

Tempo Adaptation in Non-stationary Reinforcement Learning

Hyunin Lee¹ Yuhao Ding¹ Jongmin Lee¹ Ming Jin²
Javad Lavaei¹ Somayeh Sojoudi¹

¹UC Berkeley ²Virginia Tech



Overlooked issue: Time synchronization

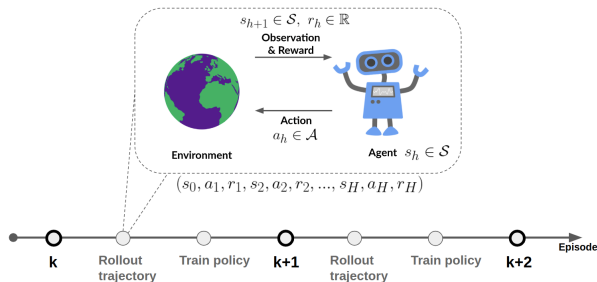


Figure 1: Conventional Non-stationary RL environment

Key observation: In reality, environmental changes occur over wall-clock time (t) rather than episode progress (k).

Existing works: episode k collect data & train policy episode $k + 1$.

In reality: time t_k spend Δt for collecting data & training time $t_{k+1} = t_k + \Delta t$.

Remove time synchronization

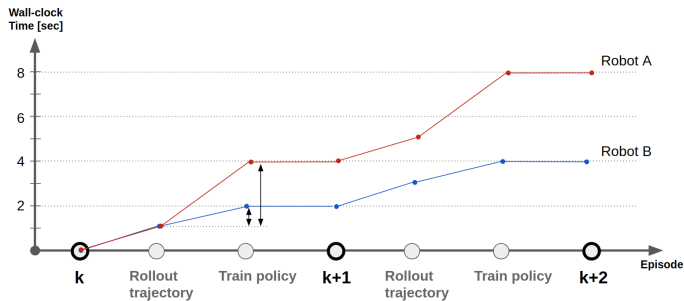


Figure 2: Different training time makes agent encounters different environment

In time-desynchronized environment, the agent should choose **when to interact** ($t_1; t_2; \dots; t_K$) additional to **how many times to interact** (K)
The choice of *interaction times* ($t_1; t_2; \dots; t_K$) significantly impacts the suboptimality gap of the policy.

Contribution

We propose a Proactively Synchronizing Tempo (ProST) framework that computes suboptimal $\{t_1; t_2; \dots; t_K\} (= \{t\}_{1:K})$.

ProST framework computes suboptimal $\{t\}_{1:K}$ by minimizing the upper bound of its performance metric, dynamic regret.

One interesting property is that we show suboptimal $\{t\}_{1:K}$ strikes a balance between the policy training time (**agent tempo**) and how fast the environment changes (**environment tempo**).

ProST framework

For given $t \in [0; T]$, ProST framework computes K , $\{t_1; t_2; \dots; t_K\}$, then $\{t_1; t_2; \dots; t_K\}$ into two components

Time optimizer

Future policy optimizer

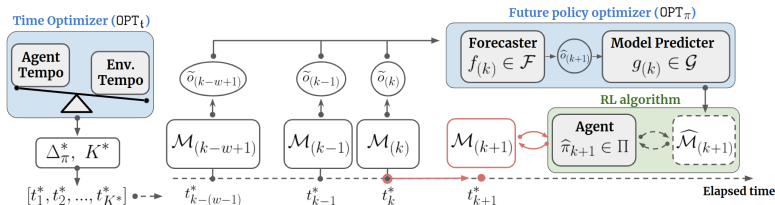


Figure 3: ProST framework

Future policy optimizer

For given t_k, t_{k+1} , it computes a near-optimal policy of t_{k+1} at time t_k

Definition (MDP forecaster g, f)

Consider two function classes F and G such that $F: \mathcal{O}^w \rightarrow \mathcal{O}$ and $G: \mathcal{S} \times \mathcal{A} \times \mathcal{O} \rightarrow \mathbb{R} \times \Delta(\mathcal{S})$, where $w \in \mathbb{N}$. Then, for $f_{(k)} \in F$ and $g_{(k)} \in G$, we define MDP forecaster at time t_k as $(g, f)_{(k)}: \mathcal{O}^w \times \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} \times \Delta(\mathcal{S})$.

