

פְּרָמֵשׁ – פְּרָמֵשׁ – MDP

פְּרָמֵשׁ 1:

בתרגול ראינו את משוואת בלמן כאשר התגמול ניתן עבור המצב הנוכחי בלבד, כלומר $R: S \rightarrow \mathbb{R}$, למתן תגמול זה נקרא "תגמול על הצמתים" מכיוון שהוא תלוי בצומת שהסוכן נמצא בו. בהתאם להגדרה זו הצגנו בתרגול את האלגוריתמים Value iteration ו-Policy Iteration למציאת המדיניות האופטימלית.

כעת, נרחיב את ההגדרה הזו, לתגמול המקבל את המצב הנוכחי והפעולה לביצוע שבה בחר הסוכן, כלומר: $R: S \times A \rightarrow \mathbb{R}$, למתן תגמול זה נקרא "תגמול על פעולה".

א. (2 נק') התאימו את הנוסחה של התוחלת של התועלת מהתרגול, עבור התוחלת של התועלת המתקבלת במקרה של "תגמול על פעולה", אין צורך לנמק.

$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(S_t, \pi(S_t)) \mid S_0 = s \right]$$

ב. (2 נק') כתבו מחדש את נוסחת משוואת בלמן עבור המקרה של "תגמול על פעולה", אין צורך לנמק.

$$U(s) = \max_{a \in A(s)} [R(s, a) + \gamma \sum_{s'} P(s'|s, a) * U(s')]$$

ג. (4 נק') נסחו את אלגוריתם Value Iteration עבור המקרה של "תגמול על פעולה". התייחסו גם למקרה בו $\gamma = 1$, והסבירו מה לדעתכם התנאים שצריכים להתקיים על הסביבה **mdp** על מנת שתמיד נצליח למצוא את המדיניות האופטימלית.

```
function VALUE-ITERATION(mdp,  $\epsilon$ ) returns a utility function
inputs: mdp, an MDP with states S, actions A(s), transition model  $P(s' \mid s, a)$ ,
rewards R(s), discount  $\gamma$ 
 $\epsilon$ , the maximum error allowed in the utility of any state
local variables: U, U', vectors of utilities for states in S, initially zero
 $\delta$ , the maximum change in the utility of any state in an iteration

repeat
  U  $\leftarrow$  U';  $\delta \leftarrow 0$ 
  for each state s in S do
    U'[s]  $\leftarrow \max_{a \in A(s)} [R(s, a) + \gamma \sum_{s'} P(s'|s, a) * U(s')]$ 
    if  $|U'[s] - U[s]| > \delta$  then  $\delta \leftarrow |U'[s] - U[s]|$ 
until  $\delta < \epsilon(1 - \gamma)/\gamma$ 
return U
```

אם $\gamma = 1$: תנאי העצירה שך value_iteration יהיה כאשר $\delta = 0$, כלומר כאשר אין הפרש בין ערכי התועלת באיטרציה הקודמת ובאיטרציה הנוכחית, במקרה זה תובטח התכנסות למדיניות אופטימלית כפי שלמדנו בקורס.

ד. (4 נק') נסחו את אלגוריתם Policy Iteration עבור המקרה של "תגמול על פעולה".

התייחסו גם למקרה בו $\gamma = 1$, והסבירו מה לדעתכם התנאים שצריכים להתקיים על הסביבה\mdp על מנת שתמיד נצליח למצוא את המדיניות האופטימלית.

