

Lecture10: Markov Decision Process for Reinforcement Learning

Notes taken by [squarezhong](#)

Repo address: [squarezhong/SDM5008-Lecture-Notes](#)

Lecture10: Markov Decision Process for Reinforcement Learning

From Classical Control to RL

What control do?

Comparison

Typical optimization problem

RL

Markov Chain

Markov Decision Process

Policy

Trajectories of MDP

Notations

Return

Bellman Equations

Value functions

Bellman Equations

Sampling

Monte Carlo Method

Intrinsic

Estimation

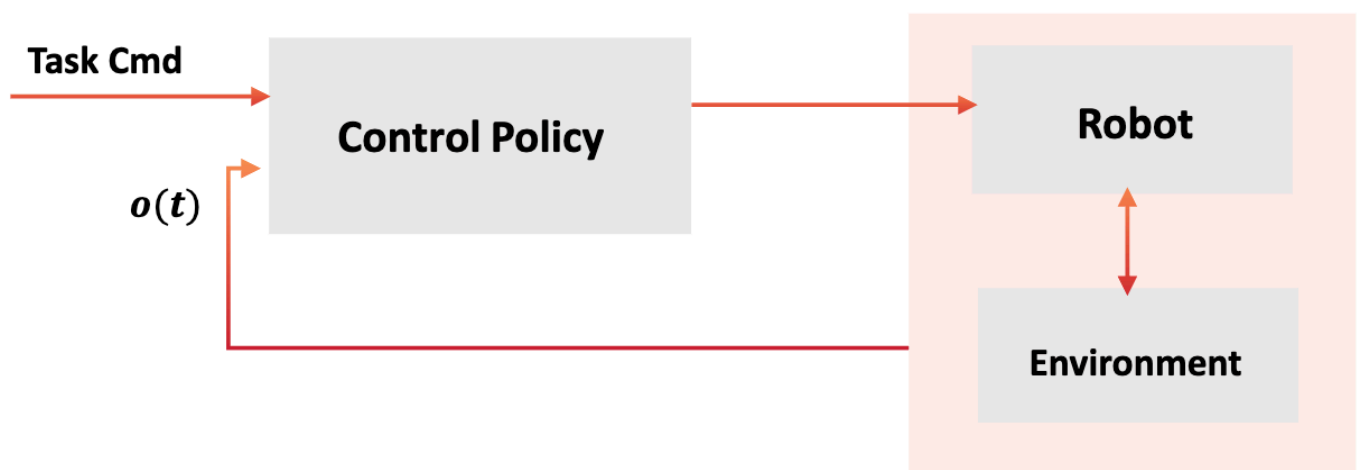
Central limit theorem (CLT)

Monte Carlo Integration

Importance sampling

From Classical Control to RL

What control do?



Comparison

Classical Control	Modern Control	MPC	Reinforcement Learning
linear system frequency domain	state space linear system limited class of nonlinear	state space linear and nonlinear (large scale)	complex model
analytical	analytical	computational	data driven

- Most above are "model" based method.
 - classical/modern: model $\xrightarrow{\text{analytical/definition}}$ policy
 - MPC: model $\xrightarrow{\text{formulate}}$ optimization \rightarrow policy
 - RL: model \rightarrow data \rightarrow policy

Typical optimization problem

$$\min f(x) \quad \text{subject to } g(x) \leq 0$$

1. typically $x \in \mathbb{R}^n$, finite dimensional optimization variable
2. objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, easy to define and compute
3. deterministic

RL

$$\max \text{Reward}(\pi) \quad \text{for all policies } \pi$$

1. optimization variable policy π infinite dimensional
2. objective function (functional 泛函) $\text{Reward} : \text{policy} \rightarrow \text{scalar}$
elevation of reward requires simulation robot \leftrightarrow env
3. probalistic/stochastic problem due to intrinsic uncertainty

Markov Chain

Markov Chain: $\text{MC} = (S, \Gamma)$

- S : state space (discrete or continuous)
- Γ : transition operator, i.e. $\Gamma(x|y) = \Pr(s_{t+1} = x|s_t = y)$
- Initial distribution $p_0(s) = \Pr(S_0 = s)$
- For discrete state space, the transition operator has a matrix representation.

- MC with P_0 specifies a way to generate sequential random samples s_0, s_1, \dots , which is called **realization/trajectory** of the MC.
- Markov chain can be seen as a stochastic dynamical system that
 - $s_{k+1} = f(s_k, w_k)$ or $s_{k+1} = f(s_k) + w_k$ where w_k is a random variable (process noise)

Markov Decision Process

System future behavior depends on state s_t and external impact (control/action)

$$\text{MDP} = (S, \mathcal{A}, \Gamma, r)$$

- S : state space (discrete or continuous)
- \mathcal{A} : action/control space (discrete or continuous)
- Γ : transition kernel/operator

$$\Gamma(s'|s, a) = \Pr(S_{t+1} = s' | S_t = s, A_t = a) = p(s'|s, a)$$
- r : reward function: $r(s, a, s')$ or typically $r(s, a)$

