

Some recent progress in continuous-time reinforcement learning

YUFEI ZHANG

Statistics Research Showcase Day

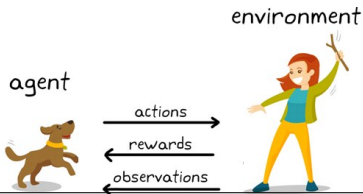
based on joint works with Matteo Basei (EDF R&D), Xin Guo (UC Berkeley), Anran Hu (UC Berkeley), Lukasz Szpruch (U of Edinburgh, Turing) and Tanut Treetanthiploet (Turing).

15 June 2022

- ▶ Stochastic control problems are ubiquitous.
- ▶ Continuous-time models well understood in this community.
 - ▶ optimal trading, dynamic hedging, autonomous driving, robots.
- ▶ Reinforcement learning (RL) methods increasingly popular.
- ▶ Analysis restricted to discrete-time models.

My research:

- ▶ systematically understand the performance of artificial agents in a continuous-time environment.



Learning algorithm focuses on policy, i.e., a function mapping system states to actions.

For a system with unknown parameter θ , issues in RL:

- ▶ Model identification: how to learn the parameter θ ?
Examples: consumer behaviour for online retailer, price impact factor in optimal execution.

Learning algorithm focuses on policy, i.e., a function mapping system states to actions.

For a system with unknown parameter θ , issues in RL:

- ▶ Model identification: how to learn the parameter θ ?
Examples: consumer behaviour for online retailer, price impact factor in optimal execution.
- ▶ Robustness of policies w.r.t. θ .

Learning algorithm focuses on policy, i.e., a function mapping system states to actions.

For a system with unknown parameter θ , issues in RL:

- ▶ Model identification: how to learn the parameter θ ?
Examples: consumer behaviour for online retailer, price impact factor in optimal execution.
- ▶ Robustness of policies w.r.t. θ .
- ▶ Convergence rate analysis:
 - ▶ critical for understanding algorithm efficiency.

- ▶ Tabular MDPs (finite states and actions): Watkins and Dayan 1992 (Q-learning), Williams 1992 (policy gradient), Jaksch, Ortner and Auer 2009, and many others.
- ▶ Infinite states and actions: (LQ-RL, $T = \infty$)
 - ▶ Sublinear regret: Abbasi-Yadkori and Szepesvari 2011, Dean et al 2018, Mania, Tu and Recht 2019, Cohen, Koren and Mansour 2019
 - ▶ Logarithmic regret (for special cases): Faradonbeh, Tewari and Michailidis 2020, Cassel, Cohen and Koren 2020, Lale, Azizzadenesheli, Hassibi and Anandkumar 2020

- ▶ Algorithm design: Modares and Lewis 2014, Doya 2000, Tallec, Blier and Ollivier 2019, Jia and Zhou 2021a, 2021b
- ▶ Asymptotic convergence analysis
 - ▶ LQ-RL ($T=\infty$): Duncan, Guo and Pasik-Duncan 1999, Rizvi and Lin 2018, Pang, Bian and Jiang 2020

— No convergence rate/regret analysis

- ▶ Algorithm design: Modares and Lewis 2014, Doya 2000, Tallec, Blier and Ollivier 2019, Jia and Zhou 2021a, 2021b
- ▶ Asymptotic convergence analysis
 - ▶ LQ-RL ($T=\infty$): Duncan, Guo and Pasik-Duncan 1999, Rizvi and Lin 2018, Pang, Bian and Jiang 2020

— No convergence rate/regret analysis

This talk:

- ▶ analyses regret of continuous-time RL over finite-time horizon.
- ▶ starts with linear-quadratic models and then extends to linear-convex models.

Fix $(A^*, B^*) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times k}$, minimise

$$J(\alpha; \theta^*) = \mathbb{E} \left[\int_0^T ((X_t^{\theta^*, \alpha})^\top Q X_t^{\theta^*, \alpha} + \alpha_t^\top R \alpha_t) dt \right],$$

where $X^{\theta^*, \alpha}$ satisfies the dynamics with parameter $\theta^* = (A^*, B^*)$:

$$dX_t = (A^* X_t + B^* \alpha_t) dt + dW_t, \quad X_0 = x_0,$$

and α is adapted to the information generated by $X^{\theta^*, \alpha}$. Here Q and R are given positive definite matrices.

