

CS7643: Deep Learning
Fall 2020
HW5 Solutions

Ruize Yang

November 5, 2020

Problem 1a

1.(a). When start at S_1 and always choose stay, it will always in S_1 . The reward of each step is -2 .

$$\sum_{t=0}^{\infty} \gamma^t r_t(s_t, a_t) = -2 \sum_{t=0}^{\infty} \gamma^t = -\frac{2}{1-\gamma}$$

Problem 1b

1.(b) If $r \geq 0.1367$, $(a_1, a_2) = (go, go)$. If $r < 0.1367$, $(a_1, a_2) = (stay, go)$ or $(stay, stay)$.

There are four possible policy:

- ① $(stay, stay) \Rightarrow (S_1, R)$ ② $(go, go) \Rightarrow (S_1) \rightarrow (S_2)$
 ③ $(stay, go) \Rightarrow (S_1, R)$ ④ $(go, stay) \Rightarrow (S_1) \rightarrow (S_2, R)$

The reward of ①, ③ is $-\frac{2}{1-r}$, and is higher than that of ④, because $(S_1) \rightarrow (S_2)$ is -3 , and stay in S_1 is -2 .

The reward of ② is $-3 + 5r$, because at the second step, the agent takes "go" from S_2 and ends the episode.

$-3 + 5r \geq -\frac{2}{1-r} \Rightarrow r \geq 0.1367$, $\Rightarrow 0.1367 \leq r < 1$, $(a_1, a_2) = (go, go)$

~~$r < 0.1367$~~ $0 \leq r < 0.1367$, $(a_1, a_2) = (stay, stay)$
 or $(a_1, a_2) = (stay, go)$

Problem 1c

1 (c). $r=1$.

$$V_1(s_1) = \max \{ r(s_1, \text{stay}) + V_0(s_1), r(s_1, \text{go}) + V_0(s_2) \}$$
$$= \max \{ -2 + 0, -3 + 0 \} = -2$$

$$V_1(s_2) = \max \{ r(s_2, \text{stay}) + V_0(s_2), r(s_2, \text{go}) \}$$
$$= \max \{ -2 + 0, 5 \} = 5$$

$$V_1 = [-2, 5].$$

$$V_2(s_1) = \max \{ r(s_1, \text{stay}) + V_1(s_1), r(s_1, \text{go}) + V_1(s_2) \}$$
$$= \max \{ -2 + (-2), -3 + 5 \} = 2.$$

$$V_2(s_2) = \max \{ r(s_2, \text{stay}) + V_1(s_2), r(s_2, \text{go}) \}$$
$$= \max \{ -2 + 5, 5 \} = 5$$

$$V_2 = [2, 5].$$

$$V_3(s_1) = \max \{ r(s_1, \text{stay}) + V_2(s_1), r(s_1, \text{go}) + V_2(s_2) \}$$
$$= \max \{ -2 + 2, -3 + 5 \} = 2.$$

$$V_3(s_2) = \max \{ r(s_2, \text{stay}) + V_2(s_2), r(s_2, \text{go}) \}$$
$$= \max \{ -2 + 5, 5 \} = 5.$$

$$V_3 = [2, 5].$$

V_2, V_3 are both optimal, but V^3 is better, because it shows V converges.

Problem 2a

$$\begin{aligned} & \geq (a). \quad V^1: \|V^1 - V^*\| = \max\{|-2-2|, |5-5|\} = 4. \\ & \quad V^2: \|V^2 - V^*\| = \max\{|2-2|, |5-5|\} = 0 \\ & \quad V^3: \|V^3 - V^*\| = \max\{|2-2|, |5-5|\} = 0. \end{aligned}$$

