

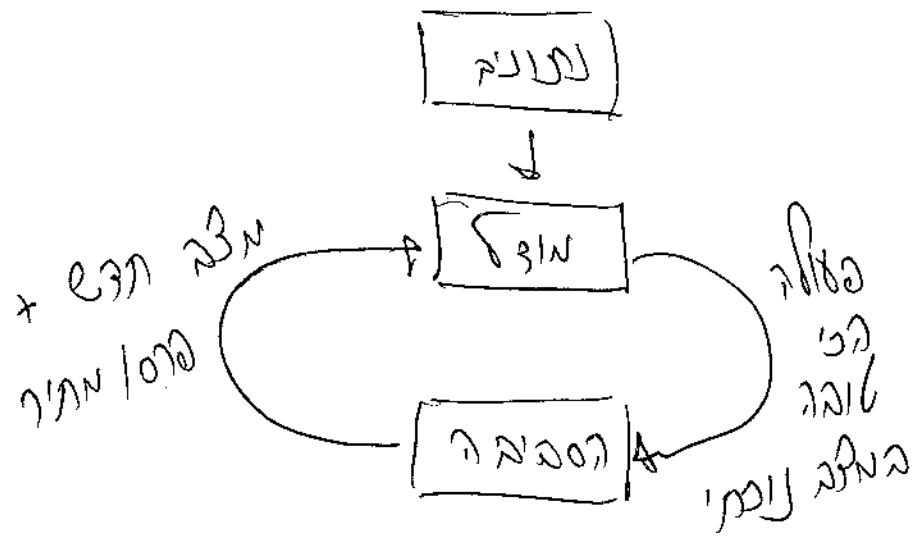
1

Reinforcement Learning



כ'מ'רה סט'ר'ט'ג'יה:

ה'עו'לה ה'א'מ'ית:



ב'מ'ה כ'ה ס'ו'ף מ'א'מ'יה א'ם ס'ו'ק'צ'י'ת מ'א'רה ס'א'מ'ת' ס'ו'ח?
 מ'ה מ'א'ס'ו'ן RL? א'מ'יה + מ'כ'ו'ן.

Markov Decision Process

- א'ו'ל מ'צ'ב'ו'ת S . S_t ה'מ'צ'ו'ת ב'מ'ן t
- א'ו'ל ס'ו'א'ת A . A_t ה'מ'צ'ו'ת ט'נ'ב'ו'ת ב'מ'ן t
- (s, a, s') מ'א'ר'צ'ו'ת מ'א'מ'יה
- R א'ו'ל ס'ו'א'ת. R_t ס'ו'א'ת ב'מ'ן t ל'מ'ר ס'מ'ח'ה ב'מ'נה
- $R(s, a)$ א'ו'ל ס'ו'א'ת מ'א'מ'יה
- $P(R, s' | s, a)$ א'מ'ק'ו'ה

(2)

מה הטענה?

לקבוע מדיניות $\pi(s) \rightarrow A$ (policy) (ייתכן סובסלו)
 כך שנמקסם את $R_t + R_{t+1} + R_{t+2} + \dots$

האם נלמדים? בתחילת הדרך, אם ייתכן כי נרצה לנסות
 פרטים חוקיים שונים. לזאת המטרה תהיה

$$\pi^* = \operatorname{argmax}_{\pi} E \left[\sum \gamma^t R_t \right]$$

נציג $\pi_{\theta}(a|s) = f(a, s; \theta)$: Policy Gradient

$$L(\theta) = E \left[\sum_{t=0}^T \gamma^t R(s_t, A_t) \right]$$

מהם טיפוסים קיימים.

הנמצא
 מידת פופולריות בו $f(a, s, \theta)$ יכול להיות רשת מעוקה
 היא עסוקה בקיום. במקרה זה, נשמקד בשיטות
 מבוססות value (עבור מצב, עבור מצב ופעולה).

$$V_{\pi}(s) = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

$$Q_{\pi}(s, a) = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

3

$$V^\pi(s) = r(s, \pi(s)) + \gamma \sum_{s'} p(s'|s, \pi(s)) V^\pi(s') \quad \text{: למשל}$$

$$V^\pi(s) = E \left[\sum_t \gamma^t R_t \mid S_0 = s \right] \quad \text{: המשל}$$

$$= E [R_0 \mid S_0 = s] + E \left[\sum_{t=1} \gamma^t R_t \mid S_0 = s \right]$$

$$= r(s, \pi(s)) + E \left[\sum_{t=1} \gamma^t R_t \mid S_0 = s \right]$$

$$= r(s, \pi(s)) + E \left[E \left[\sum_{t=1} \gamma^t R_t \mid S_0 = s, S_1 \right] \right]$$

$$= r(s, \pi(s)) + \sum_{s'} p(S_1 = s' \mid S_0 = s) E \left[\sum_{t=1} \gamma^t R_t \mid S_0 = s, S_1 = s' \right]$$

$$= r(s, \pi(s)) + \sum_{s'} p(s'|s, \pi(s)) E \left[\sum_{t=1} \gamma^t R_t \mid S_1 = s' \right]$$

