



Biologically Plausible Variational Policy Gradient with Spiking Recurrent Winner-Take-All Networks

Zhile Yang¹, Shangqi Guo^{†2}, Ying Fang³, Jian K. Liu^{†1} ¹University of Leeds, ²Tsinghua University, ³Fujian Normal University

Keywords: reinforcement learning, spiking neural networks

Message: + SVPG, a spiking-based variational policy gradient method with RWTA network and R-STDP. + Experiment results reveal its potential for solving RL tasks and to have inherent robustness.



1. Background & Problem



Background

RWTA Network



3. Policy Inference & Optimization

Policy Inference

 $p(\boldsymbol{v}_a, \boldsymbol{v}_h | s) := \frac{1}{Z(s)} \exp\{E(\boldsymbol{v})\}$ Target function Mean-field inference $\hat{p}(\boldsymbol{v}_a, \boldsymbol{v}_h | s) := \hat{p}(\boldsymbol{v}_a | s) \hat{p}(\boldsymbol{v}_{h_1} | s) \cdots \hat{p}(\boldsymbol{v}_{h_{n_h}} | s)$ $\textbf{Minimize KL-divergence } \partial D_{\mathrm{KL}}(s) / \partial q_i = 0 \blacklozenge q_i = \frac{1}{Z(q_{G(i)})} \exp\{\boldsymbol{w}_{\mathrm{row},i}^{\mathrm{T}} \boldsymbol{q} + \boldsymbol{w}_{\mathrm{col},i}^{\mathrm{T}} \boldsymbol{q} + b_i\}$ In RWTA network Let $\rho_i(l) = \hat{\rho}q_i$, then $\rho_i = \hat{\rho}\exp\{\boldsymbol{w}_{\text{row},i}^{\text{T}}\boldsymbol{q} + \boldsymbol{w}_{\text{col},i}^{\text{T}}\boldsymbol{q} + b_i - \log\sum_{j \in G(i)}\exp\{\boldsymbol{w}_{\text{row},j}^{\text{T}}\boldsymbol{q} + \boldsymbol{w}_{\text{col},j}^{\text{T}}\boldsymbol{q} + b_j\}\}$ \rightarrow RWTA net is suitable for policy inference In neuron model Let $\int_0^\infty \kappa(y) dy = 1/\hat{\rho}$, then $\Longrightarrow \begin{array}{l} \rho_{i(l)} - \rho \exp\left(w_{i(l)}\right) & = \sum_{j \in N(i)} w_{ij} \int_0^\infty \kappa(y) S_{ij}(l-y) dy + b_i. \end{array}$

Policy Optimization

Policy function $q_i = \frac{1}{Z(q_{G(i)})} \exp\{w_{\text{row},i}^{\text{T}} q + w_{\text{col},i}^{\text{T}} q + b_i\}$ (iterate until convergence)

Policy π based on an energy function of the firing states.

 $\pi(\boldsymbol{v}_a|s) = \sum_{\boldsymbol{v}_h} p(\boldsymbol{v}_a, \boldsymbol{v}_h|s) \quad p(\boldsymbol{v}_a, \boldsymbol{v}_h|s) := \frac{1}{Z(s)} \exp\{E(\boldsymbol{v})\} \quad E(\boldsymbol{v}) := \boldsymbol{v}^{\mathrm{T}} \boldsymbol{W} \boldsymbol{v} + \boldsymbol{b}^{\mathrm{T}} \boldsymbol{v}$

v: firing states of the neurons; q: firing probabilities of the neurons

4. Experiments



Precise differential X Requires pseudo-inverse of $M(W + W^{T})$ -- computationally infeasible $\frac{\partial \boldsymbol{q}_{ha}}{\partial \boldsymbol{w}_{jk}} = \boldsymbol{M} \big(\boldsymbol{U}_{jk} + \boldsymbol{U}_{kj} \big) \boldsymbol{q} + \boldsymbol{M} (\boldsymbol{W} + \boldsymbol{W}^{\mathrm{T}}) \frac{\partial \boldsymbol{q}}{\partial \boldsymbol{w}_{jk}}, \quad \frac{\partial \boldsymbol{q}_{ha}}{\partial \boldsymbol{b}_{j}} = \boldsymbol{M} \boldsymbol{b} + \boldsymbol{M} (\boldsymbol{W} + \boldsymbol{W}^{\mathrm{T}}) \frac{\partial \boldsymbol{q}}{\partial \boldsymbol{b}_{j}},$ $M = \operatorname{diag}(q_{ha})[-G_{ha}\operatorname{diag}(q) + D_{sel}]$



		21	2111	(supervised)	
MNIST	0.926 ± 0.001	0.933 ± 0.062	0.898 ± 0.067	0.971 ± 0.001	0.933 ± 0.062
GymIP	199.87 ± 0.27	199.95 ± 0.13	199.96 ± 0.12	N/A	190.79 ± 27.64

📥 BPTT

– EP

-- SVPG

BPTT

0.4

0.25 0.50 0.75

Noise Standard Deviation

0.2

Noise Amplitude

QC 9C 3

 $_{\rm L}^{\rm est}$ 0.4

2 150

Test 100 -

0.00

