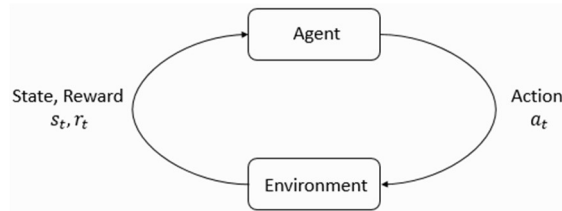


## 11. Reinforcement Learning (RL)

רוב האלגוריתמים של עולם הלמידה הינם מבוססי דאטה, כלומר, בהינתן מידע מסוים הם מנסים למצוא בו חוקיות מסוימת, ועל בסיסה לבנות מודל שיוכל להתאים למקרים נוספים. אלגוריתמים אלה מחולקים לשניים:

1. אלגוריתמים של למידה מונחית, המבוססים על דאטה  $S = \{x, y\}$ , כאשר  $x \in \mathbb{R}^{n \times d}$  הינו אוסף של אובייקטים (למשל נקודות במרחב, אוסף של תמונות וכדו'), ו- $y \in \mathbb{R}^n$  הינו אוסף של labels. לכל אובייקט  $x \in \mathbb{R}^d$  יש label מתאים  $y \in \mathbb{R}^1$ .
2. אלגוריתמים של למידה לא מונחית עבורם הדאטה  $x \in \mathbb{R}^{n \times d}$  הוא אוסף של אובייקטים ללא labels, ומנסים למצוא כללים מסוימים על דאטה זה (למשל – חלוקה לאשכולות, הורדת ממד ועוד).

למידה מבוססת חיזוקים הינה פרדיגמה נוספת תחת התחום של למידת מכונה, כאשר במקרה זה הלמידה לא מסתמכת על דאטה קיים, אלא על חקירה של הסביבה ומציאת המדיניות/האסטרטגיה הטובה ביותר לפעולה. ישנו סוכן שנמצא בסביבה שאינה מוכרת, ועליו לבצע צעדים כך שהתגמול המצטבר אותו הוא יקבל יהיה מקסימלי. בלמידה מבוססת חיזוקים, בניגוד לפרדיגמות האחרות של למידת מכונה, הסביבה לא ידועה מבעוד מועד. הסוכן נמצא באי ודאות ואינו יודע בשום שלב מה הצעד הנכון לעשות, אלא הוא רק מקבל פידבק על הצעדים שלו, וכך הוא לומד מה כדאי לעשות וממה כדאי להימנע. באופן כללי ניתן לומר שמטרת הלמידה היא לייצר אסטרטגיה כך שבכל מיני מצבים לא ידועים הסוכן יבחר בפעולות שבאופן מצטבר יהיו הכי יעילות עבורו. נתאר את תהליך הלמידה באופן גרפי:



איור 11.1 מודל של סוכן וסביבה.

בכל צעד הסוכן נמצא במצב  $s_t$  ובחר פעולה  $a_t$  המעבירה אותו למצב  $s_{t+1}$ , ובהתאם לכך הוא מקבל מהסביבה תגמול  $r_t$ . האופן בה מתבצעת הלמידה היא בעזרת התגמול, כאשר נרצה שהסוכן יבצע פעולות המזכות אותו בתגמול חיובי (חיזוק) וימנע מפעולות עבורן הוא מקבל תגמול שלילי, ובמצטבר הוא ימקסם את כלל התגמולים עבור כל הצעדים שהוא בחר לעשות. כדי להבין כיצד האלגוריתמים של למידה מבוססת חיזוקים עובדים ראשית יש להגדיר את המושגים השונים, ובנוסף יש לנסח באופן פורמלי את התיאור המתמטי של חלקי הבעיה השונים.

### 11.1 Introduction to RL

בפרק זה נגדיר באופן פורמלי תהליכי מרקוב, בעזרתם ניתן לתאר בעיות של למידה מבוססת חיזוקים, ונראה כיצד ניתן למצוא אופטימום לבעיות אלו בהינתן מודל וכל הפרמטרים שלו. לאחר מכן נדון בקצרה במספר שיטות המנסות למצוא אסטרטגיה אופטימלית עבור תהליך מרקוב כאשר לא כל הפרמטרים של המודל נתונים, ובפרקים הבאים נדבר על שיטות אלה בהרחבה. שיטות אלה הן למעשה הלב של למידה מבוססת חיזוקים, כיוון שהן מנסות למצוא אסטרטגיה אופטימלית על בסיס תגמולים ללא ידיעת הפרמטרים של המודל המרקובי עבורו רוצים למצוא אופטימום.

#### 11.1.1 Markov Decision Process (MDP) and RL

המודל המתמטי העיקרי עליו בנויים האלגוריתמים השונים של RL הינו תהליך החלטה מרקובי, כלומר תהליך שבה המעברים בין המצבים מקיים את תכונת מרקוב, לפיה ההתפלגות של מצב מסוים תלויה רק במצב הקודם לו:

$$P(s_{t+1} = j | s_1, \dots, s_t) = P(s_{t+1} = j | s_t)$$

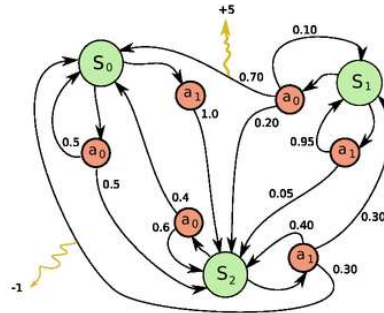
תהליך קבלת החלטות מרקובי מתואר על ידי סט הפרמטרים  $\{\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}\}$ :

- State space ( $\mathcal{S}$ ) – מרחב המצבים של המערכת. המצב ההתחלתי מסומן ב- $S_0$ .
- Action space ( $\mathcal{A}$ ) – מרחב הפעולות.  $A_s$  הוא מרחב הפעולות האפשריות במצב  $S$ .
- Transition ( $\mathcal{T}$ ) – הביטוי:  $T(s'|s, a) \rightarrow [0, 1]$  הינו פונקציית מעבר, המחשבת את ההסתברות לעבור בזמן  $t$  ממצב  $s_t$  למצב  $s_{t+1}$  על ידי הפעולה  $a: a = T(s'|s, a) = P(s_{t+1} = s' | s_t = s, a_t = a)$ . ביטוי זה למעשה מייצג את המודל – מה ההסתברות שבחירת הפעולה  $a$  במצב  $s$  תביא את הסוכן למצב  $s'$ .
- Reward ( $\mathcal{R}$ ) – הביטוי:  $\mathcal{R}_a(s, s') \rightarrow \mathbb{R}$  הינו פונקציה הנותנת תגמול/רווח לכל פעולה  $a$  הגורמת למעבר ממצב  $s$  למצב  $s'$ , כאשר בדרך כלל  $\mathcal{R}_a \in [0, 1]$ . לעיתים מסמנים את התגמול של הצעד בזמן  $t$  ב- $r_t$ .

המרקוביות של התהליך באה לידי ביטוי בכך שמצב  $s_t$  מכיל בתוכו את כל המידע הנחוץ בכדי לקבל החלטה לגבי  $a_t$ , או במילים אחרות – כל ההיסטוריה בעצם שמורה בתוך המצב  $s_t$ .

ריצה של MDP מאופיינת על ידי הרביעייה הסדורה  $\{s_t, a_t, r_t, s_{t+1}\}$  – פעולה  $a_t$  המתרחשת בזמן  $t$  וגורמת למעבר ממצב  $s_t$  למצב  $s_{t+1}$ , ובנוסף מקבלת תגמול מידי  $r_t$ , כאשר  $r_t \sim \mathcal{R}(s_t, a_t)$  ו- $s_{t+1} \sim p(\cdot | s_t, a_t)$ .

מסלול (trajectory) הינו סט של שלשות  $\tau = \{s_0, a_0, r_0, \dots, s_t, a_t, r_t\}$ , כאשר המצב התחלתי מוגדר מהתפלגות כלשהיא  $s_0 \sim \rho_0(\cdot)$ , והמעבר בין המצבים יכול להיות דטרמיניסטי  $s_{t+1} = f(s_t, a_t)$  או סטוכסטי  $s_{t+1} \sim p(\cdot | s_t, a_t)$ .



איור 11.2 תהליך קבלת החלטות מרקובי. ישנם שלושה מצבים –  $\{s_0, s_1, s_2\}$ , ובכל אחד מהם יש שתי פעולות אפשריות (עם הסתברויות מעבר שונות) –  $\{a_0, a_1\}$ . עבור חלק מהפעולות יש תגמול שונה מ-0. מסלול יהיה מעבר על אוסף של מצבים דרך אוסף של פעולות, שלכל אחד מהן יש תגמול.

אסטרטגיה של סוכן, המסומנת ב- $\pi$ , הינה בחירה של אוסף מהלכים. בבעיות של למידה מבוססת חיזוקים, נרצה למצוא **אסטרטגיה אופטימלית** (Optimal Policy)  $\pi: S \rightarrow A$  הממקסמת את התגמול המצטבר  $\sum_{t=0}^{\infty} \mathcal{R}(s_t, \pi(s_t))$ . כיוון שלא תמיד אפשרי לחשב באופן ישיר את האסטרטגיה האופטימלית, ניתן להגדיר ערך החזרה (Return) המבטא סכום של תגמולים, ומנסים למקסם את התוחלת שלו  $\mathbb{E}[Return | \mathcal{S}, \mathcal{A}]$ . ערך החזרה הכי נפוץ נקרא discount return, והוא מוגדר באופן הבא: עבור פרמטר  $\gamma \in (0, 1)$ , ה-Return הינו הסכום הבא:

$$Return = \sum_{t=1}^T \gamma^{t-1} r_t$$

אם  $\gamma = 0$ , אז מתעניינים רק בתגמול המיידי, וככל ש- $\gamma$  גדל כך נותנים יותר משמעות לתגמולים עתידיים. כיוון ש- $r_t \in [0, 1]$ , הסכום חסום על ידי  $\frac{1}{1-\gamma}$ .

