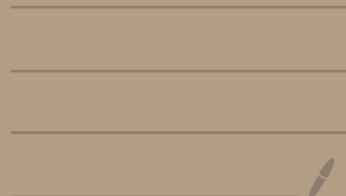


ML & SC

- ① Terms of RL & MDP
- ② DL for SC & PDEs



① Foundation of ML & SC.

① Basics of MDP $(X, A, P, r = (f, g))$

① X : state ; A : action ; $X_{t+1} \sim P_t(x_t, a_t)$;
transition probability (t, x, a)

Reward: ① running reward $f(x, a)$ ② terminal reward $g(x)$ ③ Discount factor $\gamma \in [0, 1]$

② Basics of RL

① policy: π (sequence $T\pi$) \Rightarrow ① deterministic ② stochastic
 $T\pi: X \rightarrow A$ $T\pi: X \rightarrow P(A)$ control: $a \sim \pi$

② State value function: $V_t^\pi(x) = E_{\pi} T \sum_{s=t}^{T-1} f(x_s, a_s) + g(x_T) | X_0=x$, $t \in \{0, \dots, T\}$, $x \in X$

$$V_t^\pi(x) = E[\pi Q_t^\pi(x, a)]$$

③ Q-function of π : $Q_t^\pi(x, a) = E_{\pi} T \sum_{s=t}^{T-1} f(x_s, a_s) + g(x_T) | X_0=x, a_0=a$, $t \in \{0, \dots, T\}$, $(x, a) \in X \times A$.

④ goal: find an optimal policy π^* maximize V^π $V^* = \sup_{\pi} V^\pi$, $Q^* = \sup_{\pi} Q^\pi$

③ DP: MDP is global dynamic optimisation \Rightarrow reduce the optimisation problem to local optimisation problem.

① Bellman Equation $V_T^\pi(x) = g(x)$; $V_t^\pi(x) = E[f(x, a) + \underbrace{\gamma E[V_{t+1}^\pi(x')]}_{Q_t^\pi(x, a)}]$ immediate reward
~~fix $t \in X$ to get certificate~~
~~expected future reward $E[V^\pi(x)]$~~

④ Classes of RL: ① Value-based methods: learn value & Q \rightarrow policy derived
② Policy-based methods: learn policy \rightarrow value function estimated

⑤ TD-learning: $V_t^\pi(x) = E[f(x, a) + V_{t+1}^\pi(x')]$ prediction
 $\text{TD}(t) \propto \text{TD}(t-1) + \eta \hat{V}_t^\pi(x) + \eta [f_t + \hat{V}_{t+1}^\pi(x') - \hat{V}_t^\pi(x)]$
 $\frac{\text{TD}(t)}{\text{TD}(t-1)}$ learning rate.

⑥ TD-learning & Q-learning

① Bellman Equation $V_t^\pi(x) = E[f(x, a) + V_{t+1}^\pi(x')]$ \rightarrow ① initialize ② $\pi \sim \pi_t(x, a, f_t, x')$

$$\hat{V}_t^\pi(x) \leftarrow \hat{V}_t^\pi(x) + \eta [f_t + \hat{V}_{t+1}^\pi(x') - \hat{V}_t^\pi(x)]$$

$$\text{TD } \delta_t = f_t + \hat{V}_{t+1}^\pi(x') - \hat{V}_t^\pi(x)$$

$\text{TD}(t) = \lambda = 1$ 用整系数加速的回报更新, $\lambda \in [0, 1]$.

