

סהב	7	6	5	4	3	2	1

## מבחן מועד א' – למידה מחיזוקים סמסטר ב' תשע"ח (2018/9)

בית הספר למדעי המחשב, אוניברסיטת תל-אביב

מרצה: פרופ' ישי מנצור

מתרגלת: גב' לי כהן

9.7.2019

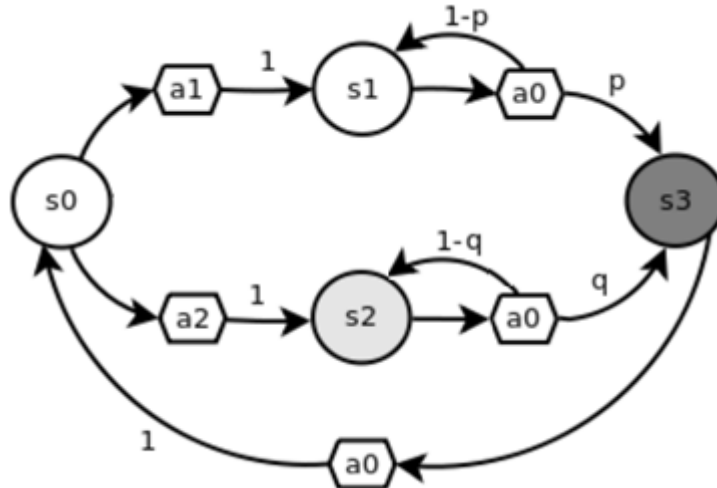
### הוראות

1. מומלץ לקרוא את כל ההנחיות והשאלות בתחילת המבחן, לפני תחילת כתיבת התשובות.
2. משך הבחינה – **שלוש שעות**. לא תינתן כל הארכה נוספת.
3. חומר עזר מותר: דף עזר אחד בגודל A4 דו-צדדי, **ומחשבוני**.
4. **יש לענות על השאלות במקום המיועד לכך בטופס השאלון (טופס זה)**. מחברות הבחינה לא ייקראו, וישמשו כטיטות בלבד.
5. יש למלא בכל דף של השאלון מספר ת.ז. ומספר מחברת.
6. במבחן 4 שאלות:
  - הניקוד לכל שאלה מופיע לידה מספר השאלה.
  - יש לענות תשובות ברורות ענייניות ותמציתיות.
7. מותר להשתמש בכל טענה שהוכחה בכיתה (בהרצאה, בתרגול, או בתרגיל בית) בתנאי שמצטטים אותה במדויק. טענות אחרות (כאלה שהוכחו בספרים, בהרצאות מהסמסטר הקודם, וכו') יש להוכיח.

בהצלחה !

**שאלה 1 (30 נקודות)**

נתון MDP שמוגדר ע"י הגרף שבציור, כאשר המצבים מסומנים במעגלים והמצב ההתחלתי הוא  $s_0$ . ישנן שלוש פעולות המתוארות באמצעות משושים. המספרים מעל החצים המחברים את הפעולות למצבים מתארים את ההסתברות לעבור מהמצב הנוכחי באמצעות הפעולה הזו למצב הבא. לדוגמא, ההסתברות לעבור למצב  $s_3$  ממצב  $s_2$  באמצעות  $a_0$  הינה  $q$ . במידה ואין חץ בין פעולה למצב, ההסתברות לעבור למצב הזה הינה 0. הרווח המיידני ממצב  $s_3$  הוא 10, ממצב  $s_2$  הוא 1 ומהשאר. ההחזר (return) הוא discounted עם פרמטר  $\gamma$ .



א. כתוב/כתבי את ה-MDP של הבעיה באופן פורמלי.