**התוחלת של ערך החזרה נקראת Value function**, והיא נותנת לכל מצב ערך מסוים המשקף את תוחלת התגמול שניתן להשיג דרך מצב זה. באופן פורמלי, כאשר מתחילים ממצב  $s$ , ה-Value function מוגדר להיות:

$$\mathcal{V}^\pi(s) = \mathbb{E}[R(\tau) | s_0 = s]$$

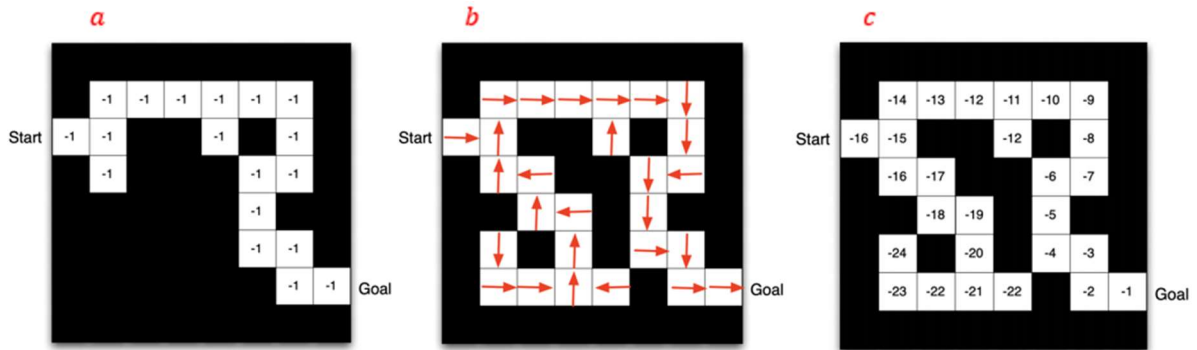
בעזרת ביטוי זה ניתן לחשב את **האסטרטגיה האופטימלית**, כאשר ניתן לנקוט בגישה ישירה ובגישה עקיפה. הגישה הישירה מנסה למצוא בכל מצב מה הפעולה הכי כדאית. בהתאם לכך, חישוב האסטרטגיה האופטימלית יעשה באופן הבא:

$$\pi(s) = \arg \max_a \sum_{s'} p_a(s, s') (\mathcal{R}_a(s, s') + \gamma \mathcal{V}(s'))$$

לעיתים החישוב הישיר מסובך, כיוון שהוא צריך לקחת בחשבון את כל הפעולות האפשריות, ולכן מסתכלים רק על ה-Value function. לאחר שלכל מצב יש ערך מסוים, בכל מצב הסוכן יעבור למצב בעל הערך הכי גדול מבין כל המצבים האפשריים אליהם ניתן לעבור. חישוב הערך של כל מצב נעשה באופן הבא:

$$\mathcal{V}(s) = \sum_{s'} p_\pi(s, s') (\mathcal{R}_\pi(s, s') + \gamma \mathcal{V}(s'))$$

ניתן לשים לב שבעוד הגישה הראשונה מתמקדת במציאת אסטרטגיה/מדיניות אופטימלית על בסיס הפעולות האפשריות בכל מצב, הגישה השנייה לא מסתכלת על הפעולות אלא על הערך של כל מצב, המשקף את תוחלת התגמול שניתן להשיג כאשר נמצאים במצב זה.



איור 11.3 (a) מודל: המצב של הסוכן הוא המשבצת בו הוא נמצא, הפעולות האפשריות הן ארבעת הכיוונים, כל פעולה גוררת תגמול של -1, והסתברויות המעבר נקבעות לפי הצבעים של המשבצות (אי אפשר ללכת למשבצות שחורות). (b) מדיניות – החלטה בכל מצב איזה צעד לבצע. (c) Value של כל משבצת.

**לסיכום**, ניתן לומר שכל התחום של RL מבוסס על שלוש אבני יסוד:

- מודל: האופן בו אנו מתארים את מרחב המצבים והפעולות. המודל יכול להיות נתון או שנצטרך לשערך אותו, והוא מורכב מהסתברויות מעבר בין מצבים ותגמול עבור כל צעד:

$$\mathcal{P}_{ss'}^a = p_\pi(s, s') = P(s_{t+1} = s' | s_t = s, a_t = a)$$

$$\mathcal{R}_{ss'}^a = \mathcal{R}_\pi(s, s') = \mathbb{E}[r_{t+1} | s_t = s, a_t = a, s_{t+1} = s']$$

- Value function – פונקציה המתארת את התוחלת של התגמולים העתידיים:

$$\mathcal{V}^\pi(s) = \mathbb{E}[R(\tau) | s_0 = s]$$

- מדיניות/אסטרטגיה (Policy) – בחירה (דטרמיניסטית או אקראית) של צעד בכל מצב נתון:  

$$\pi(s|a)$$

### 11.1.2 Bellman Equation

לאחר שהגדרנו את המטרה של למידה מבוססת חיזוקים, ניתן לדבר על שיטות לחישוב אסטרטגיה אופטימלית. בפרק זה נתייחס למקרה הספציפי בו נתון מודל מרקובי עם כל הפרמטרים שלו, כלומר אוסף המצבים, הפעולות והסתברויות המעבר ידועים. כאמור, Value function הינה התוחלת של ערך ההחזרה עבור אסטרטגיה נתונה  $\pi$ , כאשר מתחילים ממצב  $s$ :

$$\mathcal{V}^\pi(s) = \mathbb{E}[R(\tau) | s_0 = s] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s \right]$$

ביטוי זה מסתכל על הערך של כל מצב, בלי להתייחס לפעולות המעבירות את הסוכן ממצב אחד למצב אחר. נתינת ערך לכל מצב יכולה לסייע במציאת אסטרטגיה אופטימלית, כיוון שהיא מדרגת את המצבים השונים של המודל. באופן דומה, ניתן להגדיר את ה-Action-Value function – התוחלת של ערך ההחזרה עבור אסטרטגיה נתונה  $\pi$ , כאשר במצב  $s$  מבצעים את פעולה  $a$ , ולאחר מכן ממשיכים לפי האסטרטגיה  $\pi$ :

$$\mathcal{Q}^\pi(s, a) = \mathbb{E}[R(\tau) | s_0 = s, a_0 = a] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s, a_t = a \right]$$

ביטוי זה מסתכל על הזוג  $(s_t, a_t)$ , כלומר בכל מצב יש התייחסות למצב הנוכחי ולפעולות האפשריות במצב זה. בדומה ל-Value function, גם ביטוי זה יכול לסייע במציאת אסטרטגיה אופטימלית, כיוון שהוא מדרג עבור כל מצב את הפעולות האפשריות.

נוכל לסמן ב- $\mathcal{V}^*(s)$  ו- $\mathcal{Q}^*(s, a)$  את הערכים של האסטרטגיה האופטימלית  $\pi^*$  – Optimal Value function ו-Optimal Action-Value function. עבור אסטרטגיה זו מתקיים:

$$\mathcal{V}^*(s) = \max_{\pi} \mathbb{E}[R(\tau)|s_0 = s], Q^*(s, a) = \max_{\pi} \mathbb{E}[R(\tau)|s_0 = s, a_0 = a]$$

הרבה פעמים מתעניינים ביחס שבין  $\mathcal{V}$  ו- $Q$ , וניתן להיעזר במעברים הבאים:

$$\mathcal{V}^{\pi}(s) = \mathbb{E}[Q^{\pi}(s, a)]$$

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

באופן קומפקטי ניתן לרשום את  $\mathcal{V}^*(s)$  כך:

$$\forall s \in S \quad \mathcal{V}^*(s) = \max_{\pi} \mathcal{V}^{\pi}(s)$$

כלומר, האסטרטגיה  $\pi^*$  הינה האופטימלית עבור כל מצב  $s$ .

כעת נתון מודל מרקובי עם כל הפרמטרים שלו – אוסף המצבים והפעולות, הסתברויות המעבר והתגמול עבור כל פעולה, ומעוניינים למצוא דרך פעולה אופטימלית עבור מודל זה. ניתן לעשות זאת בשתי דרכים עיקריות – מציאת האסטרטגיה  $\pi(a|s)$  האופטימלית, או חישוב ה-Value של כל מצב ובחירת מצבים בהתאם לערך זה. משימות אלו יכולות להיות מסובכות מאוד עבור משימות מורכבות וגדולות, ולכן לעיתים קרובות משתמשים בשיטות איטרטיביות ובקירובים על מנת לדעת כיצד לנהוג בכל מצב. הדרך הפשוטה לחישוב  $\mathcal{V}^{\pi}(s)$  משתמשת ב-Bellman equation, המבוססת על תכונות דינמי. נפתח את הביטוי של  $\mathcal{V}^{\pi}(s)$  מתוך ההגדרה שלו:

$$\mathcal{V}^{\pi}(s) = \mathbb{E}[R(\tau)|s_0 = s] = \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s \right]$$

נפצל את הסכום שבתוחלת לשני איברים – האיבר הראשון ויתר האיברים:

$$= \mathbb{E}_{\pi} \left[ r_{t+1} + \gamma \cdot \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} | s_t = s \right]$$

כעת נשתמש בהגדרת התוחלת ונקבל:

$$\begin{aligned} &= \sum_{a, s'} \pi(a|s) p_{\pi}(s, s') \left( \mathcal{R}_{\pi}(s, s') + \gamma \cdot \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} | s_t = s \right] \right) \\ &= \sum_{a, s'} \pi(a|s) p_{\pi}(s, s') \left( \mathcal{R}_{\pi}(s, s') + \gamma \cdot \mathcal{V}^{\pi}(s') \right) \end{aligned}$$

הביטוי המתקבל הוא מערכת משוואות לינאריות הניתנות לפתרון באופן אנליטי, אם כי סיבוכיות החישוב יקרה. נסמן:

$$V = [V_1, \dots, V_n]^T, R = [r_1, \dots, r_n]^T$$

$$T = \begin{pmatrix} p_{11} & \dots & p_{1n} \\ \vdots & \ddots & \vdots \\ p_{n1} & \dots & p_{nn} \end{pmatrix}$$

ונקבל משוואה מטריציונית:

$$V = R + \gamma TV \rightarrow V = R + \gamma TV$$

$$\rightarrow \mathcal{V}^{\pi}(s) = (\mathbb{I}_n - \gamma T)^{-1} R$$

בגלל שהערכים העצמיים של  $T$  חסומים על ידי 1, בהכרח יהיה ניתן להפוך את  $\mathbb{I}_n - \gamma T$  מה שמבטיח שיהיה פתרון למשוואה, ופתרון זה הוא אף יחיד עבור  $\mathcal{V}^{\pi}$ . כשמוצאים את  $V$  ניתן למצוא גם את  $Q^{\pi}$  על ידי הקשר:

$$Q^{\pi}(s, a) = \sum_{s'} p_{\pi}(s, s') (\mathcal{R}_{\pi}(s, s') + \gamma \mathcal{V}^{\pi}(s')) = \sum_{s'} p_{\pi}(s, s') \left( \mathcal{R}_{\pi}(s, s') + \gamma \sum_{a'} \pi(a'|s') Q^{\pi}(a'|s') \right)$$