② Q-learning $Q_t(x, a) = f(x, a) + E[\sup Q_{t+1}(x', a')]$

Robbins-Monro condition: $\sum \eta_t = \infty$ $\Rightarrow \hat{Q}_t(x, a) \xrightarrow{\eta_t} Q_t(x, a)$

① SARSA: online: 同时更新 policy & action
offline: 全量高维更新 policy & action.

$$\text{① PE: } Q_t^{\pi_t^k}(x_t, a) = f(x_t, a) + E [Q_{t+1}^{\pi_t^k}(x', a')] \quad \text{② PI: } \pi_t^{k+1}(x) \in \arg\max_{a \in A} Q_t^{\pi_t^k}(x, a)$$

$$(\text{SARSA}) = \delta_t = f_t + \hat{Q}_t(x', a') - \hat{Q}_t(x, a) \quad (\text{Q-learning}) \quad \delta_t = f_t + \sup_a \hat{Q}_t(x', a') - \hat{Q}_t(x, a)$$

③ policy:

① ϵ -greedy policy: $a_t \sim \begin{cases} \arg\max_{a \in A} Q_t(x, a) \text{ with } 1-\epsilon \\ \text{uniform distribution} \end{cases}$ - exploration
 $\mathcal{U}(A)$ with ϵ - exploitation

② softmax policy: $a_t \sim \pi_t(x, a)$ $\pi_t(x, a) = \frac{\exp(Q_t(x, a)/\lambda)}{\sum_{a' \in A} \exp(Q_t(x, a')/\lambda)}$ $a \in A$ $\lambda \rightarrow \infty$ constant explore
 $\lambda \rightarrow 0^+$ extreme greedy temperature parameter.

③ parametrised randomized policy:

discrete ① Softmax: $\pi_t(x_t, a_t) = \frac{\exp(\varphi_0(t, x_t, a))}{\sum \exp(\varphi_0(t, x_t, a))} \Rightarrow \nabla_t \ln \pi_t(x_t, a) = \nabla \varphi_0(t, x_t, a) - E[\nabla \varphi_0(t, x_t, \cdot)]$

continuous ② Gaussian: $a \mapsto \pi_t(x_t, a) \sim N(\mu_t(x), \sigma^2) \Rightarrow \nabla_t \ln \pi_t(x_t, a) = \frac{(a - \mu_t(x)) \nabla \mu_t(x)}{\sigma^2}$

$$J(\theta) = E_{\pi_\theta} \left[\sum_{t=0}^{T-1} f(x_t, a_t) + g(x_T) \right] \quad R(S) = \sum_{t=0}^{T-1} f(x_t, a_t) + g(x_T) \quad J(\theta) = E[R(S)]$$

$$\nabla_\theta J(\theta) = E_{\pi_\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta \ln \pi_\theta(t, x_t, a_t) \right]$$

$$\begin{aligned} \nabla_\theta J(\theta) &= \int R(s) \nabla_\theta p_\theta(s) Q(ds) = \int R(s) [\nabla_\theta \ln p_\theta(s)] p_\theta(s) Q(ds) \\ &= \mathbb{E}_{\pi_\theta} [R(S) \nabla_\theta \ln p_\theta(S)] = \mathbb{E}_{\pi_\theta} [R(S) \sum_{t=0}^{T-1} \nabla_\theta \ln \pi_\theta(t, x_t, a_t)]. \end{aligned}$$

④ Reinforce Algorithms:

① Sampling trajectory: $S^k = (x_0^k, a_0^k, x_1^k, \dots, x_T^k)$ $a_t^k \sim \pi_{\theta_k}(t, x_t^k, \cdot)$ $x_{t+1}^k \sim P_t(x_t^k, a_t^k, \cdot)$ $f_t^k = f(x_t^k, a_t^k)$

$$R(S^k) = \sum_{t=0}^{T-1} f_t^k - f_T^k$$

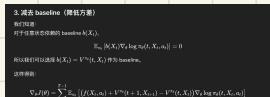
$$\text{② } \theta \leftarrow \theta + \eta \left[\frac{1}{K} \sum_{k=1}^K R(S^k) \sum_{t=0}^{T-1} \nabla_\theta \ln \pi_\theta(t, x_t^k, a_t^k) \right]$$

⑩ Actor-Critic network approach:

Value function (the critic); policy (actor) by NN.

by PE

by policy gradient in direction.



⑪ Basics of SC: dynamical system: $X(t)$; t ; $\frac{dX}{dt} = f(X, t)$, $X_{t+1} = F(X_t, u)$

optimisation of DS in a random environ.

