# `osl-dynamics`: HMM Cost Function

C. Gohil

OHBA, Department of Psychiatry, Warneford Hospital, Oxford, OX3 7JX

(February 11, 2023)

**Abstract**

We describe the calculation of the cost function used to update the observation model parameters (state means and covariances) in the `osl-dynamics` implementation of a Hidden Markov Model (HMM).

## 1  Variational Free Energy

In variational Bayesian inference we infer a posterior distribution for model parameters, $q(.)$, by minimising the *variational free energy*, $\mathcal{F}$, given some data we have observed, $x_t$.

For the HMM, our model parameters are:

- The hidden state at each time point, $s_t$.
- The state transition probability at each time point, $\pi_t$, which is dependent on $s_{t-1}$.
- The initial state probability, $\pi_0$.
- The observation model parameters, $\theta_{\mathrm{obs}}$.

Therefore, we infer our model parameters by minimising the following variational free energy[1] [1]

$$\mathcal{F} = \iiiint q(s_{1:T})q(\pi_t)q(\pi_0)q(\theta_{\mathrm{obs}}) \log \left[ \frac{q(s_{1:T})q(\pi_t)q(\pi_0)q(\theta_{\mathrm{obs}})}{p(x_{1:T}, s_{1:T}, \pi_t, \pi_0, \theta_{\mathrm{obs}})} \right] ds_{1:T} d\pi_t d\pi_0 d\theta_{\mathrm{obs}}, \quad (1)$$

where $s_{1:T}$ and $x_{1:T}$ denote $s_1, ..., s_T$ and $x_1, ..., x_T$ respectively. However, in the `osl-dynamics` implementation of an HMM, we will not be Bayesian on $\theta_{\mathrm{obs}}$, instead of learning $q(\theta_{\mathrm{obs}})$ we will learn point estimates for $\theta_{\mathrm{obs}}$. We will learn the posterior distributions $q(s_{1:T})$, $q(\pi_t)$, $q(\pi_0)$ and point estimates for $\theta_{\mathrm{obs}}$ by minimising the following variational free energy,

$$\mathcal{F} = \iiint q(s_{1:T})q(\pi_t)q(\pi_0) \log \left[ \frac{q(s_{1:T})q(\pi_t)q(\pi_0)}{p(x_{1:T}, s_{1:T}, \pi_t, \pi_0)} \right] ds_{1:T} d\pi_t d\pi_0. \quad (2)$$

We will show that Eq. (2) implicitly depends on the point estimates for $\theta_{\mathrm{obs}}$ below.

## 2  Generative Model

The term $p(x_{1:T}, s_{1:T}, \pi_t, \pi_0)$ is determined by our generative model. For the HMM, if we were being fully Bayesian this would be [1]

$$p(x_{1:T}, s_{1:T}, \pi_t, \pi_0, \theta_{\mathrm{obs}}) = p(s_0|\pi_0)p(\pi_0) \prod_{t=1}^{T} p(x_t|s_t, \theta_{\mathrm{obs}})p(s_t|s_{t-1}, \pi_t)p(\pi_t)p(\theta_{\mathrm{obs}}). \quad (3)$$

---

[1]We have used the mean field approximation.

However, because we are learning point estimates for $\theta_{\mathrm{obs}}$ we do not have the prior $p(\theta_{\mathrm{obs}})$. We will use the following generative model,

$$p(x_{1:T}, s_{1:T}, \pi_t, \pi_0) = p(s_0|\pi_0)p(\pi_0) \prod_{t=1}^{T} p(x_t|s_t, \theta_{\mathrm{obs}})p(s_t|s_{t-1}, \pi_t)p(\pi_t), \tag{4}$$

where $\theta_{\mathrm{obs}}$ are point estimates.

We assume a multivariate normal distribution for the observed data,

$$p(x_t|s_t = k, \theta_{\mathrm{obs}}) = \mathcal{N}(m_k, C_k), \tag{5}$$

where $m_k$ and $C_k$ are the mean and covariance for state $k$ respectively. Our observation model parameters $\theta_{\mathrm{obs}}$ are the set of state means and covariances, $\theta_{\mathrm{obs}} = \{m_k, C_k\}$.

# 3 Cost Function for Learning $\theta_{\mathbf{obs}} = \{m_k, C_k\}$

We update our point estimate for $\theta_{\mathrm{obs}}$ by minimising Eq. (2). We separate Eq. (2) into the following terms[2]

$$
\begin{aligned}
\mathcal{F} = & - \iiint q(s_{1:T})q(\pi_t)q(\pi_0) \log\left[p(x_{1:T}, s_{1:T}, \pi_t, \pi_0)\right] ds_{1:T} d\pi_t d\pi_0 \\
& + \iiint q(s_{1:T})q(\pi_t)q(\pi_0) \log\left[q(s_{1:T})q(\pi_t)q(\pi_0)\right] ds_{1:T} d\pi_t d\pi_0 \\
\mathcal{F} = & - \iiint q(s_{1:T})q(\pi_t)q(\pi_0) \log\left[p(x_{1:T}, s_{1:T}, \pi_t, \pi_0)\right] ds_{1:T} d\pi_t d\pi_0 \\
& + \int q(s_{1:T}) \log\left[q(s_{1:T})\right] ds_{1:T} + \int q(\pi_t) \log\left[q(\pi_t)\right] d\pi_t + \int q(\pi_0) \log\left[q(\pi_0)\right] d\pi_0
\end{aligned}
\tag{6}
$$

