

Powering Hidden Markov Model by Generative Models

Firsthand Scientists

March 5, 2019

1 Notation

Time is indexed by subscript and sequence is denoted by underline. \mathbf{x}_t is signal at time t . The sequential time is denoted by $\underline{\mathbf{x}} = [\mathbf{x}_1, \dots, \mathbf{x}_T]^\top$, where $[\cdot]^\top$ means transpose and T is the length of the sequence. Sequential signal or clip uses underline notation and is indexed by superscript, for instance $\underline{\mathbf{x}}^{(r)}$ means the r -th sequential signal, where $r = 1, 2, \dots, R$, and $\underline{\mathbf{x}}^{(r)} = [\mathbf{x}_1^{(r)}, \mathbf{x}_2^{(r)}, \dots, \mathbf{x}_{T^{(r)}}^{(r)}]$ with length $T^{(r)}$. Note different sequential signal $\underline{\mathbf{x}}^{(r)}$ could have different lengths.

The hypothesis of Hidden Markov Model (HMM): $\mathcal{H} := \{\mathbf{H} | \{\mathcal{S}, \mathbf{q}, A, p(\mathbf{x}|s; \Phi)\}\}$,

- \mathcal{S} is the set of states of HMM \mathbf{H} ;
- $\mathbf{q} = [q_1, q_2, \dots, q_{|\mathcal{S}|}]^\top$ initial distribution of HMM \mathbf{H} with $|\mathcal{S}|$ is cardinality of \mathcal{S} , $q_k = p(s = k)$ for random state variable s .
- A is the transition matrix for the HMM \mathbf{H} of size $|\mathcal{S}| \times |\mathcal{S}|$.
- Observable signal density $p(\mathbf{x}|s; \Phi)$ given hidden state sequence, where Φ is the parameter set that defines this conditional probabilistic model.

2 Problem Statement

Given a empirical distribution $\hat{p}(\underline{\mathbf{x}}) = \frac{1}{R} \sum_{r=1}^R \delta_{\underline{\mathbf{x}}^{(r)}}(\underline{\mathbf{x}})$. We want to find a probabilistic model such that:

$$\min KL(\hat{p}(\underline{\mathbf{x}}) \| p(\underline{\mathbf{x}})) \quad (1)$$

where $KL(\cdot \| \cdot)$ denotes the Kullback-Leibler divergence.

When we use HMM to model the empirical distribution and approach the unknown true distribution, the problem boils down to:

$$\operatorname{argmax}_{\mathbf{H} \in \mathcal{H}} p(\underline{\mathbf{X}}; \mathbf{H}) \quad (2)$$

where $\underline{\mathbf{X}} = [\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}, \dots, \underline{\mathbf{x}}^{(R)}]$

The problem can be reformulated as

$$\operatorname{argmax}_{\mathbf{H} \in \mathcal{H}} \sum_{r=1}^R \log p(\underline{\mathbf{x}}^{(r)}; \mathbf{H}) \quad (3)$$

for independent identical distributed assumption of $\underline{\mathbf{x}}$.

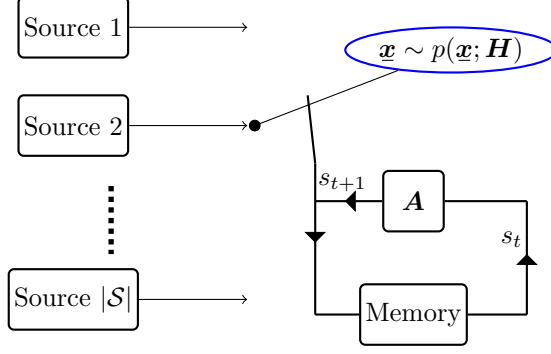


Figure 1: HMM Model defined by $\mathbf{H} = \{\mathcal{S}, \mathbf{q}, \mathbf{A}, p(\mathbf{x}|\mathbf{s}; \Phi)\}$

3 Proposal

Since model \mathbf{H} contains hidden sequential variable \mathbf{s} , we can not directly solve the maximum likelihood problem in Equation 3. We use expectation maximization (EM) to address the hidden variable problem by