$S = \{s_0, s_1, s_2, s_3\}$   
 $A = \{a_0, a_1, a_2\}$   
 $s_0 = s_0$   
 $R: R(s_3, a_0) = 10, R(s_2, a_0) = 1, \forall a \in A, s \notin \{s_2, s_3\} R(s, a) = 0$   
 $P: P(s_1|s_0, a_1) = 1, P(s_1|s_1, a_0) = 1 - p, P(s_3|s_1, a_0) = p, P(s_2|s_0, a_2) = 1,$   
 $P(s_0|s_2, a_0) = 1 - q, P(s_3|s_2, a_0) = q, P(s_0|s_3, a_0) = 1, \text{rest are } 0$

ב. כתוב/כתבי את משוואות בלמן לערך האופטימלי לכל מצב ( $V^*(s)$ ) במונחים של  $\gamma, p, q$  ו- $V^*(s')$ .

$$V^*(s_0) = \max_{a \in \{a_1, a_2\}} 0 + \gamma \sum_{s'} P(s'|s, a) V^*(s') = \gamma \max(V^*(s_1), V^*(s_2))$$

$$V^*(s_1) = \max_{a \in \{a_0\}} 0 + \gamma \sum_{s'} P(s'|s, a) V^*(s') = \gamma((1 - p)V^*(s_1) + pV^*(s_3))$$

$$V^*(s_2) = \max_{a \in \{a_0\}} 1 + \gamma \sum_{s'} P(s'|s, a) V^*(s') = 1 + \gamma((1 - q)V^*(s_2) + qV^*(s_3))$$

$$V^*(s_3) = \max_{a \in \{a_0\}} 10 + \gamma \sum_{s'} P(s'|s, a) V^*(s') = 10 + \gamma V^*(s_0)$$

ג. מעתה נניח ש  $q > p$ . מתי מתקיים  $V^*(s_2) > V^*(s_1)$ ?  
 רמז: התחילו מלהביע את  $V^*(s_1), V^*(s_2)$  באמצעות  $\gamma, p, q$  ו- $V^*(s_3)$  בלבד.

תשובה: תמיד

נמק:

$$V^*(s_1) = \frac{\gamma p V^*(s_3)}{1 - \gamma(1 - p)} < \frac{1 + \gamma q V^*(s_3)}{1 - \gamma(1 - q)} = V^*(s_2)$$

$$\begin{aligned} \gamma p V^*(s_3)(1 - \gamma(1 - q)) &< (1 + \gamma q V^*(s_3))(1 - \gamma(1 - p)) \\ \gamma p V^*(s_3) - \gamma^2 p V^*(s_3)(1 - q) &< 1 - \gamma(1 - p) + \gamma q V^*(s_3) - \gamma^2 q V^*(s_3)(1 - p) \\ \gamma(1 - p) + \gamma^2 V^*(s_3)(q(1 - p) - p(1 - q)) &< 1 + (q - p)\gamma V^*(s_3) \\ \gamma(1 - p) + \gamma^2 V^*(s_3)(q - p) &< 1 + (q - p)\gamma V^*(s_3) \\ \gamma^2 V^*(s_3)(q - p) &< (q - p)\gamma V^*(s_3), \gamma(1 - p) < 1 \end{aligned}$$

וזה מתקיים כיוון ש  $q - p > 0$

ד. האם קיימים  $\gamma, q$  שעבורם  $\pi^*(s_0) = a_1$  (זכרו ש  $q > p$ ).

תשובה: כן/לא

נמק: לא, מכיוון שעל מנת ש  $\pi^*(s_0) = a_1$ ,

צריך שיתקיים  $V^*(s_1) \geq V^*(s_2)$  והראנו ב-(ג) שזה לא מתקיים.

ה. חשבו/י לכל מצב את  $V^*(s)$  (כפונקציה של  $(p, q, \gamma, V^*(s_3))$ )

$$V^*(s_0) = \gamma V^*(s_2) = \gamma \frac{1 + \gamma q V^*(s_3)}{1 - \gamma(1 - q)}$$

$$V^*(s_1) = \frac{\gamma p V^*(s_3)}{1 - \gamma(1 - p)}$$

$$V^*(s_2) = \frac{1 + \gamma q V^*(s_3)}{1 - \gamma(1 - q)}$$

$$V^*(s_3) = 10 + \gamma V^*(s_0) = 10 + \gamma^2 \frac{1 + \gamma q V^*(s_3)}{1 - \gamma(1 - q)}$$

