

CMPUT 655: RL-1 Lecture 3

14th September 2020

Contents

We will talk about

- ▶ Small recap about what we did last time about probability
- ▶ And discuss these a little bit in the context of MDPs

Reminder about Some Probability Rules

- ▶ Conditional probability (also chain rule)

$$\underline{\underline{P(A|B) = \frac{P(A, B)}{P(B)}}}$$

$$P(X = x | Y = y) = \frac{P(X = x, Y = y)}{P(Y = y)}$$

die # (pointing to X)
coin outcome (pointing to Y)

$$P(A, B, C) = P(A, B | C) P(C)$$

$$P(X = 1 | Y = \{H\})$$

- ▶ Marginalization over a joint probability distribution

$$P(X = x) = \sum_{y \in \mathcal{Y}} P(X = x, Y = y)$$

Reminder about Some Probability Rules (contd.)

- ▶ Conditional Expectation

$$\mathbb{E}[X] = \sum_x x p(x)$$

$$\mathbb{E}_X[X|Y=y] = \sum_{x \in \mathcal{X}} x \cdot p(X=x|Y=y).$$

- ▶ ~~Conditional Expectation~~ and Law of Total Expectation

$$\mathbb{E}_X[X|Y=y] = \mathbb{E}_Y \left[\mathbb{E}_X[X|Y=y] \right].$$



RL Book Notation

- ▶ Mathcal (fancy) Symbols:

\mathcal{S} → $\{s_0, s_1, s_2$
behind a car,
speeding,
dropping, ...}

- ▶ Capital Symbols:

S_0 → r.v. stores the state
at $(t=0)$
 A_0 → action at $t=0$

- ▶ Small Symbols:

$s_0 =$ behind a car
 $A_0 =$ steer left → a_1

Example MDP

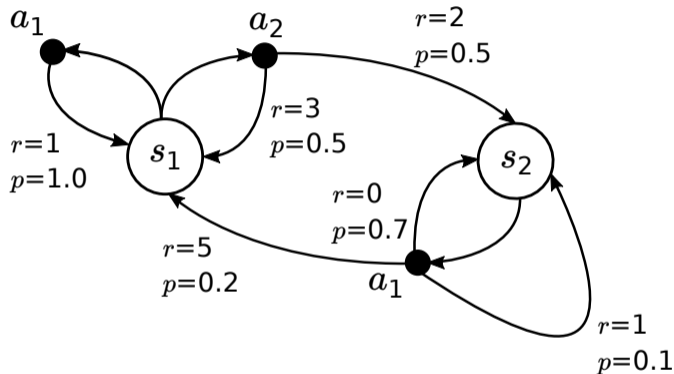


Figure: MDP $M = (\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$.

$\{s_1, s_2\}$

Agent's Interaction and the Trajectory (\sim stream of experience)

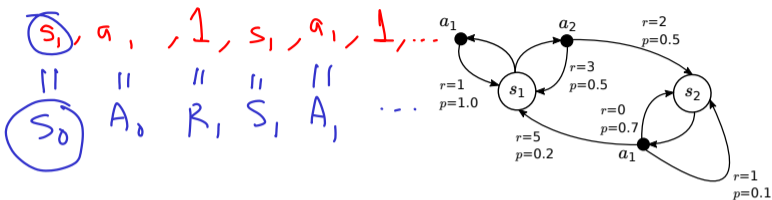


Figure: Trajectory $\tau_t = (S_0, A_0, R_1, S_1, A_1, R_2, S_2, \dots, S_T)$.

