

מבוא לבינה מלאכותית - תשפ"ד - תרגיל 5

אחמד דלאשה 324059856 | עבד נסאר 207063108

24 ביולי 2024

שאלה מספר 1.1.1 :

להלן התיאור של רכיבי ה-MDP ותהליך איטרציית הערכים כולו, כולל החישובים :

רכיבי ה-MDP :

• מצבים :

$$(S) : \{s_1, s_2, s_3, s_4, s_5, s_6, s_7, s_8, s_9\}$$

• פעולות :

$$(A) : \{Left, Right, Up, Down\}$$

• מודל המעבר (T) : דטרמיניסטי, כלומר התוצאה של כל פעולה היא ודאית. יציאה מהגרید מחזירה את הסוכן לאותו מצב.

• פונקציית התגמול (R) :

$$r(s) = -0.05 \text{ for } s \notin \{s_5, s_7, s_9\}$$

$$r(s_5) = -10$$

$$r(s_7) = 15$$

$$r(s_9) = 30$$

s_5, s_7, s_9 are absorbing states.

אתחול :

ערך התחלתי לכל המצבים $V_0(s) = 0$

$$V_{k+1}(s) = \max_a \left(\sum_{s'} T(s' | s, a) \left[R(s, a, s') + \gamma V_k(s') \right] \right)$$

תהליך האיטרציה:

נבצע איטרציית ערכים למשך 7 איטרציות, ונעדכן את הערך של כל מצב בכל שלב.

- איטרציה 1: נחשב את הערכים עבור כל מצב בהתבסס על משוואת בלמן.
- איטרציות 2 עד 7: נמשיך לעדכן את הערכים בהתבסס על תוצאות מהאיטרציה הקודמת.

להלן התוצאות של תהליך איטרציית הערכים למשך 7 איטרציות:

	iteration	s1	s2	s3	s4	s5	s6	s7	s8	s9
0	1.0	-0.05	-0.05	-0.05	-0.05	-10.0	-0.05	15.0	-0.05	30.0
1	2.0	-0.0995	-0.0995	-0.0995	14.8	-10.0	29.65	15.0	29.65	30.0
2	3.0	14.602	-0.1485	29.3035	14.8	-10.0	29.65	15.0	29.65	30.0
3	4.0	14.602	28.9605	29.3035	14.8	-10.0	29.65	15.0	29.65	30.0
4	5.0	28.6209	28.9605	29.3035	14.8	-10.0	29.65	15.0	29.65	30.0
5	6.0	28.6209	28.9605	29.3035	28.2847	-10.0	29.65	15.0	29.65	30.0
6	7.0	28.6209	28.9605	29.3035	28.2847	-10.0	29.65	15.0	29.65	30.0

איור 1:

ערכים אלו מייצגים את התגמולים המצופים המוזלים לכל מצב לאחר 7 איטרציות.

שאלה מספר 1.1.2 :

חישוב הפעולות האופטימליות והצגת הטבלה :

נחשב את הפעולות האופטימליות עבור כל מצב בכל איטרציה ונציג את התוצאות בטבלה. הטבלה תכלול 7 שורות (אחת לכל איטרציה) ו-10 עמודות (מספר האיטרציה והפעולה האופטימלית לכל מצב).

	iteration	s1	s2	s3	s4	s5	s6	s7	s8	s9
0	1	Left	Left	Left	Left	Left	Left	Left	Left	Left
1	2	Left	Left	Left	Up	Left	Up	Left	Right	Left
2	3	Up	Left	Up	Up	Left	Up	Left	Right	Left
3	4	Up	Right	Up	Up	Left	Up	Left	Right	Left
4	5	Right	Right	Up	Up	Left	Up	Left	Right	Left
5	6	Right	Right	Up	Down	Left	Up	Left	Right	Left
6	7	Right	Right	Up	Down	Left	Up	Left	Right	Left

איור 2:

The optimal policy grid correct :

←	→	←
↓	←	↑
→	→	↑

איור 3:

שאלה מספר 1.2 :

בשלב הזה, חישבנו את ערכי המדינות עבור MDP סטוכסטי בו כל פעולה מצליחה עם הסתברות $P = 0.9$ ובכשלון מעבירה את הסוכן לתא סמוך עם הסתברות אחידה.

