

CS 188 Discussion 5:

MDPs

Kenny Wang (kwkw@berkeley.edu)
Wed Sep 27, 2023

Slides based on Sashrika Pandey's

Administrivia

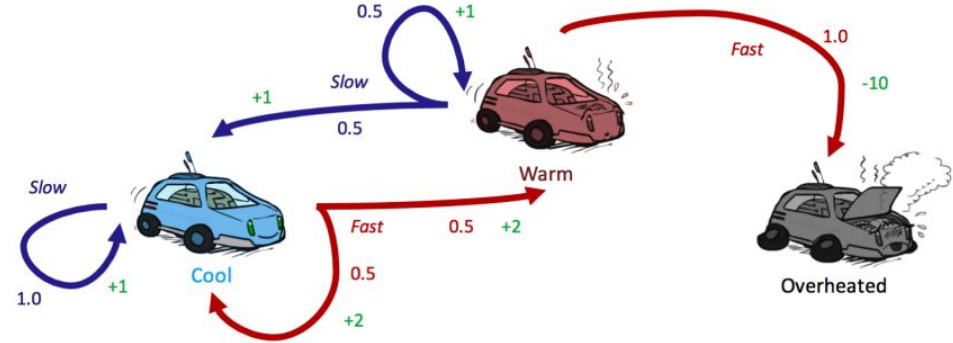
- Project 3 due on Friday, Oct 6
- Homework is due on Tuesdays
- We have office hours pretty much all day every weekday (12-7), come to Soda 341B!
- Reminder: Need extensions? We will give you extensions!
- Discussion slides are on Ed

Today's Topics

- MDPs (Markov Decision Processes)
 - Value Iteration
 - Policy Iteration

What is an MDP?

- Set of states **S**
- Set of actions **A**
- Start state
- Terminal state(s)
- Discount factor γ [gamma]
 - Rewards decay as time passes, so we prefer sooner rewards
- Transition function **$T(s, a, s')$**
 - Probability of ending up in state s' by starting in s and taking action a
- Reward function **$R(s, a, s')$**



Example:

- $S = \{\text{cool, warm, overheated}\}$
- $A = \{\text{slow, fast}\}$
- $T(\text{warm, slow, cool}) = 0.5$
- $R(\text{warm, slow, cool}) = 1$

Incorporating Discount γ [gamma]

- **Additive utility**

$$U([s_0, a_0, s_1, a_1, s_2, \dots]) = R(s_0, a_0, s_1) + R(s_1, a_1, s_2) + R(s_2, a_2, s_3) + \dots$$

- **Discounted utility** incorporates **discount factor γ**

- Reward of $\gamma^t R(s_t, a_t, s_{t+1})$ instead of $R(s_t, a_t, s_{t+1})$

$$U([s_0, a_0, s_1, a_1, s_2, \dots]) = R(s_0, a_0, s_1) + \gamma R(s_1, a_1, s_2) + \gamma^2 R(s_2, a_2, s_3) + \dots$$

Bellman Equations

- **$Q^*(s,a)$: the optimal value of (s, a) [state, action pair]**

- The expected value of the utility an agent receives after starting in s , taking a , and acting optimally

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U^*(s')]$$

- **$U^*(s)$: the optimal value of state s**

- The expected value of the utility that an agent starting in s and acting optimally will receive

$$U^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U^*(s')]$$

$$U^*(s) = \max_a Q^*(s, a)$$

Worksheet 1(a)

Value Iteration

Value Iteration

- A dynamic programming algorithm we use to compute values until convergence ($\forall s, U_{k+1}(s) = U_k(s)$)
- Algorithm
 - $\forall s \in S$, initialize $U_0(s) = 0$ [initial value estimate is 0]
 - Repeat rule until convergence (U-values stop changing)

$$\forall s \in S, U_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U_k(s')]$$

- Convergence is when
 - $\forall s \in S, U_k(s) = U_{k+1}(s) = U^*(s)$

Policy Iteration

Policy Iteration

- Issues with value iteration
 - $O(|S|^2|A|)$ runtime
 - Overcomputes, policy tends to converge faster than values
- Policy iteration: preserve the optimality from value iteration but with better performance by iterating until only the *policy* converges instead of the U-values

Policy Iteration

- Algorithm

- Define initial policy π_0 (can be arbitrary)
- Repeat until convergence:
 - **Policy evaluation:** compute expected utility of starting in state s when following policy π , for all states s

$$U^\pi(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma U^\pi(s')]$$

- **Policy improvement:** generate a better policy

$$\pi_{i+1}(s) = \operatorname{argmax}_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U^{\pi_i}(s')]$$

- Convergence when $\pi_{i+1} = \pi_i$ (policy stops changing)

Summary

- **Value Iteration**

- $\forall s \in S$, initialize $U_0(s) = 0$
[initial value estimate is 0]
- Repeat rule until convergence
(U-values stop changing)

$$\forall s \in S, U_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U_k(s')]$$

- Convergence is when
 $\forall s \in S, U_k(s) = U_{k+1}(s) = U^*(s)$

- **$Q^*(s,a)$: the optimal value of (s, a)** [state, action pair]

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U^*(s')]$$

- **$U^*(s)$: the optimal value of state s** $U^*(s) = \max_a Q^*(s, a)$

- **Policy Iteration**

- Define initial policy π_0 (can be arbitrary)
- Repeat until convergence:
 - **Policy evaluation:** compute expected utility of starting in state s when following policy π , for all states s

$$U^\pi(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma U^\pi(s')]$$

- **Policy improvement:** generate a better policy

$$\pi_{i+1}(s) = \operatorname{argmax}_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U^{\pi_i}(s')]$$

- Convergence when $\pi_{i+1} = \pi_i$ (policy stops changing)

Rest of the Worksheet

Thank you for attending!

Attendance link:

- <https://tinyurl.com/cs188fa23>

Discussion No: 5

Remember my name is Kenny

My email: kwkw@berkeley.edu