Let us calculate the probability of seeing this trajectory:

$$\begin{aligned}
 P(S_0, A_0, R_1, S_1, A_1, \dots) &= P(S_0) \times P(A_0, R_1, S_1, \dots \mid S_0) \\
 &= P(S_0) \underbrace{\pi(A_0 \mid S_0)}_{\text{action}} P(S_1, R_1 \mid A_0, S_0) \dots
 \end{aligned}$$

Different Transition Probabilities (Section 3.1 RL Book¹)

$p(s', r | s, a) \doteq \Pr(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a)$

joint (under $p(s', r | s, a)$)
conditioned (under $| S_t = s, A_t = a$)

$p(s' | s, a) \doteq \Pr(S_{t+1} = s' | S_t = s, A_t = a)$
 $= \sum_{r \in \mathcal{R}} P(s', r | s, a)$

(Dependent on the policy π)

$p(s' | s) = \sum_a \pi(a | s) p(s', a | s)$

$p(s', a | s) = P(s', r | s, a) \xrightarrow{r} p(s' | s, a)$

$P(s', a | s) = P(s', r | s, a) \times \pi(a | s)$

$p(s' | s) = \sum_a p(s', a | s) = \sum_a P(s', r | s, a) \pi(a | s)$

$p(s', a' | s, a) = \pi(a' | s')$

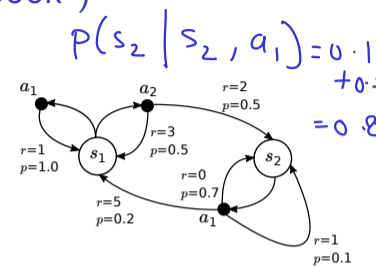


Figure: MDP

Different Reward Functions (Section 3.1 RL Book²)

- ▶ The random variable R_t

- ▶ Reward function

$$p(r|s,a)$$

$$\sum_{s'} p(s', r | s, a)$$

$$r(s, a) \doteq \mathbb{E} [R_t | S_{t-1} = s, A_{t-1} = a] \quad \leftarrow \text{cond. Expectation}$$

$$= \sum_r r \cdot p(R_t = r | S_{t-1} = s, A_{t-1} = a) \quad \text{Figure: MDP}$$

$$= \sum_r r \sum_{s'} p(s', r | s, a)$$

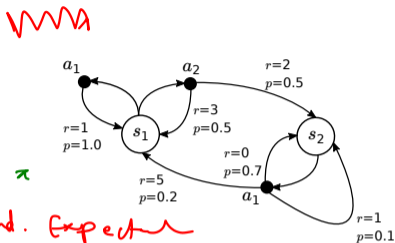
- ▶ $r(s, a, s') \doteq \mathbb{E} [R_t | s_t = s', S_{t-1} = s, A_{t-1} = a]$

$$= \sum_r r \cdot p(r | s', a, s)$$

$$p(r | s', a, s) = \frac{p(s', r | s, a)}{p(s' | s, a)}$$

- ▶ $r_\pi(s) \doteq \mathbb{E}_\pi [R_t | S_t = s]$

cond. buys, margin



²<http://incompleteideas.net/book/the-book-2nd.html>

Bellman Equation ($V^{\pi} \rightarrow V^{\pi}$)

Read from back



Bellman Equation ($V \rightarrow V$) (using $r(s, a)$)



What does \mathbb{E}_π mean?

$$\pi(a_t | s) = P(A_t = a | S_t = s)$$

$$V_\pi(s) = \mathbb{E}_\pi [R_1 + \gamma R_2 + \gamma^2 R_3 + \dots | S_0 = s]$$

$$= \sum_{a_1 \in A(s)} \pi(a_1 | s) \sum_{s_1, r_1} p(s_1, r_1 | s, a_1) \left[r_1 + \gamma \sum_{a_2} \pi(a_2 | s_1) \sum_{s_2, r_2} p(s_2, r_2 | s_1, a_2) [r_2 + \gamma \dots] \right]$$

a_0 is an instantiation at $t=0$ (A_0)

x y z
 a_0 a_1 a_2

$a_i \in A(s_i)$ all the actions at $t=i$

$$\mathbb{E}_\pi \cong \mathbb{E}_{A_0 \cup \pi(\cdot | s); S, R, \cup P(\cdot, \cdot | s, A_0); A, \cup \pi(\cdot | s) \dots}$$



$$\mathbb{E}_x[X] = \sum_{x \in \mathcal{X}} x \cdot p(X=x)$$

Bellman Equation ($V \rightarrow V$) (more thoughts)

Summary