

שאלה מס' 1

א.

$$U^\pi(s) = E_\pi \left[\sum_{t=0}^{\infty} \gamma^t R(S_t, \pi(S_t), S_{t+1}) \mid S_0 = s \right]$$

ב.

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

ג. עבור המקרה בו $\gamma = 1$ נקבל ש- $\frac{1-\gamma}{\gamma} = 0$ כלומר האלגוריתם לא יעצור לעולם כי תמיד $\delta \geq 0$

Repeat

$U \leftarrow U'; \delta \leftarrow 0$

For each state s in S do:

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

if $|U'[s] - U[s]| > \delta$ then $\delta \leftarrow |U'[s] - U[s]|$

Until $\delta < \epsilon (1 - \gamma)/\gamma$

Return U

ד. עבור המקרה בו $\gamma = 1$ ואופק אינסופי האלגו' לא יתכנס, אמנדרוש שהאופק יהיה סופי אזי שנקבל תועלת מקסימלית והאלגו יסתיים ויחזיר מדיניות אופטימלית.

Repeat

$U \leftarrow \text{POLICY-EVALUATION}(\pi, U, mdp)$

Unchanged? \leftarrow true

For each state s in S do:

$$\text{if } \max_{a \in A(s)} \left[\sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma U(s')] \right] > \sum_{s'} P(s'|s, \pi[s]) [R(s, \pi[s], s') + \gamma U(s')]$$

$$\text{then do } \pi[s] \leftarrow \operatorname{argmax}_{a \in A(s)} \left[\sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma U(s')] \right]$$

Unchanged? \leftarrow true

Until Unchanged?

Return π

שאלה 2:

חלק א' – MDP ו-RL

שאלה 2

נתונים שני אנשים – "סוחט" ו-"קורבן". בכל שלב ה"סוחט" יכול:

- (1) "לפרוש" – לפרוש עם רווחי הסחיטה.
- (2) "לסחוט" – לדרוש תשלום של 1. בהסתברות p , ה"קורבן" יענה לדרישה. ובהסתברות $1 - p$ ה"קורבן" יסרב לשלם וידווח למשטרה.

הנחות:

- לאחר שה"קורבן" דווח למשטרה, ה"סוחט" מאבד את כל הרווחים שנצברו ואינו יכול לסחוט שוב.
- לאחר שה"סוחט" מגיע לרווחים מצטברים של n , הוא פורש מיד.
- מטרת הסוחט היא למקסם את סכום הכסף שהוא מרוויח.
- אופק סופי, ניתן להניח שגדול מאוד $\gamma = 1$.

1. (4 נק') נסחו את הבעיה כבעיית MDP עם המצבים $i = 0, 1, \dots, n$ ומצב סיום T . (0 הוא מצב התחלתי) באופן ספציפי, כתבו את המצבים, הפעולות בכל מצב, ההסתברויות המעבר והתגמולים.

הערה: התגמולים חייבים להיות אי-שליליים.

2. (2 נק') האם ניתן לנסח את הבעיה כבעיית MDP עם מצב יחיד ומצב סיום? נמקו.

3. (2 נק') האם ניתן לנסח את הבעיה כבעיית MDP כאשר חלק מהתגמולים שליליים? נמקו.

4. (9 נק') נתון כי $n=3$.

כעת נרצה למצוא מדיניות אופטימליות ומה התועלת של המצב ההתחלתי כפונקציה של p .

בתשובתכם מצאו עבור אילו ערכי p מקבל כל מדיניות – מצאו את a ו- b כך שהמדיניות בטווח הנתון לא

תשתנה. מלאו את הערכים החסרים בטבלה שבעמוד הבא ונמקו היטב את תשובתכם.

הערה: כאשר המדיניות של מצב i יכולה לקבל יותר מפעולה אחת יש לציין את כל הפעולות.