## Iterative Policy Evaluation

הסיבוכיות של היפוך מטריצה הינו  $\mathcal{O}(n^3)$ , ועבור  $n$  גדול החישוב נהיה מאוד יקר ולא יעיל. כדי לחשב את הפתרון באופן יעיל, ניתן כאמור להשתמש בשיטות איטרטיביות. שיטות אלו מבוססות על אופרטור בלמן, המוגדר באופן הבא:

$$BO(V) = R^\pi + \gamma T^\pi \cdot V$$

ניתן להוכיח שאופרטור זה הינו העתקה מכווצת (contractive mapping), כלומר הוא מקיים את התנאי:

$$\forall x, y: \|f(x) - f(y)\| < \gamma \|x - y\| \text{ for } 0 < \gamma < 1$$

במילים: עבור שני וקטורים במרחב, אופרטור  $f(\cdot)$  ומספר  $\gamma$  החסום בין 0 ל-1, אם נפעיל את האופרטור על כל אחד מהווקטורים ונחשב את נורמת ההפרש, נקבל מספר קטן יותר מאשר הנורמה בין הווקטורים כפול הפקטור  $\gamma$ . אופרטור המקיים תכונה זו הינו העתקה מכווצת, כיוון שנורמת ההפרש של האופרטור על שני וקטורים קטנה מנורמת ההפרש בין הווקטורים עצמם. הוכחה:

$$\|f(u) - f(v)\|_\infty = \|R^\pi + \gamma T^\pi \cdot v - (R^\pi + \gamma T^\pi \cdot u)\|_\infty = \|\gamma T^\pi(v - u)\|_\infty$$

מטריקת אינסוף מוגדרת לפי:  $\|u - v\|_\infty = \max_{s \in \mathcal{S}} |u(s) - v(s)|$ . לכן נוכל לרשום:

$$\|\gamma T^\pi(v - u)\|_\infty \leq \|\gamma T^\pi \mathbb{1}\|_\infty \|u - v\|_\infty$$

הביטוי  $T^\pi \mathbb{1}$  למעשה סוכם את כל ערכי מטריצת המעברים, לכן הוא מסתכם ל-1, ונקבל:

$$= \gamma \|u - v\|_\infty$$

ובכך הוכחנו את הדרוש.

לפי משפט נקודת השבת של בנק, להעתקה מכווצת יש נקודת שבת (fixed point) יחידה המקיימת  $x = f(x)$  וסדרה  $x_{t+1} = f(x_t)$  המתכנסת לאותה נקודת שבת. לכן נוכל להשתמש באלגוריתם איטרטיבי עבור  $\mathcal{V}^\pi$  שיביא אותנו לנקודת שבת, ולפי המשפט זוהי נקודת השבת היחידה וממילא הגענו להתכנסות. בפועל, נשתמש באלגוריתם האיטרטיבי הבא:

$$V_{k+1} = BO(V_k) = R^\pi + \gamma T^\pi \cdot V_k$$

נסתכל על הדוגמה הבאה:

$$T^\pi = \begin{pmatrix} 0.8 & 0.1 & 0.1 & 0 & 0 \\ 0.1 & 0.8 & 0.1 & 0 & 0 \\ 0 & 0.1 & 0.8 & 0.1 & 0 \\ 0 & 0 & 0.1 & 0.8 & 0.1 \\ 0 & 0 & 0.1 & 0.1 & 0.8 \end{pmatrix}, R^\pi = \begin{pmatrix} 0.1 \\ 1.3 \\ 3.4 \\ 1.9 \\ 0.4 \end{pmatrix}, \gamma = 0.9$$

באמצעות השיטה האיטרטיבית נקבל:

$$V_0 = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, V_1 = \begin{pmatrix} 0.1 \\ 1.3 \\ 3.4 \\ 1.9 \\ 0.4 \end{pmatrix}, V_2 = \begin{pmatrix} 0.6 \\ 2.6 \\ 6.1 \\ 3.7 \\ 1.2 \end{pmatrix}, \dots, V_{10} = \begin{pmatrix} 7.6 \\ 10.8 \\ 18.2 \\ 16.0 \\ 9.8 \end{pmatrix}, \dots, V_{50} = \begin{pmatrix} 14.5 \\ 17.1 \\ 26.4 \\ 26.8 \\ 18.4 \end{pmatrix}, V^\pi = \begin{pmatrix} 14.7 \\ 17.9 \\ 26.6 \\ 27.1 \\ 18.7 \end{pmatrix}$$

ניתן לשים לב שאחרי 50 איטרציות הפתרון המתקבל בצורה האיטרטיבית קרוב מאוד לפתרון המתקבל בצורה האנליטית.

## Policy Iteration (PI)

חישוב ה-Value function מאפשר לנו לחשב את ערכו של  $\mathcal{V}^\pi(s)$  עבור כל  $s$ , אך הוא אינו מבטיח שנגיע לאסטרגיה האופטימלית. נניח והצלחנו לחשב את  $\mathcal{V}^\pi(s)$  וממנו אנו יודעים לגזור אסטרגיה, עדיין יתכן שקיימת פעולה  $a$  שיותר משתלמת מאשר הפעולה המוצעת לפי האסטרגיה הנגזרת מ- $\mathcal{V}^\pi$ . באופן פורמלי ניתן לתאר זאת בצורה פשוטה – נניח שחישבנו את  $\mathcal{V}^\pi$  ואת  $Q^\pi(s, a)$  יתכן וקיימת פעולה עבורה:

$$\text{for such } s, a: Q^\pi(s, a) > V^\pi(s)$$

אם קיימת פעולה כזו, אז ישתלם לבחור בה ולאחר מכן לחזור לפעול בהתאם לאסטרטגיה  $\pi(a|s)$  הנגזרת מחישוב ה-Value function. למעשה, ניתן לחפש את כל הפעולות עבורן כדאי לבצע פעולה מסיימת עבורה התגמול יהיה גבוה יותר מאשר האסטרטגיה של  $V^\pi(s)$ . באופן פורמלי יותר, נרצה להגדיר אסטרטגיה דטרמיניסטית, עבורה בהסתברות 1 ננקוט בכל מצב  $s$  בפעולה הכי כדאית  $a$ :

$$\pi'(a|s) = 1 \text{ for } a = \arg \max_{a'} Q^\pi(s, a')$$

נשים לב שרעיון זה הוא בעצם להשתמש באסטרטגיה גרידית – בכל מצב לנקוט בפעולה הכי משתלמת בטווח של צעד יחיד, ואז להמשיך עם האסטרטגיה הנתונה. השאלה העולה היא כמובן – מדוע זה בהכרח נכון? כלומר, האם הרעיון שאומר שלא משנה באיזה מצב אנו נמצאים, הבחירה של הפעולה האופטימלית בהכרח תוביל לקבלת אסטרטגיה יותר טובה מאשר האסטרטגיה הנוכחית? בכדי להוכיח זאת ננסח זאת כמשפט:

בהינתן 2 אסטרטגיות  $\pi, \pi'$ , כאשר  $\pi'$  דטרמיניסטית, אז כאשר  $Q^\pi(s, \pi'(s)) > V^\pi(s)$  בהכרח לכל  $s$  יתקיים:  $V^{\pi'}(s) > V^\pi(s)$ . ראשית נפתח לפי הגדרה:

$$V^\pi(s) < Q^\pi(s, \pi'(s)) = \mathbb{E}_\pi[r_{t+1} + \gamma \cdot V^\pi(s_{t+1}) | s_t = s, a_t = \pi'(s)]$$

כיוון שהאסטרטגיה הינה דטרמיניסטית, הפעולה הנבחרת אינה רנדומלית ביחס ל- $\pi'$ , ולכן נוכל לרשום:

$$= \mathbb{E}_{\pi'}[r_{t+1} + \gamma \cdot V^\pi(s_{t+1}) | s_t = s]$$

כעת לפי אותו אי שוויון שבהנחה נוכל לבצע את אותו חישוב גם לצעד הבא  $s_{t+2}$ :

$$< \mathbb{E}_{\pi'}[r_{t+1} + \gamma \cdot Q^\pi(s_{t+1}, \pi'(s_{t+1})) | s_t = s]$$

וזו שוב שווה ל:

$$= \mathbb{E}_{\pi'}[r_{t+1} + \gamma \cdot r_{t+2} + \gamma^2 \cdot V^\pi(s_{t+2}) | s_t = s]$$

וכך הלאה, ולמעשה הוכחנו את הדרוש – נקיטת הפעולה הכי יעילה בכל מצב תמיד תהיה יותר טובה מהפתרון של  $V^\pi(s)$ .

כעת יש בידינו שתי טכניקות שאנו יודעים לבצע:

**Evaluation (E)** – בהינתן אסטרטגיה מסוימת נוכל לפתור את משוואות בלמן ולקבל את  $V^\pi(s)$  ו- $Q^\pi(s, a)$ .

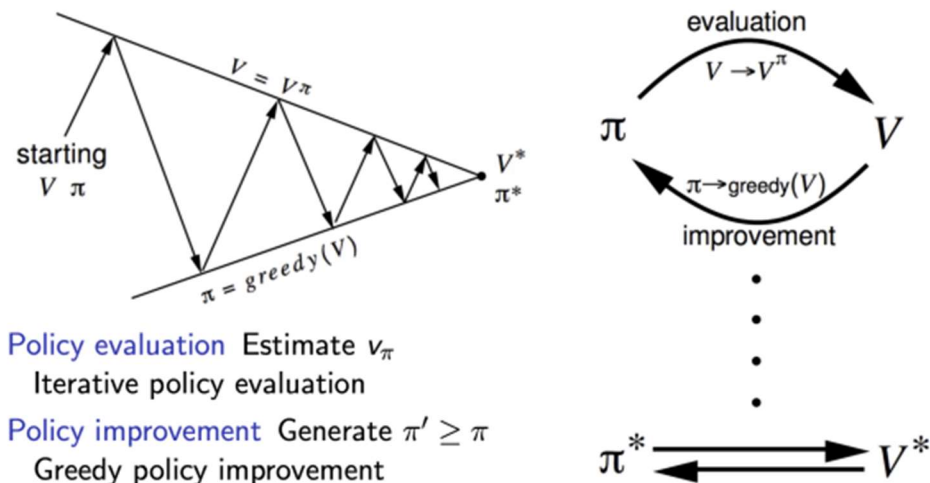
**Improve (I)** – בהינתן הערך של Value function, נוכל לשפר אותה באמצעות בחירה גרידית של פעולה.