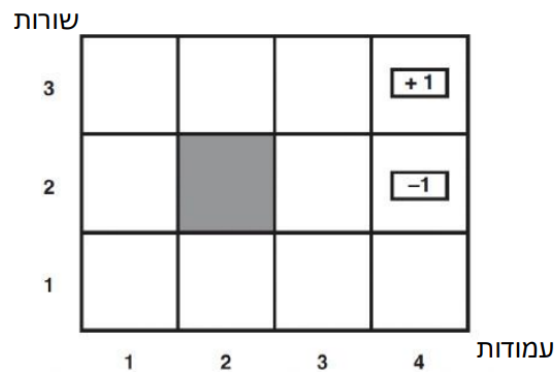
```
function POLICY-ITERATION(mdp) returns a policy
  inputs: mdp, an MDP with states  $S$ , actions  $A(s)$ , transition model  $P(s' | s, a)$ 
  local variables:  $U$ , a vector of utilities for states in  $S$ , initially zero
                   $\pi$ , a policy vector indexed by state, initially random

  repeat
     $U \leftarrow \text{POLICY-EVALUATION}(\pi, U, \text{mdp})$ 
     $\text{unchanged?} \leftarrow \text{true}$ 
    for each state  $s$  in  $S$  do
      if  $\max_{a \in A(s)} [R(s, a) + \gamma \sum_{s'} P(s'|s, a) * U(s')] > R(s, \pi(s)) + \gamma \sum_{s'} P(s'|s, \pi(s)) * U(s')$  do
         $\pi[s] \leftarrow \max_{a \in A(s)} [R(s, a) + \gamma \sum_{s'} P(s'|s, a) * U(s')]$ 
       $\text{unchanged?} \leftarrow \text{false}$ 
  until  $\text{unchanged?}$ 
  return  $\pi$ 
```

אם $\gamma = 1$: אין צורך לשנות את האלגוריתם, ההתכנסות מובטחת ב-policy_iteration אם מרחב המצבים סופי, מספר הפעולות סופי, פונקציית התגמולים חסומה, קיימים מצבים סופיים ואין מעגל תגמולים חיובי. במקרה זה גם אם $\gamma = 1$ מובטחת התכנסות למדיניות אופטימלית (כיוון שמרחב המדיניות סופי) גם ללא שינוי באלגוריתם.

תרגיל 2:

נתון ה-MDP הבא $\langle S, A, P, R, \gamma \rangle$, אופק אינסופי:



מצבים:

$$S = \{(1, 1), (1, 2), (1, 3), (1, 4), (2, 1), (2, 3), (2, 4), (3, 1), (3, 2), (3, 3), (3, 4)\}$$

$$S_G = \{(2, 4), (3, 4)\}$$

פעולות

$$\forall S \setminus S_G: A(s) = \{\text{Up, Down, Left, Right}\}$$

תגמולים:

$$R((2, 4)) = -1, R((3, 4)) = +1$$

נתונים התגמולים של המצבים הסופיים בלבד: שימו לב, התגמולים הינם תגמולים על המצבים.

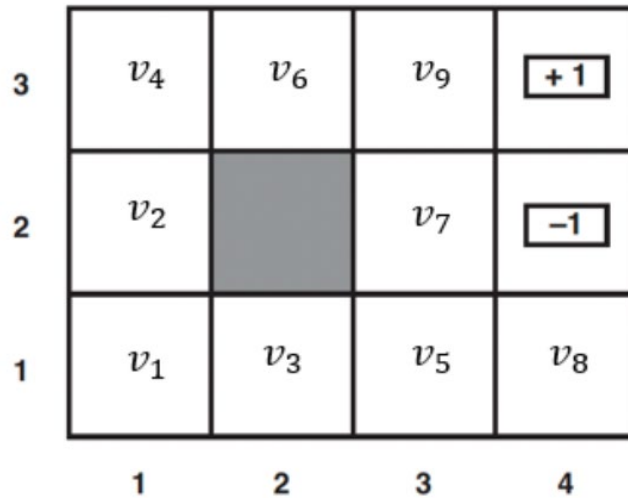
ישנם תגמולים עבור שאר המצבים, הם פשוט לא נתונים כחלק מהשאלה.

מודל מעבר:

כל פעולה "מצליחה" בהסתברות 0.8, ואם היא לא מצליחה אז בהסתברות שווה מתבצעת אחת הפעולות המאונכות לפעולה המתבקשת. כאשר הסוכן הולך לכיוון הקיר או מחוץ ללוח הוא נשאר במקום.

$$0 < \gamma < 1$$

הרצתם את האלגוריתם **value iteration** עם $\epsilon \rightarrow 0$ וקיבלתם את הפלט הבא:
 (משמעות הדבר ש- $\epsilon \rightarrow 0$ היא שתנאי העצירה קלים שבזמנה האינסוף בין ווקטורי התועלת הייתה אפסית,
 כלומר לאחר הריצה ערכי התועלת שהתקבלו מקיימים את משוואת בלמן).