רכיבי ה- MDP :

• מצבים :

$$(S) : \{s_1, s_2, s_3, s_4, s_5, s_6, s_7, s_8, s_9\}$$

• פעולות :

$$(A) : \{Left, Right, Up, Down\}$$

• מודל המעבר (T) :

- $P(s' | s, a) = 0.9$ עבור הפעולה המוצלחת

- $P(s' | s, a) = \frac{0.1}{3}$ עבור כל תא סמוך כאשר N הוא מספר התאים הסמוכים

• פונקציית התגמול (R) :

$$r(s) = -0.05 \text{ for } s \notin \{s_5, s_7, s_9\}$$

$$r(s_5) = -10$$

$$r(s_7) = 15$$

$$r(s_9) = 30$$

s_5, s_7, s_9 are absorbing states.

להלן התוצאה :

	s1	s2	s3	s4	s5	s6	s7	s8	s9
0	24.669	25.793	27.716	21.018	-10.0	28.19	15.0	27.761	30.0

איור 4:

השוואה בין הערכים של המדיניות האופטימלית בין המקרה הדטרמיניסטי והסטוכסטי

	State	Deterministic	Stochastic
0	s1	28.6209	24.011
1	s2	28.9605	24.61
2	s3	29.3035	26.904
3	s4	28.2847	21.55
4	s5	-10.0	-10.0
5	s6	29.65	27.52
6	s7	15.0	15.0
7	s8	29.65	26.93
8	s9	30.0	30.0

איור 5:

ניתוח:

במצב הסטוכסטי, הערכים לרוב נמוכים יותר מאשר במצב הדטרמיניסטי. ההסבר לכך הוא שבמצב הסטוכסטי יש סיכוי לכישלון של פעולה, מה שגורם להורדת הערכים של המדינות מכיוון שהסוכן יכול לעבור למצב פחות רצוי.

סיכום:

ההבדל בין הערכים הדטרמיניסטיים לסטוכסטיים מדגיש את החשיבות של מודל המעבר ואי הוודאות בניהול מדיניות אופטימלית בסביבות מורכבות. ההשוואה ממחישה כיצד אי ודאות משפיעה על הערכים הצפויים של המדינות ועל קבלת ההחלטות בסביבות שונות.

שאלה מספר 1.3 :

	s1	s2	s3	s4	s5	s6	s7	s8	s9
0	14.39	24.08	26.84	13.53	-10.0	27.51	15.0	26.92	30.0

איור 6:

ניתוח השוואתי

ניתן לראות כי הערכים של המדיניות הנתונה נמוכים יותר בערכים מסוימים לעומת המדיניות האופטימלית:

- s_1 : 14.39 נמוך מ- 24.011
- s_2 : 24.08 נמוך מ- 24.61
- s_4 : 13.53 נמוך מ- 21.55

המדיניות הנתונה מובילה לערכים נמוכים יותר מאשר המדיניות האופטימלית בחלק מהמצבים, דבר שמצביע על כך שהמדיניות הנתונה פחות אופטימלית בסביבה הסטוכסטית.

מסקנות:

- השפעת המדיניות: המדיניות הנתונה משפיעה באופן ישיר על הערכים הצפויים של המצבים. מדיניות פחות אופטימלית תוביל לערכים נמוכים יותר.
- סביבה סטוכסטית: בסביבה סטוכסטית, קיימת אי-וודאות בתוצאות של פעולות. המדיניות האופטימלית מצליחה להתמודד טוב יותר עם אי-הוודאות הזו ומשיגה ערכים גבוהים יותר.
- החשיבות של אופטימיזציה: על מנת להשיג ערכים מקסימליים בסביבה סטוכסטית, חשוב לבצע אופטימיזציה למדיניות תוך התחשבות באי-וודאות.

סיכום:

המדיניות הנתונה הובילה לערכים נמוכים יותר בהשוואה למדיניות האופטימלית בסביבה סטוכסטית. זה מדגיש את החשיבות של אופטימיזציה של מדיניות בסביבות עם אי-וודאות על מנת להשיג את הערכים הטובים ביותר עבור המצבים.