פתרו את  $V^*(s_3)$  (כלומר הבע אותו באמצעות  $p, q, \gamma$  בלבד)

$$V^*(s_3) = 10 + \gamma V^*(s_0) = 10 + \gamma^2 V^*(s_2)$$

$$V^*(s_3) = 10 + \gamma^2 \frac{1 + \gamma q V^*(s_3)}{1 - \gamma(1 - q)}$$

$$V^*(s_3)(1 - \gamma(1 - q)) = 10(1 - \gamma(1 - q)) + \gamma^2 + \gamma^3 q V^*(s_3)$$

$$V^*(s_3) = \frac{10(1 - \gamma(1 - q)) + \gamma^2}{1 - \gamma(1 - q) - \gamma^3 q}$$

**שאלה 2 (30 נקודות)**

בשאלה זאת נדון במודל ה-MAB. במודל ישנם  $n = 2k + 1$  פעולות. הרווח של פעולה  $a_i$  נדגם ממשתנה ברנולי עם פרמטר  $q_i$  (בהסתברות  $q_i$  זה 1, ואחרת 0) כל הפרמטרים  $q_i$  לא ידועים שונים, וכן  $|q_i - q_j| \geq \Delta$  לכל שתי פעולות  $a_i \neq a_j$ . פעולה  $a_i$  נקראת **החציון** אם ישנן  $k$  פעולות  $a_j$  עם  $q_j < q_i$  וכן  $k$  פעולות  $a_j$  עם  $q_j > q_i$ .

א. בהינתן דגימה של  $n_i$  רווחים מפעולה  $a_i$ , מתוכם  $s_i$  הם 1 (והשאר 0), הגדרי UCB ו-LCB כפונקציה של  $n_i, s_i, \delta$  עבור פעולה  $a_i$  כך שבהסתברות  $1 - \delta$  לפחות יתקיים  $q_i \in [LCB(n_i, s_i), UCB(n_i, s_i)]$ .

$$UCB(n_i, s_i) = \frac{s_i}{n_i} + \lambda$$

$$LCB(n_i, s_i) = \frac{s_i}{n_i} - \lambda$$

$$\text{Where } \lambda = \sqrt{\frac{2 \ln \frac{2}{\delta}}{n_i}}$$

מנקודה זאת והלאה, הניחו שה-UCB וה-LCB תמיד מתקיימים, וענה על הסעיפים הבאים

ב. כתוב/י תנאי מספיק לכך שפעולה  $a_i$  היא החציון (כתלות ב- $UCB(n_i, s_i)$  ו- $LCB(n_i, s_i)$  של כל הפעולות). הנח שתמיד מתקיים  $q_i \in [LCB(n_i, s_i), UCB(n_i, s_i)]$ .

קיימות  $k$  פעולות  $a_j$  שעבורן

$$UCB(n_i, s_i) < LCB(n_j, s_j)$$

ובנוסף קיימות  $k$  פעולות (אחרות)  $a_j$  שעבורן

$$LCB(n_i, s_i) > UCB(n_j, s_j)$$

ג. כתוב/י תנאי מספיק (טוב ביותר) לכך שפעולה  $a_i$  אינה החציון (כתלות ב- $UCB(n_i, s_i)$  ו- $LCB(n_i, s_i)$  של כל הפעולות). הנח שתמיד מתקיים  $q_i \in [LCB(n_i, s_i), UCB(n_i, s_i)]$ . אם יש מספר תנאים מספיקים, העדף זה שפוסל הכי הרבה.

קיימות  $k + 1$  פעולות  $a_j$  שעבורן

$$UCB(n_i, s_i) < LCB(n_j, s_j)$$

או שקיימות  $k + 1$  פעולות  $a_j$  שעבורן

$$LCB(n_i, s_i) > UCB(n_j, s_j)$$

ד. שנה/י את אלגוריתם ה- Successive Elimination כך שיחזיר את הפעולה שהיא חציון (האלגוריתם לא יודע את הפרמטר  $\Delta$ ). הדגש איזה שורה/ות משתנים, ואיך!

