



EL2805 Reinforcement Learning

Homework 1

November 6, 2023

Division of Decision and Control Systems
School of Electrical Engineering and Computer Science
KTH Royal Institute of Technology

Instructions (read carefully):

- Solve Problems 1 and 2.
- Work in groups of 2 persons.
- **Both** students in the group should upload their scanned report as a .pdf-file to Canvas before November 17, 23:59. The deadline is strict. Please mark your answers directly on this document, and **append** hand-written or typed notes justifying your answers. Reports without justification will not be graded.

Good luck!

1 Repair or replace?

You own a bike that can be in several conditions: perfect, worn and broken. You ride your bike every month, and at the beginning of the month you observe its condition and decide what to do. If the bike is broken, you have to either repair it or buy a new bike. If it is worn, you can decide to either keep it as it is, or repair it. When you repair a worn bike, its condition becomes perfect, and when you repair a broken bike, its condition becomes worn. The cost of a new bike is C_b and the cost of repairing the bike is C_r . In one month, the probability that the bike condition degrades is θ (that is, going from perfect to worn, or from worn to broken). You wish to find a strategy that minimizes your expected cost over T months. The condition of the bike at the end of the T -th month does not matter.

✗ Model the problem as an MDP, then answer the following question: What is the correct transition matrix? *Note:* The states are indexed as perfect (1), worn (2) and broken (3).

$$P(\text{keep}) = \begin{matrix} (1) \\ (2) \end{matrix} \begin{bmatrix} 1-\theta & \theta & 0 \\ 0 & 1-\theta & \theta \end{bmatrix} \quad P(\text{repair}) = \begin{matrix} (2) \\ (3) \end{matrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

$$P(\text{buy new}) = \begin{matrix} (2) \\ (3) \end{matrix} \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}.$$

b) Solve by hand the optimal control problem when there are two decisions ($T = 2$). Then provide an explicit expression of the following quantities as a function of θ , C_r and C_b .

- $u_0^*(\text{Worn}) =$
 $\theta \cdot \min \{C_r, C_b\}$
the value of the state worn at the beginning of the first month;
- $a_0^*(\text{Broken}) =$
 $\begin{cases} \text{Buy new one} & \text{if } C_r \geq C_b \\ \text{Repair it} & \text{if } C_b < C_r \end{cases}$
the best action if the bike is broken at the beginning of the first month.

What is the best action at the beginning of the first month if the bike is broken for $\theta = 0.5$, $C_b = 9$, $C_r = 6$.

c) Assume that you start with a bike in perfect condition. You decide to never repair nor to buy a new bike. How long does it take in average to get a broken bike? (Here we assume that $T = \infty$).

Part 4:

a) We model the problem as a MDP:

• State Space: $\mathcal{S} = \{1 \text{ (Perfect)}, 2 \text{ (Worn)}, 3 \text{ (Broken)}\}$

• Actions: $A_3 := \{\text{Repair, Buy new one}\}$

$A_2 := \{\text{Keep it, Repair it, Buy new one}\} \Rightarrow A = \bigcup_{i \in \mathcal{S}} A_i$

$A_1 := \{\text{Keep it}\}$

• Objective: Minimize the expected cost over T months $\min \mathbb{E} \left[\sum_{t=0}^{T-1} r_t(s_t, a_t) \right] \rightarrow \text{Finite-Horizon MDP!!}$
 \rightarrow Since the condition at the end of the T^{th} month does not matter $r_T(s_T, a_T)$ is not included!

• Reward: $r_t(s_t=3, a_t=\text{Buy new one}) = C_b$

$r_t(s_t=1, a_t=\text{Keep it}) = 0$

$r_t(s_t=3, a_t=\text{Repair}) = C_r$

$r_t(s_t=2, a_t=\text{Buy new one}) = C_b$

$r_t(s_t=2, a_t=\text{Repair}) = C_r$

$r_t(s_t=2, a_t=\text{Keep it}) = 0$

• Transition Probabilities: $P_t(s'=2 | s_t=3, a_t=\text{Repair}) = 1$

$P_t(s'=1 | s_t=3, a_t=\text{Buy new one}) = 1$

$P_t(s'=1 | s_t=2, a_t=\text{Repair}) = 1$

\rightarrow All probabilities else are 0!!