Q2: POMDP(1) נתון e - prior belief על האופן שבו p_d מתנהג:

$$p(TL) = p(TR) = 0.5$$

$$- p(TL|GL, d=1) = \frac{p(GL|TL, d=1) \cdot p(TL)}{p(GL|TL, d=1) \cdot p(TL) + p(GL|TR, d=1) \cdot p(TR)}$$

bayes rule

$$= \frac{0.9 \cdot \frac{1}{2}}{0.9 \cdot \frac{1}{2} + 0.1 \cdot \frac{1}{2}} = \frac{0.45}{0.5} = \boxed{0.9}$$

$$- p(TR|GL, d=1) = \frac{p(GL|TR, d=1) \cdot p(TR)}{p(GL|TR, d=1) \cdot p(TR) + p(GL|TL, d=1) \cdot p(TL)}$$

bayes rule

$$= \frac{0.1 \cdot \frac{1}{2}}{0.1 \cdot \frac{1}{2} + 0.9 \cdot \frac{1}{2}} = \frac{0.05}{0.5} = \boxed{0.1}$$

$$- p(TL|GL, d=2) = \frac{p(GL|TL, d=2) \cdot p(TL)}{p(GL|TL, d=2) \cdot p(TL) + p(GL|TR, d=2) \cdot p(TR)}$$

bayes rule

$$= \frac{0.6 \cdot \frac{1}{2}}{0.6 \cdot \frac{1}{2} + 0.4 \cdot \frac{1}{2}} = \frac{0.3}{0.5} = \boxed{0.6}$$

$$- p(TR|GL, d=2) = \frac{p(GL|TR, d=2) \cdot p(TR)}{p(GL|TR, d=2) \cdot p(TR) + p(GL|TL, d=2) \cdot p(TL)}$$

Bayes rule

$$\frac{0.4 \cdot \frac{1}{2}}{0.4 \cdot \frac{1}{2} + 0.6 \cdot \frac{1}{2}} = \frac{0.2}{0.5} = \boxed{0.4}$$

$$d=1 : < 0.9, 0.1 >$$

$\Rightarrow \sqrt{2} p$

$$d=2 : < 0.6, 0.4 >$$

2.) The posterior belief about the tiger's location after hearing roar from the left door calculation :-

$$- p(TL|GL) = p(TL|GL, d=1) \cdot p(d=1) + p(TL|GL, d=2) \cdot p(d=2)$$

Law of total probability

$$\Rightarrow 0.9p + 0.6(1-p) = \boxed{0.6 + 0.3p}$$

$$- p(TR|GL) = p(TR|GL, d=1) \cdot p(d=1) + p(TR|GL, d=2) \cdot p(d=2)$$

Law of total probability

$$\Rightarrow 0.1p + 0.4(1-p) = \boxed{0.4 - 0.3p}$$

$$< 0.3p + 0.6, 0.4 - 0.3p >$$

$\Rightarrow \sqrt{2} p$

3.)

a.) Move to $d=1$, listen if GR, then OL:

$$E(R) = R(\text{move}) + R(\text{listen}) + \frac{1}{2} \cdot (p(TL|GR, d=1) \cdot R(TL, OL, d=1) + p(TR|GR, d=1) \cdot R(TR, OL, d=1))$$

$$\textcircled{+} R(\text{move}) = -2, \textcircled{+} R(\text{listen}) = -1$$

$$\textcircled{+} p(TL|GR, d=1) \cdot R(TL, OL, d=1) = p(TR|GL, d=1) \cdot R(TL, OL, d=1)$$

from symmetry

b) listen. if GL, then listen. Else if GR, then ol:

$$E(R) = R(\text{listen}) \cdot \frac{1}{2} \cdot (P(TL|GL, d=2) + P(TR|GL, d=2) \cdot R(\text{listen})) \\ + \frac{1}{2} \cdot (P(TL|GR, d=2) \cdot R(TL, OL, d=2) + P(TR|GR, d=2) \cdot R(TR, OL, d=2))$$

$$- R(\text{listen}) = -1$$

$$- P(TL|GL, d=2) \cdot R(\text{listen}) = 0.6 \cdot (-1) = -0.6$$

$$- P(TL|GL, d=2) \cdot R(\text{listen}) = 0.4 \cdot (-1) = -0.4$$

$$- P(TL|GR, d=2) \cdot R(TL, OL, d=2) = P(TR|GL, d=2) \cdot R(TL, OL, d=2) = 0.4 \cdot (-5) = -2.0$$

$$- P(TR|GR, d=2) \cdot R(TR, OL, d=2) = P(TL|GL, d=2) \cdot R(TL, OR, d=2) = 0.6 \cdot 10 = 6$$