Policy

- Markov decision: agent makes decision based on current stats
- $\pi(a|s)$ is the pdf/pmf of action a given the current state s
 - i.e. $\pi(a|s) = \Pr(A = a | S = s)$
- Deterministic policy $a = \pi(s)$
- Time varying policy, $\pi_t(a|s)$
- Policy within a class of functions with certain parameters θ

Trajectories of MDP

- Given policy π , a finite horizon T
- MDP becomes a **MC** with "closed-loop" transition operator Γ_{cl}

$$\Gamma_{cl}(s'|s) = \Pr(S_{t+1} = s' | S_t = s) = \sum_a p(s'|s, a) \pi(a|s)$$
- Trajectory $\tau = (s_0, a_0, s_1, a_1, \dots, s_T, a_T)$ is a trajectory of MDP under a policy π
- Probability of a trajectory $P(\tau|\pi) = p(s_0) \prod_{k=0}^{T-1} p(s_{k+1}|s_k, a_k)$

Notations

- $a \sim \pi(\cdot | s) \Leftrightarrow a \sim \pi$
 - $s' \sim p(\cdot | s, a) \Leftrightarrow s' \sim \mathcal{P}$
- } when clear from context
- $\tau \sim \pi$ means $\tau \sim \text{pr}(\tau | \pi)$
 \nwarrow a "big" random vector
 - $E_{\tau \sim \pi}(\gamma(\tau)) = \sum_{\tau} \text{pr}(\tau | \pi) \cdot \gamma(\tau)$ or $\int \gamma(\tau) \text{pr}(\tau | \pi) d\tau$
 \nwarrow γ is func of τ , i.e. each τ gets mapped to a value/vector
 - $E_{\tau \sim \pi}(\gamma(\tau) | S_0 = s) = \sum_{\tau} \gamma(\tau) \cdot \text{pr}(\tau | \pi, S_0 = s)$ ($\tau = (S_0, A_0, S_1, A_1, \dots)$)
 - $E_{\tau \sim \pi}(\gamma(\tau)) = E_{s \sim p_0} \left(E_{\tau \sim \pi}(\gamma(\tau) | S_0 = s) \right) = \sum_s \text{pr}(S_0 = s) \cdot \left(\sum_{\tau} \text{pr}(\tau | S_0 = s, \pi) \cdot \gamma(\tau) \right)$
 \nwarrow tower property.

For infinite sequence $\tau = (S_0, A_0, S_1, A_1, \dots)$
 $\tau_{\geq 1} = (S_1, A_1, S_2, A_2, \dots)$ } both under π .

$$\tau \sim \pi | S_0 = s \Leftrightarrow \tau_{\geq 1} \sim \pi | S_1 = s$$

same distribution, only different notation.

Return

Cumulative rewards over a trajectory, which may take several different forms.

- Finite-horizon (undiscounted) return $R(\tau) = \sum_{t=0}^T r(s_t, a_t)$
- Infinite-horizon discounted return $R(\tau) = \sum_{t=0}^{\infty} \alpha^t r_t$

where $\alpha \in (0, 1)$ is discount factor, future reward is less important than immediate reward

MDP (RL) Problem: $\max E_{\tau \sim \pi}[R(\tau)]$

Bellman Equations

Value functions

- On-policy (state)-value function:

$$V_{\pi}(s) \triangleq E_{\tau \sim \pi}(R(\tau) | S_0 = s)$$

evaluate the "performance" of a given policy π

- On-policy action-value function (**Q-function**)

$$Q_{\pi}(s, a) = E_{\tau \sim \pi}[R(\tau) | S_0 = s, A_0 = a]$$

- Optimal value function

$$V^*(s) = \max_{\pi} (E_{\tau \sim \pi}(R(\tau) | S_0 = s))$$

- Optimal action-value function

$$Q^*(s, a) = \max_{\pi} (E_{\tau \sim \pi}[R(\tau) | S_0 = s, A_0 = a])$$

By definition, we have the following deduction

- $V_{\pi}(s) = E_{a \sim \pi}[Q_{\pi}(s, a)]$
- $V^*(s) = \max_a Q^*(s, a)$

Bellman Equations

Bellman equation is a **necessary condition** for optimality associated with the mathematical optimization method known as dynamic programming. It writes the "value" of a decision problem at a certain point in time in terms of the payoff from some initial choices and the "value" of the **remaining decision problem** that results from those initial choices. [\[source\]](#)

