# A Local Temporal Difference Code for Distributional Reinforcement Learning

Brabeeba Wang
4/30/2021

# Traditional

# Distributional

$$V(s_t) \leftarrow V(s_t) + \alpha\delta(t)$$
$$\delta(t) = r_t + \gamma V(s_{t+1}) - V(s_t)$$

$$V_i(s_t) \leftarrow V_i(s_t) + \alpha_i^+ \delta_i(t) \ \text{ if } \delta_i(t) > 0$$
$$V_i(s_t) \leftarrow V_i(s_t) + \alpha_i^- \delta_i(t) \ \text{ if } \delta_i(t) < 0$$
$$\delta_i(t) = r_t + \gamma \tilde{V}(s_{t+1}) - V_i(s_t)$$

$$V(s) \rightarrow E\left[\sum_{\tau=0}^{\infty} \gamma^\tau r_{t+\tau} \Big| s_t = s\right]$$

**But sampling from value distribution is not local.**

(a) $V$ $\delta$

(b) $V_i$ $\delta_i$

# Laplace code on discount factor

$$V_\gamma(s_t) \leftarrow V_\gamma(s_t) + \alpha\delta_\gamma(t)$$

$$\delta_\gamma(t) = r_t + \gamma V_\gamma(s_{t+1}) - V_\gamma(s_t)$$

**Discrete:**
$$V_\gamma(s_t) \rightarrow E\left[\sum_{\tau=0}^{\infty} \gamma^\tau r_{t+\tau}\Big|s_t\right] = \sum_{\tau=0}^{\infty} \gamma^\tau E[r_{t+\tau}|s_t]$$

$$Z^{-1}\{V_\gamma(s_t)\}_{\gamma\in(0,1)} = \{E[r_{t+\tau}|s_t]\}_{\tau=0}^{\infty}$$

**Continuous:**
$$V_\gamma(s_t) \rightarrow \int_0^{\infty} e^{-\tau(-\log\gamma)} E[r_{t+\tau}|s_t]\, d\tau$$

$$\mathcal{L}^{-1}\{V_\gamma(s_t)\}_{\gamma\in(0,1)} = \{E[r_{t+\tau}|s_t]\}_{\tau>0}$$

**Linear readout:** $\mathbf{L}^{-1}[V_{\gamma_1}(s_t),\ldots,V_{\gamma_N}(s_t)] = [E[r_{t+0}|s_t],\ldots,E[r_{t+T}|s_t]]$

# Laplace code on reward sensitivity

$$V_{h,\gamma}(s_t) \leftarrow V_{h,\gamma}(s_t) + \alpha \delta_{h,\gamma}(t)$$

$$\delta_{h,\gamma}(t) = f_h(r_t) + \gamma V_{h,\gamma}(s_{t+1}) - V_{h,\gamma}(s_t)$$



(a)

(b)

$$V_{h,\gamma}(s_t) \rightarrow E\left[\sum_{\tau=0}^{\infty} \gamma^\tau f_h(r_{t+\tau}) \Big| s_t\right] = \sum_{\tau=0}^{\infty} \gamma^\tau E\left[f_h(r_{t+\tau}) | s_t\right]$$

$$= \sum_{\tau=0}^{\infty} \gamma^\tau E\left[H(r_{t+\tau} - \theta_h) | s_t\right] = \sum_{\tau=0}^{\infty} \gamma^\tau P(r_{t+\tau} > \theta_h | s_t)$$

$$\mathbf{L}^{-1}[V_{h,\gamma_1}(s_t), \ldots, V_{h,\gamma_N}(s_t)] = [P(r_{t+0} > \theta_h | s_t), \ldots, P(r_{t+T} > \theta_h | s_t)]$$

TD learning (Eq. 10)

$\tau = 0$    $\tau = 1$    $\tau = 2$    $\tau = 3$

$(s)$ → ○ → ○ → □

$r = 2$    $r = 2$    $r = 1$    $r = 1$
$r = -2$    $r = -2$    $r = -1$    $r = -1$