Algorithm Successive Elimination

1: Initially  $S = A$

2: While  $|S| \geq 2$ :

3: Try each action  $a_i \in S$  once, and update  $n_i$  and  $s_i$

4: For each  $a_i \in S$ , such that  $\exists a_j \in S: UCB(n_i, s_i) < LCB(n_j, s_j)$  THEN  $S \leftarrow S - \{a_i\}$

5: END WHILE

6: Return  $S$

שנה את שורה 4 באופן הבא:

For each  $a_i \in S$ , such that  $\exists a_{j_1}, \dots, a_{j_{k+1}} \in A: UCB(n_i, s_i) < LCB(n_{j_k}, s_{j_k})$   
THEN  $S \leftarrow S - \{a_i\}$

For each  $a_i \in S$ , such that  $\exists a_{j_1}, \dots, a_{j_{k+1}} \in A: LCB(n_i, s_i) > UCB(n_{j_k}, s_{j_k})$   
THEN  $S \leftarrow S - \{a_i\}$

**שאלה 3 (20 נקודות)**

התפלגות Weibull מוגדרת עם פרמטרים  $\lambda, 0 < k$ , כאשר פונקציית הצפיפות היא  $\frac{k}{\lambda} \left(\frac{x}{\lambda}\right)^{k-1} \exp\left(-\left(\frac{x}{\lambda}\right)^k\right)$  לכל  $x \geq 0$ .

משתמשים בהתפלגות הנ"ל על מנת לדגום פעולות באופן הבא:  
לכל מצב  $s$  יש וקטור  $\phi(s) \in \mathbb{R}^d$  שמאפיין אותו.

מדיניות  $\pi$  מאופיינת על ידי שני וקטורים  $\theta_\lambda, \theta_k \in \mathbb{R}^d$ .

במצב  $s$  דוגמים פעולה  $a \geq 0$  לפי התפלגות המוגדרת לעיל עם פרמטרים:

$$\lambda = \exp(\theta_\lambda^\top \phi(s)), k = \theta_k^\top \phi(s)$$

א. מה ההסתברות שדגום במצב  $s$  פעולה  $a$  כך ש  $a \in [0,1]$  כפונקציה של  $\theta_\lambda, \theta_k$  ו-  $\phi(s)$ ?

$$\int_0^1 \frac{k}{\lambda} \left(\frac{x}{\lambda}\right)^{k-1} \exp\left(-\left(\frac{x}{\lambda}\right)^k\right) dx = [-\exp\left(-\left(\frac{x}{\lambda}\right)^k\right)]_0^1 =$$

$$1 - \exp\left(-\left(\frac{1}{\lambda}\right)^k\right) = 1 - \exp\left(-\exp(-\theta_\lambda^\top \phi(s)) \cdot \theta_k^\top \phi(s)\right)$$

ב. חשבי את  $\nabla_\theta \log \pi(a|s; \theta)$ .

$$\log \pi(a|s; \theta) = \log \left( \frac{k}{\lambda} \left(\frac{a}{\lambda}\right)^{k-1} \exp\left(-\left(\frac{a}{\lambda}\right)^k\right) \right) =$$

$$= \log k - k \log \lambda + (k-1) \log a - \left(\frac{a}{\lambda}\right)^k$$

נשתמש ב:  $\frac{\partial}{\partial \theta_\lambda} \log \lambda = \phi(s), \frac{\partial \lambda}{\partial \theta_\lambda} = \phi(s) \lambda$

$$\nabla_{\theta_\lambda} \log \pi(a|s; \theta) =$$

$$-\phi(s) - (\theta_k^\top \phi(s)) \phi(s) + (\theta_k^\top \phi(s)) a^{\theta_k^\top \phi(s)} \exp\left(-\theta_\lambda^\top \phi(s) \cdot \theta_k^\top \phi(s)\right) \phi(s)$$

נשתמש ב:  $\frac{\partial}{\partial \theta_k} k = \phi(s)$

$$\nabla_{\theta_k} \log \pi(a|s; \theta) = \frac{\phi(s)}{\theta_k^\top \phi(s)} - \phi(s) \cdot \theta_\lambda^\top \phi(s) + \phi(s) \log a - (\log a - \theta_\lambda^\top \phi(s)) \left(\frac{a}{\exp(\theta_\lambda^\top \phi(s))}\right)^{\theta_k^\top \phi(s)} \phi(s)$$

תעודת זהות:

מספר מחברת:

ג. כתוב/י את העדכון של REINFORCE עבור ההתפלגות Weibull:

$$\begin{aligned}\theta_\lambda^{t+1} &= \theta_\lambda^t + \alpha G \phi(s) D_\lambda \\ \theta_k^{t+1} &= \theta_k^t + \alpha G \phi(s) D_k\end{aligned}$$