→ 100% 100% 100%

$$E(R) = 0.5 \cdot (-1) + 0.5 \cdot (-14) = -8.5$$

• observations → 100% 100% 100% 100%

{GL, GL} {GL, GR} {GR, GL} {GR, GR}

→ {GL, GL} 100% (1)

$$P(TL|GL, d=1) = 0.9 \quad \rightarrow 100\% 100\%$$

$$P(TR|GR, d=1) = 0.1$$

$$p(T_L | G_L, G_L, d=1) = \frac{p(G_L | T_L, G_L, d=1) \cdot p(T_L | G_L, d=1)}{p(G_L | G_L, d=1)} \quad \text{--- } \int_{\Omega} p \Omega$$

$$\hookrightarrow p(G_L | G_L, d=1) = p(G_L | G_L, T_L, d=1) \cdot p(T_L | G_L, d=1) + p(G_L | G_L, T_R, d=1) \cdot p(T_R | G_L, d=1)$$

$$\Rightarrow p(G_L | T_L, d=1) \cdot p(T_L | G_L, d=1) +$$

$$p(G_L | T_R, d=1) \cdot p(T_R | G_L, d=1) = 0.9 \cdot 0.9 + 0.1 \cdot 0.1 = 0.82$$

$$= \frac{p(G_L | T_L, d=1) \cdot p(T_L | G_L, d=1)}{0.82} = \frac{0.9 \cdot 0.9}{0.82} = 0.988$$

--- beliefs agent pos

$$d=1 : 0.988$$

$$d=2 : 1 - 0.988 = 0.012$$

$\{GL, GR\}$ 128 (II)

$$P(TL|GL, d=1) = 0.9$$

\therefore 90% נכון

$$P(TR|GR, d=1) = 0.1$$

\therefore 10% טעות

$$P(TL|GL, GR, d=1) = \frac{P(GR|TL, GL, d=1) \cdot P(TL|GL, d=1)}{P(GR|GL, d=1)}$$

$$\begin{aligned} \rightarrow P(GR|GL, d=1) &= P(GR|GL, TL, d=1) \cdot P(TL|GL, d=1) \\ &\quad + P(GR|GL, TR, d=1) \cdot P(TR|GL, d=1) \\ &\Rightarrow P(GR|TL, d=1) \cdot P(TL|GL, d=1) \\ &\quad + P(GR|TR, d=1) \cdot P(TR|GL, d=1) \\ &= 0.1 \cdot 0.9 + 0.9 \cdot 0.1 = \boxed{0.18} \end{aligned}$$

\therefore 18% נכון

$$\frac{P(GR|TL, d=1) \cdot P(TL|GL, d=1)}{0.18} = \frac{0.1 \cdot 0.9}{0.18} = \boxed{0.5}$$

\therefore 50% נכון

$$d=1 : 0.5$$

$$d=2 : 1 - 0.5 = 0.5$$

הסתברות 4 הוצאות והסתברות

$$\{GL, GL\} = \langle 0.988, 0.012 \rangle$$

$$\{GR, GL\} = \langle 0.5, 0.5 \rangle$$

$$\{GL, GR\} = \langle 0.5, 0.5 \rangle$$

$$\{GR, GR\} = \langle 0.012, 0.988 \rangle$$

הסתברות 4 הוצאות והסתברות

הסתברות 4 הוצאות והסתברות

Question 3:

1.) we know that policy determines how the agent behaves

a.) μ^t - policy parameter, i.e. the location the robot will pick

b.) y^t - the action of the robot $\Rightarrow y^t = \mu^t + z, z \sim N(0, \sigma)$