המשל
?/כן /כן

$$= r(s, \pi(s)) + \gamma \sum_{s'} p(s'|s, \pi(s)) V^\pi(s')$$

Q^* אף בצורה דומה: אף Q אולי נראה שזה בהמשך אף Q

משל: V^π אולי Q^π בהמשך π (Policy Evaluation)

$$V_{t+1}^*(s) = r(s, \pi(s)) + \gamma \sum_{s'} p(s'|s, \pi(s)) V_t^*(s') \quad \forall s$$

למשל: המהלך מתקדם $V^\pi(s)$

4

$V_{t+1} = T(V_t)$ / μ / הוכחה:
 (contracting) γ / T
 $0 \leq c \leq 1$ $d(T(x), T(x')) \leq c \cdot d(x, x')$

Banach: d של T נפילתו של T על X נכנסת לנקודה יחידה

$$d(x_{t+1}, x^*) \leq c \cdot d(x_t, x^*)$$

$$\Downarrow$$

$$d(x_{t+1}, x^*) = \gamma^t d(x_0, x^*)$$

$$\Downarrow$$

$$t = \log \epsilon / d(x_0, x^*)$$

$$d(T[V], T[\tilde{V}]) = \|T(V) - T(\tilde{V})\|_\infty$$

$$= \max_s \left| p(s, \pi(s)) + \gamma \sum_{s'} p(s'|s, \pi(s)) V(s') - \tilde{V}(s) \right|$$

$$\leq \gamma \max_s \left| \sum_{s'} p(s'|s, \pi(s)) [V(s') - \tilde{V}(s')] \right|$$

$$\leq \gamma \max_s \sum_{s'} p(s'|s, \pi(s)) |V(s') - \tilde{V}(s')|$$

$$\leq \gamma \max_s \max_{s'} |V(s') - \tilde{V}(s')| = \gamma \|V - \tilde{V}\|_\infty$$

(הוכחה) γ של T הוא $\max_{s,s'} \sum_{s'} p(s'|s, \pi(s)) |V(s') - \tilde{V}(s')|$
 וזהו γ של T .

$Q^\pi(s, a)$ / μ / γ / T

5

Optimal Policy

π בוד $V_{\pi}(s) \equiv V_{\pi}(s)$ א'ק'ן מרובו π_* - א ש'י'ן מ'ן מ'ן
 פ'ת'ן מ'ן מ'ן π_* ו'ת'ן s פ'ת'ן מ'ן

$$Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

$$V_*(s) = \max_{\pi} V_{\pi}(s)$$

ש'ד s - א ר'ת'ן א ר'ת'ן מ'ן מ'ן
 π^* - א מ'ת'ן מ'ן

$$V_*(s) = \max_a Q^*(s, a)$$

י'ת'ן א π פ'ת'ן מ'ן מ'ן $Q^*(s, a)$ מ'ת'ן מ'ן

$$\textcircled{A} \quad \boxed{\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)}$$

$$Q^*(s, a) = \max_{\pi} E \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid S_0 = s, A_0 = a \right] \quad \text{--- מ'ת'ן מ'ן מ'ן}$$

$$= \max_{\pi} r(s, a) + \sum_{s'} p(s' | s, a) E \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid S_1 = s' \right]$$

$$= r(s, a) + \gamma \max_a \sum_{s'} p(s' | s, a) E \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid S_1 = s' \right]$$

$V^*(s')$ י'ת'ן מ'ן מ'ן s ו'ת'ן

$$= r(s, a) + \gamma \sum_{s'} p(s' | s, a) V^*(s')$$

$$= r(s, a) + \gamma \sum_{s'} p(s' | s, a) \max_a Q^*(a, s')$$

$$\textcircled{B} \quad \boxed{V^*(s) = \max_a r(s, a) + \gamma \sum_{s'} p(s' | s, a) V^*(s')}$$

י'ת'ן מ'ן מ'ן מ'ן
 $V^*(s) = \max_a Q^*(s, a)$