כאשר v_i הוא ערך התועלת למצב ה- i כפי שניתן לראות בתרשים. בנוסף נסמן את התגמול למצב ה- i ב- r_i .
 ענו נכון \ לא נכון, וספקו הסבר קצר או דוגמה נגדית מפורטת.

א. (3 נק') אם $v_9 > 1$, אז בהכרח מתקיים ש- $r_9 > 1$. נכון \ לא נכון.

דוגמה נגדית:

נגדיר $\gamma = 0.9$ ובביט בלוח התגמולים:

-0.04	-0.04	0.7	+1
-0.04	WALL	0.9	-1
-0.04	-0.04	-0.04	-0.04

נריץ את אלגוריתם value_iteration עם אפסילון ששוואף ל-0 ונקבל לוח התועלות:

4.881	5.695	6.542	1
4.237	WALL	7.277	-1
4.573	5.305	6.098	4.682

והמדיניות:

RIGHT	RIGHT	DOWN	1
UP	WALL	LEFT	-1
RIGHT	RIGHT	UP	LEFT

נראה כי התועלות שמתקבלות מקיימות את משוואת בלמן:

$$v_1 = -0.04 + 0.9$$

$$\begin{aligned} & * \max(0.8v_2 + 0.1v_1 + 0.1v_3, 0.8v_1 + 0.1v_1 + 0.1v_3, \\ & 0.8v_3 + 0.1v_1 + 0.1v_2, 0.8v_1 + 0.1v_1 + 0.1v_2) = 4.573 \end{aligned}$$

$$v_2 = -0.04 + 0.9$$

$$\begin{aligned} & * \max(0.8v_4 + 0.1v_2 + 0.1v_2, 0.8v_1 + 0.1v_2 + 0.1v_2, \\ & 0.8v_2 + 0.1v_1 + 0.1v_4, 0.8v_2 + 0.1v_1 + 0.1v_4) = 4.237 \end{aligned}$$

$$v_3 = -0.04 + 0.9$$

$$\begin{aligned} & * \max(0.8v_3 + 0.1v_1 + 0.1v_5, 0.8v_3 + 0.1v_1 + 0.1v_5, \\ & 0.8v_1 + 0.1v_3 + 0.1v_3, 0.8v_5 + 0.1v_3 + 0.1v_3) = 5.305 \end{aligned}$$

$$v_4 = -0.04 + 0.9$$

$$\begin{aligned} & * \max(0.8v_4 + 0.1v_4 + 0.1v_6, 0.8v_2 + 0.1v_4 + 0.1v_6, \\ & 0.8v_4 + 0.1v_4 + 0.1v_2, 0.8v_6 + 0.1v_4 + 0.1v_2) = 4.881 \end{aligned}$$

$$v_5 = -0.04 + 0.9$$

$$\begin{aligned} & * \max(0.8v_7 + 0.1v_3 + 0.1v_8, 0.8v_5 + 0.1v_3 + 0.1v_8, \\ & 0.8v_3 + 0.1v_5 + 0.1v_7, 0.8v_8 + 0.1v_5 + 0.1v_7) = 6.098 \end{aligned}$$

$$v_6 = -0.04 + 0.9$$

$$\begin{aligned} & * \max(0.8v_6 + 0.1v_4 + 0.1v_9, 0.8v_6 + 0.1v_4 + 0.1v_9, \\ & 0.8v_4 + 0.1v_6 + 0.1v_6, 0.8v_9 + 0.1v_6 + 0.1v_6) = 5.695 \end{aligned}$$

$$\begin{aligned} v_7 = 0.9 + 0.9 * \max(0.8v_9 + 0.1 * -1 + 0.1v_7, 0.8v_5 + 0.1v_7 + 0.1 * -1, \\ 0.8v_7 + 0.1v_9 + 0.1v_5, 0.8 * -1 + 0.1v_9 + 0.1v_5) = 7.277 \end{aligned}$$

$$\begin{aligned} v_9 = 0.7 + 0.9 * \max(0.8v_9 + 0.1 * 1 + 0.1v_6, 0.8v_7 + 0.1v_6 + 0.1 * 1, \\ 0.8v_6 + 0.1v_9 + 0.1v_7, 0.8 * 1 + 0.1v_9 + 0.1v_7) = 6.542 \end{aligned}$$

מתקיים כי משוואת בלמן אכן מתקיימת עבור 9 המצבים אבל $v_9 = 0.7$ קטן מ-1 אע"פ ש- $v_9 = 6.542$ שזה גדול מ-1.