$P_t(s'=3 | s_t=2, a_t=\text{Keep it}) = 0$

$P_t(s'=2 | s_t=1, a_t=\text{Keep it}) = 0$

$P_t(s'=2 | s_t=2, a_t=\text{Keep it}) = 1-0$

$P_t(s'=1 | s_t=1, a_t=\text{Keep it}) = 1-0$

\rightarrow The transition probabilities are used to fill in the transition matrices above!

b) We want to solve the following problem for $T=2$:

$$\min_{\pi} \mathbb{E} \left[\sum_{t=0}^1 r_t(s_t, a_t) \right]$$

This can be solved using Dynamic Programming

$$t=2: \quad u_2^B(s_2) = 0 \quad \forall s_2 \in \mathcal{S}$$

$$t=1: \quad u_1^B(s_1) = \min_{a \in A_{s_1}} \left[r_1(s_1, a) + \sum_{j \in \mathcal{S}} P_1(j | s_1, a) \cdot u_2^B(j) \right]$$

$$\forall s_t \in \mathcal{S}: \quad u_t^B(s_t) = \min_{a \in A_{s_t}} \left[r_t(s_t, a) \right]$$

$$① \quad U_1^B(s_1=1) = r_1(s_1=1, a_1=\text{keep it}) = 0$$

$$U_1^B(s_1=2) = \min_{a_1 \in \{\text{keep it, Repair it, Buy new one}\}} \{ r_1(s_1=2, a_1) \}$$

$$= \min \{ r_1(s_1=2, a_1=\text{keep it}), r_1(s_1=2, a_1=\text{Repair it}) \}$$

$$= \min \{ 0, C_r, C_b \} \stackrel{\uparrow \text{assuming } C_r, C_b > 0}{=} 0$$

$$U_1^B(s_1=3) = \min_{a_1 \in \{\text{Repair it, Buy new one}\}} \{ r_1(s_1=3, a_1) \}$$

$$= \min \{ r_1(s_1=3, a_1=\text{Repair it}), r_1(s_1=3, a_1=\text{Buy new one}) \}$$

$$= \min \{ C_r, C_b \}$$

$$t=0: \quad U_0^B(s_0) = \min_{a_0 \in A_{s_0}} \left\{ r_0(s_0, a_0) + \sum_{j \in S} P(j|s_0, a_0) \cdot U_1^B(j) \right\}$$

$$U_0^B(s_0=1) = \min_{a_0 \in \{\text{keep it}\}} \left\{ r_0(s_0=1, a_0) + p(s'=1|s_0=1, a_0) \cdot U_1^B(s'=1) + p(s'=2|s_0=1, a_0) \cdot U_1^B(s'=2) \right\}$$

$$= r_0(s_0=1, a_0=\text{keep it}) + p(s'=1|s_0=1, a_0=\text{keep it}) \cdot U_1^B(s'=1) + p(s'=2|s_0=1, a_0=\text{keep it}) \cdot U_1^B(s'=2)$$

$$= (1-\theta) \cdot 0 + \theta \cdot 0 = 0$$

$$U_0^B(s_0=2) = \min_{a_0 \in \{\text{keep it, Repair it, Buy new one}\}} \left\{ r_0(s_0=2, a_0) + \sum_j p(s'=j|s_0=2, a_0) \cdot U_1^B(s'=j) \right\}$$

$$= \min \{ 0 + (1-\theta) \cdot 0 + \theta \cdot \min \{ C_r, C_b \}, C_r + 1 \cdot 0, C_b \}$$

$$= \min_{c \in [0,1]} \{ \theta \cdot \min \{ C_r, C_b \}, C_r, C_b \} = \underline{\theta \cdot \min \{ C_r, C_b \}}$$

$$U_0^B(s_0=3) = \min_{a_0 \in \{\text{Buy new one, Repair it}\}} \left\{ r_0(s_0=3, a_0) + \sum_j p(s'=j|s_0=3, a_0) \cdot U_1^B(j) \right\}$$

$$= \min \{ C_b + 1 \cdot 0, C_r + 1 \cdot 0 \} = \underline{\min \{ C_b, C_r \}}$$