נגדיר MDP בצורה הבא : $\langle S, A, P, R, \gamma \rangle$

קבוצת המצבים S : $S = \{0, 1, \dots, n\} \cup \{T\}$

Set of actions A:

$$\forall 0 \leq i \leq n-1: A(i) = \{\text{לסחות, לפרוש}\}, A(n) = \{\text{לפרוש}\}, A(T) = \emptyset$$

Rewards R: נבחר בתגמול על הקשתות $R(s, a, s')$

הסבר הבחירה שלי :

שאנחנו נמצאים במצב i נוכל לבחור לסחות או לפרוש והבחרה מכן תגמול על המצווים לא יעבוד, עבור הבחרה ב לסחות קיים שני מצביים להצליח ולעבור למצב $i+1$ ו לא להצליח ולעבר ל T ולכן בחרנו ב תגמול על הקשתות .

יהי מודר בצורה הבא :

- $\forall 0 \leq i \leq n: R(i, \text{לפרוש}, T) = i$
- $\forall 0 \leq i \leq n-1: R(i, \text{לסחות}, i+1) = 0$
- $\forall 0 \leq i \leq n-1: R(i, \text{לסחות}, T) = 0$

Transitions P: $P(s'|s, a)$

- $\forall 0 \leq i \leq n: P(T|i, \text{לפרוש}) = 1$
- $\forall 0 \leq i \leq n: P(T|i, \text{לסחות}) = 1 - p$
- $\forall 0 \leq i \leq n: P(i+1|i, \text{לסחות}) = p$

Discount factor: $\gamma = 1$

start state: $s_{init} = 0$

terminal state: $s_{final} = T$

2. (2 נק') האם ניתן לנסח את הבעיה בבעיית MDP עם מצב יחיד ומצב סיום? נמקו.

לא ניתן

נפקע בעקרון Markov

Action outcomes depend only on the current state.

כי עבור עובר 2 מצבים לא נוכל לדעת עבור הצעד ה- i את התגמול אז נצטרך לחשיב התקמול עד הצעד הנוכי אז עבורו בעולת לסחוט לא מוצלחית נצטרך לתקמול אי שליל לא קבוע והתקמולים הם קבועים עבור כל מצב ולא יכול להיות תלוי בהסטוריה ולכן לא ניתן.

3. (2 נק') האם ניתן לנסח את הבעיה בבעיית MDP כאשר חלק מהתגמולים שלילים? נמקו.

בדומה לסעיף 1 אבל שעוברים ממצב i ל $i+1$ נוסף אחד עבור פעולת לסחות מוצלחת

עבור פעולת לסחות לא מוצלחת ממצב i נוריד i ונעבור ל T

ועבור פרישה נוסף 0 .

4. (9 נק') נתון כי $n=3$.

כעת נרצה למצוא מדיניות אופטימליות ומה התועלת של המצב ההתחלתי כפונקציה של p .
בתשובתכם מצאו עבור אילו ערכי p נקבל כל מדיניות – מצאו את a ו- b כך שהמדיניות בטווח הנתון לא תשתנה. מלאו את הערכים החסרים בטבלה שבעמוד הבא ונמקו היטב את תשובתכם.
הערה: כאשר המדיניות של מצב i יכולה לקבל יותר מפעולה אחת יש לציין את כל הפעולות.

מאופן ההגדרה $A(3)=$ לפרוש, עבור מדוניות אופטימלית תמיד קדי לנו בהתחלה לסחות לכן $A(0) =$ לסחוט ולכן קיים רק שלוש בסלסות אפשריות וכול ש p יותר גדול יודר לנו יותר לסחוט, יכולים לרות בטבלה שמיליתי את שלושת הבלת האפשריות וחישוב התועלת בכול בולסיה נשאר לחשיב את a ו b לפי מה הספרנו קודם עבור b יהי יותר קידי לשתמש בה כול שי עירך p יותר גדול אז עבור תועלת יותר גדולה יותר קדי לשמשם בה מלשתמש מ a זה מתקיים קאשר $2p^2 > p$ מכן נקביל כי $a=0.5$ באופן דומה עבור חישוב b