- E-step: The “expected likelihood” function:

$$\mathcal{Q}(\mathbf{H}; \mathbf{H}^{\text{old}}) = \mathbb{E}_{p(\mathbf{s}^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}})} \left[\sum_{r=1}^R \log p(\mathbf{x}^{(r)}, \mathbf{s}^{(r)}; \mathbf{H}) \right] \quad (4)$$

- M-step: the optimization step:

$$\max_{\mathbf{H}} \mathcal{Q}(\mathbf{H}; \mathbf{H}^{\text{old}}) \quad (5)$$

The Equation 5 can be reformulated as:

$$\max_{\mathbf{H}} \mathcal{Q}(\mathbf{H}; \mathbf{H}^{\text{old}}) = \max_{\mathbf{q}} \mathcal{Q}(\mathbf{q}; \mathbf{H}^{\text{old}}) + \max_{\mathbf{A}} \mathcal{Q}(\mathbf{A}; \mathbf{H}^{\text{old}}) + \max_{\Phi} \mathcal{Q}(\Phi; \mathbf{H}^{\text{old}}) \quad (6)$$

where

$$\mathcal{Q}(\mathbf{q}; \mathbf{H}^{\text{old}}) = \sum_{r=1}^R \mathbb{E}_{p(\mathbf{s}^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}})} \left[\log p(\mathbf{s}_1^{(r)}; \mathbf{q}) \right] \quad (7)$$

$$\mathcal{Q}(\mathbf{A}; \mathbf{H}^{\text{old}}) = \sum_{r=1}^R \mathbb{E}_{p(\mathbf{s}^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}})} \left[\log \sum_{t=1}^{T^{(r)}-1} p(s_{t+1}^{(r)}|s_t^{(r)}; \mathbf{A}) \right] \quad (8)$$

$$\mathcal{Q}(\Phi; \mathbf{H}^{\text{old}}) = \sum_{r=1}^R \mathbb{E}_{p(\mathbf{s}^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}})} \left[\log p(\mathbf{x}^{(r)}|\mathbf{s}^{(r)}; \Phi) \right] \quad (9)$$

We can see that the solution of \mathbf{H} depends on the posterior probability $p(\mathbf{s}|\mathbf{x}; \mathbf{H})$. Though the evaluation of posterior according to Bayesian theorem is simple, the computation complexity of $p(\mathbf{s}|\mathbf{x}; \mathbf{H})$ grows exponentially with the length of \mathbf{s} . Therefore, we would employ Forward/Backward algorithm [1] to do the posterior computation efficiently. The marginal $p(s_t|\mathbf{x}; \mathbf{H})$ is also efficiently computed as the joint posterior.

We summarize the optimization algorithm as:

Algorithm 1 Meta algorithm for HMM powered by Generative Models

```

1: Input: Building  $\mathbf{H}^{\text{old}}$ ,  $\mathbf{H} \in \mathcal{H}$  gives:
    $\mathbf{H}^{\text{old}} = \{\mathcal{S}, \mathbf{q}^{\text{old}}, \mathbf{A}^{\text{old}}, p(\mathbf{x}|s; \Phi^{\text{old}})\}$ ,  $\mathbf{H} = \{\mathcal{S}, \mathbf{q}, \mathbf{A}, p(\mathbf{x}|s; \Phi)\}$  ,
2: Initialize  $\mathbf{H}$ 
3:  $\mathbf{H}^{\text{old}} \leftarrow \mathbf{H}$ 
4: for  $\mathbf{H}$  not converge do
5:   Sample a batch of data  $\{\mathbf{x}^{(r)}\}_{r=1}^{R_b}$  from the dataset  $\hat{p}(\mathbf{x})$ 
6:   Compute  $p(s_t^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}})$ ,  $p(s_t^{(r)}, s_{t+1}^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}})$  by forward/backward algorithm;
7:    $\mathbf{q} \leftarrow \underset{\mathbf{q}}{\text{argmin}} \mathcal{Q}(\mathbf{q}; \mathbf{H}^{\text{old}})$  by Equation 13;
8:    $\mathbf{A} \leftarrow \underset{\mathbf{A}}{\text{argmin}} \mathcal{Q}(\mathbf{A}; \mathbf{H}^{\text{old}})$  by Equation 15;
9:    $\Phi \leftarrow \underset{\Phi}{\text{argmin}} \mathcal{Q}(\Phi; \mathbf{H}^{\text{old}})$  by calling Algorithm 2 or 3;
10:   $\mathbf{H}^{\text{old}} \leftarrow \mathbf{H}$ 
11: end for

