

Towards Stochastic Recurrent Neural Control Policy

Sequential Monte Carlo variant of Generalized EM algorithm Training Objective: variational lower bound of data log-likelihood $l(heta) = \log p_{ heta}(u|x) = \log \sum_{h} p_{ heta}(u,h|x) = \log \mathbb{E}_{h \sim q} \frac{p_{ heta}(u,h|x)}{q(h)}$ $\geq \mathbb{E}_{h \sim q} \log \frac{p_{\theta}(u, h|x)}{q(h)} = Q(\theta) - \mathcal{H}(q) = \mathcal{L}(q, \theta)$ (3) <u>Approximate Distribution</u>: the tighest choice is $q^*(h) = p(h|u,x)$ $\mathcal{L}(q, heta) = \mathbb{E} \log rac{p(h|u, x)p_{ heta}(u|x)}{q(h)}$ $= \log p_{\theta}(u|x) - D_{KL}(q(h) \parallel p_{\theta}(h|u,x))$ (4) EM Algorithm:

> $q^* = \max_q \mathcal{L} = p(h|u, x)$ (5a)E-step: $\theta^* = \max_{\theta} \mathcal{L} = \max_{\theta} Q(\theta)$ (5b) M-step:

Sequential Monte Carlo (E-step): importance sampling is needed because the true posterior distribution is hard to compute

$$w(\mathbf{h}_{0:T}^{(i)}) = \frac{p(\mathbf{h}_{0:T}^{(i)} | \mathbf{u}_{0:T}, \mathbf{x}_{0:T})}{p(\mathbf{h}_{0:T}^{(i)} | \mathbf{x}_{0:T})}$$

$$= \frac{p(h_0^{(i)} | x_0, u_0)}{p(h_0^{(i)} | x_0)} \prod_{t=1}^T \frac{p(h_t^{(i)} | \mathbf{h}_{0:t-1}^{(i)}, \mathbf{u}_{0:t}, \mathbf{x}_{0:t})}{p(h_t^{(i)} | \mathbf{h}_{0:t-1}^{(i)}, \mathbf{x}_{0:t})} \quad (6)$$

$$w_T^i \triangleq \zeta_0^i \prod_{t=1}^T \zeta_t^i = w_{T-1}^i \zeta_T^i \qquad (7)$$

The last equation suggests that the importance weights can be computed recursively by:

$$\begin{aligned} \zeta_{t}^{i} &= \frac{p(h_{t}^{(i)} | \mathbf{h}_{0:t-1}^{(i)}, \mathbf{u}_{0:t}, \mathbf{x}_{0:t})}{p(h_{t}^{(i)} | \mathbf{h}_{0:t-1}^{(i)}, \mathbf{x}_{0:t})} \\ &= \frac{p(\mathbf{h}_{0:t}^{(i)}, \mathbf{u}_{0:t}, \mathbf{x}_{0:t})}{p(\mathbf{h}_{0:t}^{(i)}, \mathbf{x}_{0:t})p(\mathbf{u}_{0:t} | \mathbf{h}_{0:t-1}^{(i)}, \mathbf{x}_{0:t})} = \frac{p(\mathbf{u}_{0:t} | \mathbf{h}_{0:t}^{(i)}, \mathbf{x}_{0:t})}{p(\mathbf{u}_{0:t} | \mathbf{h}_{0:t-1}^{(i)}, \mathbf{x}_{0:t})} \\ &= \frac{p(u_{0} | h_{0}, x_{0})p(u_{t} | \mathbf{h}_{0:t-1}^{(i)}, \mathbf{x}_{0:t})\prod_{\tau=1}^{t-1}p(u_{\tau} | \mathbf{h}_{0:\tau}, \mathbf{x}_{0:\tau})}{p(u_{0} | h_{0,\tau}, \mathbf{x}_{0:t})p(u_{t} | \mathbf{h}_{0:t-1}, \mathbf{x}_{0:t})\prod_{\tau=1}^{t-1}p(u_{\tau} | \mathbf{h}_{0:\tau}, \mathbf{x}_{0:\tau})} \\ &= \frac{p(u_{t} | \mathbf{h}_{0:t}^{(i)}, \mathbf{x}_{0:t})}{p(u_{t} | \mathbf{h}_{0:t-1}, \mathbf{x}_{0:t})\prod_{\tau=1}^{t-1}p(u_{\tau} | \mathbf{h}_{0:\tau}, \mathbf{x}_{0:\tau})} \end{aligned}$$
(8)

where M is the total number of samples from q_t for each t and the last equation follows from a standard Monte Carlo:

$$p(u_t | \mathbf{h}_{0:t-1}, \mathbf{x}_{0:t}) = \mathbb{E}_{h_t \sim q_t} p(u_t | h_t, \mathbf{h}_{0:t-1}, \mathbf{x}_{0:t})$$

By now we should have collected:

$$q_t^i \triangleq p(h_t^{(i)} | \mathbf{h}_{0:t-1}, \mathbf{x}_{0:t}) \qquad l_t^i \triangleq p(u_t | \mathbf{h}_{0:t}^{(i)}, \mathbf{x}_{0:t}) \quad (9)$$