$$b=2/3 \text{ מכן נקביל כי } 3p^3 > 2p^2$$

$$a = \underline{0.5}$$

$$b = \underline{\frac{2}{3}}$$

ערכי p	מדיניות	תועלות
$0 < p < a$	$\pi_1(0) = \underline{\text{לסמוך}}$ $\pi_1(1) = \underline{\text{לברוח}}$ $\pi_1(2) = \underline{\text{לברוח}}$ $\pi_1(3) = \underline{\text{לברוח}}$	$V^{\pi_1}(0) = \underline{p}$ $(1-p) \cdot 0 + p \cdot 1 = p$
$a < p < b$	$\pi_2(0) = \underline{\text{לסמוך}}$ $\pi_2(1) = \underline{\text{לסמוך}}$ $\pi_2(2) = \underline{\text{לברוח}}$ $\pi_2(3) = \underline{\text{לברוח}}$	$V^{\pi_2}(0) = \underline{2p^2}$ $(1-p) \cdot 0 + p \cdot ((1-p) \cdot 0 + p \cdot 2) = 2p^2$
$b < p < 1$	$\pi_3(0) = \underline{\text{לסמוך}}$ $\pi_3(1) = \underline{\text{לסמוך}}$ $\pi_3(2) = \underline{\text{לסמוך}}$ $\pi_3(3) = \underline{\text{לברוח}}$	$V^{\pi_3}(0) = \underline{3p^3}$ $(1-p) \cdot 0 + p \cdot ((1-p) \cdot 0 + p \cdot ((1-p) \cdot 0 + p \cdot 3)) = 3p^3$

היכן חסמה: יש

$$= p[p(p(3))] = 3p^3$$

חלק ג:

- Done
- Done
- Done
- Done
- Done
- גרסת anytime:

```
1  function anytime:
2  k=1
3  episodes=1
4  lastResult=None
5  while time > 0
6      result=adp_algorithm(episodes)
7      if (time < 0)
8          return lastResult
9
10     result = lastResult
11     episodes *= 10
12 return lastResult
```

חלק ב' – מבוא ללמידה

חלק א' – יבש

א.1.

עבור $d=1$ נקבל שאין תלות בבחירת פונק' המרחק עבור על ערך K כי שתיהן ייתנו אותו ערך:

$$\text{Euclid}(x, y) = \sqrt{(x_1 - y_1)^2} = |x_1 - y_1|$$

$$\text{Manhattan}(x, y) = |x_1 - y_1|$$

א.2.

עבור בעיית קלסיפיקציה נבחר $d=2, k=1$, נקבע את הסיווג של הדוג' לפי סיווג השכן הקרוב ביותר מאוסף דוגמאות האימון.

הסט של האימון יהיה: $D = \{d_1 = <(0,0), ->, d_2 = <(1,1.25), +>\}$

הדוג' תהיה: $t_1 = (1,1)$

לפי האוקלידי: נקבל שהסיווג של t_1 הוא $-$ כי d_1 הוא השכן הקרוב ביותר לכן סיווגו יהיה כמוהו:

$$\text{Euclid}(d_1, t_1) = \sqrt{(0-1)^2 + (0-1)^2} = \sqrt{2} \cong 1.4$$

$$\text{Euclid}(d_2, t_1) = \sqrt{(1-1)^2 + (2.5-1)^2} = \sqrt{(1.5)^2} \cong 1.5$$

לפי מנהטן: נקבל שהסיווג של t_1 הוא $+$ כי d_2 הוא השכן הקרוב ביותר לכן סיווגו יהיה כמוהו:

$$\text{Manhattan}(d_1, t_1) = |0-1| + |0-1| = 2$$

$$\text{Manhattan}(d_2, t_1) = |1-1| + |2.5-1| = 1.5$$

א.3.

עבור $k=5$ ו- $k=7$ הדיוק מרבי על קב' האימון \leftarrow 10 מתוך 14 סיווגים נכונים.

א.4.