Only the first term depends on $\theta_{\mathrm{obs}}$ so the rest can be ignored. Substituting Eq. (4) into the first term, we have

$$
\begin{aligned}
\mathcal{F} \propto & - \iiint q(s_{1:T})q(\pi_t)q(\pi_0) \log\left[p(x_{1:T}, s_{1:T}, \pi_t, \pi_0)\right] ds_{1:T} d\pi_t d\pi_0 \\
\propto & - \iiint q(s_{1:T})q(\pi_t)q(\pi_0) \log\left[p(s_0|\pi_0)p(\pi_0) \prod_{t=1}^{T} p(x_t|s_t, \theta_{\mathrm{obs}})p(s_t|s_{t-1}, \pi_t)p(\pi_t)\right] ds_{1:T} d\pi_t d\pi_0.
\end{aligned}
\tag{7}
$$

Again, only retaining the factors that depend on $\theta_{\mathrm{obs}}$, we have

$$
\begin{aligned}
\mathcal{F} \propto & - \iint q(s_{1:T})q(\pi_t) \log\left[\prod_{t=1}^{T} p(x_t|s_t, \theta_{\mathrm{obs}})p(s_t|s_{t-1}, \pi_t)p(\pi_t)\right] ds_{1:T} d\pi_t \\
\propto & - \sum_{t=1}^{T} \iint q(s_{1:T})q(\pi_t) \log\left[p(x_t|s_t, \theta_{\mathrm{obs}})p(s_t|s_{t-1}, \pi_t)p(\pi_t)\right] ds_{1:T} d\pi_t \\
\propto & - \sum_{t=1}^{T} \iint q(s_{1:T})q(\pi_t) \left\{\log\left[p(x_t|s_t, \theta_{\mathrm{obs}})\right] + \log\left[p(s_t|s_{t-1}, \pi_t)p(\pi_t)\right]\right\} ds_{1:T} d\pi_t
\end{aligned}
\tag{8}
$$

---

[2] We have used $\int q(\xi)d\xi = 1$ to evaluate some of the integrals.

Only retaining the term that depends on $\theta_{\text{obs}}$, we have

$$
\begin{aligned}
\mathcal{F} &\propto -\sum_{t=1}^{T} \iint q(s_{1:T})q(\pi_t) \log\left[p(x_t|s_t, \theta_{\text{obs}})\right] ds_{1:T}d\pi_t \\
&\propto -\sum_{t=1}^{T} \int q(s_{1:T}) \log\left[p(x_t|s_t, \theta_{\text{obs}})\right] ds_{1:T} \\
&\propto -\sum_{t=1}^{T} \int \ldots \int q(s_1)\ldots q(s_T) \log\left[p(x_t|s_t, \theta_{\text{obs}})\right] ds_1 \ldots ds_T \\
&\propto -\sum_{t=1}^{T} \int q(s_t) \log\left[p(x_t|s_t, \theta_{\text{obs}})\right] ds_t = \mathcal{L}.
\end{aligned}
\tag{9}
$$

Here, we have defined the negative log-likelihood loss, $\mathcal{L}$, which is minimised via stochastic gradient descent to learn the parameters $\theta_{\text{obs}}$. As $q(s_t)$ is a discrete probability distribution for the state, we can evaluate the integral as

$$
\begin{aligned}
\mathcal{L} &= -\sum_{t=1}^{T}\sum_{k=1}^{K} q(s_t = k) \log\left[p(x_t|s_t = k, \theta_{\text{obs}})\right] \\
&= -\sum_{t=1}^{T}\sum_{k=1}^{K} \gamma_{kt} \log\left[p(x_t|s_t = k, \theta_{\text{obs}})\right],
\end{aligned}
\tag{10}
$$

where $K$ is the number of states and $q(s_t = k) = \gamma_{kt}$ is the probability of state $k$ at time $t$. Substituting Eq. (5) into this we have

$$
\mathcal{L} = -\sum_{t=1}^{T}\sum_{k=1}^{K} \gamma_{kt} \log\left[\mathcal{N}(m_k, C_k)\right],
\tag{11}
$$

which is the log-likelihood loss function implemented in `osl-dynamics` for inferring the point estimates for the observation model parameters $\theta_{\text{obs}} = \{m_k, C_k\}$.

# References

[1] I. Rezek and S. Roberts, Ensemble hidden Markov models with extended observation densities for biosignal analysis. Probabilistic modeling in bioinformatics and medical informatics. Springer, London, 419-450 (2005).