Estimate the future MDP model and optimize.

At $t = t_k$

During $t \in (t_k, t_{k+1})$

$$\begin{aligned}\hat{\theta}_{(k+1)} &= f_{(k)}(\{\theta\}_{(k-w+1:k)}) \\ (\hat{R}_{(k+1)}(s, a), \hat{P}_{(k+1)}(\cdot|s, a)) &= g_{(k)}(s, a, \hat{\theta}_{(k+1)}) \\ \hat{\pi}_{(k+1)} &\leftarrow \hat{\mathcal{M}}_{(k+1)} = \langle \mathcal{S}, \mathcal{A}, H, \hat{P}_{(k+1)}, \hat{R}_{(k+1)}, \gamma \rangle\end{aligned}$$

At $t = t_{k+1}$

Time optimizer

Strategy: Δ is a minimizer of the dynamic regret's upper bound

Analysis on finite space S ; $A < \text{ProST-T}$

Theorem (ProST-T dynamic regret R)

Let $\epsilon_h^{(k)} = \sum_{h=0}^{H-1} \sum_{k=1}^{K-1} (s_h^{(k+1)}; a_h^{(k+1)})$ and $\epsilon^{(k)} = \sum_{k=1}^{K-1} \epsilon^{(k+1)}$, where $\epsilon_h^{(k)}$ is a data-dependent error. For a given $p \in (0, 1)$, the dynamic regret of the forecasted policies $\{\pi_h^{(k+1)}\}_{1 \leq k \leq K-1}$ of ProST-T is upper bounded with probability at least $1 - p$ as follows:

$$R(\{\pi_h^{(k+1)}\}_{1 \leq k \leq K-1}; K) \leq R_I + R_{II}$$

where $R_I = \sum_{k=1}^{K-1} (1 - \epsilon) - \frac{K}{H} + C_p \sqrt{K-1}$; $R_{II} = C_{II}[\Delta] (K-1)$, and $C_p, C_{II}[\Delta]$ are some functions of p, Δ , respectively.

R_I Forecasting model error $B(\Delta)$ (rate of environment's change)

R_{II} Policy optimization error Δ (rate of agent's adaption)

Δ strikes a balance between R_I and R_{II}

Δ bounds for sublinear R_{II}

Δ should satisfy sublinear dynamic regret to K

ϵ : approximation gap

η : entropy regularization parameter

α : learning rate

Proposition (Δ bounds for sublinear R_{II})

A total step H is given by MDP. For a number $\epsilon > 0$ such that

$H = \Omega(\log((r_{\max} - r_{\min})))$, we choose ϵ ; η ; α to satisfy

$\epsilon = O(\frac{1}{\sqrt{H}})$; $\eta = \Omega(\frac{1}{\log A})$ and $\alpha = (1 - \epsilon)$, where r_{\max} and r_{\min} are the maximum reward of the forecasted model and the maximum reward of the environment, respectively. Define $N_{II} = \{n \mid n > \frac{1}{\epsilon} \log \frac{C_1(\epsilon + 2)}{\epsilon}; n \leq N\}$, where C_1 is a constant. Then $R_{II} \leq 4(K - 1)$ for all $\Delta \in N_{II}$.

R_T Forecasting model error $B(\Delta)$

SW-LSE : Sliding window regularized LSE

Theorem (Dynamic regret R_T when $f = \text{SW-LSE}$)

For given $p \in (0; 1)$, if the exploration bonus constant β and regularization parameter λ satisfy $\beta = \Omega\left(\frac{S H}{\log(H/p)}\right)$; $\lambda \geq 1$, then the R_T is bounded with probability $1 - p$,

$$R_T \leq C_l[B(\Delta)] w + C_k \sqrt{\frac{1}{w} \log \left(1 + \frac{H}{w}\right)}$$

where $C_l[B(\Delta)] = (1 - \beta) + H) B_r(\Delta) + (1 + H\hat{p}_{\max}) (1 - \beta) B_p(\Delta)$, and C_k is a constant on the order of $\mathcal{O}(K)$.