$\gamma$-space

$V_{h-1,\gamma}(s)$
$-V_{h,\gamma}(s)$

Sensitivity $(\theta_h)$

Temporal discount $(\gamma)$

$\mathbf{L}^{-1}$

$\tau$-space

$P(r_\tau | s)$

Sensitivity $(\theta_h)$

Future time $(\tau)$

# What can this code recover?

$$E\left[\sum_{\tau=0}^{\infty} \tilde{\gamma}^{\tau} r_{t+\tau} \Big| s_t\right] = \sum_{\tau=0}^{\infty} \tilde{\gamma}^{\tau} \sum_{r} r P(r_{t+\tau} = r | s_t)$$

$$E\left[\sum_{\tau=0}^{\infty} \tilde{\gamma}^{\tau} r_{t+\tau} \Big| s_t\right] = \sum_{h=1}^{H} \left(V_{h-1,\tilde{\gamma}}(s_t) - V_{h,\tilde{\gamma}}(s_t)\right)\theta_h$$

$$P\left(\sum_{\tau=0}^{T} \tilde{\gamma}^{\tau} r_{t+\tau} = V \Big| s_t\right) = (1 - \delta_{(V,0)}) \sum_{\tau=0}^{T} P\left(r_{t+\tau} = \tilde{\gamma}^{-\tau} V \Big| s_t\right)$$

**However, notice that this approach recovers the reward distribution evolution but not the value distribution since reward might be correlated across time.**

# Laplace code on temporal discount

$$R_t^n = r_t + \ldots + \tilde{\gamma}^n r_{t+n}$$

$$V_{h,\gamma,n}(s_t) \leftarrow V_{h,\gamma,n}(s_t) + \alpha \delta_{h,\gamma,n}(t+n)$$

$$\delta_{h,\gamma,n}(t+n) = f_h(R_t^n) + \gamma V_{h,\gamma,n}(s_{t+1}) - V_{h,\gamma,n}(s_t)$$

$(a)$

$0.5 \rightarrow \boxed{r_1}$

$s$

$\tau = T$

$0.5 \rightarrow \bigcirc \rightarrow \cdots \rightarrow \bigcirc \rightarrow \boxed{r_2}$

$(b)$

$P(V \mid s)$

Probability

$r_1 \quad \tilde{\gamma}^T r_2$

$(c)$

$\tau$-space

$P(r_\tau \mid s)$

Sensitivity $(\theta_h)$

$r_1 \quad r_2$

Future time $(\tau)$

$(d)$

Error

$T \quad T'$

Trial

$(a)$ Expectile / Laplace density plots (Reward vs Density).

$(b)$ $V_h = E_{r \sim D}[f_h(r)] = P_D(r > \theta_h)$

$\delta_h(R) = f_h(R) - V_h$

$(c)$ Expectile / Laplace plots (Reward asymm. vs Rev. point), with $T$ and $T - \tau$.

# Connection to successor representation

$$[SR^\gamma(s)]_{s'} = E\left[\sum_{\tau=0}^{\infty} \gamma^\tau \delta_{(s',s_{t+\tau})} \middle| s_t = s\right] = \sum_{\tau=0}^{\infty} \gamma^\tau P(s_{t+\tau} = s' | s_t = s).$$

$$V_{h,\gamma}(s_t) \rightarrow \sum_{\tau=0}^{\infty} \gamma^\tau \sum_s P(s_{t+\tau} = s | s_t) P(r_{t+\tau} > \theta_h | s_t, s_{t+\tau} = s)$$

$$V_{h,\gamma}(s_t) \rightarrow \sum_s \left(\sum_{\tau=0}^{\infty} \gamma^\tau P(s_{t+\tau} = s | s_t)\right) P(r > \theta_h | s) = SR^\gamma(s_t) \cdot \boldsymbol{r}_h$$