- ▶ When θ^* is **known**, the optimal control α^{θ^*} is given by:

$$\alpha_t^{\theta^*} = \phi_t^{\theta^*}(X_t^{\theta^*}),$$

where

- ▶ X^{θ^*} is the optimal state process satisfying

$$dX_t = (A^*X_t + B^*\phi_t^{\theta^*}(X_t))dt + dW_t,$$

- ▶ $\phi_t^{\theta^*}(x) = K_t^{\theta^*}x$, where K^{θ^*} solves a Riccati ODE associated with θ^* .

- ▶ When θ^* is **known**, the optimal control α^{θ^*} is given by:

$$\alpha_t^{\theta^*} = \phi_t^{\theta^*}(X_t^{\theta^*}),$$

where

- ▶ X^{θ^*} is the optimal state process satisfying

$$dX_t = (A^* X_t + B^* \phi_t^{\theta^*}(X_t))dt + dW_t,$$

- ▶ $\phi_t^{\theta^*}(x) = K_t^{\theta^*} x$, where K^{θ^*} solves a Riccati ODE associated with θ^* .
- ▶ When θ^* is **unknown**, one needs to balance **exploration** (learning via interactions with the random environment) and **exploitation** (optimal control).

Learning via trial and error

Episodic setting



- ▶ Let $\hat{\theta}^{(m-1)}$ be the estimated parameter before m -th episode.

- ▶ Let $\hat{\theta}^{(m-1)}$ be the estimated parameter before m -th episode.
- ▶ Given $\hat{\theta}^{(m-1)}$, agent exercises a policy $\phi^{(m)}$ (which may depend on $\hat{\theta}^{(m-1)}$ or not) and observes a trajectory of

$$dX_t^m = (A^* X_t^m + B^* \phi_t^{(m)}(X_t^m))dt + dW_t^m,$$

Cost for the m -th episode is

$$J(\phi^{(m)}; \theta^*) = \mathbb{E} \left[\int_0^T \left((X_t^m)^\top Q X_t^m + \phi_t^{(m)}(X_t^m)^\top R \phi_t^{(m)}(X_t^m) \right) dt \right].$$

- ▶ Let $\hat{\theta}^{(m-1)}$ be the estimated parameter before m -th episode.
- ▶ Given $\hat{\theta}^{(m-1)}$, agent exercises a policy $\phi^{(m)}$ (which may depend on $\hat{\theta}^{(m-1)}$ or not) and observes a trajectory of

$$dX_t^m = (A^* X_t^m + B^* \phi_t^{(m)}(X_t^m))dt + dW_t^m,$$

Cost for the m -th episode is

$$J(\phi^{(m)}; \theta^*) = \mathbb{E} \left[\int_0^T \left((X_t^m)^\top Q X_t^m + \phi_t^{(m)}(X_t^m)^\top R \phi_t^{(m)}(X_t^m) \right) dt \right].$$

- ▶ Agent constructs $\hat{\theta}^{(m)}$ using observed trajectories of $(X^i)_{i=1}^m$.

- ▶ Let $\hat{\theta}^{(m-1)}$ be the estimated parameter before m -th episode.
- ▶ Given $\hat{\theta}^{(m-1)}$, agent exercises a policy $\phi^{(m)}$ (which may depend on $\hat{\theta}^{(m-1)}$ or not) and observes a trajectory of

$$dX_t^m = (A^* X_t^m + B^* \phi_t^{(m)}(X_t^m))dt + dW_t^m,$$

Cost for the m -th episode is

$$J(\phi^{(m)}; \theta^*) = \mathbb{E} \left[\int_0^T \left((X_t^m)^\top Q X_t^m + \phi_t^{(m)}(X_t^m)^\top R \phi_t^{(m)}(X_t^m) \right) dt \right].$$

- ▶ Agent constructs $\hat{\theta}^{(m)}$ using observed trajectories of $(X^i)_{i=1}^m$.
- ▶ Performance criteria – regret up to episode $N \in \mathbb{N}$:

$$\mathcal{R}(N) = \sum_{m=1}^N (J(\phi^{(m)}; \theta^*) - J(\phi_{\theta^*}; \theta^*)).$$

- ▶ Take actions $(\phi^{(1)}, \phi^{(2)}, \dots)$ to learn (A^*, B^*) and minimise \mathcal{R} .

Theorem

θ^* is identifiable under ϕ^e iff

- ▶ if $u \in \mathbb{R}^d$ and $v \in \mathbb{R}^p$ satisfy $u^\top x + v^\top \phi^e(t, x) = 0$ for all $(t, x) \in [0, T] \times \mathbb{R}^d$, then u and v are zero vectors.

