# DM887 Assignment2 Q1

Jiawei Zhao

19.03.2024

---

**Algorithm 1** Least-Squares Temporal Differences (LSTD) Deep Q-Learning with Nonlinear Feature Extraction that maps a state to lower-dimensional latent embedding. While the action-value function should be a linear function of the output of the feature extractor.

---

1: Initialize a variational autoencoder $\mathbf{A}$ as the feature extraction network with initial weights $\theta_0$ [$\theta_0$ are drawn from a zero-mean Gaussian distribution]

2: Initialize the number of maximum episode $N_0 = 10$, $N_1 = 10$, and $N_2 = 30$ for each phase of training [It would be a wise choice to seperate the warm-up phase, the antoencoder improvement phase, and the LSTD improvement phase] [Initialization of training hyperparameters]

3: Initialize the number of maximum time step per episode $t = 5000$

4: Initialize the weights $w_0$ for linear approximation of the action-value function $y = f(phi(s))$ which will be used for LSTD

5: Initialize a replay memory buffer $\mathcal{D}$ with a capacity $N = t(N_0 + N_1 + N_2)$

6: Initialize a discount factor $\gamma = 0.9$ [A relatively large discount factor encourages long-term planning and faster convergence during training]

   [The warm-up phase starts here]

7: Freeze the weights of $\mathbf{A}$ as $\theta_0$

8: Freeze the weights of the $f(phi(s))$ as $w_0$

9: **for** *episode* $e = 1$ to $N_0$ **do**

10:    Initialize state $s_0$

11:    Preprocess $s$ into $s_1$ to adapt it as input of $\mathbf{A}$ [preprocessing of high-dimensional states is necessary w.r.t. autoenconders]

12:    **for** each time step $t = 1$ to $T$ **do**

13:       Encode state $s_t$ using $\mathbf{A}$ to get latent embedding $\phi(s_t)$

14:       Select action $a_t$ using $\epsilon$-greedy policy with $\epsilon = 0.3$ [The learning curves of Game B of Assignment 1 corroborate that a relatively high $\epsilon$ could improve $Q$ more efficiently when an $e$ does not end prematurely in the majority of cases]

15:       Execute $a_t$ and obtain reward $r_t$ and new state $s_{t+1}$

16:       Store the transition of the current $t$, i.e. $(s_t, a_t, r_t, s_{t+1})$ in $\mathcal{D}$

17:       Encode state $s_t$ using $\mathbf{A}$ to get latent embedding $\phi(s_t)$

18:       Select action $a_t$ using $\epsilon$-greedy policy with $\epsilon = 0.3$ [The learning curves of Game B of Assignment 1 corroborate that a relatively high $\epsilon$ could improve $Q$ more efficiently when an $e$ does not end prematurely in the majority of cases]

19:       Execute $a_t$ and obtain reward $r_t$ and new state $s_{t+1}$

20:       Store the transition of the current $t$, i.e. $(s_t, a_t, r_t, s_{t+1})$ in $\mathcal{D}$

21:    **end for**

22: **end for**

[1]

---

**Algorithm 2** The LSTD update starts here

1: Unfreeze the weights of $\mathbf{A}$, i.e. $\theta$
2: **for** *episode* $e = 1$ to $N_1$ **do**
3:     Repeat the same steps at phase 1
4:     **for** each time step $t = 1$ to $T$ **do**
5:         Repeat the same steps at phase 1
6: [Start using the autoencoder algorithm to update $w$]
7:         Sample a minibatch of transitions from $\mathcal{D}$ with a batchsize with a randomly selected
8:     **end for**
9: **end for**

[1]

**Algorithm 3** The autoencoder + LSTD training starts here

1: Freeze the weights of $\mathbf{A}$, i.e. $\theta$
2: Unfreeze the weights of $f(phi(s))$
3: **for** $episode\ e = 1$ to $N_2$ **do**
4:     Repeat the same steps at phase 1
5:     **for** each time step $t = 1$ to $T$ **do**
6:         Repeat the same steps at phase 1
7: [Start using the online LSTD algorithm to update $w$]
8:         Sample a minibatch of transitions from $\mathcal{D}$ with a batchsize
9:     **end for**
10: **end for**

[1]