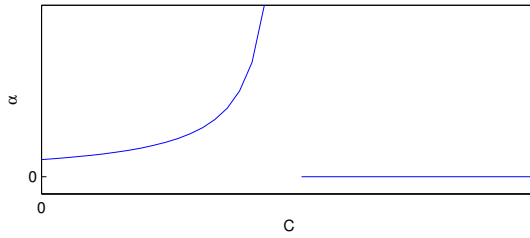


Figure 1: Regularization path.

$$\min \mathcal{J}(\alpha) - C \sum_{i,j} |\alpha_{ij}^2|$$

yields

$$\alpha_{ij} = \frac{-B + \sqrt{B^2 + 4C \sum_{n=1}^N \sum_{l=1}^{n-1} \delta_{i_l, i} \delta_{i_n, j} \eta_{ln}}}{2C}$$

where  $B = \sum_{n=1}^N \delta_{i_n, i} K(T - t_n)$ . Unfortunately, this solution is not sparse either (note that  $\eta_{ln} = 0$  iff the corresponding  $\alpha_{ij} = 0$ ). Also, this solution is not well-posed when the regularization strength is small (e.g.,  $C \rightarrow 0$ ).

**Infeasibility of sparseness in MHPs: the reason.** What happened? Why the sparsity encouraging penalties cannot yield solutions that are supposed to be sparse? The key reason has to do with the way we lower-bounding the log-likelihood. Particular, by introducing the branchings and breaking down the log-sum in the log-likelihood  $\mathcal{L}$ , the infectivities  $\alpha$ s become singular points in the lower-bound  $\mathcal{J}$ , hence, sparsity is infeasible if  $\mathcal{J}$  is optimized.

As such, in order to obtain sparsity, we need to optimize the log-likelihood  $\mathcal{L}$  directly instead of its lower-bound  $\mathcal{J}$ . For a plain MHP model, this is not a big deal as the optimization can be solved with any convex optimization algorithm. However, for the mixture of MHPs model, optimizing  $\mathcal{L}$  directly is intractable. For example, if we didn't introduce the branchings to break down the log-sum, each meme-identity variable  $Z$  would be coupled with all  $Z$ s in the past and all  $Z$ s in the future in a very complicated nonlinear manner, inference of which is unimaginably troublesome.