בעזרת טכניקות אלו ניתן להתחיל מאסטרטגיה רנדומלית, ואז לבצע איטרציות המורכבות משתי הטכניקות האלה באופן הבא:

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots$$

תהליך זה נקרא Policy iteration – בכל צעד בו יש לנו אסטרטגיה נפתור עבורה משוואות בלמן ובכך נחשב את ה-Value function שלה, ולאחר מכן נשפר את האסטרטגיה באמצעות policy improvement, שכאמור מבצע בחירה גרידית שבטווח הקצר טובה יותר מאשר ה-Value function שחישבנו. ניתן להוכיח שאחרי מספר סופי של איטרציות האסטרטגיה תתכנס לנקודת שבת (fixed point), ואז הפעולה הבאה לפי האסטרטגיה תהיה זהה לבחירה הגרידית:

$$\pi(s) = \arg \max_a Q^\pi(s, a) = \pi'(s)$$



איור 11.4 Policy iteration – ביצוע איטרציות של Policy evaluation ו-Policy improvement על מנת למצוא בכל שלב את ה-value function ולשפר אותו באמצעות בחירה גרידית.

### Bellman optimality equations

השלב הבא בשימוש ב-Policy iteration הוא **להוכיח** שהאסטרטגיה אליה מתכנסים הינה אופטימלית. נסמן את נקודת השבת ב- $\pi^*$  ונקבל את הקשר הבא:

$$V^{\pi^*}(s) \equiv V^*(s) = \max_a Q^*(s, a) = \max_a \sum_{s'} p_\pi(s, s') (\mathcal{R}_\pi(s, s') + \gamma \cdot V^*(s'))$$

ובאופן דומה:

$$Q^*(s, a) = \sum_{s'} p_\pi(s, s') (\mathcal{R}_\pi(s, s') + \gamma \cdot \max_{a'} Q^*(s', a'))$$

משוואות אלה נקראות Bellman optimality equation. ניתן לשים לב שהן מאוד דומות למשוואות בלמן מהן יצאנו, אך במקום התוחלת שהייתה לנו בהתחלה, כעת יש  $\max$ . נרצה להראות שהפתרון של משוואות אלה הוא ה-Value של **האסטרטגיה האופטימלית**. ננסח את הטענה באופן הבא:

אסטרטגיה הינה אופטימלית אם ורק אם היא מקיימת את Bellman optimality equation. כיוון אחד להוכחה הוא טריוויאלי – אם האסטרטגיה הינה אופטימלית אז היא בהכרח מקיימת את משוואות האופטימליות, כיוון שהראינו שהן מתקבלות מנקודת השבת אליה האיטרציות מתכנסות. אם האסטרטגיה לא הייתה אופטימלית אז היה ניתן לשפר עוד את האסטרטגיה ולא היינו מגיעים עדיין לנקודת השבת. בשביל להוכיח את הכיוון השני נשתמש שוב ברעיון של העתקה מכווצת. נגדיר את האופרטור הבא:

$$BV(s) = \max_a \sum_{s'} p_\pi(s, s') (\mathcal{R}_\pi(s, s') + \gamma \cdot V(s))$$

ניתן להראות שאופרטור זה הינו העתקה מכווצת, וממילא לפי המשפט של בנק יש לו נקודת שבת יחידה. כיוון שהראינו ששימוש ב-Policy iteration מביא את האסטרטגיה לנקודת שבת מסוימת, נוכל לצרף לכך את העובדה שהאופרטור שהגדרנו הינו העתקה מכווצת וממילא נקבל שאותה נקודת שבת הינה יחידה, וממילא אופטימלית.

### Value Iteration

הראנו שבעזרת שיטת Policy iteration ניתן להגיע לאסטרטגיה אופטימלית, אך התהליך יכול להיות איטי. ניתן לנקוט גם בגישה יותר ישירה ולנסות לחשב באופן ישיר את הפתרון של משוואות האופטימליות של בלמן (ופתרון הינו אופטימלי כיוון שהראינו שהפתרון הוא נקודת שבת יחידה). נתחיל עם פתרון רנדומלי  $V_0$  ולאחר מכן נצבע איטרציות באופן הבא עד שנגיע להתכנסות:

$$\mathcal{V}_{k+1} = \max_a \sum_{s'} p_\pi(s, s') (\mathcal{R}_\pi(s, s') + \gamma \cdot \mathcal{V}_k(s'))$$

נשים לב שבשיטה זו אין לנו מידע לגבי האסטרטגיה אלא רק חישבנו את ה-Value function, אך ממנה ניתן לגזור את  $Q$  ואז לבחור באסטרטגיה גרידית, שהינה במקרה זה גם אופטימלית:

$$\pi(s) = \arg \max_a Q^\pi(s, a)$$

ניתן להראות כי בשיטה זו ההתכנסות מהירה יותר ודרושות פחות איטרציות מהשיטה הקודמת, אך כל איטרציה יותר מורכבת.

## Limitations

לשתי השיטות – Policy iteration ו-Value iteration – יש שני חסרונות מרכזיים:

1. הן דורשות לדעת את המודל והסביבה באופן שלם ומדויק.
2. הן דורשות לעדכן בכל שלב את כל המצבים בו זמנית. עבור מערכות עם הרבה מצבים, זה לא מעשי.

### 11.1.3 Learning Algorithms

בפרק הקודם הוסבר כיצד ניתן לחשב את האסטרטגיה האופטימלית וערך ההחזרה **בהינתן** מודל מרקובי. השתמשנו בשתי הנחות עיקריות על מנת להתמודד עם הבעיה:

1. Tabular MDP – הנחנו שהבעיה סופית ולא גדולה מדי, כך שנוכל לייצג אותה בזיכרון ולפתור אותה.
  2. Known environment – הנחנו שהמודל ידוע לנו, כלומר נתונה לנו מטריצת המעברים שקובעת מה הסיכוי לעבור ממצב  $s$  למצב  $s'$  כשנוקטים בפעולה  $a$  (סימנו את זה בתור  $\mathcal{P}_{ss'}^a = p_\pi(s, s')$ , ובנוסף נתון לנו מה ה-reward המתקבל עבור כל action (סימנו את זה בתור  $\mathcal{R}_{ss'}^a = \mathcal{R}_\pi(s, s')$ ).
- בעזרת שתי ההנחות פיתחנו את משוואות בלמן, כאשר היו לנו שני צמדים של משוואות. משוואות בלמן עבור אסטרטגיה נתונה נכתבות באופן הבא:

$$\mathcal{V}^\pi(s) = \sum_{a, s'} \pi(a|s) \mathcal{P}_{ss'}^a (\mathcal{R}_{ss'}^a + \gamma \cdot \mathcal{V}^\pi(s'))$$

$$Q^\pi(s, a) = \sum_{s'} \mathcal{P}_{ss'}^a \left( \mathcal{R}_{ss'}^a + \gamma \sum_{a'} Q^\pi(s', a') \right)$$

ובנוסף פיתחנו את המשוואות עבור הפתרון האופטימלי:

$$\mathcal{V}^*(s) = \max_a \sum_{s'} \mathcal{P}_{ss'}^a (\mathcal{R}_{ss'}^a + \gamma \cdot \mathcal{V}^*(s'))$$

$$Q^*(s, a) = \sum_{s'} \mathcal{P}_{ss'}^a \left( \mathcal{R}_{ss'}^a + \gamma \max_{a'} Q^*(s', a') \right)$$

הראינו שתי דרכים להגיע לפתרון האופטימלי:

1. Policy iteration המורכב מ-Policy evaluation ולאחריו Policy improvement.

2. Value iteration – פתרון משוואות בלמן באופן ישיר בעזרת איטרציות על ה-Value function.

כאמור, דרכי פתרון אלו מניחים שהמודל ידוע, ובנוסף שמרחב המצבים אינו גדול מדי ויכול להיות מיוצג בזיכרון. האתגר האמיתי מתחיל בנקודה בה לפחות אחת מהנחות אלה אינה תקפה, ולמעשה פה מתחיל התפקיד של אלגוריתמי RL. עיקר ההתמקדות של אלגוריתמים אלו יהיה למצוא באופן יעיל את האסטרטגיה האופטימלית כאשר לא נתונים הפרמטרים של המודל, ואז צריך לשערך אותם (Model-based learning) או למצוא דרך אחרת לחישוב האסטרטגיה האופטימלית ללא שימוש במודל (Model free learning). אם למשל יש משחק בין משתמש לבין המחשב, אלגוריתמים השייכים ל-Model based learning ינסו ללמוד את המודל של המשחק או להשתמש במודל



קיים, ובעזרת המודל הם ינסו לבחון כיצד יגיב המשתמש לכל תור שהמחשב יבחר. לעומת זאת אלגוריתמים מסוג Model free learning לא יתעניינו בכך, אלא ינסו ללמוד ישירות את האסטרטגיה הטובה ביותר עבור המחשב.

היתרון המשמעותי של אלגוריתמים המסתכלים על המודל של הבעיה (Model-based) נובע מהיכולת לתכנן מספר צעדים קדימה, כאשר עבור כל בחירה של פעולה המודל בוחן את התגובות האפשריות, את הפעולות המתאימות לכל תגובה, וכך הלאה. דוגמא מפורסמת לכך היא תוכנת המחשב AlphaZero שאומנה לשחק משחקי לוח כגון שחמט או גו. במקרים אלו המודל הוא המשחק והחוקים שלו, והתוכנה משתמשת בידע הזה בכדי לבחון את כל הפעולות והתגובות למשך מספר צעדים רב ובחירה של הצעד הטוב ביותר.

עם זאת, בדרך כלל אף בשלב האימון אין לסוכן מידע חיצוני מהו הצעד הנכון באופן אולטימטיבי, ועליו ללמוד רק מהניסיון. עובדה זו מציבה כמה אתגרים, כאשר העיקרי ביניהם הוא הסכנה שהאסטרטגיה הנלמדת תהיה טובה רק עבור המקרים אותם ראה הסוכן, אך לא תתאים למקרים חדשים שיבואו. אלגוריתמים שמחפשים באופן ישיר את האסטרטגיה האופטימלית אמנם לא משתמשים בידע שיכול להגיע מבחינת צעדים עתידיים, אך הם הרבה יותר פשוטים למימוש ולאמון.