ב. (3 נק') אם $v_i > 0, \forall i \in [9]$, אז בהכרח $r_i > 0, \exists i \in [9]$. נכון \(\square\) לא נכון.

דוגמה נגדית:

לא נכון, עבור אותה דוגמה מהסעיף הקודם נראה כי $\forall v_i \in [1,9]: v_i > 0$ אבל לא מתקיים שכל התגמולים חיוביים.

ג. (3 נק') אם $r_1 = r_2 = \dots = r_9 < 0$, אז בהכרח $v_1 = \min\{v_i | i \in [9]\}$. נכון \ לא נכון.

דוגמה נגדית:

נשתמש באותה דוגמה שראינו בתרגול:

טבלת התגמולים:

-0.04	-0.04	-0.04	1
-0.04	WALL	-0.04	-1
-0.04	-0.04	-0.04	-0.04

עבורה קיבלנו טבלת התועלות:

0.812	0.868	0.918	1
0.762	WALL	0.660	-1
0.705	0.655	0.611	0.388

קיבלנו כי $\min\{v_i | i \in [1,9]\} = 0.388$ אבל מתקיים כי $v_1 = 0.705 > 0.388$ לכן הטענה לא נכונה.

ד. (3 נק') אם $v_1 > v_2 > v_3 > 0$, אז בהכרח $U^p((1,1)) = \pi^*$. נכון \ לא נכון.

דוגמה נגדית:

ניקח את הערכים הבאים: $v_1 = 3, v_2 = 2, v_3 = 1$, נראה כי בכדי שמשוואת בלמן תתקיים צריך לעשות את הפעולה LEFT :

$$v_1 = r_1 + \gamma * \max(0.8v_2 + 0.1v_1 + 0.1v_3, 0.8v_1 + 0.1v_1 + 0.1v_3, 0.8v_3 + 0.1v_1 + 0.1v_2, 0.8v_1 + 0.1v_1 + 0.1v_2)$$

כלומר הפעולה שתבחר במדיניות האופטימלית היא זאת שמקיימת:

$$\max(0.8v_2 + 0.1v_1 + 0.1v_3, 0.8v_1 + 0.1v_1 + 0.1v_3, 0.8v_3 + 0.1v_1 + 0.1v_2, 0.8v_1 + 0.1v_1 + 0.1v_2)$$

נמצא איזה פעולה מהווה את המקסימים:

$$up \rightarrow 0.8v_2 + 0.1v_1 + 0.1v_3 = 0.8 * 2 + 0.1 * 3 + 0.1 * 1 = 2$$

$$\text{down} \rightarrow 0.8v_1 + 0.1v_1 + 0.1v_3 = 2.8$$

$$\text{left} \rightarrow 0.8v_1 + 0.1v_1 + 0.1v_2 = 2.9$$

$$\text{right} \rightarrow 0.8v_3 + 0.1v_1 + 0.1v_2 = 1.3$$

כלומר הפעולה שממקסמת את v_1 היא LEFT ולא UP.

ה. (2 נק') אם $\gamma = 0$, מה מספר המדיניות האופטימליות הקיימות? נמקו.

אם $\gamma = 0$ זאת אומרת שהביטוי שנקבל עבור התועלת הוא:

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) * U(s') = R(s)$$

כלומר עבור כל מצב לא משנה לנו מהי התועלת העתידית אלא רק מה שמרוויחים מיידי, ולכן בהתאם לתועלת הבאה המקסימלית נבחר הפעולה הבאה, ישנם 12 מצב בלוח שלנו, 2 מהם מצבי מטה ואחד קיר שאי אפשר לגשת אליו, נשארים 9 מצבים.