Theorem

θ^* is identifiable under ϕ^e iff

- ▶ if $u \in \mathbb{R}^d$ and $v \in \mathbb{R}^p$ satisfy $u^\top x + v^\top \phi^e(t, x) = 0$ for all $(t, x) \in [0, T] \times \mathbb{R}^d$, then u and v are zero vectors.

For any given $\theta = (A, B)$, the greedy policy $\phi_\theta(t, x) = K_t^\theta x$ identifies θ^* iff B is full column rank.

- ▶ **Self-exploration:** if B^* is full column rank, the greedy policy ϕ_{θ^*} explores the environment, and hence explicit exploration is not required.

- ▶ Given fixed policy ϕ , estimate $\theta^* = (A^*, B^*)$ via

$$dX_t = \theta^* Z_t^\phi dt + dW_t, \quad Z_t^\phi = (X_t, \phi(t, X_t))^\top.$$

- ▶ Agent only observes the state process X (also Z^ϕ), but not the corresponding Brownian path.
- ▶ View unknown θ^* as a **hidden** random variable, and observe the log-likelihood of θ^* is quadratic.

- ▶ Discrete time approximation is given by

$$X_{k+1} - X_k = \theta^* Z_k^\phi \Delta t + \sqrt{\Delta t} \xi_k, \quad \text{and} \quad \xi_k \sim_{IID} N(0, 1).$$

- ▶ Discrete time approximation is given by

$$X_{k+1} - X_k = \theta^* Z_k^\phi \Delta t + \sqrt{\Delta t} \xi_k, \quad \text{and} \quad \xi_k \sim_{IID} N(0, 1).$$

- ▶ One can compute likelihood function

$$\ell(\theta^* \mid X_1, \dots, X_n) \propto$$

$$\exp \left(-\frac{1}{2} \theta^* \left(\sum_{k=0}^{n-1} Z_k^\phi (Z_k^\phi)^\top \Delta t \right) (\theta^*)^\top + \theta^* \sum_{k=0}^{n-1} Z_k^\phi (X_{k+1} - X_k) \right),$$

- ▶ Discrete time approximation is given by

$$X_{k+1} - X_k = \theta^* Z_k^\phi \Delta t + \sqrt{\Delta t} \xi_k, \quad \text{and} \quad \xi_k \sim_{IID} N(0, 1).$$

- ▶ One can compute likelihood function

$$\ell(\theta^* \mid X_1, \dots, X_n) \propto$$

$$\exp\left(-\frac{1}{2}\theta^* \left(\sum_{k=0}^{n-1} Z_k^\phi (Z_k^\phi)^\top \Delta t\right) (\theta^*)^\top + \theta^* \sum_{k=0}^{n-1} Z_k^\phi (X_{k+1} - X_k)\right),$$

which leads to

$$\ell(\theta^* \mid X) \propto \exp\left(-\frac{1}{2}\theta^* \left(\int_0^t (Z_s^\phi)(Z_s^\phi)^\top ds\right) (\theta^*)^\top + \theta^* \int_0^t Z_s^\phi dX_s\right).$$

Given a prior $\theta^* \sim N(\hat{\theta}^{(m-1)}, V^{(m-1)})$, the posterior of θ^* based on observation Z^ϕ is also Gaussian $N(\hat{\theta}^{(m)}, V^{(m)})$.

Given a prior $\theta^* \sim N(\hat{\theta}^{(m-1)}, V^{(m-1)})$, the posterior of θ^* based on observation Z^ϕ is also Gaussian $N(\hat{\theta}^{(m)}, V^{(m)})$.

Theorem

For all $\delta > 0$,

$$|\hat{\theta}^{(m)} - \theta^*|^2 \leq \frac{1}{\lambda_{\min}((V^{(m)})^{-1})} \text{poly}(\ln m, \ln(\frac{1}{\delta})), \quad \forall m \geq 2.$$

- ▶ $\lambda_{\min}((V^{(m)})^{-1})$ increases if ϕ^e is executed.
- ▶ Sub-exponential concentration: Z_t^ϕ is Gaussian, and hence $\int_0^T Z_t^\phi (Z_t^\phi)^\top dt$ and $\int_0^T Z_t^\phi dX_t$ are sub-exponential.