```

3.1 Initial Probability Update

Equation 7 can be written as:

$$\begin{aligned}
\mathcal{Q}(\mathbf{q}; \mathbf{H}^{\text{old}}) &= \sum_{r=1}^R \sum_{\mathbf{s}^{(r)}} p(\mathbf{s}^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}}) \log p(s_1^{(r)}; \mathbf{q}) \\
&= \sum_{r=1}^R \sum_{s_1^{(r)}=1}^{|\mathcal{S}|} \sum_{s_2^{(r)}=1}^{|\mathcal{S}|} \cdots \sum_{s_{T^r}^{(r)}=1}^{|\mathcal{S}|} p(s_1^{(r)}, s_2^{(r)}, \dots, s_{T^r}^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}}) \log p(s_1^{(r)}; \mathbf{q}) \tag{10}
\end{aligned}$$

$$= \sum_{r=1}^R \sum_{s_1^{(r)}=1}^{|\mathcal{S}|} p(s_1^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}}) \log p(s_1^{(r)}; \mathbf{q}) \tag{11}$$

Since $p(s_1^{(r)}; \mathbf{H})$ is the probability of initial state of HMM \mathbf{H} for r -th sequential, actually $q_i = p(s_1^{(r)} = i; \mathbf{H})$ for $i = 1, 2, \dots, |\mathcal{S}|$. Solution to problem:

$$\begin{aligned}
\mathbf{q}^{\text{new}} &= \underset{\mathbf{q}}{\text{argmax}} \mathcal{Q}(\mathbf{q}; \mathbf{H}^{\text{old}}), \\
&\text{s.t. } \sum_{i=1}^{|\mathcal{S}|} q_i = 1 \\
&q_i \geq 0, \forall i. \tag{12}
\end{aligned}$$

is

$$q_i = \frac{1}{R} \sum_{r=1}^R p(s_1^{(r)} = i|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}}), \forall i = 1, 2, \dots, |\mathcal{S}|. \tag{13}$$

3.2 Transition Probability Update

Equation 8 can be written as

$$\begin{aligned}
\mathcal{Q}(\mathbf{A}; \mathbf{H}^{\text{old}}) &= \sum_{r=1}^R \mathbb{E}_{p(\mathbf{s}^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}})} \left[\log \sum_{t=1}^{T^{(r)}-1} p(s_{t+1}^{(r)}|s_t^{(r)}; A) \right] \\
&= \sum_{r=1}^R \sum_{\mathbf{s}^{(r)}} p(\mathbf{s}^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}}) \log \sum_{t=1}^{T^{(r)}-1} p(s_{t+1}^{(r)}|s_t^{(r)}; A) \\
&= \sum_{r=1}^R \sum_{t=1}^{T^{(r)}-1} \sum_{s_t^{(r)}=1}^{|S|} \sum_{s_{t+1}^{(r)}=1}^{|S|} p(s_t^{(r)}, s_{t+1}^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}}) \log p(s_{t+1}^{(r)}|s_t^{(r)}; A)
\end{aligned} \tag{14}$$

Since $\mathbf{A}_{i,j} = p(s_{t+1}^{(r)} = j|s_t^{(r)} = i; A)$ where $A_{i,j}$ is the element of transition matrix A , the solution to problem:

$$\begin{aligned}
\mathbf{A}^{\text{new}} &= \underset{\mathbf{A}}{\text{argmax}} \mathcal{Q}(\mathbf{A}; \mathbf{H}^{\text{old}}), \\
\text{s.t. } &\mathbf{A} \cdot \mathbf{1} = \mathbf{1} \\
&\mathbf{A}^\top \cdot \mathbf{1} = \mathbf{1} \\
&\mathbf{A}_{i,j} \geq 0.
\end{aligned} \tag{15}$$

is

$$\mathbf{A}_{i,j}^{\text{new}} = \frac{\bar{\xi}_{i,j}}{\sum_{k=1}^{|S|} \bar{\xi}_{i,k}}, \tag{16}$$

where

$$\bar{\xi}_{i,j} = \sum_{r=1}^R \sum_{t=1}^{T^{(r)}-1} p(s_t^{(r)} = i, s_{t+1}^{(r)} = j|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}}) \tag{17}$$