Gradient Asent (M-step):

$$Q(\theta) = \mathbb{E}_{\mathbf{h}_{0:T} \sim q^{*}} \log p_{\theta}(\mathbf{u}_{0:T}, \mathbf{h}_{0:T} | \mathbf{x}_{0:T})$$

$$\simeq \sum_{i=1}^{M} \tilde{w}_{T}^{i} \Big[\log p_{\theta}(\mathbf{u}_{0:T} | \mathbf{h}_{0:T}^{(i)}, \mathbf{x}_{0:T}) + \log p_{\theta}(\mathbf{h}_{0:T}^{(i)} | \mathbf{x}_{0:T}) \Big]$$

$$= \sum_{i=1}^{M} \Big(\prod_{t=0}^{T} \frac{q_{t}^{i}}{\sum_{j} q_{t}^{j}} \Big) \Big[\sum_{t=0}^{T} \log q_{t}^{i} + \log l_{t}^{i} \Big]$$

$$\frac{\partial Q}{\partial \theta} \simeq \sum_{i=1}^{M} \Big(\prod_{t=0}^{T} \frac{q_{t}^{i}}{\sum_{j} q_{t}^{j}} \Big) \frac{\partial}{\partial \theta} \Big[\sum_{t=0}^{T} \log q_{t}^{i} + \log l_{t}^{i} \Big]$$
(10)

Tianhao Zhang

Abstract

<u>Goal</u>: MLPs are popular control policies for, and achieve state-of-theart results in, many robotics tasks. However, the feedforward architecture of MLPs limits their ability to extend further than a sophisticated reflexive agent. We introduce stochastic units into the neural network control policies, with an ultimate goal of achieving a stochastic recurrent neural network.

<u>Method</u>: We derive a Sequential Monte Carlo variant of Generalized EM algorithm for learning stochastic recurrent neural policy.

Graphical Models

<u>Hybrid Stochastic Feed-forward Network (SFNN)</u>: Shares same structure with MLPs, except that some hidden units are replaced by Bernoulli stochastic neurons. In the extreme case, this becomes a Sigmoid Belief Network.

<u>Bernoulli units</u>: $p(h_{out} = 1 | h_{in}) = \sigma(h_{in})$

Recurrent Stocastic units: We can readily add additional dependency passed through time.



The below algorithm summarizes the overall learning procedure of a hybrid stochastic recurrent neural control policy:

Algorithm 1 SIS Generalized EM

```
Input: training trajectories \{(\mathbf{x}_{0:T}^{(n)}, \mathbf{u}_{0:T}^{(n)}, \boldsymbol{\Sigma}_{\mathbf{u}_{0:T}}^{(n)})\}_{n=1}^{N},
hidden stochastic units for each step unrolled \{\mathbf{h}_t\}_{t=1}^T,
network parameter \theta, number of samples M
repeat
     {E-step}
    for n = 1 to N do
        for i = 1 to M do
           h_1^{(i)} \sim p(h_1|h_0, x_1)
           l_1^i = p(h_1^{(i)}|h_0, x_1)
           w_1^i = p(u_1 | \mathbf{h}_{0:1}^{(i)}, x_1)
        end for
        \tilde{w}_1^i = w_1^i / \sum_i w_1^i
        for t = 2 to T do
           for i = 1 to M do
               h_t^{(i)} \sim p(\mathbf{h}_{0:t} | \mathbf{h}_{0:t-1}, \mathbf{x}_{0:t})
               l_t^i = p(\mathbf{h}_{0:t}^{(i)} | \mathbf{h}_{0:t-1}, \mathbf{x}_{0:t})
               w_t^i = p(u_t | \mathbf{h}_{0:t}^{(i)}, \mathbf{x}_{0:t})
            end for
           \tilde{w}_t^i = w_t^i / \sum_i w_t^i
        end for
    end for
     {M-step}
    for i = 1 to M do
       compute \Delta Q^{(i)} \triangleq \frac{\partial}{\partial \theta} \left[ \sum_{t=0}^{T} \log q_t^i + \log l_t^i \right] using
       standard back-propagation
   end for
   sga_update(\theta, \sum_{i=1}^{M} \tilde{w}_T^i \Delta Q^{(i)}) {any stochastic gradi-
   ent ascent algorithm such as ADAM or RMSProp}
until convergence
```



Experimental Evaluation

Implementation: We independently implemented harness3, based on Computation Graph Toolkit, that allows convenient construction of arbitrary structure of neural network. A highlight of our harness is that it supports layer with mixed types of units, including standard fully-connected units, Bernoulli stochastic units, long-short-term-memory block, and com- mon nonlinearities.

Synthetic Datasets: Here we review the capabilities of SFNN on 1D dataset



Right: optimal solution (conditional average) of a deterministic network Left: error bars represent the range of samples from the learnt stochastic^{0.5} network

* Horizontal/vertical axis is inputs/outputs



<u>Quadrotor Control Task:</u> cylindrical obstacle avoidance

Data Generation:

- Total of 15 randomized initial positions roughly in front of the obstacle
- For each scenario, a trajectory optimization algorithm (iLQG) is used to generate optimal (time-variant Gaussian) control policy given perfect state information (i.e. position, etc)
- Rollout each optimal policy 5 times in simulation and collect the raw observations (i.e. laser ranges) along the resulting trajectories as our supervised training set





Set-up: we trained MLP and SFNN on same training set with conflicting examples which often arise in typical data collection process.

For SFNN, use M = 30 samples for gradient estimate. For both, apply weight decay (5e-3) and allow objective to plateau with proper number of epochs.

Both have fully-connected layers of size 40 (inputs) -> 40 -> 20 -> 15 -> 4 (outputs); SFNN has additional one Bernoulli node at each of the last two layers.

Result: MLP control policy is permanently biased towards dodging right and leads to frequent crashes, whereas SNN retains the variety and keeps a lower crash rate.

Table 1. Comparison of stochastic and deterministic neural network control policies (numbers in parenthesis represent training epochs

Method	%CRASH	%LEFT
SNN(5) SNN(10) SNN(15) MLP(15)	32.0 49.3 33.3	22.9 15.2 11.3