Δ bounds for sublinear R_I

Proposition (Δ bounds for sublinear R_I)

Denote $B(1)$ as the environment tempo when $\Delta = 1$, which is a summation over all time steps. Assume that the environment satisfies $B_r(1) + B_p(1)\hat{r}_{max} (1 - \gamma) = o(K)$ and we choose $w = O((K-1)^{2/3} (C_I[B(\Delta)])^{2/3})$. Define the set N_I to be $\{n \mid n < K; n \in \mathbb{N}\}$. Then R_I is upper-bounded as $R_I = O(C_I[B(\Delta)]^{1/3} (K-1)^{2/3} \sqrt{\log((K-1) C_I[B(\Delta)])})$ and also satisfies a sublinear upper bound, provided that $\Delta \in N_I$.

Δ strikes a balance between R_I and R_{II}

R_I upperbound is increasing on a interval $N_I \quad N_{II}$

R_{II} upperbound is decreasing on a interval $N_I \quad N_{II}$

Theorem (Suboptimal tempo Δ^*)

Let $k_{Env} = (r - p)^2 C_I[B(1)]$, $k_{Agent} = \log(1 - (1 - \rho)) C_1(K - 1)(\rho + 2)$.

Consider three cases: **case1:** $r - p = 0$, **case2:** $r - p = 1$, **case3:**

$0 < r - p < 1$ or $r - p > 1$. Then Δ^* depends on the environment's drifting constants as follows:

- Case1: $\Delta^* = T$.
- Case2: $\Delta^* = \log_{1-\rho}(k_{Env} k_{Agent}) + 1$.
- Case3: $\Delta^* = \exp(-W) \frac{\log(1-\rho)}{\max(r, p)-1}$, provided that the parameters are chosen so that $k_{Agent} = (1 - \rho) k_{Env}$.

Performance

Benchmark methods

MBPO : state of the art model-based policy optimization.

Pro-OLS : policy optimization algorithm that predicts future V .

ONPG : adaptive algorithm that fine-tunes the policy on current data.

FTRL : adaptive algorithm that maximizes the performance on all previous data.

Table 1: Average reward returns

Speed	$B(G)$	Swimmer-v2					Halfcheetah-v2					Hopper-v2				
		Pro-OLS	ONPG	FTML	MBPO	ProST-G	Pro-OLS	ONPG	FTML	MBPO	ProST-G	Pro-OLS	ONPG	FTML	MBPO	ProST-G
1	16.14	-0.40	-0.26	-0.08	-0.08	0.57	-83.79	-85.33	-85.17	-24.89	-19.69	98.38	95.39	97.18	92.88	92.77
2	32.15	0.20	-0.12	0.14	-0.01	1.04	-83.79	-85.63	-86.46	-22.19	-20.21	98.78	97.34	99.02	96.55	98.13
3	47.86	-0.13	0.05	-0.15	-0.64	1.52	-83.27	-85.97	-86.26	-21.65	-21.04	97.70	98.18	98.60	95.08	100.42
4	63.14	-0.22	-0.09	-0.11	-0.04	2.01	-82.92	-84.37	-85.11	-21.40	-19.55	98.89	97.43	97.94	97.86	100.68
5	77.88	-0.23	-0.42	-0.27	0.10	2.81	-84.73	-85.42	-87.02	-20.50	-20.52	97.63	99.64	99.40	96.86	102.48
A	8.34	1.46	2.10	2.37	-0.08	0.57	-76.67	-85.38	-83.83	-40.67	83.74	104.72	118.97	115.21	100.29	111.36
B	4.68	1.79	-0.72	-1.20	0.19	0.20	-80.46	-86.96	-85.59	-29.28	76.56	80.83	131.23	110.09	100.29	127.74

*Whole training procedure is in Appendix

Ablation study

Figure 4: (a) Optimal . (b-1) Different forecaster f (ARIMA, SA). (b-2) The Mean squared Error (MSE) model loss of four ProST-G with different forecasters (ARIMA and three SA) and the MBPO. x -axis are all episodes.