# ELEN E6885: Reinforcement Learning
# Homework 1

Josh Rutta

September 28, 2018

Josh Rutta
Prof. Li
Reinforcement Learning

HW #1   3.1, 3.2, 3.4, 4.2, 4.4, 4.6, 4.8

3.1 a) $Pr(A_t = j) = \dfrac{e^{\mu_j/\tau}}{\sum\limits_{i=1}^{k} e^{\mu_i/\tau}} = \dfrac{1}{\sum\limits_{i=1}^{k} e^{(\mu_i - \mu_j)/\tau}}$ ; for largest $\mu_j$, $\mu_i - \mu_j < 0$ except $i=j$.

$\lim\limits_{\tau \to 0} e^{(\mu_i - \mu_j)/\tau} = 0$ except $i=j$ $\Rightarrow \lim\limits_{\tau \to 0} \dfrac{1}{\sum\limits_{i=1}^{k} e^{(\mu_i - \mu_j)/\tau}} = 1$ for largest $\mu_j$

for all other $\mu_j$, there will be some $\mu_i - \mu_j > 0 \Rightarrow \lim\limits_{\tau \to 0} \dfrac{1}{\sum\limits_{i=1}^{k} e^{(\mu_i - \mu_j)/\tau}} = 0$. Thus softmax becomes same as greedy action selection, always choosing largest $\mu_j$.

b) $\dfrac{e^{\mu_1/\tau}}{e^{\mu_1/\tau} + e^{\mu_2/\tau}} \cdot \dfrac{e^{-\mu_1/\tau}}{e^{-\mu_1/\tau}} = \dfrac{1}{1 + e^{(\mu_2 - \mu_1)/\tau}}$

3.2 The greedy method could perform significantly worse than the $\epsilon$-greedy method in the long run, because in the greedy case, the agent picks the action with the current largest estimated expected return $\mu_j$, however it may not have explored other actions enough for its estimated expected returns to be accurate, thus it could get stuck doing an action which does not actually have the largest expected return. $\epsilon$-greedy mitigates this risk by having the agent do a random action with probability $\epsilon$, so that the agent won't get stuck with a suboptimal action.

3.4

K=0

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

k=1

| 0 | -1 | -1 | -1 |
|---|---|---|---|
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 0 |

$V(1) = 4\left(\frac{1}{4}(-1 + 1(0))\right)$

$V(2) = 4\left(\frac{1}{4}(-1 + 1(0))\right)$

k=2

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |

$V(1) = \frac{1}{4}(-1+0) + \frac{1}{4}(-1 + 1(-1)) + \frac{1}{4}(-1 + 1(-1)) + \frac{1}{4}(-1 + 1(-1))$

$= -\frac{1}{4} - 3\left(\frac{3}{4}\right) = -\frac{7}{4} = -1.75$

$V(2) = \frac{1}{4}(-1 + 1(-1)) + \frac{1}{4}(-1 + 1(-1)) + \frac{1}{4}(-1 + 1(-1)) + \frac{1}{4}(-1 + 1(-1))$

$= 4\left(\frac{1}{4}(-2)\right) = -2$

4.2, 4.4, 4.6, 4.8

**4.2** The value of the discount factor in infinite horizon problems is a measure of how much the agent is interested in long term vs. short term rewards. For small discount factor, agent will care more about immediate rewards. For a large discount factor, the agent will take long term rewards more into account

**4.8** We should model this as a POMDP because with car automation, much data is often noisy.

**4.4** Example: immediate reward of all states is -1 except terminal state

$\gamma = 1$ (val. iteration:

$$V(s) = \max_a \sum_{s'} p(s'|s,a)\left[r(s,a,s') + \gamma V(s')\right])$$

| | | | |
|---|---|---|---|
| ▨ | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

$V_1$

| | | | |
|---|---|---|---|
| ▨ | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |

$V_2$

| | | | |
|---|---|---|---|
| ▨ | -1 | -2 | -2 |
| -1 | -2 | -2 | -2 |
| -2 | -2 | -2 | -2 |
| -2 | -2 | -2 | -2 |

$V_3$

| | | | |
|---|---|---|---|
| ▨ | -1 | -2 | -3 |
| -1 | -2 | -3 | -3 |
| -2 | -3 | -3 | -3 |
| -3 | -3 | -3 | -3 |

$V_4$

| | | | |
|---|---|---|---|
| ▨ | -1 | -2 | -3 |
| -1 | -2 | -3 | -4 |
| -2 | -3 | -4 | -4 |
| -3 | -4 | -4 | -4 |

$V_5$

| | | | |
|---|---|---|---|
| ▨ | -1 | -2 | -3 |
| -1 | -2 | -3 | -4 |
| -2 | -3 | -4 | -5 |
| -3 | -4 | -5 | -5 |

$V_6$

| | | | |
|---|---|---|---|
| ▨ | -1 | -2 | -3 |
| -1 | -2 | -3 | -4 |
| -2 | -3 | -4 | -5 |
| -3 | -4 | -5 | -6 |

$V_7$

| | | | |
|---|---|---|---|
| ▨ | -1 | -2 | -3 |
| -1 | -2 | -3 | -4 |
| -2 | -3 | -4 | -5 |
| -3 | -4 | -5 | -6 |

$V_8$

$V_7 = V_8$

In such a case, policy can't ever change again because value doesn't change again and $\pi(s) = \arg\max_a \sum_{s'} p(s'|s,a)\left[r(s,a,s') + \gamma V(s')\right]$ (greedy) so $\pi(s)$ for $V_7$ will equal $\pi(s)$ for $V_8, V_9, ..., V_\infty (= V_7)$

**4.6** Synchronous value iteration updates every single state in one iteration, asynchronous value iteration does not, but rather updates the state values individually in some order. The ordering of how we update state value is important in asynchronous value iteration because if the state space is very large, you'll want to prioritize updating states which are important by some measure such as those with largest Bellman error $\left|\max_{a \in A}\left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V(s')\right) - V(s)\right|$ or states which are relevant to agent