3.3 Generative Model Update

Equation 9 can be rewritten as

$$\begin{aligned}
\mathcal{Q}(\Phi; \mathbf{H}^{\text{old}}) &= \sum_{r=1}^R \sum_{\mathbf{s}^{(r)}} p(\mathbf{s}^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}}) \log p(\mathbf{x}^{(r)}|\mathbf{s}^{(r)}; \Phi) \\
&= \sum_{r=1}^R \sum_{t=1}^{T^{(r)}-1} \sum_{s_t^{(r)}=1}^{|S|} p(s_t^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}}) \log p(\mathbf{x}_t^{(r)}|s_t^{(r)}; \Phi).
\end{aligned} \tag{18}$$

Then the third subproblem of Equation 6 becomes:

$$\begin{aligned}
&\underset{\Phi}{\text{argmax}} \mathcal{Q}(\Phi; \mathbf{H}^{\text{old}}), \\
\text{s.t. } &p(\mathbf{x}|s; \Phi) \text{ is our general model}
\end{aligned} \tag{19}$$

It could be seen from Equation 18 that the key to update generate model is to evaluate $p(\mathbf{x}|s; \Phi)$ for all $s \in \mathcal{S}$. In Forward/Backward algorithm, evaluation of $p(\mathbf{x}|s; \Phi)$ is also all what is needed to compute $p(s|\mathbf{x}; \Phi)$. In the following two subsections, we will provide two neural network based generative models that fulfill this requirement and also have high capability for complex signal modeling.

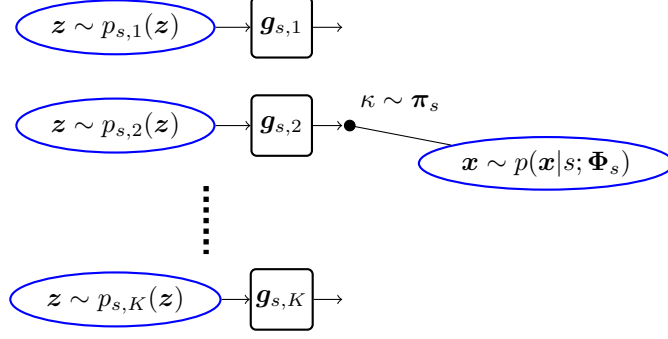


Figure 2: Source s powered by generator mixed generative model (GenHMM).

3.3.1 Generator Mixed HMM (GenHMM)

Base on the generative model assumption:

$$p(\mathbf{x}|s; \Phi_s) = \sum_{\kappa=1}^K \pi_{s,\kappa} p(\mathbf{x}|s, \kappa; \Phi_{s,\kappa}) \quad (20)$$

where $\Phi_s = \{\Phi_{s,\kappa} | \kappa = 1, 2, \dots, K\}$, and

$$\sum_{\kappa=1}^K \pi_{s,\kappa} = 1 \quad (21)$$

The help function should be revised to deal with the new latent variable κ into:

$$\begin{aligned} \mathcal{Q}(\Phi; \mathbf{H}^{\text{old}}) &= \sum_{r=1}^R \sum_{t=1}^{T^{(r)}-1} \sum_{s_t^{(r)}=1}^{|S|} \sum_{\kappa_t^{(r)}=1}^K p(s_t^{(r)}, \kappa_t^{(r)} | \underline{\mathbf{x}}^{(r)}; \mathbf{H}^{\text{old}}) \log p(\kappa_t^{(r)}, \mathbf{x}_t^{(r)} | s_t^{(r)}; \Phi) \\ &= \sum_{r=1}^R \sum_{t=1}^{T^{(r)}-1} \sum_{s_t^{(r)}=1}^{|S|} \sum_{\kappa_t^{(r)}=1}^K p(s_t^{(r)} | \underline{\mathbf{x}}^{(r)}; \mathbf{H}^{\text{old}}) p(\kappa_t^{(r)} | s_t^{(r)}, \underline{\mathbf{x}}^{(r)}; \mathbf{H}^{\text{old}}) \left[\log \pi_{s_t^{(r)}, \kappa_t^{(r)}} \right. \\ &\quad \left. + \log p(\mathbf{x}_t^{(r)} | s_t^{(r)}, \kappa_t^{(r)}; \Phi) \right] \end{aligned} \quad (22)$$