Hence we have obtained the value function: $U_2^*(s) = U_0^B(s), \forall s \in S$

The above expressions can easily be read off, $U_0^B(\text{Worn}) = U_0^B(s_0=2) = \theta \cdot \min \{ C_r, C_b \}$

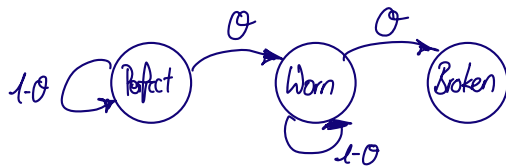
$U_0^B(\text{Broken}) = \arg \min (U_0(s_0=1)) = \begin{cases} \text{Buy new one}, & C_r \geq C_b \\ \text{Repair}, & C_b > C_r \end{cases}$

Assuming: $\theta = 0.5, C_b = 9, C_r = 6$

$U_0^B(s_0=\text{Broken}) = \arg \min_{a_0 \in \{\text{Buy new one, Repair it}\}} \{ U_0(s_0=\text{Broken}) \} = \underline{\text{Repair it}}$

→ Best action is to Repair it!

c) We start in state "Perfect" = ① and can only take one action "Keep it"! The state space is viewed below:



Let T define the event of the no. of stages until we reach the state Broken.

$$P[T=0] = 0$$

Probability to reach the state Broken in 0-steps

$$P[T=1] = 0$$

$$P[T=2] = 0^2$$

$$P[T=3] = 2 \cdot 0^2 (1-0)$$

$$P[T=4] = 3 \cdot 0^2 (1-0)^2$$

⋮

$$P[T=k] = \begin{cases} (k-1) \cdot 0^2 \cdot (1-0)^{k-2} & , \text{ for } k \geq 1 \\ 0 & , \text{ otherwise} \end{cases}$$

$$E[T] = \sum_{k=1}^{\infty} k \cdot P[T=k]$$

$$\begin{aligned}
 &= \sum_{k=1}^{\infty} k \cdot (k-1) \cdot 0^2 \cdot (1-0)^{k-2} = \frac{0^2}{(1-0)^2} \cdot \sum_{k=1}^{\infty} k(k-1)(1-0)^k \\
 &= \frac{0^2}{(1-0)^2} \cdot \sum_{k=2}^{\infty} k(k-1)(1-0)^k = \frac{0^2}{(1-0)^2} \cdot \sum_{k=1}^{\infty} (k+1) \cdot k (1-0)^{k+1} \\
 &= \frac{0^2}{(1-0)} \cdot \left[\sum_{k=1}^{\infty} k(1-0)^k + \sum_{k=1}^{\infty} k^2 (1-0)^k \right]
 \end{aligned}$$

⊗ We use two results from Geometric series:

$$\sum_{k=1}^{\infty} k \cdot x^k = \frac{x}{(1-x)^2}$$

$$\begin{aligned}
 \sum_{k=1}^{\infty} k^2 x^k &= (1-x)^{-3} \cdot [x^1 + x^2] \\
 &= \frac{x^1 + x^2}{(1-x)^2}
 \end{aligned}$$

$$\begin{aligned}
 &\left. \begin{aligned}
 &\otimes \quad \left[\frac{0^2}{(1-0)} \cdot \left[\frac{1-0}{0^2} + \frac{(1-0) + (1-0)^2}{0^3} \right] \right] \\
 &= \frac{1}{(1-0)} \cdot \left[1 + \frac{1+1-0}{0} \right] = \left[\frac{2}{0} \right] = \underline{\underline{0}}
 \end{aligned} \right|
 \end{aligned}$$

2 Optimal Stopping

You observe a fair coin being tossed T times. You may stop observing at any time, and when you do you receive as a reward the proportion of heads observed. For example, if the first toss is head, you should stop immediately. Your problem is to identify a stopping rule maximizing the average reward.

a) Model the problem as an MDP. How many states will you use? $T(T+1)+1$
Justify your answer and write Bellman's equations.

b) Establish by induction one of the following statement. Which one is true? A
Let $V_t(n)$ denote the maximal average reward if after t tosses, we got n heads.