The error decreases monotonically

Problem 2b

2(b). For state s :

$$T(V)(s) = \max_a \sum_{s'} p(s'|s, a) [r(s, a) + \gamma V(s')]$$

$$T(V')(s) = \max_a \sum_{s'} p(s'|s, a) [r(s, a) + \gamma V'(s')]$$

$$\begin{aligned} \|T(V) - T(V')\| &\leq \max_a \left| \sum_{s'} p(s'|s, a) [r(s, a) + \gamma V(s')] - \sum_{s'} p(s'|s, a) [r(s, a) + \gamma V'(s')] \right| \\ &= \max_a \left| \sum_{s'} p(s'|s, a) [\gamma V(s') - \gamma V'(s')] \right| \\ &= \gamma \max_a \sum_{s'} p(s'|s, a) |V(s') - V'(s')|. \end{aligned}$$

$$\|T(V) - T(V')\|_\infty = \gamma \max_{a, s} \sum_{s'} p(s'|s, a) |V(s') - V'(s')|$$

$$\leq \gamma \max_{a, s} \sum_{s'} p(s'|s, a) \|V - V'\|_\infty$$

$$= \gamma \|V - V'\|_\infty \max_{a, s} \sum_{s'} p(s'|s, a)$$

$$= \gamma \|V - V'\|_\infty$$

$$\Rightarrow \|T(V) - T(V')\|_\infty \leq \gamma \|V - V'\|_\infty, \quad \forall V, V'.$$

Problem 2c

$$\begin{aligned}
 \geq (c). \quad \|V^{n+1} - V^*\|_\infty &= \|V^{n+1} - T(V^{n+1}) + T(V^{n+1}) - V^*\|_\infty \\
 &\leq \|V^{n+1} - T(V^{n+1})\|_\infty + \|T(V^{n+1}) - V^*\|_\infty \\
 &= \|T(V^n) - T(V^{n+1})\|_\infty + \|T(V^{n+1}) - T(V^*)\|_\infty \\
 &\leq r \|V^n - V^{n+1}\|_\infty + r \|V^{n+1} - V^*\|_\infty.
 \end{aligned}$$

$$\begin{aligned}
 (1-r) \|V^{n+1} - V^*\|_\infty &\leq r \|V^n - V^{n+1}\|_\infty \\
 \|V^{n+1} - V^*\|_\infty &\leq \frac{r}{1-r} \|V^n - V^{n+1}\|_\infty.
 \end{aligned}$$

As $\{V^n\}$ is a Cauchy sequence, then $\forall \varepsilon, \exists n$ s.t. $\|V^n - V^{n+1}\|_\infty < \varepsilon$.

$$\Rightarrow \|V^{n+1} - V^*\|_\infty \leq \frac{r}{1-r} \varepsilon.$$

Problem 2d

> (d) Define the sequence $\{x_n\}$ as $x_n = T(x_{n-1})$. ~~and for given~~
 $\|T(x_{n+1}) - T(x_n)\|_\infty \leq \alpha \|x_{n+1} - x_n\|_\infty \stackrel{=}{=} \alpha \|T(x_n) - T(x_{n-1})\|_\infty \leq \alpha^n \|x_1 - x_0\|_\infty$
 $0 \leq \alpha < 1$, $\|x_1 - x_0\|_\infty$ is a constant.

$$\begin{aligned} \text{For } m > n, \quad \|T(x_m) - T(x_n)\|_\infty &\leq \|T(x_m) - T(x_{m-1})\|_\infty + \dots + \|T(x_{n+1}) - T(x_n)\|_\infty \\ &\leq \alpha^{m-1} \|x_1 - x_0\|_\infty + \dots + \alpha^n \|x_1 - x_0\|_\infty \\ &= \alpha^n \|x_1 - x_0\|_\infty \sum_{k=0}^{m-n-1} \alpha^k \\ &\leq \alpha^n \|x_1 - x_0\|_\infty \sum_{k=0}^{\infty} \alpha^k \\ &= \alpha^n \|x_1 - x_0\|_\infty \frac{1}{1-\alpha} \end{aligned}$$

For any ϵ , $\exists n$ s.t. $\alpha^n \|x_1 - x_0\|_\infty \frac{1}{1-\alpha} < \epsilon$, $\{x_n\}$ is Cauchy sequence.
 So $\{x_n\}$ will converge, let it converges to x^* .

$$x^* = \lim_{n \rightarrow \infty} x_n = \lim_{n \rightarrow \infty} T(x_{n-1}) = T(\lim_{n \rightarrow \infty} x_{n-1}) = T(x^*)$$

x^* is a fixed point of $\{x_n\}$.