באופן מעט יותר פורמלי ניתן לנסח את ההבדל בין הגישות כך: גישת Model-based learning מנסה למצוא את הפרמטרים המגדירים את המודל  $\{\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}\}$  ואז בעזרתם לחשב את האסטרטגיה האופטימלית (למשל בעזרת משוואות בלמן). הגישה השנייה לעומת זאת לא מעוניינת לחשב במפורש את הפרמטרים של המודל אלא למצוא באופן ישיר את האסטרטגיה האופטימלית  $\pi(a_t|s_t)$  שעבור כל מצב קובעת באיזה פעולה לנקוט. ההבדל בין הגישות נוגע גם לפונקציית המחיר לה נרצה למצוא אופטימום.

בכל אחד משני סוגי הלמידה יש אלגוריתמים שונים, כאשר הם נבדלים אחד מהשני בשאלה מהו האובייקט אותו מעוניינים ללמוד.

## Model-free learning

בגישה זו יש שתי קטגוריות מרכזיות של אלגוריתמים:

- א. Policy Optimization – ניסוח האסטרטגיה כבעיית אופטימיזציה של מציאת סט הפרמטרים  $\theta$  המקסם את  $\pi_\theta(a|s)$ . פתרון בעיה זו יכול להיעשות באופן ישיר על ידי שיטת Gradient Ascent עבור פונקציית המחיר  $\mathcal{J}(\pi_\theta) = \mathbb{E}[R(\tau)]$ , או בעזרת קירוב פונקציה זו ומציאת מקסימום עבורה.
- ב. Q-learning – שערך  $Q^*(s, a)$  על ידי  $Q_\theta(s, a)$ . מציאת המשערך האופטימלי יכולה להתבצע על ידי חיפוש  $\theta$  שיספק את השערך הטוב ביותר שניתן למצוא, או על ידי מציאת הפעולה שתמקסם את המשערך:  
$$a(s) = \arg \max_a Q_\theta(s, a)$$

השיטות המנסות למצוא אופטימום לאסטרטגיה הן לרוב on-policy, כלומר כל פעולה נקבעת על בסיס האסטרטגיה המעודכנת לפי הפעולה הקודמת. Q-learning לעומת זאת הוא לרוב אלגוריתם off-policy, כלומר בכל פעולה ניתן להשתמש בכל המידע שנצבר עד כה. היתרון של שיטות האופטימיזציה נובע מכך שהן מנסות למצוא באופן ישיר את האסטרטגיה הטובה ביותר, בעוד שאלגוריתם Q-learning רק משערך את  $Q^*(s, a)$ , ולעיתים השערך לא מספיק ואז התוצאה המתקבלת אינה מספיק טובה. מצד שני, כאשר השערך מוצלח, הביצועים של Q-learning טובים יותר, כיוון שהשימוש במידע על העבר מנוצל בצורה יעילה יותר מאשר באלגוריתמים המבצעים אופטימיזציה של האסטרטגיה. שתי הגישות האלה אינן זרות לחלוטין, וישנם אלגוריתמים שמנסים לשלב בין הרעיונות ולנצל את החוזקות והיתרונות שיש לכל גישה.

## Model-based learning

גם בגישה זו יש שתי קטגוריות מרכזיות של אלגוריתמים:

- א. Model-based RL with a learned model – אלגוריתמים המנסים ללמוד הן את המודל עצמו והן את ה-Value function או את האסטרטגיה  $\pi$ .
- ב. Model-based RL with a known model – אלגוריתמים המנסים למצוא את ה-Value function ו/או את האסטרטגיה כאשר המודל עצמו נתון.

ההבדל בין הקטגוריות טמון באתגר איתו מנסים להתמודד. במקרים בהם המודל ידוע, הממד של אי הוודאות לא קיים, ולכן ניתן להתמקד בביצועים אסימפטוטיים. במקרים בהם המודל אינו ידוע, הדגש העיקרי הוא על למידת המודל.

## 11.2 Model Free Prediction

לאחר שסקרנו בפרק המבוא את הבסיס המתמטי של בעיות RL והצגנו את משוואות בלמן ופתרוןן, בפרקים הבאים נציג גישות שונות להתמודדות עם בעיות RL עבורן פתרונות אלה אינם מספיקים – או מפני שהמודל אינו ידוע או מפני שהן Scale גדול יותר מזה שניתן לפתור באמצעות משוואות בלמן. בפרק זה נציג שתי שיטות הבאות להתמודד עם מקרים בהם המודל אינו ידוע (כלומר האלגוריתם הינו Model-Free), והדרך שלהן להתמודד עם אתגר זה הינו לשערך את המודל.

### 11.2.1 Monte-Carlo (MC) Policy Evaluation

האלגוריתם הראשון אותו נציג הינו Monte Carlo, והוא מציע דרך לשערך את ה-Value function בלי לדעת את המודל. ראשית נסביר בקצרה מהו אלגוריתם Monte Carlo ואז נראה כיצד ניתן ליישם אותו בבעיות RL.

נניח ונרצה לשערך תוחלת של פונקציית התפלגות כלשהיא  $\mathbb{E}_p[f(x)]$ . התוחלת יכולה להיות סכום או אינטגרל שקשה מאוד לחשב. ניתן לשערך את התוחלת על ידי דגימות רנדומליות מההתפלגות וחישוב הממוצע של הדגימות:

$$x_1, \dots, x_n \sim p(x)$$

$$\mathbb{E}_p[f(x)] \approx \frac{1}{n} \sum_i f(x_i)$$

לפי חוק המספרים הגדולים הממוצע של הדגימות מתכנס לתוחלת. משערך זה הינו חסר הטיה, ובנוסף השונות שלה קטנה ביחס לינארי לכמות הדגימות:

$$\mathbb{E} \left[ \frac{1}{n} \sum_i f(x_i) \right] = \mathbb{E}[f(x)]$$

$$\text{Var} \left[ \frac{1}{n} \sum_i f(x_i) \right] = \frac{\text{Var}[f(x)]}{n}$$

כאמור לעיל, ה-Value function הינה תוחלת עבור אסטרטגיה נתונה  $\pi$ , כאשר מתחילים ממצב  $s$ :

$$V^\pi(s) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s \right]$$

אמנם התוחלת היא על סכום אינסופי, אך עבור מקרים רבים נוכל להניח שהריצה היא סופית (Episodic MDP). בפועל נבצע את הפעולה הבאה: עבור אסטרטגיה נתונה  $\pi$ , נייצר ממנה ריצה של  $T$  צעדים, ואז ניקח מצב מסוים ונסתכל על כל ה-rewards שמתקבלים בריצה זו החל ממנו. באופן הזה קיבלנו value עבור הערך של אותו מצב. חזרה על אותה פעולה שוב ושוב תייצר ריצות שונות וממילא ערכים שונים למצב מסוים, ומיצוע על פני הערכים ייתן לנו שערך לערך האמיתי של אותו מצב. נעיר בהערת אגב שכיוון שמצבים יכולים לחזור על עצמם, נוצרת בעיה שבריצה כזו המצבים אינם בלתי תלויים. בכדי להתגבר על כך, אם מצב חוזר על עצמו יותר מפעם אחת מאפשרים לדגום רק את המופע הראשון של אותו מצב ולא את יתר המופעים (ישנן עוד דרכים להתגבר על כך, אך זוהי הדרך פשוטה ביותר). באופן פורמלי, נניח ויש לנו ריצה של  $T$  צעדים:

$$S_1, A_1, R_1, \dots, S_{T-1}, A_{T-1}, R_{T-1}, S_T$$

אז דגימה אחת מתוך ההתפלגות  $p_\pi(\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = S_t)$  תראה באופן הבא:

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$$

למשל, הדגימה  $G_1$  תהיה מורכבת מכל ה-rewards שהגיעו לאחר הצעד הראשון:  $R_1 + \gamma R_2 + \dots + \gamma^{T-1} R_T$ . הסכום הזה ייתן לנו value עבור אותו מצב ממנו התחלנו ( $S_1$ ), ועל ידי שערך אותו מצב שוב ושוב ביחס לריצות שונות נוכל לקבל ערכים שונים, ולאחר מכן למצע אותם בכדי לשערך את ה-value האמיתי של אותו מצב  $S_1$ .

באופן פורמלי, העדכון של מצב לאחר כל דגימה נראה כך:

#sample of  $s_t$ :  $N(S_t) = N(S_t) + 1$

update the value of  $s_t$ :  $V(s_t) = V(s_t) + \frac{1}{N(S_t)}(G_t - V(s_t))$

ולאחר הרבה דגימות השערוך מתקבל על ידי התחלת שלהן:

$$V(s) = \mathbb{E}_\pi[G_t | s_t = s]$$

באופן סכמתי ניתן לתאר את האלגוריתם באופן הבא:

#### First-visit MC prediction, for estimating $V \approx v_\pi$

Initialize:

$\pi \leftarrow$  policy to be evaluated  
 $V \leftarrow$  an arbitrary state-value function  
 $Returns(s) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$

Repeat forever:

Generate an episode using  $\pi$   
 For each state  $s$  appearing in the episode:  
 $G \leftarrow$  the return that follows the first occurrence of  $s$   
 Append  $G$  to  $Returns(s)$   
 $V(s) \leftarrow \text{average}(Returns(s))$

איור 11.5 אלגוריתם MC עבור שערוך ה-Value function בהינתן Policy. עבור כל מצב מאתחלים רשימה ריקה, ולאחר מכן מייצרים המון ריצות שונות. עבור כל ריצה, עוברים על כל המצבים ובודקים מה ה- $G$  שלהם, ומוסיפים אותו לרשימה של אותו מצב. לבסוף מחשבים את ה-value של כל מצב על ידי מיצוע הרשימה (=ערכי  $G$  השונים) של אותו מצב.

יש מספר יתרונות לשיטה זו: היא מספקת משערוך חסר הטיה עבור ה-Value function, מבטיחה התכנסות אחרי מספיק איטרציות, ובנוסף ניתן לשערך באמצעותה את השגיאה. אפשר גם באותה דרך לשערך גם את  $Q^\pi(s, a)$  אך זה יהיה יותר רועש וידרוש יותר דגימות. מן הצד השני יש גם חסרונות לשיטה זו: ראשית, היא מתאימה רק ל-Episodic MDP ולא לריצות אינסופיות. עם בעיה זו ניתן להתמודד בקלות כיוון שעבור בעיה אינסופית ניתן לקחת ריצה סופית ולחסום את השגיאה באמצעות  $\gamma$ . שנית, ההתכנסות יכולה להיות מאוד איטית, ובנוסף השונות יחסית גבוהה.