$y^t \sim N(\mu^t, \sigma^2)$ מכיוון $z \sim N(0, \sigma^2)$ ו- μ^t הוא policy של agent

$$\pi(y|\mu) = \mathcal{L}(\mu|y) = f_{\mu}(y=y) = \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{\frac{-(y-\mu)^2}{2\sigma^2}}$$

2) update rule :

לפי הסבר הסרט - G_t הוא ה-gain

$$\mu = \mu + \alpha \cdot G_t \cdot \nabla_{\mu} \log \pi_{\mu}(y^t)$$

$$G_t = \sum_{i=t+1}^T r_i \quad \text{ה-Return}$$

שאלה מספר 3.3 :

(3) נתון $z \sim N(0, \sigma^2)$ ו- $y = \mu + z$ $y \sim N(\mu, \sigma^2)$ μ ו- σ^2 הם פרמטרים של y .

נניח שיש לנו m תצפיות של y ונרצה להעריך את μ . נכתוב את הפונקציה של התוחלת:

$$E(r) = E[-(m-y)^2] = -E[(m-y)^2] = -E[m^2 - 2my + y^2]$$

$$= -E[m^2] + E[2my] - E[y^2] = -m^2 + 2mE[y] - E[y^2]$$

אנחנו יודעים ש- $E[y] = \mu$ ו- $E[y^2] = \mu^2 + \sigma^2$.

$$E(r) = -m^2 + 2m\mu - (\mu^2 + \sigma^2)$$

נאזן את $E(r)$ ל-0 כדי למצוא את הערך של μ שממזער את הפונקציה:

$$(-m^2 + 2m\mu - \mu^2 - \sigma^2)' = 2m - 2\mu = 0$$

$$\Rightarrow \mu = m$$

שאלה מספר 3.4 :

כדי למצוא את שיעור השינוי הממוצע של μ עלינו לחשב את התוחלת של שלב העדכון:

$$\mu = \mu + \alpha \cdot \left(-(m - y_t)^2 \right) \cdot \frac{y_t - \mu}{\sigma^2}$$

פתרון שלב אחר שלב:

הגדרת שלב העדכון:

$$\Delta\mu = \alpha \cdot \left(-(m - y_t)^2 \right) \cdot \frac{y_t - \mu}{\sigma^2}$$

לקיחת התוחלת:

$$E[\Delta\mu] = E \left[\alpha \cdot \left(-(m - y_t)^2 \right) \cdot \frac{y_t - \mu}{\sigma^2} \right]$$

החלפת $y_t = \mu + z$ כך ש- $z \sim N(0, \sigma^2)$:

$$E \left[\alpha \cdot \left(-(m - (\mu + z))^2 \right) \cdot \frac{(\mu + z) - \mu}{\sigma^2} \right]$$

$$= \alpha \cdot E \left[-(m - \mu - z)^2 \cdot \frac{z}{\sigma^2} \right]$$

פיתוח הריבוע:

$$= \alpha \cdot E \left[-\frac{(m - \mu - z)^2 \cdot z}{\sigma^2} \right]$$

$$= \alpha \cdot \mathbb{E} \left[- \frac{\left((m - \mu)^2 - 2(m - \mu) + z + z^2 \right) \cdot z}{\sigma^2} \right]$$

הפרדת התוחלת:

$$= \alpha \cdot \left[- \frac{(m - \mu)^2 \cdot \mathbb{E}[z]}{\sigma^2} + \frac{2(m - \mu) \cdot \mathbb{E}[z^2]}{\sigma^2} + \frac{\mathbb{E}[z^3]}{\sigma^2} \right]$$

הערכת התוחלות:

$$\mathbb{E}[z] = 0 \text{ (mean of a zero - mean Gaussian)}$$

$$\mathbb{E}[z^2] \text{ (variance of Gaussian)}$$

$$\mathbb{E}[z^3] \text{ (skewness of Gaussian is zero)}$$

לכן :

$$\begin{aligned} \mathbb{E}[\Delta\mu] &= \alpha \cdot \left[0 + \frac{2(m - \mu)\sigma^2}{\sigma^2} - 0 \right] \\ &= 2\alpha \cdot (m - \mu) \end{aligned}$$