Assume there is another fixed point x' , $T(x') = x'$.

$$\|x^* - x'\|_\infty = \|T(x^*) - T(x')\|_\infty \leq \alpha \|x^* - x'\|_\infty.$$

$$\alpha < 1 \Rightarrow \|x^* - x'\|_\infty = 0 \Rightarrow x^* = x'.$$

So the fixed point x^* exists and is unique.

Problem 3a

Problem 3b

Problem 3c

Problem 3d

Problem 4a

$$\begin{aligned}
 4(a). \quad \nabla_{\theta} J(\theta) &= \nabla_{\theta} \mathbb{E}_{z \sim \pi_{\theta}} [R(z)] \\
 &\Rightarrow \frac{1}{N} \sum_{i=1}^N (R(z_i) - b) \nabla_{\theta} \log \pi_{\theta}(z_i) \\
 &= \frac{1}{N} \sum_{i=1}^N R(z_i) \nabla_{\theta} \log \pi_{\theta}(z_i) - \frac{1}{N} \sum_{i=1}^N b \nabla_{\theta} \log \pi_{\theta}(z_i)
 \end{aligned}$$

$$\begin{aligned}
 \frac{1}{N} \sum_{i=1}^N b \nabla_{\theta} \log \pi_{\theta}(z_i) &= \mathbb{E}_{z \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(z_i) \cdot b] \\
 &= \mathbb{E}_{s_{0:t}, a_{0:t-1}} [\mathbb{E}_{s_{t+1:T}, a_{t:T-1}} [\nabla_{\theta} \log \pi_{\theta}(z_i) \cdot b]] \\
 &= \mathbb{E}_{s_{0:t}, a_{0:t-1}} [b \cdot \mathbb{E}_{s_{t+1:T}, a_{t:T-1}} [\nabla_{\theta} \log \pi_{\theta}(z_i)]] \\
 &= \mathbb{E}_{s_{0:t}, a_{0:t-1}} [b \cdot \mathbb{E}_{a_t} [\nabla_{\theta} \log \pi_{\theta}(z_i)]]
 \end{aligned}$$

$$= \mathbb{E}_{s_{0:t}, a_{0:t-1}} [b \cdot \int \frac{\nabla_{\theta} \pi_{\theta}(a_t | s_t)}{\pi_{\theta}(a_t | s_t)} \pi_{\theta}(a_t | s_t) da_t]$$

$$= \mathbb{E}_{s_{0:t}, a_{0:t-1}} [b \cdot \nabla_{\theta} \int \pi_{\theta}(a_t | s_t) da_t]$$

$$= \mathbb{E}_{s_{0:t}, a_{0:t-1}} [b \cdot \nabla_{\theta} 1]$$

$$= \mathbb{E}_{s_{0:t}, a_{0:t-1}} [b \cdot 0]$$

$$= 0$$

So $\frac{1}{N} \sum_{i=1}^N (R(z_i) - b) \nabla_{\theta} \log \pi_{\theta}(z_i) = \frac{1}{N} \sum_{i=1}^N (R(z_i) \nabla_{\theta} \log \pi_{\theta}(z_i))$,
 doesn't change the gradient estimate.