① By SDEs: $dX_t = b(X_t, \alpha_t) dt + \sigma(X_t, \alpha_t) dW_t \rightarrow \text{SDS}$

State dynamics: drift term: 方向速度 diffusion term: 扩散速度 Brownian motion 随机游走 $dW_t \sim N(0, dt)$

② objective: $J(\alpha) = E \left[\int_0^T f(X_t, \alpha_t) dt + g(X_T) \right]$ max/min over $\alpha = (\alpha_t)_t$

average running cost sum up terminal reward

② DL for SC & PDEs

① For MDP

Algo NN Conv P: control learning by policy iteration.

Algo Action / Critic

$$\mathbb{E}_M |V_t(X_t) - \hat{V}_t^M(X_t)| = \mathcal{O}_{\mathbb{P}} \left(\underbrace{\left(\gamma^t K \frac{\ln M}{M} \right)^{\frac{1}{2(t-n)}}}_{\text{statistical error}} + \underbrace{\varepsilon_{t,n}^{\text{NN}}(\pi^*)}_{\text{policy approximation error}} + \underbrace{\varepsilon_{t,n}^{\text{NN}}(V)}_{\text{value approximation error}} \right)$$

DNN + SGD Backward induction

$$\hat{\pi}_{\pi^*} \in \underset{\pi}{\operatorname{argmax}} \left[\mathbb{E} \left[f(X_t, \pi(X_t)) + \sum_{s=t+1}^{T-1} f(\hat{X}_s, \hat{\pi}_{\pi^*}(\hat{X}_s)) + g(\hat{X}_T) \right] \right]$$

$$\hat{\pi}_{\pi^*} \in \underset{\pi}{\operatorname{argmax}} \left[\mathbb{E} \left[f(X_t, \pi(X_t)) + V_{\pi^*}(X_{t+1}) \right] \right]$$

$$\mathbb{E}_M |V_t(X_t) - \hat{V}_t^M(X_t)| = \mathcal{O}_{\mathbb{P}} \left(\underbrace{\gamma^{T-t} \sqrt{\frac{\ln M}{M}}}_{\text{statistical error}} + \underbrace{\sup_{s=t, \dots, T-1} \inf_{p \in \mathcal{A}_K} \|p(X_s) - \pi_s^*(X_s)\|_{L^1}}_{\varepsilon_t^{\text{NN}}(\pi^*) \text{ approximation error}} \right)$$

② For PDE

$$\frac{\partial v(t,x)}{\partial t} = \mathcal{N}(v)(t,x) = f(t,x, v, \nabla_x v, \dots), \quad (+ \text{ boundary conditions})$$

diff w.r.t. t

PDE

初值条件 + 边界条件

① PDE $\mathcal{E}[0,T] \times \mathbb{R}^d$; objective function $v(t,x)$

DNN approximates $v(t,x) / \nabla_x v$

① loss function \Rightarrow ② SGD on X (domain) ③ TensorFlow autocompute differentiation

$$\begin{cases} \frac{\partial v(t,x)}{\partial t} = \mathcal{N}(v)(t,x) \quad (t,x) \in \mathcal{E}[0,T] \times \mathbb{R}^d \\ v(T,x) = g(x) \quad x \in \mathbb{R}^d \end{cases}$$

\rightarrow NN \Rightarrow PDE

$$\text{loss function: } L(v_0, t, x) = \|v_0(t,x) - g(x)\|^2 + \|\partial_t v_0(t,x) - \mathcal{N}(v_0)(t,x)\|^2$$

① why PDE: independent variables: PDEs - $t \& x + \partial_t V / \nabla_x V$

PDE + terminal condition / boundary condition \Rightarrow unique solution

② NN approximation: NN $v_0(t,x) \rightarrow V(t,x)$

input $t, x \rightarrow$ output $V(t,x)$

auto-diff: time derivative $\Delta t \rightarrow$ space derivative $\Delta x \rightarrow$