הסבר :

כלל העדכון מראה ששיעור השינוי הממוצע של μ הוא יחסי ישיר להבדל בין היעד m והערכת הנוכחית μ . קצב הלמידה α מגדיר את מהירות השינוי. זה אומר שככל שהסוכן מעדכן שוב ושוב את μ , הוא יכוון בהדרגה את μ לכיוון היעד m , והמהירות של התאמה זו מושפעת מה- α . פתרון זה מראה כיצד אלגוריתם *REINFORCE* מתאים את פרמטר המדיניות μ עם הזמן, ומבטיח שהרובוט משתפר בקליעת המטרה על ידי הפחתת ההבדל בין μ ל- m .

$$\nabla_{\sigma} \log \pi(y|\mu, \sigma) \quad (5) \quad \text{צריך לחשב :-}$$

$$\begin{aligned} \nabla_{\sigma} \log \pi(y|\mu, \sigma) &= \nabla_{\sigma} \left(\log \frac{1}{\sigma \sqrt{2\pi}} \right) + \log \left(e^{-\frac{(y-\mu)^2}{2\sigma^2}} \right) \\ &= \nabla_{\sigma} \left(-\log(\sigma \sqrt{2\pi}) - \frac{(y-\mu)^2}{2\sigma^2} \right) = \frac{-\sqrt{2\pi}}{\sigma \sqrt{2\pi}} + \frac{(y-\mu)^2}{\sigma^2} = \frac{(y-\mu)^2}{\sigma^2} - \frac{1}{\sigma} \\ &= \frac{(y-\mu)^2 - \sigma^2}{\sigma^2} \end{aligned}$$

$$\begin{bmatrix} \mu^t \\ \sigma^t \end{bmatrix} = \theta^t = \theta^{t-1} - \alpha (m - y)^2 \begin{bmatrix} \frac{(y-\mu)^{t-1}}{(\sigma^{t-1})^2} \\ \frac{(y-\mu^{t-1})^2 - (\sigma^{t-1})^2}{(\sigma^{t-1})^3} \end{bmatrix} \quad \text{סך הכל הסתכלו}$$

מסתכל בקורבס קיבלנו כי $\mathbb{E}[\Delta \mu] = 2\alpha(m - \mu)$, כי זה נמצא על סף כן
 והסתכלנו ה'טו' פרוגרדנציה' להיבטיל בין $\mu - \sigma$ ו- μ .

$$\mathbb{E}[\Delta \sigma] = \mathbb{E} \left[-\alpha (m - y)^2 \cdot \frac{2(y - \mu)^2 - \sigma^2}{\sigma^3} \right] = \frac{-\alpha}{\sigma^3} \left(\mathbb{E}[(m - y)^2 (y - \mu)^2] - \sigma^2 \mathbb{E}[(m - y)^2] \right)$$

$$y \sim N(\mu, \sigma) \quad \leftarrow \quad \frac{-\alpha}{\sigma^3} \left(\mathbb{E}[(m - y)^2 (y - \mu)^2] - \sigma^2 \mathbb{E}[(m - y)^2] \right) =$$

$$(m - y) \sim N(m - \mu, \sigma) \quad = \quad \frac{-\alpha}{\sigma^3} \left((m - \mu)^2 \sigma^2 + 3\sigma^4 - \sigma^2 (\sigma^2 (m - \mu)^2) \right) = \boxed{-2\alpha\sigma}$$

* אלו 'בול'ים לנאות הם שיש להם $n-m$, היעדרם 'כדי' את m , ואם m גדול
 $n-m$, היעדרם 'בול'ים אף m וכל $n-m$ מהקבוצה היעדרם 'בול'ים
 זהו לרוב (מ- μ) 'כדי' לראות.

[illegible]

לכן צ"ע לה' ז' - ב' קטן עם הסמן בלתי נחשב עם שיתאם.

לפי הסיפור, מבין e-agent בוסס במטרה, כלומר הוא 'עובד' על המטרה!
מבחינת מומחיות - יבנה מחדש 'עובד' מומחיות