### 11.2.2 Temporal Difference (TD) – Bootstrapping

במקום להסתכל על ריצה שלמה, ניתן אחרי כל צעד לעדכן את ה-Value. ממשוואת בלמן ניתן להראות שהביטוי  $R_{t+1} + \gamma V^\pi(S_{t+1})$  הינו משערוך חסר הטיה עבור  $V^\pi(S_t)$ . אי אפשר להשתמש במשערוך זה כמו שהוא, כיוון שאנחנו לא יודעים את ה-Value function, וממילא הביטוי  $V^\pi(S_{t+1})$  לא ידוע. בשביל בכל זאת לשערך את  $V^\pi(S)$  באמצעות אותו משערוך, ניתן להחליף את  $V^\pi(S_{t+1})$  ב- $V(S_{t+1})$ , ולבצע את השערוך באופן הבא:

$$V^\pi(S_t) = R_{t+1} + \gamma V(S_{t+1})$$

הרעיון מאחורי השימוש הזה הוא להיעזר במידע שיש לנו מ- $R_t$ . נניח וננחש ערך כלשהוא עבור  $V^\pi(S_t)$  וניחוש זה יהיה גרוע. אפשר מעט לשפר את הניחוש באמצעות ניחוש  $V(S_{t+1})$  ושימוש ב-reward  $R_{t+1}$  שהתקבל עבור אותו מצב  $S_t$ , שבעצם מספק מידע כלשהוא על מצב זה. שערוך זה עדיף על ניחוש מוחלט כיוון שהאלמנט של הניחוש מקבל משקל נמוך יותר עקב המכפלה ב- $\gamma$ , ויש יותר משקל ל- $R_{t+1}$  שמספק מידע אמיתי על המצב  $S_t$ . צריך לשים לב שאנו מנסים לשערך את ה-Value function מתוך הערכים שלה בעצמה. מאתחלים את כל הערכים במספרים כלשהם (למשל – וקטור של 0), ואז עוברים צעד צעד ומעדכנים את הניחושים בעזרת התגמולים. אינטואיטיבית זה נראה מעט משונה, אך מסתבר שפתרון זה הוא אחד הכלים החזקים בבעיות RL. באופן פורמלי האלגוריתם נראה כך:

### Tabular TD(0) for estimating $v_\pi$

Input: the policy  $\pi$  to be evaluated  
Initialize  $V(s)$  arbitrarily (e.g.,  $V(s) = 0$ , for all  $s \in \mathcal{S}^+$ )  
Repeat (for each episode):  
    Initialize  $S$   
    Repeat (for each step of episode):  
         $A \leftarrow$  action given by  $\pi$  for  $S$   
        Take action  $A$ , observe  $R, S'$   
         $V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$   
         $S \leftarrow S'$   
    until  $S$  is terminal

איור 11.6 אלגוריתם Temporal Difference (TD) עבור שערך ה-Value function בהינתן Policy. מנחשים ערך עבור כל מצב ואז באופן איטרטיבי משפרים את הניחושים בעזרת שערך התלוי ב-reward ובערך המצב הבא.

מגדירים את השגיאה של ה-TD כהפרש שבין הניחוש עבור ערך המצב לבין השערך שלו:

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

בשונה משערך MC, השערך בשיטת TD הוא בעל הטיה, אך השונות קטנה יותר. בנוסף, כיוון שבכל צעד מבצעים שיפור לערך של מצב, תהליך השערך יותר מהיר מאשר ב-MC. הרבה פעמים מוסיפים פרמטר  $\alpha$  לשערך (כפי שמופיע באיור 11.6):

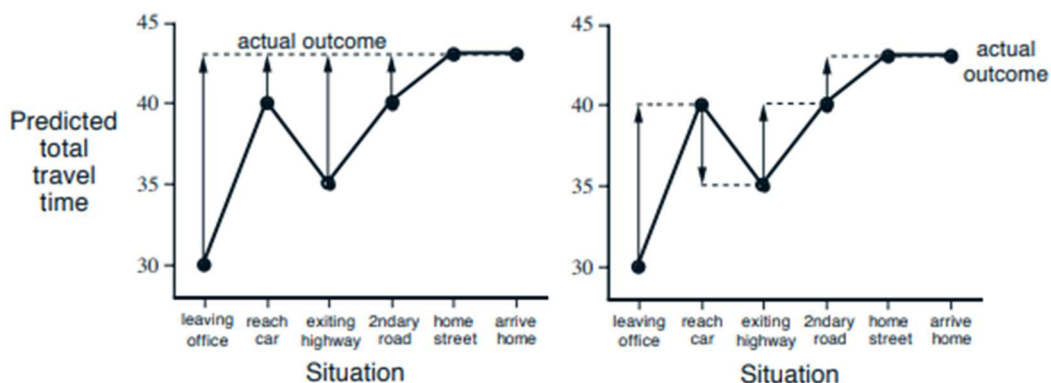
$$V(S_t) = V(S_t) + \alpha \cdot [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

עבור ערכי  $\alpha$  מתאימים, המשערך מתכנס לתוחלת האמיתית.

ניתן דוגמה שתמחיש את שיטת MC ושיטת TD ואת היחס ביניהם: נהג יצא מהמשרד שלו ונסע לביתו, ובדרך הוא ניסה לשערך את זמן ההגעה שלו וקיבל את הזמנים הבאים:

שערך זמן הנסיעה הכולל	שערך הזמן שנותר לנסיעה	כמה זמן עבר	מצב
30	30	0	יציאה מהמשרד
40	35	5	הליכה למכונית תחת גשם
35	15	20	הגעה לכביש מהיר
40	10	30	נסיעה מאחורי משאית
43	3	40	הגעה לרחוב של הבית
43	0	43	הגעה הביתה

נשרטט את שני המשערכים:



איור 11.7 שערך MC (משמאל) ושערך TD (מימין) ביחס לתצפית הנתונה.

כיוון ששיטת MC מספקת ערך לאחר ריצה שלמה, אז ניקח את זמן ההגעה בפועל של הנהג וניתן את הערך הזה לכל מצבי הביניים. שיטת TD לעומת זאת מעדכנת את הערך בכל מצב בהתאם למצב הבא. ניתן לראות שהשונות בשערך TD קטנה מזו שהתקבלה בשערך MC.

ניתן להסתכל על שיטת TD כבעיית רגרסיה "דינמית":

עבור כל צעד נרצה ש- $\mathcal{V}(S_t)$  יהיה שווה למשוואות בלמן, כלומר נרצה לשערך את  $\mathcal{V}(S_t)$  כך שיתקיים:

$$\mathcal{V}(S_t) = \mathbb{E}_\pi[R_{t+1} + \gamma\mathcal{V}(S_{t+1})|s_t = S_t]$$

עבור בעיה זו נוכל להגדיר פונקציית מחיר (Loss) מתאימה:

$$L = \frac{1}{2} (\mathbb{E}_\pi[R_{t+1} + \gamma\mathcal{V}(S_{t+1})|s_t = S_t] - \mathcal{V}(S_t))^2$$

את הפונקציה הזו ניתן למזער על ידי דגימות סטוכסטיות של  $R_{t+1} + \gamma\mathcal{V}(S_{t+1})$ . נשים לב לדבר חשוב – המטרה שלנו היא לשערך את ההווה באמצעות העתיד ולהיפך, כיוון שהעתיד הוא בעל יותר מידע – הוא ראה reward ומצב חדש. הבחנה זו משפיעה על איך שאנחנו רוצים שפונקציית המחיר תתנהג – אנחנו רוצים ש- $\mathcal{V}(S_t)$  יתקרב לערך של  $\mathbb{E}_\pi[R_{t+1} + \gamma\mathcal{V}(S_{t+1})|s_t = S_t]$  ולא להיפך. בדוגמה של הנהג שמשערך את זמן הנסיעה – הוא רוצה לעדכן את השערוך של כל מצב בהתאם למצב הבא. אם למשל ההערכה שלו בזמן  $t$  הינה 30 דקות ואז הוא מגיע לפקק ומעדכן את ההערכה ל-35 דקות, אז הוא ירצה לתקן את השערוך הקודם ( $\mathcal{V}(S_t)$ ) כך שיהיה דומה לשערוך הנוכחי ( $R_{t+1} + \gamma\mathcal{V}(S_{t+1})$ ), ולא לעדכן את השערוך הנוכחי כך שיהיה דומה לקודם. בכדי לדאוג לכך, נתייחס לעתיד כקבוע ולא נחשב עבורו גרדיאנט (למרות ש- $\mathcal{V}$  מופיע בו). באופן פורמלי נוכל לנסח זאת כך:

$$T(S_t) = \mathbb{E}_\pi[R_{t+1} + \gamma\mathcal{V}(S_{t+1})|s_t = S_t]$$

$$L = \frac{1}{2} (T(S_t) - \mathcal{V}(S_t))^2$$

ואז הגרדיאנט יהיה:

$$\frac{\partial L}{\partial \mathcal{V}} = T(S_t) - \mathcal{V}(S_t)$$

בהתאם לכך, הדגימה  $R_{t+1} + \gamma\mathcal{V}(S_{t+1}) - \mathcal{V}(S_t)$  הינה שערוך סטוכסטי לגרדיאנט. כיוון שכל צעד תלוי בצעד הבא, אז המטרה  $T(S_t)$  משתנה בכל צעד (אמנם קיבענו אותה בכל צעד יחיד, אך היא עדיין תלויה ב- $\mathcal{V}(S_t)$ ). במובן הזה בעיית הרגרסיה שהגדרנו הינה "דינמית", כיוון שהמטרה משתנה בכל צעד.

### 11.3.2 TD( $\lambda$ )

ננסה לבחון את הקשר בין שתי השיטות שראינו. שערוך TD מבצע דגימה של ריצה ואז מעדכן את הערך של כל מצב בהתאם למצב הבא בלבד:

$$\mathcal{V}(S_t) = \mathcal{V}(S_t) + \alpha \cdot [R_{t+1} + \gamma\mathcal{V}(S_{t+1}) - \mathcal{V}(S_t)]$$

בגלל ההתייחסות למצב אחד בכל פעם השונות של המשערך נמוכה, אך יש הטיה.

שיטת MC לעומת זאת דוגמת ריצה ומעדכנת את הערך של כל המצב בהתאם לכל המצבים שבאים לאחר מכן:

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$$

$$\mathcal{V}(s_t) = \mathcal{V}(s_t) + \frac{1}{N(S_t)} (G_t - \mathcal{V}(s_t))$$

משערך זה הינו חסר הטיה, אך עם זאת השונות שלו גבוהה.

