Jack's car rental

Corrado Possieri

Machine and Reinforcement Learning in Control Applications

1/12

Problem



Manage cars between two locations.

Problem formulation

- Jack manages two locations for a car rental company.
- Each day, some customers arrive at each location to rent cars.
- If Jack has a car available, he rents it out and is credited \$10.
- Cars are available for renting the day after they are returned.
- Jack can move cars between the two locations overnight.
- The cost of moving a car is \$2.
- Each location is capable of accommodating 30 cars.

Problem data

 Cars requested and returned at each location are Poisson random variables

$$\mathbb{P}[\mathsf{cars} = n] = \frac{\lambda^n}{n!} \exp(-\lambda).$$

- ullet λ is 3 and 4 for rental requests.
- \bullet λ is 3 and 2 for returns.
- Jack's foresight modeled with discount $\gamma = 0.9$.
- Jack can move up to 7 cars between the two locations.

Model

- We can model the process as an MDP.
- The state is the number of car at each location
 - is it a Markov state?
 - we have $\#1 \cdot \#2$ states.
- The action is the number of cars moved
 - we have 2#c+1 actions.

5/12

State update

- Jack reintroduces car returned at previous day.
- Jack rents available cars.
- 3 Jack moves cars between the two locations.

6/12

Transition and reward

- Let S and A be the number of states and actions.
- ullet Transition probabilities can be stored in a $S \times S \times A$ matrix P

$$P_{s,s',a} = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a].$$

• Rewards can be stored in a $S \times A$ matrix R

$$R_{s,a} = \mathbb{E}[R_{t+1}|S_t = s, A_t = a].$$

Transition probabilities

- Probability of returns do not depend on actions
 - lacktriangle define the return probability matrix P_{ret}

$$[P_{\mathsf{ret}}]_{s,s'} = \mathbb{P}[S_{\mathsf{after return}} = s' | S_t = s].$$

- Probability of rentals do not depend on actions
 - define the rental probability matrix P_{ren}

$$[P_{\text{ren}}]_{s,s'} = \mathbb{P}[S_{\text{after rental}} = s' | S_{\text{after return}} = s].$$

- Probability of movement depend on actions
 - lacktriangle define the movement probability matrix P_{mov}

$$[P_{\mathsf{mov}}]_{s,s',a} = \mathbb{P}[S_{\mathsf{after movement}} = s' | S_{\mathsf{after rental}} = s, A_t = a].$$

By the law of total probability

$$P = P_{\text{ret}} \cdot P_{\text{ren}} \cdot P_{\text{mov}}$$
.

Expected rewards

Expected earning do not depend on action

$$\mathbb{E}\left[\mathsf{earning}_{t+1}|S_{\mathsf{after\ return}} = s\right] = \sum_{r} r \mathbb{P}\left[r|S_{\mathsf{after\ return}} = s\right].$$

By the law of total probability

$$\begin{split} & \mathbb{E}\left[\mathsf{earning}_{t+1}|S_t = s\right] \\ & = \sum_{s'} \mathbb{P}[S_{\mathsf{after \, return}} = s'|S_t = s] \sum_{r} r \mathbb{P}\left[r|S_{\mathsf{after \, return}} = s'\right]. \end{split}$$

• The expected reward is given by

$$R_{s,a} = [P_{\mathsf{ret}} \cdot \mathsf{earning}]_s - \mathsf{cost}_a.$$

PI and VI revisited

 \bullet Given a deterministic policy $\pi,$ define

Bellman expectation update can be rewritten as

$$v \leftarrow R^{\pi} + \gamma P^{\pi} v \qquad (v^{\pi} = (I - \gamma P^{\pi})R^{\pi}).$$

Bellman optimality update can be rewritten as

$$v \leftarrow \max_{\pi} \left\{ R^{\pi} + \gamma P^{\pi} v \right\}$$

Matrix Formulation

Recall classical Bellman update

$$\begin{split} v(s) &\leftarrow \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,a) \left(r + \gamma v(s')\right) \\ &= \sum_{s',r} p(s',r|s,\pi(s)) \left(r + \gamma v(s')\right) \\ &= \sum_{s',r} p(s',r|s,\pi(s))r + \gamma \sum_{s',r} p(s',r|s,\pi(s))v(s') \\ &= \sum_{s',r} r \sum_{s'} (s',r|s,\pi(s)) + \gamma \sum_{s'} v(s') \sum_{r} p(s',r|s,\pi(s)) \\ &= r(s,\pi(s)) + \gamma \sum_{s'} p(s'|s,\pi(s))v(s'). \end{split}$$

A similar relation holds for Bellman optimality update

$$v(s) \leftarrow \max_{a} \left\{ r(s, a) + \gamma \sum_{s'} p(s'|s, a) v(s') \right\}.$$

Assignment #2

- Write a code for PI (in class).
- Write a code for VI (in class).
- Model gambler's problem
 - a gambler has the opportunity to make bets on the outcomes of a sequence of coin flips;
 - if the coin comes up heads, he wins as many dollars as he has staked on that flip; if it is tails, he loses his stake;
 - the game ends when the gambler wins by reaching his goal of \$100, or loses by running out of money;
 - \blacksquare reward +1 for winning.