- (A) For all t and n , $V_t(n+1) \geq V_t(n)$
- (B) For all t and n , $V_t(n+1) \leq V_t(n)$
- (C) For all t and n , $V_t(n+1) = V_t(n)$

c) One of the following policies is optimal. Which one? Justify your choice. Hint: proceed by elimination (justify why 2 of 3 strategies are not optimal). C

- (A) After the second toss, stop only if the number of heads reaches $T/2$
- (B) Never stop, except when the first toss is head
- (C) After t tosses and n observed heads, stop if and only if $n > \frac{t}{2}$

d) The coin is biased, with an unknown bias. We are using an off-policy RL algorithm converging to the optimal policy. The algorithm works with one of the following behavior policies. Which one? A

- (A) After t tosses and n observed heads, stop if and only if $n > t/2$
- (B) Never stop, i.e., always select the same action

a) We model the MDP (similar to Ex 2.1a)

- Statespace: + We define the state (t, n) where t : current time
 n : no. of times head has been observed

+ We define a further state: E indicating we have stopped playing

$$\mathcal{S} = \left(\underbrace{\{1, \dots, T\}}_{\text{dom}(t)} \times \underbrace{\{0, \dots, T\}}_{\text{dom}(n)} \right) \cup \{E\}$$

- Actions: + Continue (C)
 + Stop (S)

$$A = \{C, S\}$$

- Reward Function:

$$\forall t: 1 \leq t \leq T: r((t, n), S) = \left(\frac{n}{t}\right)$$

$$r((t, n), C) = 0$$

$$r(E, a) = 0 \quad \forall a \in A$$

- Transition probabilities:

$$p_t(s' = F | s_t = (t, n), a_t = S) = 1$$

$$p_t(s' = F | s_t = (t, n), a_t = C) = 0$$

$$p_t(s' = E | s_t = E, a_t = S) = 1$$

$$p_t(s' = F | s_t = E, a_t = C) = 1$$

$$p_t(s' = (t+1, n) | s_t = (t, n), a_t = C) = \frac{1}{2}$$

$$p_t(s' = (t+1, n+1) | s_t = (t, n), a_t = C) = \frac{1}{2}$$

$$p_t(s' = (t+1, n+1) | s_t = (t, n), a_t = S) = 0$$

$$p_t(s' = (t+1, n) | s_t = (t, n), a_t = S) = 0$$

+ All other transition probabilities are zero!

→ Objective: $\max_{\pi} \mathbb{E} \left[\sum_{t=1}^T r_t(s_t, a_t) \right]$

→ Finite-horizon Problem

It is obvious that the cardinality of the statespace is: $|\mathcal{S}| = T \cdot (T+1) + 1$

- Bellman's Equations: $\forall s_t \in \mathcal{S}: U_T^B(s_T) = \max_{a_T \in A} r_T(s_T, a_T): U_T^B(F) = 0; U_T^B(s_T = (T, n)) = \frac{n}{T} \quad \forall n \in \text{dom}(n)$

$$\underline{t \in [2, \dots, T]}: \underline{\forall s_{t-1} \in \mathcal{S}}: U_{t-1}^B(s_t) = \max_{a_{t-1} \in A} \left[r_{t-1}(s_{t-1}, a_{t-1}) + \sum_j p_t(s' = j | s_{t-1} = j, a_{t-1}) \cdot U_t^B(s' = j) \right]$$

→ ① $U_{t-1}^B(s_{t-1} = F) = 0$ ② $U_{t-1}^B(s_{t-1} = (t-1, n)) = \max \left\{ \frac{n}{t-1}, \frac{1}{2} U_t^B(s' = (t, n)) + \frac{1}{2} U_t^B(s' = (t, n+1)) \right\}, n = 0, 1, \dots, t-1$

f) $V_t(n)$ corresponds to the maximal average reward after t tosses and observing n heads. This means we are at stage t and state $s_t = (t, n)$. It is clear that $V_t(n) = U_t^B(s_t = (t, n))$!

We prove the following statement $V_t(n) \leq V_{t+1}(n+1)$ via induction:

+ Induction Start: $k=T$: $V_T(n) = \frac{n}{T} < V_T(n+1) = \frac{n+1}{T} \quad n = 0, 1, \dots, T-1$ which satisfies the above statement!