נראה כיצד ניתן לחבר בין שני המשערכים ולמצוא שיטה שתהיה אופטימלית מבחינת היחס שבין ההטיה לשונות. ניתן להכליל את שני המשערכים לנוסחה כללית באופן הבא:

$$\mathcal{V}(s_t) = \mathcal{V}(s_t) + \alpha (G_t^{(n)} - \mathcal{V}(s_t))$$

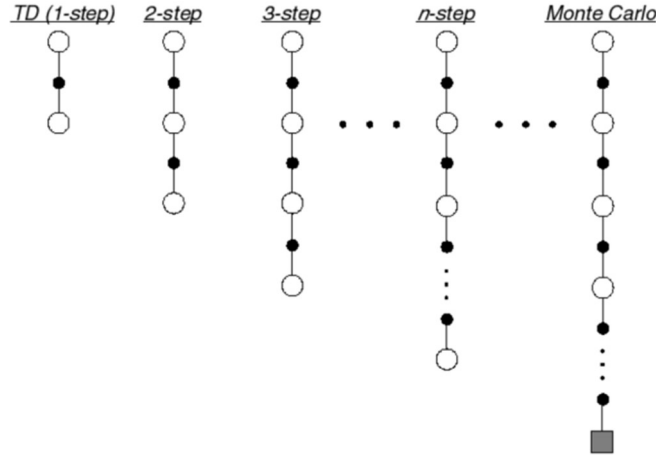
כאשר  $G_t^{(n)}$  הוא הסכום של  $n$  ה-rewards הבאים:

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_n = \sum_{k=0}^{n-t-1} \gamma^k R_{t+k+1}$$

כעת נוכל להגדיר:

$$TD(n) \rightarrow \mathcal{V}(s_t) = \mathcal{V}(s_t) + \alpha (G_t^{(n)} - \mathcal{V}(s_t))$$

ולפי סימון זה נוכל לשים לב שמשערך MC הוא למעשה  $TD(n \rightarrow \infty)$  ואילו משערך TD שווה ל- $TD(1)$ .



איור 11.8 משערך  $TD(n)$  MC הינו המקרה הפרטי עבורו  $n \rightarrow \infty$  ואילו משערך TD הינו המקרה הפרטי בו  $TD(n = 1)$ .

אם ניקח  $1 < n < \infty$  נקבל משערך "ממוצע" בין MC לבין TD. ניתן להציע גרסה יותר טובה למשערך זה תחת ההנחה שכלל שמאורעות סמוכים אחד לשני כך יש להם יותר השפעה. בדומה ל-discount factor שמוריד את ההשפעה של תגמול ככל שהוא יותר רחוק מהמצב הנוכחי, כך גם כאן ניתן לכלל מאורע עתידי משקל הולך וקטן. נדאג שכל המשקלים יסתכמו ל-1 ונקבל את השערך הבא, הידוע גם בכינוי  $TD(\lambda)$ :

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

$$\mathcal{V}(s_t) = \mathcal{V}(s_t) + \alpha (G_t^\lambda - \mathcal{V}(s_t))$$

שערך זה הוא בעל שונות גדולה יותר מאשר TD, כיוון שהוא מתחשב ביותר צעדים, אך ההטיה שלו קטנה יותר. כאמור, הוא מנסה למצב בין שני המשערכים שראינו ולאזן בין השונות להטיה.

נסכם בקצרה את מה שראינו בפרק זה. התייחסנו למקרים בהם המודל אינו ידוע ואנו רוצים לשערך אותו, והצגנו שלוש גישות המנסות לשערך את המודל בעזרת דגימות סטוכסטיות: TD, MC, (שזהו מקרה פרטי של  $TD(n)$ ) וגישה המנסה לשלב בין השתיים –  $TD(\lambda)$ .

לסיום נעיר ששערך המודל כשלעצמו אינו מספיק, כיוון שהוא מספק מודל עבור הבעיה, אך הוא אינו בהכרח אופטימלי. בפרק הקודם ראינו כיצד בהינתן המודל ניתן לשפר אותו ולמצוא את האסטרטגיה האופטימלית. יכולנו לעשות זאת כיוון שידענו את המודל במלואו, מה שאיפשר לבצע בכל צעד מהלך חמדני ובכך להתכנס לבסוף לאסטרטגיה האופטימלית. במקרים בהם אנו רק משערכים את המודל, אין לנו בהכרח מידע מספיק טוב עבור כל המצבים. מצבים ופעולות שלא נוסו כלל (או נוסו רק בפעמים נדירות), השערך עבורם יכול להיות די גרוע, ואז השיפור באמצעות בחירה גרידית לא בהכרח יכול להביא את האסטרטגיה להיות אופטימלית.

### 11.3 Model Free Control

בפרק הקודם ראינו כיצד ניתן לשערך את המודל באמצעות דגימות סטוכסטיות. שיטות השערך אפשרו לנו לקבל מידע על המודל, אך הן אינן התייחסו לשאלה האם הוא אופטימלי. בפרק הראשון ראינו כיצד ניתן לקחת מודל או אסטרטגיה ולהביא אותם לאופטימליות, אך אי אפשר להשתמש בשיטה זו עבור מודל משוערך. הסיבה לכך נעוצה

שמודל זה לא בהכרח ראה את כל הזוגות האפשריים של המצבים והפעולות, וממילא הוא לא יכול לשפר את הבחירות שלו על סמך אסטרטגיה גרידית. ניקח לדוגמה מקרה בו עבור מצב מסוים האסטרטגיה הינה דטרמיניסטית והפעולה שנבחרת הינה תמיד ללכת למעלה. במקרה כזה אין לנו שום מידע על יתר הפעולות האפשריות במצב זה ולכן המודל שלנו לא שלם, וממילא לא נוכל להשתמש ב-Policy improvement המבוסס על כך שיש לנו מידע על כל המצבים והפעולות.

בכדי להתמודד עם בעיה זו נהיה חייבים לדאוג לכך שנבקר בכל המצבים. כמובן שנוכל לנקוט בגישה פשטנית של אסטרטגיה אקראית, שאחרי מספר מספיק גדול של פעולות קרוב לוודאי שנבקר בכל המצבים. אסטרטגיה זו אמנם טובה לבדיקת מצבים ופעולות חדשים, אך כמובן שהיא רחוקה מלהיות אופטימלית, לכן נרצה להשתמש באסטרטגיה שמצד אחד מנסה להיות אופטימלית ומצד שני יש בה ממד של אקראיות המביא לכך שנבקר גם במצבים שלא היינו מגיעים אליהם לפי האסטרטגיה הנוכחית. בחירת פעולות בהתאם לאסטרטגיה הנוכחית נקראת **Exploitation** ואילו בחירה של מצבים חדשים נקראת **Exploration**, והמטרה שלנו תהיה לאזן בין השניים תחת הדרישות הבאות:

1. ביקור בכל הזוגות של המצבים והפעולות אינסוף פעמים.
2. ככל שמספר הצעדים גדל, כך האסטרטגיה מתכנסת לאסטרטגיה הגרידית:

$$\lim_{k \rightarrow \infty} \pi_k(a|s) = 1 \left[ a = \arg \max_{a'} Q(s, a') \right]$$

מצד אחד נרצה לבקר בכל המצבים ומצד שני נרצה שכל שנעשה יותר צעדים ככה נשפר את האסטרטגיה שלנו. אסטרטגיה המקיימת דרישות אלה נקראת **Greedy in the Limit of Infinite Exploration (GLIE)**, וניתן להוכיח שקיום דרישות אלה מביא את האסטרטגיה להיות אופטימלית. דוגמה לאסטרטגיה פשוטה העונה על הדרישות הינה  $\epsilon - greedy$  – בחירה של האסטרטגיה האופטימלית בהסתברות  $(1 - \epsilon)$  ובחירה מצב אחר בהסתברות  $\epsilon$ . עבור  $\epsilon$  שהולך וקטן עד ל-0 בקצב שאינו מהיר מדי, אסטרטגיה זו הינה GLIE.

לעיל הראינו כיצד בהינתן מודל ניתן לשפר את האסטרטגיה (Policy improvement) על ידי כך שבחרנו באופן גרידי פעולות. כעת נרצה להראות שבאופן דומה מתקיים אותו רעיון גם עבור  $\epsilon - greedy$ , כלומר שאם השתמשנו בשיטה זו והגענו ל-Value function, אז ניתן בצעד הבא לשפר את ה-Value על ידי אסטרטגיה  $\epsilon$ -גרידית. אסטרטגיה זו בוחרת באופן גרידי את הפעולה האופטימלית לפי האסטרטגיה הנתונה, אך נותנת לכל יתר הפעולות הסתברות הגדולה או שווה להסתברות שהייתה לפעולה זו בשימוש באסטרטגיה  $\epsilon$ -גרידית. ננסח את המשפט באופן פורמלי:

$$\text{If Policy } \pi_i \text{ has } \forall s, a: \pi_i(s, a) \geq \frac{\epsilon}{|A|} \text{ and is } \epsilon - \text{greedy w.r.t } Q^{\pi_i}, \text{ then } \mathcal{V}^{\pi_{i+1}} \geq \mathcal{V}^{\pi_i}$$

הוכחה: נסתכל על המקרה בו נוקטים צעד אחד גרידי לפי  $\pi_{i+1}$  ואז ממשיכים לפי האסטרטגיה הקודמת  $\pi_i$ :

$$\sum_a \pi_{i+1}(a|s) Q^{\pi_i}(s, a) = \sum_a \frac{\epsilon}{|A|} Q^{\pi_i}(s, a) + (1 - \epsilon) \max_a Q^{\pi_i}(s, a)$$

נרשום את  $(1 - \epsilon)$  בצורה אחרת: במקום  $\epsilon$  נרשום  $\sum_a \frac{\epsilon}{|A|}$ , ובמקום 1 נרשום  $\sum_a \pi_i(a|s)$  (הסכום אכן שווה ל-1 כי סוכמים את כל האפשרויות עבור התפלגות נתונה). נקבל:

$$= \sum_a \frac{\epsilon}{|A|} Q^{\pi_i}(s, a) + \max_a Q^{\pi_i}(s, a) \sum_a \left( \pi_i(a|s) + \frac{\epsilon}{|A|} \right)$$

כעת נחליף את הביטוי  $\max_a Q^{\pi_i}(s, a)$  באסטרטגיה עצמה  $Q^{\pi_i}(s, a)$ . כיוון שמתקיים  $Q^{\pi_i}(s, a) \leq \max_a Q^{\pi_i}(s, a)$ , נקבל:

$$\geq \sum_a \frac{\epsilon}{|A|} Q^{\pi_i}(s, a) + Q^{\pi_i}(s, a) \sum_a \left( \pi_i(a|s) + \frac{\epsilon}{|A|} \right)$$

כעת יש שני איברים זהים שמצטמצמים, ונשאר עם:

$$= \sum_a \pi_i(a|s) Q^{\pi_i}(s, a)$$

ובסך הכל קיבלנו:

$$\sum_a \pi_{i+1}(a|s) Q^{\pi_i}(s, a) \geq \sum_a \pi_i(a|s) Q^{\pi_i}(s, a)$$

כלומר, לעשות צעד אחד לפי האסטרטגיה  $\pi_{i+1}$  ואז להמשיך לפי  $\pi_i$  בהכרח יותר טוב מאשר גם את הצעד הראשון לעשות לפי  $\pi_i$ . כעת בדומה להוכחה שהראינו לעיל, ניתן להוכיח שמתקיים  $V^{\pi_i}(s) \leq V^{\pi_{i+1}}(s)$  וביוצא השיפור שוב ושוב יביא את האסטרטגיה להתכנס לזו האופטימלית. הבעיה בשיטה זו היא חוסר היעילות שבה, כיוון שהיא דורשת המון דגימות והמון איטרציות. עקב כך שיטה זו לא פרקטית, ובמקומה נציג כעת שיטות אחרות המאפשרות לשפר את האסטרטגיה עבור מודל משוער. פורמלית, שיטות אלה נכנסות תחת קטגוריה הנקראת Model Free Control – שליטה (ושיפור) אסטרטגיה שאינה מתבססת על מודל ידוע מראש.

