

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) .

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 $\theta \equiv (W_{xx}^L, \cancel{W_{xu}^L}, b_x^L, W_y, b_y)$. We can define a new matrix x_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 , with covariance matrices Σ_y and Σ_x . We can now write our model in terms of two functions f and g as:

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

In the following section, we will focus on maximizing the joint log likelihood:

$$\log p_\theta(X_{0:T}, y_{0:T}, u_{0:T}) \quad (2)$$

3 Minimization

We can start by developing the log likelihood:

$$\begin{aligned}
\log p_\theta(X_{0:T}, y_{0:T}, u_{0:T}) &= \frac{1}{T} \log \left(p_\theta(x_0) p_\theta(y_0|x_0) \prod_{k=1}^T p_\theta(x_k|x_{k-1}, u_k) p_\theta(y_k|x_k) \right) \\
&= \frac{1}{T} \log p_\theta(x_0) \\
&\quad + \frac{1}{T} \sum_{k=1}^T \log p_\theta(x_k|x_{k-1}, u_k) + \frac{1}{T} \sum_{k=0}^T \log p_\theta(y_k|x_k) \\
&= \frac{1}{T} \log p_\theta(x_0) - \frac{1}{2} \log |\Sigma_x| - \frac{1}{2} \log |\Sigma_y| + Cst \\
&\quad - \frac{1}{2T} \sum_{k=1}^T (x_k - g_\theta(x_{k-1}, u_k))' \Sigma_x^{-1} (x_k - g_\theta(x_{k-1}, u_k)) \\
&\quad - \frac{1}{2T} \sum_{k=0}^T (y_k - f_\theta(x_k))' \Sigma_y^{-1} (y_k - f_\theta(x_k))
\end{aligned}$$

We aim at maximizing 2 by gradient descent, by leveraging fisher's identity:

$$\nabla \log p_\theta(x_{0:T}, y_{0:T}, u_{0:T}) = \mathbb{E}_\theta [\nabla \log p_\theta(x_{0:T}, y_{0:T}, u_{0:T}) | Y_{0:T}]$$

In order to approximate this expectation, we need to sample from the posterior distribution $p_\theta(x|y)$. In Section 4, we detail a Sequential Monte Carlo approach to this end. In Section 5, we describe the algorithm to train our model through gradient descent.

4 Sequential Monte Carlo Approach

4.1 Filter

In order to compute the conditional expectations in the previous expressions, we will iteratively sample trajectories $\xi_{1:T}^i$ associated with weights ω^i with respect to the density $p_\theta(x|y)$, using a sequential Monte Carlo particle filter.

At time step $k = 0$, $(\xi_0^l)_{l=1}^N$ are sampled independently from the first hidden state, and associated with sampling weights proportional to the observation density q_θ :

$$\begin{aligned}
\xi_0^i &\sim \mathcal{N}(x_0, \Sigma_x) \\
\omega_0^i &\sim q_\theta(\xi_0^i)
\end{aligned}$$

At time step $k + 1$, we sample indices I of the particles to propagate, based on their previous weights. After propagation, particles weights are computed following the observation density function:

$$\begin{aligned}
\mathbb{P}(I_{k+1}^i = j) &= \omega_k^j \quad \forall 1 \leq j \leq N \\
\omega_{k+1}^i &\sim q_\theta(\xi_{k+1}^i)
\end{aligned}$$

4.2 Smoother

Using the poor man filter, we get N trajectories:

$$\xi_{1:k+1}^i = (\xi_{1:k}^{I_{k+1}^i}, \xi_{k+1}^i)$$

4.3 Approximation

We can now approximate this conditional expectation for any measurable bounded function h :

$$\mathbb{E}_\theta [h(x)|y_{1:T}] = \sum_{i=1}^N \omega_T^i h(\xi_{0:T}^i)$$

5 Gradient descent

5.1 Forward pass

During the forward pass, we generate a set of N particles under the law $p(x|y)$ for fixed values of θ , Σ_x and Σ_y . We initialize the sequence with a initial hidden state sampled from the noise's distribution.

$$x_0^i \sim \mathcal{N}(0, \Sigma_x)$$

In order to predict each new time step $k + 1$, particles from the previous step are attributed weights ω_k^i proportionally to the density probability around the targeted value y_k .

$$\omega_k^i \sim \exp\left(-\frac{1}{2}(y_k - f_{\theta_p}(x_k^i))' \Sigma_{y,p}^{-1} (y_k - f_{\theta_p}(x_k^i))\right)$$

We then select a new population from these particles indexed by I_{k+1}^i , based on their weights.

$$\mathbb{P}(I_{k+1}^i = j) = \omega_k^j \quad \forall 1 \leq j \leq N$$

The current hidden state is computed for the selected particles.

$$x_{k+1}^i = g_\theta(x_k^{I_{k+1}^i}, u_{k+1}) + \eta_{k+1}^i$$

5.2 Loss function

Considering that we have computed a set of N trajectories $(\xi_{0:T}^i)$, $1 \leq i \leq N$, associated with weights (ω^i) , we can approximate the gradient of the log likelihood by computing the gradient of:

$$\begin{aligned} \mathbb{J}(\theta) &= \log |\Sigma_x| + \log |\Sigma_y| \\ &+ \frac{1}{T} \sum_{k=1}^T \sum_{i=1}^N \omega^i (y_k - f_\theta(\xi_k^i))' \Sigma_y^{-1} (y_k - f_\theta(\xi_k^i)) \\ &+ \frac{1}{T} \sum_{k=0}^T \sum_{i=1}^N \omega^i (\xi_k^i - g_\theta(\xi_{k-1}^i, u_k))' \Sigma_x^{-1} (\xi_k^i - g_\theta(\xi_{k-1}^i, u_k)) \end{aligned}$$

5.3 Backward pass

During this step, each parameter of the model is updated by gradient descent.

6 Finetuning

In this section, we aim at improving the weights of a model already trained in a traditional fashion. Given a dataset of samples $(u_{0:T}^{(i)}, y_{0:T}^{(i)})_{i=1}^m$, we consider a training of our model as defined in 1 without the added weights, resulting in an initial set of weights θ_0 . Because we trust the initial training has successfully extracted relevant information from the command u , and reached satisfactory parameters for the hidden state model, we will only adapt the parameters of the observation model. Σ_x is set to a small value.