Problem 4b

$$4(b). \quad x = R(z) \nabla_{\theta} \log \pi_{\theta}(z), \quad x' = (R(z) - b) \nabla_{\theta} \log \pi_{\theta}(z)$$

$$\text{Var}(x') = E([(R(z) - b) \nabla_{\theta} \log \pi_{\theta}(z)]^2) - [E((R(z) - b) \nabla_{\theta} \log \pi_{\theta}(z))]^2$$

~~The second term~~

$$\text{Var}(x) = E([R(z) \nabla_{\theta} \log \pi_{\theta}(z)]^2) - [E(R(z) \nabla_{\theta} \log \pi_{\theta}(z))]^2$$

The second ~~the~~ term in $\text{Var}(x)$ and $\text{Var}(x')$ are the same, only the first term is different.

$$\begin{aligned} E([(R(z) - b) \nabla_{\theta} \log \pi_{\theta}(z)]^2) &\leq E([R(z) \nabla_{\theta} \log \pi_{\theta}(z)]^2) \\ &\leq E([R(z)]^2) E([\nabla_{\theta} \log \pi_{\theta}(z)]^2) \\ &\leq E([R(z)]^2) E([\nabla_{\theta} \log \pi_{\theta}(z)]^2) \end{aligned}$$

$$\Rightarrow \text{Var}(x') \leq \text{Var}(x)$$

$$\frac{\partial \text{Var}(x')}{\partial b} = \frac{\partial}{\partial b} [E([(R(z) - b) \nabla_{\theta} \log \pi_{\theta}(z)]^2) - [E((R(z) - b) \nabla_{\theta} \log \pi_{\theta}(z))]^2]$$

$$= \frac{\partial}{\partial b} E([(R(z) - b) \nabla_{\theta} \log \pi_{\theta}(z)]^2)$$

$$= \frac{\partial}{\partial b} E([R(z) - b]^2 [\nabla_{\theta} \log \pi_{\theta}(z)]^2)$$

$$= \frac{\partial}{\partial b} E([R^2(z) - 2bR(z) + b^2] [\nabla_{\theta} \log \pi_{\theta}(z)]^2)$$

$$= \frac{\partial}{\partial b} [E(R^2(z) [\nabla_{\theta} \log \pi_{\theta}(z)]^2) - 2b E(R(z) [\nabla_{\theta} \log \pi_{\theta}(z)]^2) + b^2 E([\nabla_{\theta} \log \pi_{\theta}(z)]^2)]$$

$$= \frac{\partial}{\partial b} [-2b E(R(z) [\nabla_{\theta} \log \pi_{\theta}(z)]^2) + b^2 E([\nabla_{\theta} \log \pi_{\theta}(z)]^2)]$$

$$= -2 E(R(z) [\nabla_{\theta} \log \pi_{\theta}(z)]^2) + 2b E([\nabla_{\theta} \log \pi_{\theta}(z)]^2)$$

Set $\frac{\partial \text{Var}(x')}{\partial b}$ to 0.

$$\Rightarrow -2 E(R(z) [\nabla_{\theta} \log \pi_{\theta}(z)]^2) + 2b E([\nabla_{\theta} \log \pi_{\theta}(z)]^2) = 0$$

$$\Rightarrow b = \frac{E(R(z) [\nabla_{\theta} \log \pi_{\theta}(z)]^2)}{E([\nabla_{\theta} \log \pi_{\theta}(z)]^2)}$$

$$E([\nabla_{\theta} \log \pi_{\theta}(z)]^2)$$

Problem 5

The authors show an algorithm to formulate RL as supervised learning problem, called Upside-Down RL. They use the desired reward and desired horizon as input command, the model then learns the actions to achieve that reward in that amount of steps. The training is split into Training and Gathering which update each other. They test the algorithm on several tasks, which shows Upside-Down RL can solve RL problems and has better performance when the reward is sparse.

To make such a input command, need to check each possible reward value and the actions to get the reward. This is fine with small number of possible rewards, e.g. the reward is given at the end of episode and have a few known value. This is the case of LunarLanderSparse, where UDRL outperforms other algorithm. But a task with continuously valued rewards observed at arbitrary time steps would be difficult for UDRL.