+ Induction Hypothesis: Assume $V_{t+1}(n) \leq V_{t+1}(n+1) \quad n = 0, 1, \dots, T-1$

+ Induction Step: $t+1 \rightarrow t$: We use the Bellman Equation to obtain

$$V_t(n) = \max \left\{ \frac{n}{t}, \frac{V_{t+1}(n) + V_{t+1}(n+1)}{2} \right\} \quad \text{and} \quad V_t(n+1) = \max \left\{ \frac{n+1}{t}, \frac{V_{t+1}(n+1) + V_{t+1}(n+2)}{2} \right\}$$

Since $\frac{n+1}{t} > \frac{n}{t}$ it remains to show $\frac{V_{t+1}(n) + V_{t+1}(n+1)}{2} \leq \frac{V_{t+1}(n+1) + V_{t+1}(n+2)}{2}$

which simplifies to $V_{t+1}(n) \leq V_{t+1}(n+2)$ which follows from the induction hypothesis!

Hence, we've shown that $V_t(n) \geq V_{t+1}(n+1)$

c) Let's solve the problem for $T=4$:

$$\underline{t=4}: U_4^B(s_4=(4,0)) = 1$$

$$U_4^B(s_4=(4,1)) = \frac{3}{4}$$

$$U_4^B(s_4=(4,2)) = \frac{1}{2}$$

$$U_4^B(s_4=(4,3)) = \frac{1}{4}$$

$$U_4^B(s_4=(4,4)) = 0$$

$$U_4^B(s_4=E) = 0$$

$$\underline{t=3}: U_3^B(s_3=(3,3)) = \max \left\{ 1, \frac{1}{2} \cdot \frac{3}{4} + \frac{1}{2} \cdot 1 \right\}$$

$$= \max \left\{ 1, \frac{7}{8} \right\} = 1$$

$$U_3^B(s_3=(3,2)) = \max \left\{ \frac{3}{4}, \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{3}{4} \right\}$$

$$= \max \left\{ \frac{3}{4}, \frac{5}{8} \right\} = \frac{3}{4}$$

$$U_3^B(s_3=(3,1)) = \max \left\{ \frac{1}{2}, \frac{1}{2} \cdot \frac{1}{4} + \frac{1}{2} \cdot \frac{1}{2} \right\} = \frac{3}{8}$$

$$U_3^B(s_3=(3,0)) = \max \left\{ 0, \frac{1}{2} \cdot 0 + \frac{1}{2} \cdot \frac{1}{4} \right\} = \frac{1}{8}$$

$$U_3^B(s_3=E) = 0$$

$$\underline{t=2}: U_2^B(s_2=(2,2)) = \max \left\{ 1, \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{3}{4} \right\} = 1$$

$$U_2^B(s_2=(2,1)) = \max \left\{ \frac{1}{2}, \frac{1}{2} \cdot \frac{1}{4} + \frac{1}{2} \cdot \frac{3}{8} \right\} = \frac{25}{48}$$

$$U_2^B(s_2=(2,0)) = \max \left\{ 0, \frac{1}{2} \cdot \frac{1}{8} \right\} = \frac{3}{16}$$

$$U_2^B(s_2=E) = 0$$

$$\underline{t=1}: U_1^B(s=(1,1)) = 1$$

$$U_1^B(s=(1,0)) = \max \left\{ 0, \frac{1}{2} \cdot \frac{3}{16} + \frac{1}{2} \cdot \frac{25}{48} \right\} = \frac{17}{48}$$

$$U_1^B(s=E) = 0$$

A) is clearly false! Let's say T is fixed. The policy tells us we can only stop after the 2^{nd} toss if we've reached $\frac{T}{2}$ heads. However, we can stop before reaching $\frac{T}{2}$ [for example if we've observed only heads] and obtain a higher reward on average!

B) is also false! The above example serves as a counter example: If we've observed 2 heads in 3 tosses the optimal policy will be to stop!!

⇒ Hence C) is the optimal policy!

d) We know from Part 3 of the lecture that for an off-policy algorithm to converge, the behaviour policy should visit every state, action pair.

Hence, if we choose B) we will always select the same action and no exploring is taking place.

⇒ However, if we select A) we are guaranteed to explore all actions!