



Computer Engineering Department

Reinforcement Learning: Model Based RL

Mohammad Hossein Rohban, Ph.D.

Spring 2025

Courtesy: Most of slides are adopted from CS 285 Berkeley.

Lecture 11 - 1

Overview

- Introduction to model-based reinforcement learning
- What if we know the dynamics? How can we make decisions?
- Stochastic optimization methods
- Monte Carlo tree search (MCTS)
- Trajectory optimization
- Goal: Understand how we can perform planning with known dynamics models in discrete and continuous spaces

Recap: Model-Free RL



Recap: Model-Free RL



 $\pi_{\theta}(\tau)$

assume this is unknown don't even attempt to learn it

$$\theta^{\star} = \arg\max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$

What if we knew the transition dynamics?

- Often we do know the dynamics
 - Games (e.g., Atari games, chess, Go)
 - Easily modeled systems (e.g., navigating a car)
 - Simulated environments (e.g., simulated robots, video games)
- Often we can learn the dynamics
 - System identification fit unknown parameters of a known model
 - Learning fit a general-purpose model to observed transition data

Does knowing the dynamics make things easier?

Often, yes!

Model-based RL

- Model-based reinforcement learning: learn the transition dynamics, then figure out how to choose actions.
- Today: how can we make decisions if we know the dynamics?
 - a. How can we choose actions under perfect knowledge of the system dynamics?
 - b. Optimal control, trajectory optimization, planning

The deterministic case



$$\mathbf{a}_1, \dots, \mathbf{a}_T = \arg \max_{\mathbf{a}_1, \dots, \mathbf{a}_T} \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \text{ s.t. } \mathbf{a}_{t+1} = f(\mathbf{s}_t, \mathbf{a}_t)$$

Lecture 11 - 7

The stochastic open-loop case



The stochastic open-loop case

کری می خواست به عیادت بیماری برود.اندیشید که هنگام احوال پرسی ممکن است صدای اورانشنوم وپاسخی ناشایسته بدهم.ازاین رودرپی چاره برآمدوبالاخره باخود گفت:بهتراست پرسشهارا پیش ازرفتن بسنجم وپاسخ رانیزبرآورد کنم تادچاراشتباه نشوم. بنابراین پرسشهای خودراچنین پیش بینی کرد: -ابتداازاومی پرسم حالت بهتراست؟ اوخواهد گفت "آری" من درجواب می گویم:خدا را شکر -بعدازاومی پرسم چه خورده ای؟ لابد نام غذایی راخواهد آورد.من می گویم گوارا باد. -درپایان می پرسم پزشکت کیست؟ نام پزشکی رامی گویدومن پاسخ می دهم:مقدمش مبارک باد.

> چون به خانه ی بیماررسید همان گونه که ازپیش آماده شده بودبه احوال پرسی پرداخت: -کر گفت: "چگونه ای؟" کر گفت: خدارا شکر بیمارازاین سخن بیجا برآشفت. بیمارازاین سخن بیجا برآشفت. -بعدازآن پرسید: "چه خورده ای؟" بیمارگفت: زهر بیمار ازاین پاسخ نیزبیشتربه خود پیچید. بیمار ازاین پاسخ نیزبیشتربه خود پیچید. بیمار که آشفتگی وناراحتی اش به نهایت رسیده بود در پاسخ گفت: عزرائیل می آید, برو.

Lecture 11 - 9

open-loop vs. closed-loop case



The stochastic open-loop case



$$p(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T) = p(\mathbf{s}_1) \prod_{t=1}^T \pi(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

$$\pi = \arg\max_{\pi} E_{\tau \sim p(\tau)} \left[\sum_{t} r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

form of π ? neural net

time-varying linear
$$\mathbf{K}_t \mathbf{s}_t + \mathbf{k}_t$$

Stochastic optimization

abstract away optimal control/planning:

$$\mathbf{a}_1, \dots, \mathbf{a}_T = \arg \max_{\mathbf{a}_1, \dots, \mathbf{a}_T} J(\mathbf{a}_1, \dots, \mathbf{a}_T)$$

$$\mathbf{A} = \arg \max_{\mathbf{A}} J(\mathbf{A})$$

don't care what this is

simplest method: guess & check "random shooting method"

1. pick $\mathbf{A}_1, \ldots, \mathbf{A}_N$ from some distribution (e.g., uniform)

2. choose \mathbf{A}_i based on $\arg \max_i J(\mathbf{A}_i)$

Cross-entropy Method (CEM)

1. pick $\mathbf{A}_1, \ldots, \mathbf{A}_N$ from some distribution (e.g., uniform)

2. choose \mathbf{A}_i based on $\arg \max_i J(\mathbf{A}_i)$ can we do better?



cross-entropy method with continuous-valued inputs:

- 1. sample $\mathbf{A}_1, \ldots, \mathbf{A}_N$ from $p(\mathbf{A})$
- 2. evaluate $J(\mathbf{A}_1), \ldots, J(\mathbf{A}_N)$
- 3. pick the *elites* $\mathbf{A}_{i_1}, \ldots, \mathbf{A}_{i_M}$ with the highest value, where M < N
- 4. refit $p(\mathbf{A})$ to the elites $\mathbf{A}_{i_1}, \ldots, \mathbf{A}_{i_M}$

Pros and Cons

• Pros

- Could be very fast (Parallelizable)
- Extremely simple

• Cons

- Very harsh dimensionality limit
- Only open-loop planning