

Periodic agent-state based Q-learning (PASQL) for POMDPs

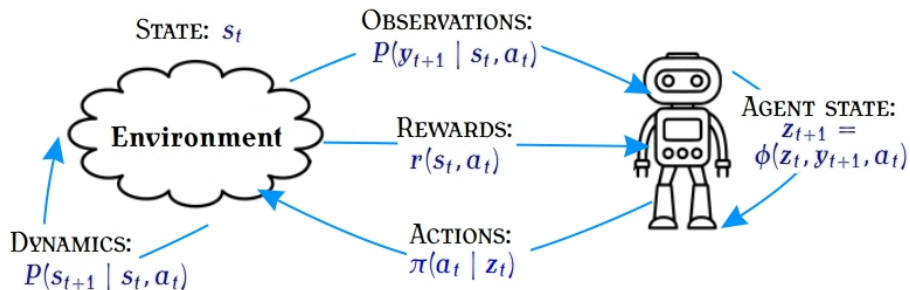
Amit Sinha¹ Matthieu Geist² Aditya Mahajan¹

¹McGill University, Mila

²Cohere

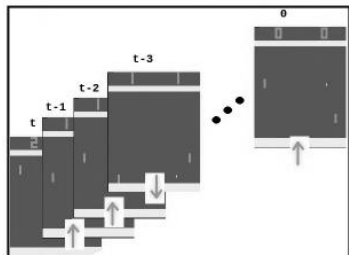
November 13, 2024

Partially observable MDP (POMDP) with agent states



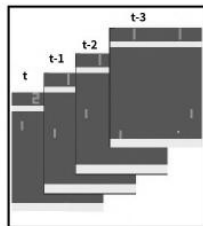
Example: Agent states in Atari-pong

Observation:  **Action:** \uparrow/\downarrow



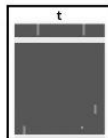
History

$$Z_t = \mathcal{Y}_{0:t}, a_{0,t-1}$$



**Finite Window
(Size 4)**

$$Z_t = \mathcal{Y}_{t-4:t}$$



Abstract Features

$$Z_t =$$

$$(x,y) = (0.1, 0.1)$$

$$(x,y) = (0.9, 0.2)$$

$$(x,y) = (0.8, 0.1)$$

Main problem

- ① Agent state Z_t may be **non-Markovian** (e.g. due to function approx. in neural networks)

Main problem

- ① Agent state Z_t may be **non-Markovian** (e.g. due to function approx. in neural networks)
- ② **Non-stationary** policies **will do better** than **stationary** ones

Why periodic policies?

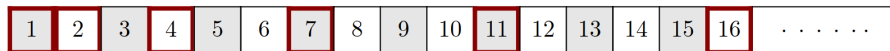
Definition (period L): $\pi = (\pi_1, \dots, \pi_L, \pi_1, \dots, \pi_L, \dots)$

Why periodic policies?

Definition (period L): $\pi = (\pi_1, \dots, \pi_L, \pi_1, \dots, \pi_L, \dots)$

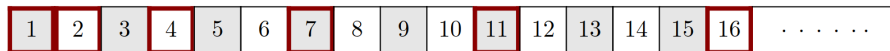
Reason : Non-stationary policies are **not realizable** in practice for infinite horizon problems.

Periodic performance



Period (L)	1	2	3	4	5	6	7	8	9	10
J_L^*	4.022	4.022	7.479	6.184	8.810	7.479	9.340	8.488	9.607	8.810

Periodic performance



Period (L)	1	2	3	4	5	6	7	8	9	10
J_L^*	4.022	4.022	7.479	6.184	8.810	7.479	9.340	8.488	9.607	8.810

Periodic policies perform better than **stationary policies** !

Periodic agent-state based Q-learning (PASQL) algorithm

Regular Q-learning:

$$Q_{t+1}(z, a) = Q_t(z, a) + \alpha_t(z, a) \left[R_t + \gamma \max_{a' \in A} Q_t(Z', a') - Q_t(z, a) \right].$$

Periodic agent-state based Q-learning (PASQL) algorithm

Regular Q-learning:

$$Q_{t+1}(z, a) = Q_t(z, a) + \alpha_t(z, a) \left[R_t + \gamma \max_{a' \in A} Q_t(Z', a') - Q_t(z, a) \right].$$

PASQL is just **Q-learning** for **periodic policies**!

$$Q_{t+1}^{\ell}(z, a) = Q_t^{\ell}(z, a) + \alpha_t^{\ell}(z, a) \left[R_t + \gamma \max_{a' \in A} Q_t^{\llbracket \ell+1 \rrbracket}(z_{t+1}, a') - Q_t^{\ell}(z, a) \right]$$
$$\forall \ell \in L, \quad \llbracket \ell \rrbracket = \ell \pmod{L}. \quad (\text{PASQL})$$