In Equation 22, $p(s_t | \underline{\mathbf{x}}, \mathbf{H}^{\text{old}})$ is computed by forward/backward algorithm. The posterior of κ is:

$$\begin{aligned} p(\kappa | s, \underline{\mathbf{x}}; \mathbf{H}^{\text{old}}) &= \frac{p(\kappa, \underline{\mathbf{x}} | s; \mathbf{H}^{\text{old}})}{p(\underline{\mathbf{x}} | s, \mathbf{H}^{\text{old}})} \\ &= \frac{\pi_{s,\kappa} p(\underline{\mathbf{x}} | s, \kappa, \mathbf{H}^{\text{old}})}{\sum_{\kappa=1}^K \pi_{s,\kappa} p(\underline{\mathbf{x}} | s, \kappa, \mathbf{H}^{\text{old}})} \\ &= \frac{\pi_{s,\kappa} p(\mathbf{x} | s, \kappa, \mathbf{H}^{\text{old}})}{\sum_{\kappa=1}^K \pi_{s,\kappa} p(\mathbf{x} | s, \kappa, \mathbf{H}^{\text{old}})} \end{aligned} \quad (23)$$

where the last equation is due to the fact that only \mathbf{x}_t among sequence $\underline{\mathbf{x}}$ depends on s_t, κ_t .

The latent prior for mixture of each source s is obtained by solving the following problem:

$$\begin{aligned} \pi_{s,\kappa} &= \underset{\pi_{s,\kappa}}{\operatorname{argmax}} \mathcal{Q}(\Phi; \mathbf{H}^{\text{old}}) \\ \text{s.t. } &\sum_{\kappa=1}^K \pi_{s,\kappa} = 1 \end{aligned} \quad (24)$$

which gives the solution:

$$\pi_{s,\kappa} = \frac{\sum_{r=1}^R \sum_{t=1}^{T^{(r)}-1} p(s_t^{(r)} = s, \kappa_t^{(r)} = \kappa | \underline{x}^{(r)}; \mathbf{H}^{\text{old}})}{\sum_{k=1}^K \sum_{r=1}^R \sum_{t=1}^{T^{(r)}-1} p(s_t^{(r)} = s, \kappa_t^{(r)} = k | \underline{x}^{(r)}; \mathbf{H}^{\text{old}})} \quad (25)$$

where $p(s, \kappa | \underline{x}; \mathbf{H}^{\text{old}}) = p(s | \underline{x}; \mathbf{H}^{\text{old}}) p(\kappa | s, \underline{x}; \mathbf{H}^{\text{old}})$ that can be computed by results of forward/backward and Equation 23.

In implementation, GenHMM uses the generator mixed emission model with κ -th component as:

$$\begin{aligned} & p(\mathbf{x} | s, \kappa; \Phi_{s,\kappa}) \\ &= p_{s,\kappa}(\mathbf{z}) \left| \det \left(\frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right) \right| \\ &= \mathcal{N}(\mathbf{f}_{s,\kappa}(\mathbf{x}); \boldsymbol{\mu}_{s,\kappa}, \mathbf{C}_{s,\kappa}) \left| \det \left(\frac{\partial \mathbf{f}_{s,\kappa}(\mathbf{x})}{\partial \mathbf{x}} \right) \right| \end{aligned} \quad (26)$$

where $\mathbf{f}_{s,\kappa} = \mathbf{g}_{s,\kappa}^{-1}$ is defined by parameter set $\boldsymbol{\theta}_{s,\kappa}$, and $\Phi_{s,\kappa} = \{\boldsymbol{\mu}_{s,\kappa}, \mathbf{C}_{s,\kappa}, \boldsymbol{\theta}_{s,\kappa}\}$, $\kappa = 1, 2, \dots, K$. We can start by setting $\boldsymbol{\mu}_{s,\kappa} = \mathbf{0}$ and $\mathbf{C}_{s,\kappa} = \text{diag}(\mathbf{1})$.

