

Chapter 1

Model definition

1.1 Notations

We consider the prediction task of a set of observations (y^1, \dots, y^T) given a set of input (u^1, \dots, u^T) .

1.2 Model

We define a L layer RNN followed by a fully connected layer. At time step t ,

$$\begin{cases} y_{t+1} = \tanh(W_y x_{t+1}^L + b_y) \\ x_{t+1}^l = \tanh(W_{xx}^l x_t^l + W_{xu}^l x_{t+1}^{l-1} + b_x^l) \quad \forall 1 \leq l \leq L \end{cases}$$

with $x_t^0 \equiv u_t \forall t$ and $x_0^l \equiv 0 \forall 1 \leq l \leq L$.

Let's consider the weights of the last RNN and fully connected layers as $\theta \equiv (W_{xx}^L, W_{xu}^L, b_x^L, W_y, b_y)$. We can define a new matrix y_t at each time step corresponding to the concatenation of all RNN layers: $x_t \equiv (x_t^1 \dots x_t^L)$. We also introduce two sequences of random noises as i.i.d real valued random variables ϵ and η . We can now write our model in terms of two functions f and g as:

$$\begin{cases} y_{t+1} = f_\theta(x_{t+1}) + \epsilon_{t+1} & \text{observation model} \\ x_{t+1} = g_\theta(x_t, u_{t+1}) + \eta_{t+1} & \text{state model} \end{cases} \quad (1.1)$$

In the following section, we will focus on minimizing the log likelihood

$$\log p_\theta(X_{1:T}, y_{1:T}, u_{1:T}) \quad (1.2)$$

1.3 Minimization

In order to minimize 1.2, we apply an EM strategy. Let

$$Q(\hat{\theta}_p, \theta) = \mathbb{E}_{\hat{\theta}_p} [\log p_\theta(X_{1:T}, y_{1:T}, u_{1:T}) | y_{1:T}] \quad (1.3)$$

We can start by developing the log likelihood:

$$\begin{aligned} \log p_\theta(X_{1:T}, y_{1:T}, u_{1:T}) &= \frac{1}{T} \log \left(\prod_{k=1}^T p_\theta(x_k | x_{k-1}, u_k) p_\theta(y_k | x_k) \right) \\ &= \frac{1}{T} \sum_{k=1}^T \log p_\theta(x_k | x_{k-1}, u_k) + \log p_\theta(y_k | x_k) \\ &= \frac{1}{T} \sum_{k=1}^T \log \left(\det(2\pi\Sigma_x)^{-1/2} \exp\left(-\frac{1}{2}(x_k - g_\theta(x_{k-1}, u_k))^T \Sigma_x^{-1} (x_k - g_\theta(x_{k-1}, u_k))\right) \right) \\ &\quad + \log \left(\det(2\pi\Sigma_y)^{-1/2} \exp\left(-\frac{1}{2}(y_k - f_\theta(x_k))^T \Sigma_y^{-1} (y_k - f_\theta(x_k))\right) \right) \\ &= -\frac{1}{2} \log |\Sigma_x| - \frac{1}{2} \log |\Sigma_y| \\ &\quad - \frac{1}{2T} \sum_{k=1}^T (x_k - g_\theta(x_{k-1}, u_k))^T \Sigma_x^{-1} (x_k - g_\theta(x_{k-1}, u_k)) \\ &\quad - \frac{1}{2T} \sum_{k=1}^T (y_k - f_\theta(x_k))^T \Sigma_y^{-1} (y_k - f_\theta(x_k)) \end{aligned}$$

We will jointly update Σ_x , Σ_y and θ iteratively. We can start by computing the explicit form of the minimum of both covariance matrices.

For Σ_y , we search the zeros of the derivate of the convex function $\Sigma_y \mapsto p_\theta(X_{1:T}, y_{1:T}, u_{1:T})$.

$$\frac{\partial p_\theta(X_{1:T}, y_{1:T}, u_{1:T})}{\partial \Sigma_y^{-1}} = \frac{1}{2} \Sigma_y - \frac{1}{2T} \sum_{k=1}^T (x_k - f_\theta(x_k)) \cdot (x_k - f_\theta(x_k))'$$

$$\frac{\partial p_\theta(X_{1:T}, y_{1:T}, u_{1:T})}{\partial \Sigma_y^{-1}} = 0 \implies \Sigma_y = \frac{1}{T} \sum_{k=1}^T (y_k - f_\theta(x_k))(y_k - f_\theta(x_k))'$$

$$\frac{\partial p_\theta(X_{1:T}, y_{1:T}, u_{1:T})}{\partial \Sigma_x^{-1}} = 0 \implies \Sigma_x = \frac{1}{T} \sum_{k=1}^T (x_k - g_\theta(x_{k-1}, u_k))(x_k - g_\theta(x_{k-1}, u_k))'$$

Since we can't compute an explicit form for θ , we will compute the argmin using a gradient descent. The EM algorithm develops is as follows:

$$\begin{aligned}
\Sigma_{y,p+1} &= \frac{1}{T} \sum_{k=1}^T \mathbb{E}_{\hat{\theta}_p} [(y_k - f_{\theta}(x_k))(y_k - f_{\theta}(x_k))' | y_{1:T}] \\
\Sigma_{x,p+1} &= \frac{1}{T} \sum_{k=1}^T \mathbb{E}_{\hat{\theta}_p} [(x_k - g_{\theta}(x_{k-1}, u_k))(x_k - g_{\theta}(x_{k-1}, u_k))' | y_{1:T}] \\
\theta_{p+1} &= \underset{\theta}{\operatorname{argmin}} \frac{1}{T} \sum_{k=1}^T \mathbb{E}_{\hat{\theta}_p} [(y_k - f_{\theta}(x_k))' \Sigma_{y,p} (y_k - f_{\theta}(x_k)) | y_{1:T}]
\end{aligned}$$

1.4 Sequential Monte Carlo Approach

We approximate the expectation conditional to the observations by using a particle filter. Let M be the number of particles, ξ_k^m and ω_k^m the m -th particle associated to it's weight for time step k .

$$\begin{aligned}
\Phi_k^M &= \mathbb{E}_{\hat{\theta}_p} [(y_k - f_{\theta}(x_k))(y_k - f_{\theta}(x_k))' | y_{1:T}] \\
&= \frac{1}{\Omega_k^M} \sum_{m=1}^M \omega_k^m (y_k - f_{\theta}(\xi_k^m))(y_k - f_{\theta}(\xi_k^m))'
\end{aligned}$$

where $\Omega_k^M = \sum_{m=1}^M \omega_k^m$. In the following sections, we consider that the particles weights sum at 1.

1.5 Gradient descent

At each iteration of the EM algorithm, we start by generating a set of particles under the law $p(x|y)$, that allows us to compute a explicit value for the expectation. We can then minimize this expectation, in order to approximate the new θ candidate, using a gradient descent.

$$\theta_{p+1} = \underset{\theta}{\operatorname{argmin}} \frac{1}{T} \sum_{k=1}^T \sum_{m=1}^M (y_k - f_{\theta}(\omega_k^m))' \Sigma_{y,p} (y_k - f_{\theta}(\omega_k^m))$$