

Supplementary Material for "Fully Neural Network based Model for General Temporal Point Processes"

Evaluating the $\partial Z_i(\tau)/\partial\tau$

In our method, we need to calculate $\partial Z_i(\tau)/\partial\tau$, the derivative of the output $Z_i(\tau)$ of the cumulative hazard function network with respect to the input τ . The derivative $\partial Z_i(\tau)/\partial\tau$ can be computed using automatic differentiation. Automatic differentiation is a method to evaluate the derivative of an arbitrary function [13]. We here employ backpropagation which is a special case of automatic differentiation. Backpropagation is sometimes considered as a specific method to evaluate the gradient of the loss function for a multi-layer perceptron, however backpropagation is a fairly general method to compute the gradient of an arbitrary function. Please refer to [14] for more details on backpropagation.

We first formulate the operation of the cumulative hazard function network. We denote the output of the j th layer in the cumulative hazard function network by $\mathbf{Y}^{(j)} \in \mathbb{R}^{m_j}$ and write the feedforward operation as follows:

$$\mathbf{Y}^{(j)} = f^{(j)}(W^{(j)}\mathbf{Y}^{(j-1)} + \mathbf{b}^{(j)}) \quad (j = 1, 2, \dots, L), \quad (15)$$

where $f^{(j)}$, $W^{(j)}$, and $\mathbf{b}^{(j)}$ represent the activation function, the weight matrix, and the bias term for the j th layer, respectively. The input to the cumulative hazard function network is given as $\mathbf{Y}^{(0)} = (\tau, \mathbf{h}_i^T)^T$, where \mathbf{h}_i is the RNN output. The L th layer is the output layer, and we have $Z_i(\tau) = \mathbf{Y}^{(L)}$.

To calculate $\partial Z_i(\tau)/\partial\tau$, we introduce a notation

$$\mathbf{y}^{(j)} = \frac{\partial Z_i(\tau)}{\partial \mathbf{Y}^{(j)}}, \quad (16)$$

and employ a chain rule to obtain a backward operation for $\mathbf{y}^{(j)}$ as

$$\mathbf{y}^{(j-1)} = \left(\frac{\partial \mathbf{Y}^{(j)}}{\partial \mathbf{Y}^{(j-1)}} \right)^\top \mathbf{y}^{(j)} \quad (17)$$

$$= W^{(j)\top} \text{diag}(f'^{(j)}(W^{(j)}\mathbf{Y}^{(j-1)} + \mathbf{b}^{(j)})) \mathbf{y}^{(j)} \quad (j = 1, 2, \dots, L), \quad (18)$$

where we have $\mathbf{y}^{(L)} = 1$. By recursively applying the backward operation, we finally obtain the desired derivative as

$$\partial Z_i(\tau)/\partial\tau = [\mathbf{y}^{(0)}]_1, \quad (19)$$

where $[\mathbf{y}^{(0)}]_1$ represents the first element of the vector $\mathbf{y}^{(0)}$. We can then construct an extended network that produces both $Z_i(\tau)$ and $\partial Z_i(\tau)/\partial\tau$ by concatenating the original feedforward network of eq. (15) and the backward network of eq. (18), as shown in Fig. S1.

Evaluating the gradient of the loss function

The loss function of our model, the negative log-likelihood function, depends on both $Z_i(\tau)$ and $\partial Z_i(\tau)/\partial\tau$. The derivative $\partial Z_i(\tau)/\partial\tau$ is computed using backpropagation, as seen in the previous section. Then the computational graph to evaluate the loss function is then obtained as in Fig. S1. The gradient of the loss function with respect to the parameters can be evaluated by applying backpropagation to the computational graph.

Here we use backpropagation twice for training the model, the first for calculating $\partial Z_i(\tau)/\partial\tau$ and the second for calculating the gradient of the loss function. This kind of procedure is sometimes called as double backpropagation [15].

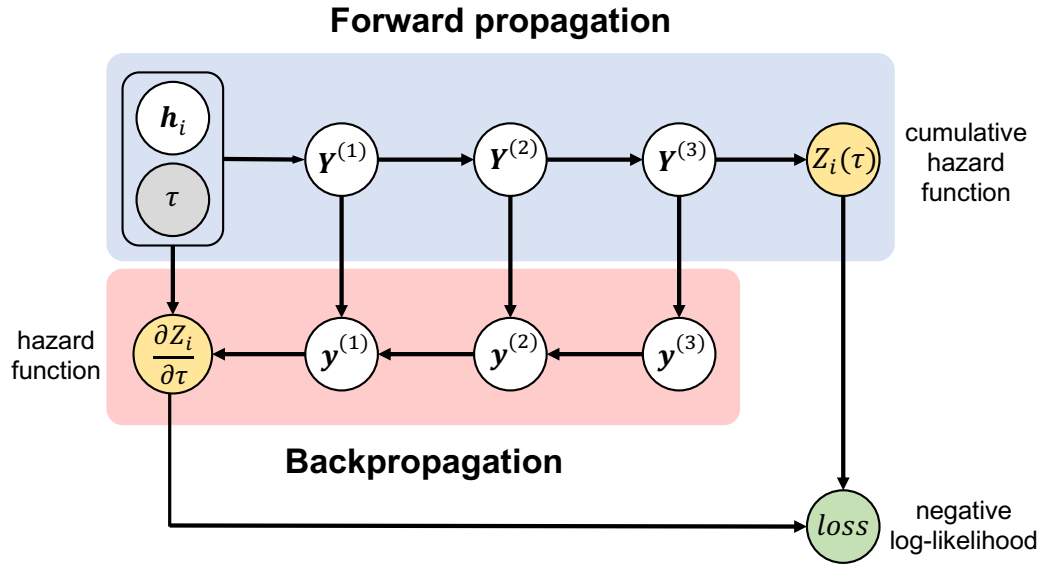


Figure S1: A network structure of our model ($L = 4$). We omitted the computational graph for the operation of the RNN for simplicity.