מכל מצב ניתן לבצע 4 פעולות (למעלה, למטה, ימינה או שמאלה) ולפי immediate reward של כל אחד מהמצבים העוקבים לאחר הפעלת הפעולות הנ"ל, נקבל את הפעולה האופטימלית, כלומר אם:

immediate reward of next state if we choose to go up > immediate reward of next state if we choose to go left

אז נעדיף פעולת up על left, והפעולה האופטימלית תהיה זאת בעלת הreward המקסימלית, ולכן מכל מצב אפשר לעשות 4 פעולות שכל אחת מהן יכולה להיות אופטימלית בהתאם לreward שלה.



סה"כ נקבל כי מספר המדיניות האופטימליות הוא 9^4 .

ו. (2 נק') לסעיף זה בלבד נתון כי $\mathbf{r}_8 = \mathbf{0}$, $\mathbf{v}_5 = -1$. מהו $\pi^*((1, 4))$? ציינו את כל האפשרויות ונמקו.

$$\pi(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a) * U(s')$$

נמצא עבור כל אחד מהפעולות $a \in A(s)$ את $a_{ut} = \sum_{s'} P(s'|s, a) * U(s')$

$$a = \text{up} \rightarrow \text{up}_{ut} = 0.8 * -1 + 0.1 * v_5 + 0.1 * v_8 = -0.9 + 0.1v_8$$

$$a = \text{down} \rightarrow \text{down}_{ut} = 0.8 * v_8 + 0.1 * v_5 + 0.1 * v_8 = -0.1 + 0.9v_8$$

$$a = \text{left} \rightarrow \text{left}_{ut} = 0.8 * v_5 + 0.1 * -1 + 0.1 * v_8 = -0.9 + 0.1v_8$$

$$a = \text{right} \rightarrow \text{right}_{ut} = 0.8 * v_8 + 0.1 * -1 + 0.1 * v_8 = -0.1 + 0.9v_8$$

$$\rightarrow \quad \text{up}_{ut} = \text{left}_{ut} \quad , \quad \text{down}_{ut} = \text{right}_{ut}$$

עבור $v_8 < -1$ הצעד האופטימלי יהיה ללכת ימינה או למטה, אחרת הצעד האופטימלי יהיה ללכת שמאלה או למעלה.

נראה עכשיו לפי משוואת בלמן כי אי אפשר שיתקיים $v_8 < -1$:

$$\begin{aligned} v_8 &= r_8 + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) * U(s') \\ &= r_8 + \gamma * \max(-0.9 + 0.1v_8, -0.1 + 0.9v_8) \\ &\geq r_8 + \gamma * (-0.9 + 0.1v_8) = 0 - 0.9\gamma + 0.1\gamma v_8 \end{aligned}$$

עבור $v_8 \leq -1$ נקבל:

$$v_8 \geq -0.9\gamma + 0.1\gamma v_8 \geq -0.9\gamma - 0.1\gamma$$

$$v_8 \leq -1 \rightarrow -1 \geq -0.9\gamma - 0.1\gamma \rightarrow \gamma \geq 1$$

בסתירה לנתון ש- $\gamma \in (0,1)$, ולכן נקבל כי הפעולות האופטימליות הן רק right, down עבור מצב (1,4).

ז. (2 נק') נתון כי $v_1 > v_2 > v_3 > 0$, מצאו חסמים צמודים, עליון ותחתון ל- r_1 כפונקציה של v_i (ולא כפונקציה של γ).

$$v_1 = r_1 + \gamma * \max(0.8v_2 + 0.1v_1 + 0.1v_3, 0.8v_1 + 0.1v_1 + 0.1v_3, 0.8v_3 + 0.1v_1 + 0.1v_2, 0.8v_1 + 0.1v_1 + 0.1v_2)$$

נתון כי $v_1 > v_2 > v_3$ לכן :

$$\max(0.8v_2 + 0.1v_1 + 0.1v_3, 0.8v_1 + 0.1v_1 + 0.1v_3, 0.8v_3 + 0.1v_1 + 0.1v_2, 0.8v_1 + 0.1v_1 + 0.1v_2) = 0.9v_1 + 0.1v_2$$