Where G is the return and

$$D_\lambda = -(\theta_k^\top \phi(s)) + (\theta_k^\top \phi(s)) a^{\theta_k^\top \phi(s)} \exp(-\theta_\lambda^\top \phi(s) \cdot \theta_k^\top \phi(s))$$

$$D_k = -\frac{1}{\theta_k^\top \phi(s)} + \log a - \theta_\lambda^\top \phi(s) - (\log a - \theta_\lambda^\top \phi(s)) \left( \frac{a}{\exp(\theta_\lambda^\top \phi(s))} \right)^{\theta_k^\top \phi(s)}$$



**שאלה 4 (20 נקודות)**

נתון MDP  $M$  המוגדר ע"י  $(S, A, p, s_0, R)$ , ומדיניות  $\pi$   
 לכל מצב  $s$  נגדיר וקטור יחידה  $e_s$  כך ש-  $e_s(s) = 1$  ו-  $e_s(s') = 0$   $\forall s' \neq s$   
 לכל פעולה  $a$  נגדיר מטריצה  $P^a[s_i, s_j] = p(s_i | s_j, a)$   
 למדיניות  $\pi$  נגדיר וקטור רווחים  $r[s] = r(s, \pi(s))$  ומטריצה סטוכסטית  $P^\pi[s_i, s_j] = p(s_i | s_j, \pi(s_j))$   
 מגדירים את האופרטורים הבאים עבור  $V \in \mathbb{R}^{|S|}$  ו-  $Q \in \mathbb{R}^{|S| \times |A|}$  בהתאמה:

$$(1) (R^{(3)}V)(s) = r^\top e_s + \gamma r^\top P^\pi e_s + \gamma^2 r^\top (P^\pi)^2 e_s + \gamma^3 V^\top (P^\pi)^3 e_s$$

$$(2) (G^{(3)}Q)(s, a) = r(s, a) + \gamma r^\top P^a e_s + \gamma^2 r^\top P^\pi P^a e_s + \max_a \gamma^3 Q(\cdot, a)^\top (P^\pi)^2 P^a e_s$$

א. כתוב/י את ההגדרה מתי אופרטור  $H$  כלשהו הוא  $\gamma$ -contracting, עבור נורמה  $\|\cdot\|_\infty$ .

הגדרה:

$$\forall V_1, V_2: \|HV_1 - HV_2\|_\infty \leq \gamma \|V_1 - V_2\|_\infty$$

ב. לכל אחד מהאופרטורים, ציין למה הוא מתכנס

תשובה:

$R^{(3)}: V^\pi$   
 $G^{(3)}: Q'$  such that  $\pi^* = greedy(Q')$

הסבר (למה הפתרון המוצע מקיים את הזהות של האופרטור):

ג. הוכח/י שהאופרטורים  $R^{(3)}$  ו-  $G^{(3)}$  הינם  $\gamma^3$ -contracting.

הוכחה:

$$\begin{aligned} & \left\| (R^{(3)}V_1)(s) - (R^{(3)}V_2)(s) \right\|_{\infty} = \\ & \left\| \gamma^3 V_1^T (P^\pi)^3 e_s - \gamma^3 V_2^T (P^\pi)^3 e_s \right\|_{\infty} = \gamma^3 \left\| (V_1 - V_2)^T (P^\pi)^3 e_s \right\| \\ & \leq \gamma^3 \|V_1 - V_2\|_{\infty} \end{aligned}$$

$$\begin{aligned} & \left\| (G^{(3)}Q_1)(s, a) - (G^{(3)}Q_2)(s, a) \right\|_{\infty} = \\ & \left\| \max_a \gamma^3 Q_1(\cdot, a)^T (P^\pi)^2 P^a e_s - \max_a \gamma^3 Q_2(\cdot, a)^T (P^\pi)^2 P^a e_s \right\|_{\infty} \\ & \leq \max_a \gamma^3 \left\| \max_a \gamma^3 (Q_1(\cdot, a) - Q_2(\cdot, a))^T (P^\pi)^2 P^a e_s \right\| \\ & \leq \gamma^3 \|Q_1 - Q_2\|_{\infty} \end{aligned}$$