Algorithm 1: PEGE Algorithm

Input: $m : \mathbb{N} \rightarrow \mathbb{N}$.

```
1 Initialize  $m = 0$ .
2 for  $k = 1, 2, \dots$  do
3   | Execute the exploration policy  $\phi^e$  for one episode, and  $m \leftarrow m + 1$ .
4   | Update the estimate  $\hat{\theta}^{(m)}$  and set  $\bar{\theta} = \hat{\theta}^{(m)}$ .
5   | for  $l = 1, 2, \dots, m(k)$  do
6   |   | Execute the greedy policy  $\phi_{\bar{\theta}}$  for one episode, and  $m \leftarrow m + 1$ .
7   | end
8 end
```

-
- ▶ $m : \mathbb{N} \rightarrow \mathbb{N}$ determines the exploration and exploitation trade-off.

Let $\mathcal{E}^\Phi = \{m \in \mathbb{N} \mid \phi^{(m)} = \phi^e\}$ and consider

$$\begin{aligned} \mathcal{R}(N) &= \sum_{m=1}^N (J(\phi^{(m)}; \theta^*) - J(\phi_{\theta^*}; \theta^*)) \\ &= \sum_{m \in [1, N] \cap \mathcal{E}^\Phi} (J(\phi^e, \theta^*) - J(\phi_{\theta^*}; \theta^*)) + \sum_{m \in [1, N] \setminus \mathcal{E}^\Phi} (J(\phi_{\hat{\theta}^{(m-1)}}; \theta^*) - J(\phi_{\theta^*}; \theta^*)) \end{aligned}$$

Let $\mathcal{E}^\Phi = \{m \in \mathbb{N} \mid \phi^{(m)} = \phi^e\}$ and consider

$$\begin{aligned} \mathcal{R}(N) &= \sum_{m=1}^N (J(\phi^{(m)}; \theta^*) - J(\phi_{\theta^*}; \theta^*)) \\ &= \sum_{m \in [1, N] \cap \mathcal{E}^\Phi} (J(\phi^e; \theta^*) - J(\phi_\theta; \theta^*)) + \sum_{m \in [1, N] \setminus \mathcal{E}^\Phi} (J(\phi_{\hat{\theta}^{(m-1)}}; \theta^*) - J(\phi_{\theta^*}; \theta^*)) \end{aligned}$$

Theorem (Performance gap)

For all $\varepsilon > 0$, $\exists C_\varepsilon \geq 0$,

$$|J(\phi_\theta; \theta^*) - J(\phi_{\theta^*}; \theta^*)| \leq C_\varepsilon |\theta - \theta^*|^2, \quad \forall \theta \in \mathbb{B}_\varepsilon(\theta^*).$$

Let $\mathcal{E}^\Phi = \{m \in \mathbb{N} \mid \phi^{(m)} = \phi^e\}$ and consider

$$\begin{aligned} \mathcal{R}(N) &= \sum_{m=1}^N (J(\phi^{(m)}; \theta^*) - J(\phi_{\theta^*}; \theta^*)) \\ &= \sum_{m \in [1, N] \cap \mathcal{E}^\Phi} (J(\phi^e, \theta^*) - J(\phi_\theta; \theta^*)) + \sum_{m \in [1, N] \setminus \mathcal{E}^\Phi} (J(\phi_{\hat{\theta}^{(m-1)}}; \theta^*) - J(\phi_{\theta^*}; \theta^*)) \end{aligned}$$

Theorem (Performance gap)

For all $\varepsilon > 0$, $\exists C_\varepsilon \geq 0$,

$$|J(\phi_\theta; \theta^*) - J(\phi_{\theta^*}; \theta^*)| \leq C_\varepsilon |\theta - \theta^*|^2, \quad \forall \theta \in \mathbb{B}_\varepsilon(\theta^*).$$

This gives $\mathcal{R}(N) \lesssim \left| [1, N] \cap \mathcal{E}^\Phi \right| + \sum_{m \in [1, N] \setminus \mathcal{E}^\Phi} |\hat{\theta}^{(m-1)} - \theta^*|^2$.