ולכן נקבל:

$$v_1 = r_1 + \gamma(0.9v_1 + 0.1v_2)$$

$$r_1 = v_1 - \gamma(0.9v_1 + 0.1v_2)$$

נתון כי $\gamma \in (0,1)$ ולכן נקבל:

$$r_{1\max} = v_1 - \gamma_{\min}(0.9v_1 + 0.1v_2) < v_1 - 0 = v_1$$

$$r_{1\min} = v_1 - \gamma_{\max}(0.9v_1 + 0.1v_2) > v_1 - 1 * (0.9v_1 + 0.1v_2) = 0.1(v_1 - v_2)$$

כלומר:

$$0.1(v_1 - v_2) < r_1 < v_1$$

חלק למידה

- א. ניקח את הנקודות הבאות : נקודות האימון $\{(9,1), (7,5)\}$, נק 2 $(-7,5)$ $(K=1, d=2)$ דוגמת מבחן : $(1,1)$
כך שבמרחק אוקלידי נקודה 2 היא הקרובה ולפיכך התוצאה תהיה - , ובמרחק מנהטן נקודה 1 היא הקרובה כלומר התוצאה תהיה +
- ב. אם נבחר $K=1$ נקבל ששגיאת האימון תהיה 0 כי הנקודה הנוכחית תהיה הקרובה ביותר וכך שכל נקודה אחרת המרחק המינימלי מהנקודה הנוכחית יהיה $\sqrt{2}$
- ג. אם נבחר $K=14$ (כלומר במספר נקודות סט האימון) אנחנו תמיד ניקח את הסיווג הנפוץ ביותר מקבוצת האימון
- ד. כפי שנלמד בתרגול: אם K יהיה קטן מדי יהיה לנו מאוד קשה להתמודד עם רעש (overfitting) ואם K גדול מדי אנחנו נתחשב גם בשכנים שלא באמת קרובים
- ה. הפרכה :
- נניח שיש לנו N נקודות של סט אימון, יהי d and k כלשהם ונניח נקודת מבחן (x_l, y_l) כך ש x_l שייך למרחב R^D
ונניח שאין רדיוס שייביא את אותו הסיווג כמו באלגו המקורי , כלומר בגרסה החדשה האלגו לא יצליח להביא את הסיווג של האלגוריתם המקורי, כלומר מספר השכנים שלנו יביא שיהיה לנו תיקו בין שתי התוצאות, ולפי נתוני השאלה, דבר זה מחייב שהסיווג יהיה חיובי, לפי נתוני השאלה כך שאין שתי נקודות בעלות אותו מרחק מנקודה המבחן אז יהיה לנו רדיוס מסויים כך שנפריד בין הנקודה ה- K לנקודה שאחריה, ולפיכך סתירה !
- ו. הפרכה:
- נניח שיש לנו N נקודות של סט אימון, יהי d and r כלשהם ונניח נקודת מבחן (x_l, y_l) כך ש x_l שייך למרחב R^D
נניח שאין k שניב את אותה תוצאה כמו האלגוריתם החדש כלומר אין מספר שכנים קבוע שניב את אותה תוצאה כמו האלגוריתם החדש, אבל עם הנחת השאלה נקבל סתירה מידית, כי אם אין שתי נקודות בעלות אותו מרחק ויש רדיוס קבוע שמניב תוצאה מסויימת אז בהכרח יש מספר מסויים וקבוע של שכנים שגם יניבו את אותה תוצאה מקורית) בגלל שבאלגוריתם החדש אנחנו לוקחים גם את הקרובים ביותר על הרדיוס הנון, אז אותם שכנים באלגו המקורי יילקחו)

מתפצלים ונהנים :

$(\epsilon' = -10\epsilon)$ $(d=1)$ אם ניקח את ערך האפסילון הנל מתקיים שהתנאי אף פעם לא מתקיים בעץ T כך שהעץ עם ערך x תמיד יונב כ $true$ בעץ T , אבל בעץ T' תמיד נקבל ערך $false$, סתירה

$\epsilon = -10\epsilon$

T

X 70

False

True

• true

15 false

' True

0 false

T'

1 True

15 false

Test Accuracy: 94.69%

2b

3a המטרה היא להחליש את התאמת היתר על ידי ולעשות עצים לא עקביים במטרה להתמודד עם הרעש של הנתונים

3d כן שיפר

Test Accuracy: 97.35%