

TEDnet Documentation

1 Model Specification

1.1 Mixture Density Outputs

$$x_t \in \mathbb{R}^{1 \times k} \quad (1)$$

$$\pi_t \in \mathbb{R} \quad (2)$$

$$\mu_t \in \mathbb{R}^{1 \times k} \quad (3)$$

$$\Sigma_t \in \mathbb{R}^{k \times k} \quad (4)$$

$$y_t = (\{\pi_t^j, \mu_t^j, \Sigma_t^j\}_{j=1}^M) \quad (5)$$

From Inputs to Mixture Gaussian Components

$$\hat{y}_t = (\{\hat{\pi}_t^j, \hat{\mu}_t^j, \hat{\Sigma}_t^j\}_{j=1}^M) = b_y + \sum_{n=1}^N W_{h^n} h_t^n \quad (6)$$

$$\pi_t^j = \frac{\exp(\hat{\pi}_t^j)}{\sum_{j'=1}^M \exp(\hat{\pi}_t^{j'})} \implies \pi_t^j \in (0, 1), \sum_{j=1}^M \pi_t^j = 1 \quad (7)$$

$$\mu_t^j = \hat{\mu}_t^j \quad (8)$$

There are two options for the number of input values necessary for describing the covariance matrix values. We have one option where the covariance is treated normally, and since the covariance matrix is symmetric we only need to account for $k(k-1)/2$ value (either upper or lower triangular value). The second option is to treat the covariance matrix as a strickly diagonal matrix and therefore we only account for k values.

Option 1: Full Covariance Matrix

$$\hat{\Sigma}_t \in \mathbb{R}^{1 \times (k(k-1)/2)} \quad (9)$$

$$\Sigma_{t,(l,m)} = \begin{cases} \hat{\Sigma}_{t,(1, l \cdot (l-1)/2+m)}, & \text{if } m \leq l \\ \hat{\Sigma}_{t,(1, m \cdot (m-1)/2+l)}, & \text{if } m > l \end{cases} \quad (10)$$

Option 2: Diagonal Covariance Matrix

$$\hat{\Sigma}_t \in \mathbb{R}^{1 \times k} \quad (11)$$

$$\Sigma_{t,(l,m)} = \begin{cases} \hat{\Sigma}_{t,(1,l)}, & \text{if } l = m \\ 0, & \text{otherwise} \end{cases} \quad (12)$$

Probability and Loss

The probability density $\Pr(x_{t+1} \mid y_t)$ of the next input x_{t+1} given the output y_t is defined as follows:

$$\Pr(x_{t+1} \mid y_t) = \sum_{j=1}^M \pi_t^j \mathcal{N}(x_{t+1} \mid \mu_t^j, \Sigma_t^j) \quad (13)$$

where

$$\mathcal{N}(x \mid \mu, \Sigma) = \frac{1}{\sqrt{2\pi |\Sigma|}} \exp\left[-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right] \quad (14)$$

Using eq 13 we can determine sequence loss as follows:

$$\mathcal{L}(x) = \sum_{t=1}^T -\log\left(\sum_{j=1}^M \pi_t^j \mathcal{N}(x_{t+1} \mid \mu_t^j, \Sigma_t^j)\right) \quad (15)$$

Deriving Gradient of Loss¹

Using responsibilities γ_t^j will make derivations clearer:

$$\hat{\gamma}_t^j = \pi_t^j \mathcal{N}(x_{t+1} \mid \mu_t^j, \Sigma_t^j) \quad (16)$$

$$\gamma_t^j = \frac{\hat{\gamma}_t^j}{\sum_{j'}^M \hat{\gamma}_t^{j'}} \quad (17)$$

$$\frac{\delta \mathcal{L}(x)}{\delta \pi_t^j} = -\gamma_t^j \frac{1}{\pi_t^j} \quad (18)$$

$$\frac{\delta \mathcal{L}(x)}{\delta \mu_t^j} = -\gamma_t^j ([\Sigma_t^j]^{-1}(x_{t+1} - \mu_t^j)) \quad (19)$$

$$\frac{\delta \mathcal{L}(x)}{\delta \Sigma_t^j} = -\gamma_t^j \left(\frac{1}{2} [-[\Sigma_t^j]^{-1} + [\Sigma_t^j]^{-1}(x_{t+1} - \mu_t^j)(x_{t+1} - \mu_t^j)^T [\Sigma_t^j]^{-1}] \right) \quad (20)$$

References

1. Petersen, K., Pedersen, M.: The Matrix Cookbook. Eq. 390–396, 44–45 (2008)