Algorithm 2 M-step w.r.t. Φ powered by GenMM

- 1: **Input:** Latent mixture distribution: $\mathcal{N}(\mathbf{z}; \mathbf{0}, \text{diag}(\mathbf{1}))$, $\forall s \in \mathcal{S}, \kappa = 1, 2, \dots, K$;
Empirical distribution $\hat{p}(\mathbf{x})$ of dataset;
 - 2: Set a total number of epochs T of training as stop criterion. A learning rate η .
 - 3: **for** epoch $t < T$ **do**
 - 4: Sample a batch of data $\{\underline{x}^{(r)}\}_{r=1}^{R_b}$ from dataset $P_d(\underline{x})$
 - 5: Compute $p(s_t^{(r)}, \kappa_t^{(r)} | \underline{x}^{(r)}; \mathbf{H}^{\text{old}})$ by Forward/backward and Equation 23.
 - 6: Compute loss $\mathcal{Q}(\Phi, \mathbf{H}^{\text{old}})$ in Equation 22
 - 7: $\partial \boldsymbol{\theta}_s \leftarrow \nabla_{\boldsymbol{\theta}} - \frac{1}{R_b} \mathcal{Q}(\Phi, \mathbf{H}^{\text{old}})$
 - 8: $\boldsymbol{\theta}_s \leftarrow \boldsymbol{\theta}_s - \eta \cdot \partial \boldsymbol{\theta}_s$, $\forall s \in \mathcal{S}$
 - 9: **end for**
 - 10: Update $\pi_{s,\kappa}$, $\forall s \in \mathcal{S}, \kappa = 1, 2, \dots, K$, according to Equation 25
 - 11: Assemble $\Phi = \{\Phi_s | s \in \mathcal{S}\}$ for $\Phi_s = \{\boldsymbol{\theta}_{s,\kappa}, \boldsymbol{\mu}_{s,\kappa}, \mathbf{C}_{s,\kappa} | \kappa = 1, 2, \dots, K\}$.
-

3.3.2 Latent-source Mixed HMM (LatMM)

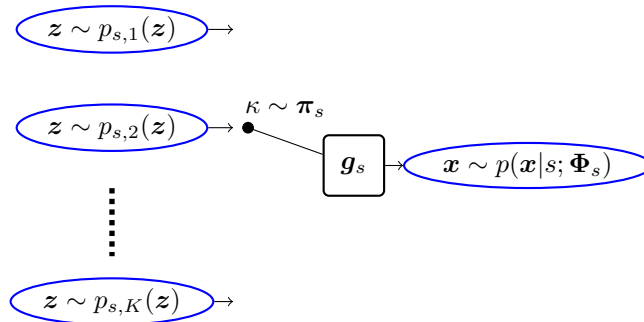


Figure 3: Source s powered by Latent mixed generative model (LatHMM).

LatHMM, (**benching marking**), the latent source mixed emission model is built by allowing different latent source sharing the same generator \mathbf{g}_s :

$$\begin{aligned} & p(\mathbf{x}|s, \kappa; \Phi_s) \\ &= p_{s,\kappa}(\mathbf{z}) \left| \det \left(\frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right) \right| \\ &= \mathcal{N}(\mathbf{f}_s(\mathbf{x}); \boldsymbol{\mu}_{s,\kappa}, \mathbf{C}_{s,\kappa}) \left| \det \left(\frac{\partial \mathbf{f}_s(\mathbf{x})}{\partial \mathbf{x}} \right) \right| \end{aligned} \quad (27)$$

where $\mathbf{f}_s = \mathbf{g}_s^{-1}$ is defined by $\boldsymbol{\theta}$, and $\Phi_s = \{\boldsymbol{\theta}_s, \boldsymbol{\mu}_{s,\kappa}, \mathbf{C}_{s,\kappa} | \kappa = 1, 2, \dots, K\}$.