Theorem

For $m(k) = k$, $k \in \mathbb{N}$, with high probability,

$$\mathcal{R}(N) \leq CN^{\frac{1}{2}}(\log N)^2, \quad \forall N \geq 2.$$

Theorem

For $m(k) = k$, $k \in \mathbb{N}$, with high probability,

$$\mathcal{R}(N) \leq CN^{\frac{1}{2}}(\log N)^2, \quad \forall N \geq 2.$$

If *self-exploration* property holds, then by setting $m(k) = 2^k$, $k \in \mathbb{N}$, with high probability,

$$\mathcal{R}(N) \leq C(\log N)^2, \quad \forall N \geq 2.$$

Theorem

For $m(k) = k$, $k \in \mathbb{N}$, with high probability,

$$\mathcal{R}(N) \leq CN^{\frac{1}{2}}(\log N)^2, \quad \forall N \geq 2.$$

If *self-exploration* property holds, then by setting $m(k) = 2^k$, $k \in \mathbb{N}$, with high probability,

$$\mathcal{R}(N) \leq C(\log N)^2, \quad \forall N \geq 2.$$

- ▶ Discrete time observations and actions have the same regret order, with an additional discretization error.

Let $\theta^* = (A^*, B^*) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times k}$ be fixed but unknown, minimise

$$J(\alpha; \theta^*) = \mathbb{E} \left[\int_0^T f_t(X_t^{\theta^*, \alpha}, \alpha_t) dt + g(X_T^{\theta^*, \alpha}) \right],$$

over stochastic processes α , where $X^{\theta^*, \alpha}$ satisfies the dynamics with θ^* :

$$dX_t = (A^* X_t + B^* \alpha_t) dt + \sigma dW_t + \int_{\mathbb{R}^p \setminus \{0\}} \gamma(u) \tilde{N}(dt, du),$$

f and g are convex in state and strongly convex in control.

- ▶ f can be nonsmooth, and includes action constraints, ℓ^1 -norm or entropy regularisers.

	LQ-RL	LC-RL
Greedy policy	Linear	Nonlinear
Policy characterisation	Riccati ODE	BSDE
Performance gap	Quadratic	Linear/quadratic
Parameter estimation	Bayesian	Least-squares
Estimation error	Sub-exponential r.v.	Sub-Weibull r.v.

- ▶ Consider the 3d controlled SDE:

$$dX_t = (AX_t + B\alpha_t) dt + dW_t, \quad t \in [0, 1.5].$$

with unknowns A, B from Dean et al. 2018, and a given cost

$$J(\alpha) = \mathbb{E} \left[\int_0^T (0.1|X_t^\alpha|^2 + |\alpha_t|^2) dt \right].$$

- ▶ Consider the 3d controlled SDE:

$$dX_t = (AX_t + B\alpha_t) dt + dW_t, \quad t \in [0, 1.5].$$

with unknowns A, B from Dean et al. 2018, and a given cost

$$J(\alpha) = \mathbb{E} \left[\int_0^T (0.1|X_t^\alpha|^2 + |\alpha_t|^2) dt \right].$$

- ▶ Run PEGE algorithm with $m(k) = 2^k$, $k \in \mathbb{N}$.
- ▶ Perform 100 independent executions to estimate statistical properties of the algorithm.

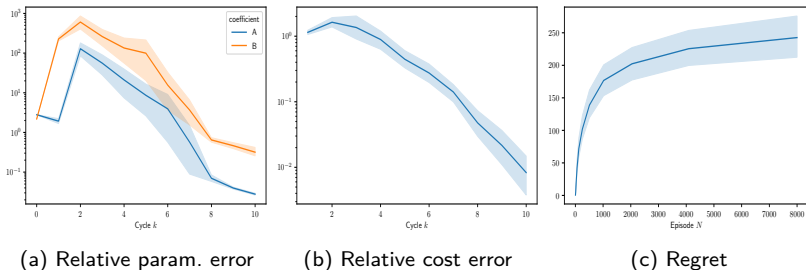


Figure: Numerical results from 100 repeated experiments; solid lines are sample means and shallow areas are 95% confidence intervals.

Two complimentary aspects on model-based RL:

- ▶ Finite-sample analysis of parameter estimation ([statistical learning theory](#)) and performance gap analysis of greedy policy ([control theory](#)).
- ▶ A phase-based learning algorithm with optimal regrets for linear-convex models.

Two complimentary aspects on model-based RL:

- ▶ Finite-sample analysis of parameter estimation ([statistical learning theory](#)) and performance gap analysis of greedy policy ([control theory](#)).
 - ▶ A phase-based learning algorithm with optimal regrets for linear-convex models.
- (1) Basei, Guo, Hu, Zhang, *Logarithmic regret for episodic continuous-time linear-quadratic reinforcement learning over a finite-time horizon*, JMLR, to appear, 2020.
 - (2) Szpruch, Treetanhiplot and Zhang, *Exploration-exploitation trade-off for continuous-time episodic reinforcement learning with linear-convex models*, arXiv preprint, 2021.