Periodic agent-state based Q-learning (PASQL) algorithm

Regular Q-learning:

$$Q_{t+1}(z, a) = Q_t(z, a) + \alpha_t(z, a) \left[R_t + \gamma \max_{a' \in A} Q_t(Z', a') - Q_t(z, a) \right].$$

PASQL is just **Q-learning** for **periodic policies**!

$$Q_{t+1}^{\ell}(z, a) = Q_t^{\ell}(z, a) + \alpha_t^{\ell}(z, a) \left[R_t + \gamma \max_{a' \in A} Q_t^{\llbracket \ell+1 \rrbracket}(z_{t+1}, a') - Q_t^{\ell}(z, a) \right]$$
$$\forall \ell \in L, \quad \llbracket \ell \rrbracket = \ell \pmod{L}. \quad (\text{PASQL})$$

Problem:

Standard Q-learning convergence cannot apply (since Z_t is non-Markovian)

Main result: Convergence theorem

$$\{(Q_t^0, \dots, Q_t^{L-1})\}_{t \geq 1} \rightarrow (Q_\mu^0, \dots, Q_\mu^{L-1}) \quad \text{a.s.}$$

Main result: Convergence theorem

$$\{(Q_t^0, \dots, Q_t^{L-1})\}_{t \geq 1} \rightarrow (Q_\mu^0, \dots, Q_\mu^{L-1}) \quad \text{a.s.}$$

$$Q_\mu^\ell(z, a) = r_\mu^\ell(z, a) + \gamma \sum_{z' \in Z} P_\mu^\ell(z'|z, a) \max_{a' \in A} Q_\mu^{\llbracket \ell+1 \rrbracket}(z', a')$$

$$\forall \ell \in L, \forall (z, a) \in Z \times A$$

Main result: Convergence theorem

$$\{(Q_t^0, \dots, Q_t^{L-1})\}_{t \geq 1} \rightarrow (Q_\mu^0, \dots, Q_\mu^{L-1}) \quad \text{a.s.}$$

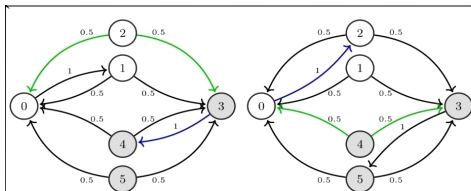
$$Q_\mu^\ell(z, a) = r_\mu^\ell(z, a) + \gamma \sum_{z' \in Z} P_\mu^\ell(z'|z, a) \max_{a' \in A} Q_\mu^{\llbracket \ell+1 \rrbracket}(z', a')$$

$$\forall \ell \in L, \forall (z, a) \in Z \times A$$

$$r_\mu^\ell(z, a) := \sum_{s \in S} r(s, a) \zeta_\mu^\ell(s | z)$$

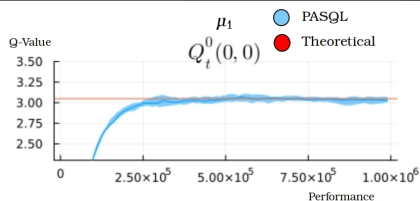
$$P_\mu^\ell(z'|z, a) := \sum_{(s, y') \in S \times Y} \mathbb{1}_{\{z' = \phi(z, y', a)\}} P(y'|s, a) \zeta_\mu^\ell(s | z)$$

PASQL Numerical



(a) Dynamics under action 0.

(b) Dynamics under action 1.



$$\mu_1 := \begin{bmatrix} 0.2 & 0.8 \\ 0.8 & 0.2 \end{bmatrix}, \mu_2 := \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}, \mu_3 := \begin{bmatrix} 0.8 & 0.2 \\ 0.2 & 0.8 \end{bmatrix}$$

$$\bar{\mu}_1 := [0.2; 0.8], \bar{\mu}_2 := [0.5; 0.5], \bar{\mu}_3 := [0.8; 0.2]$$

Performance of converged *periodic* ($L=2$) policies.

J_2^*	$J^{\pi_{\mu_1}}$	$J^{\pi_{\mu_2}}$	$J^{\pi_{\mu_3}}$
6.793	6.793	1.064	0.532

Performance of converged *stationary* policies.

J_1^*	$J^{\bar{\pi}_{\mu_1}}$	$J^{\bar{\pi}_{\mu_2}}$	$J^{\bar{\pi}_{\mu_3}}$
2.633	0.0	1.064	2.633

- 1 We show that considering **periodic policies** over **stationary policies** can be beneficial

Conclusions

- 1 We show that considering **periodic policies** over **stationary policies** can be beneficial
- 2 We provide an algorithm **PASQL** with several **useful theoretical properties**

The End