2nd derivative $\nabla^2 v_0$

③ loss function: NN satisfy: ① PDE ② boundary / terminal condition.

$$\text{PINN loss } L(v_0) = E \left[\underbrace{\|v_0(t,x) - g(x)\|^2}_{\text{terminal condition err.}} + \underbrace{\|\partial_t v_0 - \mathcal{N}[v_0]\|^2}_{\text{PDE err.}} \right]$$

$\rightarrow f \Rightarrow g(x)$

即 $\nabla_x \nabla_t \text{NN err.}$

$\rightarrow f \Rightarrow g(x)$

即 $\nabla_x \nabla_t \text{NN err.}$

DGNN loss \rightarrow DGNN + Monte Carlo Sampling

④ Optimization: domain $\mathcal{E}[0,T] \times \mathbb{R}^d \rightarrow v(t,x)$ Sampling randomly

Calculate PDE err + terminal condition err $\Rightarrow \partial_t v_0 - \partial_t v_0 - \mathcal{N}[v_0] \text{ (loss gradient)}$

5. DGM vs PINN

	DGM	PINN
采样方式	Monte Carlo 随机采样	固定网格或采样点
优点	高维问题不需要网格, 灵活	容易融合物理定律约束
典型应用	金融、控制	物理模拟、工程 PDE

⑤ with control problems:

- ① In stochastic control, optimal function V satisfy PDE (HJB function)
- ② traditional method: PDE \Rightarrow curse of dimensionality
- ③ DGMs / PINNs NN approximates V via sampling

③ RL methods in continuous times.

① Problem Setting:

- ① continuous MDP: $dX_t = b(X_t, \alpha_t) dt + \sigma(X_t, \alpha_t) dW_t$
 X_t : State, α_t : action / controls (continuous), b, σ
- ② objective: $J(\alpha) = \mathbb{E} \left[\int_0^T f(X_s, \alpha_s) ds + g(X_T) \right]$

② Policy:

- ① Deterministic policy: $\alpha_t = \pi_\theta(X_t)$
- ② Randomized policy: $\alpha_t \sim \pi_\theta(x)$

③ Value function & HJB PDE:

$$V^\pi(x) = \mathbb{E}^\pi \left[\int_x^T f(X_s, \alpha_s) ds + g(X_T) \mid X_0 = x \right]$$

$$\text{HJB PDE: } \partial_t V + \sup_{\alpha} \left[b(x, \alpha) \cdot \nabla_x V + \frac{1}{2} \text{Tr} (\sigma \sigma^T \nabla_x^2 V) + f(x, \alpha) \right] = 0 \quad V(T, x) = g(x)$$

④ Bellman Principle in Continuous Time

任意时刻的最优策略 = 从当前状态往后看仍然要遵循最优策略 \Rightarrow HJB form

⑤ RL approaches in continuous time.

- ① Policy Gradient Method: $\theta \leftarrow \theta + \eta \nabla \mathbb{E}[J(\theta)] \Rightarrow$ Ito formula + Stochastic Calculus

- ② Action-Value Method: Action policy + Value Function \Rightarrow on-policy

- ③ Q-learning in Continuous time: ① Define $Q(x, \alpha) \Rightarrow$ Bellman Equations \Rightarrow ② Sweepplay + Approximate NN

⑥ Function Approximation

- ① NN approximates $V(x)$ / $Q(x, \alpha)$

HJB PDEs:

System: $dX_t = b(X_t, \alpha_t) dt + \sigma(X_t, \alpha_t) dW_t$

Bellman: \mathbb{E}^π optimal + \mathbb{E}^π optimal = global optimal, $V(x) = \sup_{\alpha} \mathbb{E} \left[\int_x^T f(X_s, \alpha_s) ds + g(X_T) \mid X_0 = x \right]$

HJB PDE 就是一个方程，它的解。 $V(x)$ 告诉你在任意时间、任意状态下能得到的最优总收益；解这个 PDE，就等于解决了整个最优控制问题。