This emission generator sharing makes the posterior computation w.r.t. of κ easier:

$$p(\kappa|s, \mathbf{x}; \mathbf{H}^{\text{old}}) = \frac{\pi_{s,\kappa} p(\mathbf{x}|s, \kappa, \mathbf{H}^{\text{old}})}{\sum_{\kappa=1}^K \pi_{s,\kappa} p(\mathbf{x}|s, \kappa, \mathbf{H}^{\text{old}})} = \frac{\pi_{s,\kappa} \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{s,\kappa}, \mathbf{C}_{s,\kappa})}{\sum_{\kappa=1}^K \pi_{s,\kappa} \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{s,\kappa}, \mathbf{C}_{s,\kappa})} \Big|_{\mathbf{z}=\mathbf{f}_s(\mathbf{x})}. \quad (28)$$

Computation of $p(s_t|\mathbf{x}, \mathbf{H}^{\text{old}})$ remains the same, relying on forward/backward algorithm.

In implementation, $\mathbf{C}_{s,\kappa} = \text{diag}(\boldsymbol{\sigma}_{s,\kappa})$ for simplicity. To avoid the singularity problem of Gaussian mixture, we put a Gamma distribution as prior for $\boldsymbol{\sigma}_{s,\kappa}$, i.e. $\Gamma(\boldsymbol{\sigma}_{s,\kappa}^{-1}; a, b)$ where a and b are hyperparameter for Gamma distribution. Then the problem can be reformulated as:

$$\underset{\Phi}{\text{argmax}} \mathcal{Q}(\Phi; \mathbf{H}^{\text{old}}) + \frac{1}{K} \log \prod_{k=1}^K \Gamma(\boldsymbol{\sigma}_{s,k}^{-1}; a, b) \quad (29)$$

Apart from the emission probability model difference, rest computation is the same as [subsubsection 3.3.1](#). We summarize the algorithm as:

Algorithm 3 M-step w.r.t. Φ powered by LatMM

- 1: **Input:** Latent mixture distribution: $\sum_{k=1}^K \pi_{s,\kappa} \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{s,\kappa}, \text{diag}(\boldsymbol{\sigma}_{s,\kappa}^2))$
Empirical distribution $\hat{p}(\mathbf{x})$ of dataset;
 - 2: Set a total number of epochs T of training as stop criterion. A learning rate η . Set hyperparameter a , b for prior of $\boldsymbol{\sigma}_{s,\kappa}^{-1}, \forall k$.
 - 3: **for** epoch $t < T$ **do**
 - 4: Sample a batch of data $\{\mathbf{x}^{(r)}\}_{r=1}^{R_b}$ from dataset $P_d(\mathbf{x})$
 - 5: Compute $p(s_t^{(r)}, \kappa_t^{(r)} | \mathbf{x}^{(r)}; \mathbf{H}^{\text{old}})$ by Forward/backward and [Equation 28](#).
 - 6: Compute loss in [Equation 29](#)
 - 7: $\partial \boldsymbol{\theta}_s, \partial \boldsymbol{\mu}_{s,\kappa}, \partial \boldsymbol{\sigma}_{s,\kappa} \leftarrow \nabla_{\boldsymbol{\theta}, \boldsymbol{\mu}_{s,\kappa}, \boldsymbol{\sigma}_{s,\kappa}} - \frac{1}{R_b} \mathcal{Q}(\Phi, \mathbf{H}^{\text{old}}) - \frac{1}{K} \sum_{k=1}^K \log \Gamma(\boldsymbol{\sigma}_{s,k}^{-1}; a, b)$
 - 8: $\boldsymbol{\theta}_s \leftarrow \boldsymbol{\theta}_s - \eta \cdot \partial \boldsymbol{\theta}_s, \forall s \in \mathcal{S}$
 - 9: $\boldsymbol{\mu}_{s,\kappa} \leftarrow \boldsymbol{\mu}_{s,\kappa} - \eta \cdot \partial \boldsymbol{\mu}_{s,\kappa}, \forall \kappa, s$
 - 10: $\boldsymbol{\sigma}_{s,\kappa} \leftarrow \boldsymbol{\sigma}_{s,\kappa} - \eta \cdot \partial \boldsymbol{\sigma}_{s,\kappa}, \forall \kappa, s$
 - 11: **end for**
 - 12: Update π_k according to [Equation 25](#)
 - 13: Assemble $\Phi = \{\Phi_s | s \in \mathcal{S}\}$ for $\Phi_s = \{\boldsymbol{\theta}_s, \boldsymbol{\mu}_{s,\kappa}, \mathbf{C}_{s,\kappa} | \kappa = 1, 2, \dots, K\}$.
-