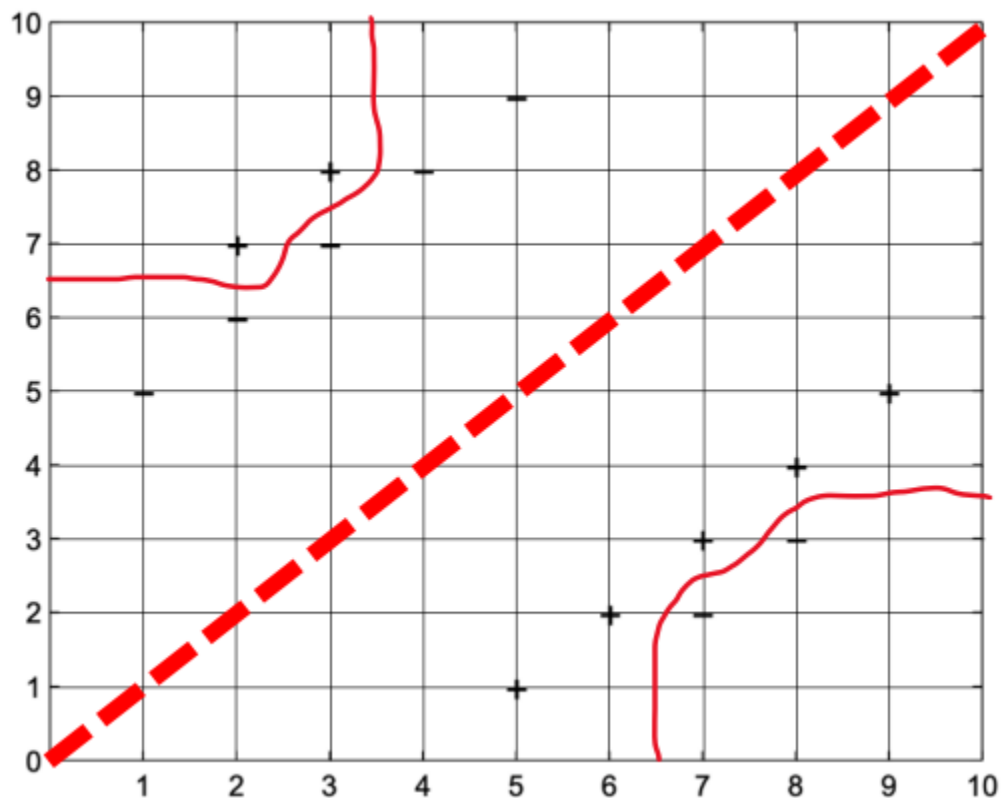
עבור $k=14$ הסיווג של קב' האימון יהיה majority כי כאשר מסווגים דוג' K השכנים הקרובים ביותר שלה הם כל דוגמאות האימון הקיימות ו- 14 הוא בעצם מס' דוגמאות האיון הקיימות. כלומר, הסיווג הוא הסיווג שיש לרוב דוגמאות האימון; majority

5) (2 נק') נמקו מדוע שימוש בערכי k גדולים או קטנים מדי יכול להיות גרוע עבור קבוצת הדגימות הנ"ל.

תשובה: עבור ערכי k שקטנים מדי דוגמה יכולה להיספג בצורה שגויה בגלל - overfitting המודל יהיה ממש רגיש וקשה להתמודד עם רעש, שזה עלול לפגוע ברמת הדיוק המירבי. לדוגמה עבור $k=3$ לפי נקודות המבחן הנתונות אם נתבונן בדוגמת מבחן שנמצאת קרוב לקצוות של אחד הצדדים (למשל בנקודה (6,1) (תסווג בהתאם לערך השגוי שלהן. עבור ערכי k שהם גדולים מדי המודל יתחשב בדוגמאות מבחן שהן לא ממש רלוונטיות לדוגמת המבחן וזה יפגע ברמת הדיוק. למשל עבור k גדול נקודת מבחן אם תיוג שלילי שנופלת באיזור הימני התחתון של הגרף תקבל סיווג חיובי בגלל כמות הפלוסים הגדולה יחסית שנמצאת סביבה, באותו אופן עבור נקודת מבחן עם ציוג חיובי שהסיווג לה נמצא באישור העליון השמאלי, היא תקבל סיווג שלילי מכיוון שיש הרבה יותר נקודות מבחן עם סיווג שלילי מאשר חיובי. (לעומת המקרה בו היינו בוחרים k קטן יותר כמו 4 למשל שהיה נותן סיווג תואם עבור מקרים אלה) .

6) (2 נק') שרטט את גבול ההחלטה של 1-nearest neighbor עבור הגרף

תשובה:



גבול ההחלטה מפריד בין איזורים בהם מסווגים דוגמאות בסימונים שונים, הסיווג לפי נקודה קרובה ביותר.