Infinite-horizon discounted return case

- $V_{\pi}(s)$

$$\begin{aligned}
 V_{\pi}(s) &= E_{\tau \sim \pi}(R(\tau) | S_0 = s) = E_{\tau \sim \pi}(r(S_0, A_0) + \alpha r(S_1, A_1) + \alpha^2 r(S_2, A_2) + \dots | S_0 = s) \\
 &\stackrel{\text{towering property}}{=} E_{a \sim \pi(\cdot | s)} \left[E_{\tau \sim \pi}(r(S_0, A_0) + \alpha r(S_1, A_1) + \dots | S_0 = s, A_0 = a) \right] \\
 &= E_{a \sim \pi(\cdot | s)} \left(r(s, a) + \alpha E_{\tau \sim \pi}(r(S_1, A_1) + \alpha r(S_2, A_2) + \dots | S_0 = s, A_0 = a) \right) \\
 &= E_{a \sim \pi(\cdot | s)} \left(r(s, a) + \alpha E_{\tau \sim \pi} \left(E_{\tau' \sim \pi}(r(S_1, A_1) + \alpha r(S_2, A_2) + \dots | S_0 = s, A_0 = a, S_1 = s') \right) \right) \\
 &= E_{a \sim \pi(\cdot | s)} \left(r(s, a) + \alpha E_{s' \sim p(\cdot | s, a)} (V_{\pi}(s')) \right)
 \end{aligned}$$

- $Q_\pi(s, a)$

$$Q_\pi(s, a) = E_{\tau \sim \pi} [r(S_0, A_0) + \alpha r(S_1, A_1) + \dots \mid S_0 = s, A_0 = a]$$

$$= r(s, a) + \alpha E_{\tau \sim \pi} [r(S_1, A_1) + \alpha r(S_2, A_2) + \dots \mid S_0 = s, A_0 = a]$$

$$= r(s, a) + \alpha E_{\substack{s' \sim p(\cdot | s, a) \\ a' \sim \pi(\cdot | s')}} \left[E_{\tau' \sim \pi} [r(S_1, A_1) + \alpha r(S_2, A_2) + \dots \mid S_1 = s', A_1 = a'] \right]$$

$$= r(s, a) + \alpha E_{\substack{s' \sim p(\cdot | s, a) \\ a' \sim \pi(\cdot | s')}} [Q_\pi(s', a')]$$

$$= r(s, a) + E_{s' \sim p(\cdot | s, a)} (V_\pi(s'))$$

- Summary

$$\blacksquare V_\pi(s) = E_{a \sim \pi} [r(s, a) + \alpha E_{s' \sim p(\cdot | s, a)} [V_\pi(s')]]$$

$$\blacksquare Q_\pi(s, a) = E_{s' \sim p} [r(s, a) + \alpha E_{a' \sim \pi} [Q_\pi(s', a')]]$$

$$\blacksquare V^*(s) = \max_a E_{s' \sim p} [r(s, a) + \alpha V^*(s')]$$

$$\blacksquare Q^*(s, a) = E_{s' \sim p} [r(s, a) + \alpha \max_{a'} [Q_\pi(s', a')]]$$

Sampling

Use random sampling to

- simulate a MC or MDP

- evaluate high-dim expectations in RL

Monte Carlo Method

Intrinsic

- X_1, X_2, \dots, X_n i.i.d random vectors
- $E(X_i) = \mu_X, Cov(X_i) = Q_X$

Estimation

- Sample mean: $\bar{X}_n = \frac{1}{n} \sum X_i \rightarrow \mu_X$
- Sample covariance: $\bar{Q}_n = \frac{1}{n-1} \sum_i (X_i - \bar{X}_n)(X_i - \bar{X}_n)^T \rightarrow Q_X$
- unbiased estimate: $E(\bar{X}_n) = \mu_X, E(\bar{Q}_n) = Q_X$

Central limit theorem (CLT)

$\sqrt{n}(\bar{X}_n - \mu_X) \rightarrow \mathcal{N}(0, Q_X)$ in distribution

- \bar{X}_n can be approximated by gaussian distribution $\mathcal{N}(\mu_X, \frac{Q_X}{n})$
- Covariance $E[(\bar{X}_n - \mu_X)(\bar{X}_n - \mu_X)^T] \approx \frac{Q_X}{n}$
- MSE of the estimate \bar{X}_n is $trace(\frac{Q_X}{n})$

Monte Carlo Integration

- $E(\phi(X)) = \frac{1}{n} \sum_i \phi(X_i)$
- $P(X \in A) = E(1_A(X)), 1_A(X) = \begin{cases} 1, & \text{if } X \in A \\ 0, & \text{otherwise} \end{cases}$

$X_i \sim f_X(x)$ are i.i.d samples

Importance sampling

Estimate $E_g(X) = \sum_x xg(x)$ with sample from $f(x)$ distribution

$$E_g(\phi(X)) = \sum_x \phi(x)g(x) = \sum_x (\phi(x)\frac{g(x)}{f(x)})f(x) = E_f(\phi(x)\frac{g(x)}{f(x)}) \approx \frac{1}{N} \sum_i \frac{g(X_i)}{f(X_i)} \phi(X_i)$$

$f(x)$ and $g(x)$ are both known

Benefits:

- possibly reduce final **sample variance**
- f is easier to evaluate and sample than g

Application

$$\int g(x)dx = \int g(x)\frac{1}{f(x)}f(x)dx = E_f(g(x)\frac{1}{f(x)}) \approx \frac{1}{N} \sum_i g(X_i)\frac{1}{f(X_i)}$$