3.3.3 Generator-shared HMM (GSHMM)...to be continued

[to be continued...](#)

Alternatively, we can use a latent-source mixed HMM (LatM-HMM) where different latent source share the same generator functioning as feature mapping. Then the generator of the LatM-HMM is defined as

$$\{g|g : \mathbf{z} \rightarrow \mathbf{x}, s \in \mathcal{S}, \mathbf{z} \sim p_s(\mathbf{z})\}. \quad (30)$$

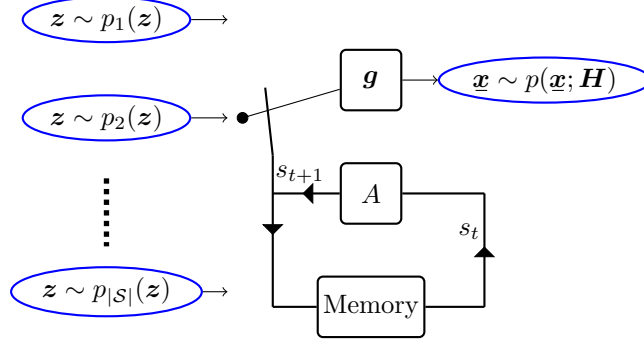


Figure 4: LatM-HMM Model defined by $\mathbf{H} = \{\mathcal{S}, \mathbf{q}, A, p(\mathbf{x}|\mathbf{s}; \Phi)\}$

We use $\mathbf{f} = \mathbf{g}^{-1}$ to denote inverse of \mathbf{g} and use θ to denote the parameter set of \mathbf{g} . Then the conditional probability for LatM-HMM is modeled as

$$\begin{aligned} p(\mathbf{x}|\mathbf{s}; \Phi) &= p_s(\mathbf{z}) \left| \det \left(\frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right) \right| \\ &= p_s(\mathbf{f}(\mathbf{x})) \left| \det \left(\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \right) \right| \end{aligned} \quad (31)$$

The parameter set for this model to be decide is $\Phi = \{\theta, \omega_s, \forall s \in \mathcal{S}\}$. Then the problem in Equation 19 can be reformulated as:

$$\begin{aligned} &\max_{\Phi} \mathcal{Q}(\Phi; \mathbf{H}^{\text{old}}) \\ &= \max_{\theta, \omega_s, \forall s \in \mathcal{S}} \sum_{r=1}^R \sum_{t=1}^{T^{(r)}-1} \sum_{s_t^{(r)}=1}^{|\mathcal{S}|} p(s_t^{(r)}|\mathbf{x}^{(r)}; \mathbf{H}^{\text{old}}) \left[\log p_{s_t^{(r)}}(\mathbf{f}(\mathbf{x}_t^{(r)})) + \log \left| \det \left(\frac{\partial \mathbf{f}(\mathbf{x}_t^{(r)})}{\partial \mathbf{x}_t^{(r)}} \right) \right| \right]. \end{aligned} \quad (32)$$

3.3.4 Generator Shared HMM (GSHMM)

To be continued...

4 On Implementation of acoustic signal

Found a HMM python lib that basics provide needed API for us, see [hmmlearn](#). Saikat also has suggestion.

For problem Equation 19 we are going to use our generative models to solve. I have the following consideration to revised our LatMM and GenMM for this application:

- Use factorized model instead of additive mixture model, to make likelihood computation logarithm domain compatible;
- Use full EM fashion instead of mini-batch fashion for training: store generative model as old for EM, there are always two neural networks working, one old for probability evaluation and one new for optimization.

5 On Implementation of Planning

Refer to [1] and its [experiments](#).

References

- [1] Thanard Kurutach, Aviv Tamar, Ge Yang, Stuart J. Russell, and Pieter Abbeel. Learning plannable representations with causal infogan. *CoRR*, abs/1807.09341, 2018.