### 11.3.1 SARSA – On-Policy TD control

בפרק הקודם ראינו כיצד ניתן באמצעות שערך TD לשפר את ה-Value function, כאשר העדכון בכל צעד מתקיים באופן הבא:

$$V(S_t) = V(S_t) + \alpha \cdot [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

כעת אנחנו מתעניינים באסטרטגיה ולא רק ב-Value function, לכן במקום לעדכן את  $V(S_t)$  נעדכן את פונקציית ה-state-action  $Q(S_t, A_t)$ . באופן הזה נקבל טבלה בגודל  $|S| \times |A|$ : המכילה מידע על כל הזוגות האפשריים של  $(S, A)$ , ועבור כל זוג יש ערך מסוים (טבלה זו נקראת  $Q$ -table). בהתחלה הערכים בטבלה לא יתקנו את הערכים האמיתיים, אך עם התקדמות הלמידה הטבלה תשתפר ואולי אף תתכנס לאופטימליות. העדכון מתבצע באופן הבא:

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha \cdot [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

העדכון בכל צעד מתבצע על סמך המצבים והפעולות של שתי יחידות זמן:  $\{S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}\}$  ולכן האלגוריתם נקרא SARSA. אלגוריתם זה הינו **On-Policy Learning**, כלומר, העדכון בכל צעד נעשה על סמך מידע המגיע מהאסטרטגיה הידועה באותו זמן: בוחרים לבצע פעולה  $A_t$  במצב  $S_t$  ( $Q(S_t, A_t)$ ) בהתאם לאסטרטגיה, ואז מעדכנים אותה על סמך התגמול  $R_{t+1}$  שהתקבל בעקבות הפעולה  $A_t$ . באופן סכמתי איטרציה אחת של האלגוריתם מתוארת באופן הבא:

```

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat (for each episode):
  Initialize  $S$ 
  Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
  Repeat (for each step of episode):
    Take action  $A$ , observe  $R, S'$ 
    Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
     $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$ 
     $S \leftarrow S'; A \leftarrow A'$ 
  until  $S$  is terminal
  
```

איור 11.7 אלגוריתם SARSA. שערך on-policy של  $Q(s, a)$  בעזרת דגימות מאסטרטגיה  $\epsilon$ -גרדית ביחס ל- $Q$  הנוכחי.

בכל מצב הפעולה שנבחרת הינה  $\epsilon$ -גרדית ביחס ל- $Q(s, a)$  הנוכחי. כלומר: אם במצב  $S_t$  יש לנקוט לפי האסטרטגיה את המצב  $A_t = \bar{A}$ , אזי האסטרטגיה ה- $\epsilon$ -גרדית ביחס לכך הינה:

$$A_t = \begin{cases} \bar{A} & \text{w.p. } 1 - \epsilon \\ A \neq \bar{A} & \text{w.p. } \epsilon \end{cases}$$

ניתן להוכיח שאלגוריתם SARSA מביא את האסטרטגיה לאופטימליות תחת שני תנאים:

א. שהאסטרטגיה תהיה GLIE.

ב. שיתקיים תנאי Robbins-Monroe עבור  $\alpha$  (תנאי זה דואג לכך שנגיע בהכרח לכל המצבים ומצד שני  $\alpha_i \rightarrow 0$ ):

$$\sum_{i=0}^{\infty} \alpha_i = \infty, \sum_{i=0}^{\infty} \alpha_i^2 < \infty$$

לגישה זו, הפועלת בגישת on-policy learning, יש מספר חסרונות:



1. המטרה היא ללמוד את האסטרטגיה האופטימלית אבל בפועל ה-exploration הוא ביחס לאסטרטגיה הנתונה בכל מצב.
2. לא ניתן להשתמש בצעדים ישנים, כיוון שהם מתייחסים לאסטרטגיה שכבר לא רלוונטית.
3. לא ניתן להשתמש במידע שמגיע מבחוץ.

### 11.3.2 Q-Learning

ניתן להפוך את אלגוריתם SARSA להיות off-policy, והאלגוריתם המתקבל, שנקרא Q-Learning, הוא אחד האלגוריתמים השימושיים בתחום של RL. נתבונן בפונקציית העדכון של SARSA:

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha \cdot [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

האלמנט שתלוי באסטרטגיה הינו  $Q(S_{t+1}, A_{t+1})$ , כיוון שבו אנחנו נוקטים בפעולה  $A_{t+1}$  בהתאם לאסטרטגיה. במקום לבחור בפעולה זו, ניתן להיות גרידי ולקחת את הפעולה בעלת הערך הכי גדול בצעד הקרוב, ועל פיה לעדכן את ה- $Q$  - value:

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha \cdot [R_{t+1} + \gamma \max_A Q(S_{t+1}, A) - Q(S_t, A_t)]$$

כעת האסטרטגיה אינה משפיעה על פונקציית העדכון, וממילא היא off-policy.

```

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat (for each episode):
  Initialize  $S$ 
  Repeat (for each step of episode):
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Take action  $A$ , observe  $R, S'$ 
     $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
     $S \leftarrow S'$ ;
  until  $S$  is terminal
  
```

איור 11.8 אלגוריתם Q-Learning. שערור off-policy של  $Q(s, a)$  בעזרת דגימות מאסטרטגיה  $\epsilon$ -גרידית ועדכון ה- $Q$  - value בהתאם לפעולה בעלת הערך הכי גדול בצעד הקרוב.

האסטרטגיה לא משפיעה על עדכון ה- $Q$  - value, אך כן יש לה השפעה על כמות הפעמים שנבקר בכל מצב. בניגוד לאלגוריתם SARSA בו דרשנו שהאסטרטגיה תהיה GLIE, באלגוריתם Q-Learning הדרישה הינה רק שנבקר אינסוף פעמים בכל מצב. הבדל זה הוא משמעותי, כיוון ש-GLIE דורש שהאסטרטגיה תתכנס לאסטרטגיה הגרידית (כלומר, הדרישה היא ש- $\epsilon$  ילך ל-0). דרישה זו מקשה על הלמידה כיוון שצעדים גרידים מביאים מעט מאוד מידע חדש. הסרת הדרישה על ה- $\epsilon$  מאפשרת exploration בצורה יותר חופשית, והלמידה נעשית בצורה הרבה יותר מהירה. עם זאת, יותר מדי exploration זה גם לא טוב, כיוון שמבקרים בהרבה מצבים לא רלוונטיים מספר רב של פעמים.

נתבונן על האלגוריתמים שראינו עד כה ונשווה בין הפתרונות שהיו מבוססים על ידיעת המודל לבין פתרונות משוערכים:

Full Backup (Dynamic Programming)	Sample Backup (TD)
Iterative Policy Evaluation: $V(S) = \mathbb{E}[R + \gamma \cdot V(S')   S]$	TD Learning: $V(S) = V(S) + \alpha \cdot [R_{t+1} + \gamma V(S') - V(S)]$
Q-Policy Iteration: $Q(S, A) = \mathbb{E}[R + \gamma \cdot Q(S', A')   S, A]$	SARSA: $Q(S, A) = Q(S, A) + \alpha \cdot [R + \gamma Q(S', A') - Q(S, A)]$
Q-Value Iteration: $Q(S, A) = \mathbb{E}[R + \gamma \cdot \max_A Q(S, A)   S, A]$	Q-Value Iteration: $Q(S, A) = Q(S, A) + \alpha \cdot [R + \gamma \max_A Q(S', A) - Q(S, A)]$

הפתרונות המשוערכים מצליחים להתמודד עם בעיות בינוניות, אך הצורך בדגימות יכול להיות בעייתי משתי סיבות:  
א. האלגוריתמים נדרשים להמון דגימות בשביל להצליח. ב. חסרון נוסף בגישות שאינן model-based קשור ל-

exploration עבור בעיות של העולם האמיתי, לא תמיד אפשרי לראות דוגמאות שליליות כדי ללמוד שהן לא טובות. רכב אוטונומי למשל, אם רוצים שהוא ילמד לא לנסוע ברמזור אדום, אי אפשר לאמן אותו על ידי זה שניתן לו יד חופשית לבצע exploration וכך הוא יגיע למצבים בהם הוא ייסע באור אדום ויקבל על כך תגמול שלילי. לעומת זאת, בעיות הכרוכות בסימולציה הן יותר מתאימות לאלגוריתמים שהצגנו המבוססים על דגימות ושערוך, כיוון שניתן בצורה יחסית זולה לבצע המון דגימות, ובנוסף אין בהן בעיות בטיחות בביצוע ה-exploration.

מלבד בעיית הדגימות והשערוך, האתגר העיקרי של גישות אלה נעוץ ביכולת ההכללה של המודלים הנלמדים. התוצאה של SARSA ו-Learning  $Q$  – הינה כאמור טבלה של ה-value  $Q$ , ובטבלה זו אין שום קשר בין הערכים השונים בטבלה. כל איבר בטבלה עומד בפני עצמו, ואי אפשר ללמוד ממנו על שאר האיברים. אם למשל הגענו למצב חדש שעוד לא ראינו אך ראינו הרבה מצבים דומים לו, לא נוכל ללמוד מהם שום דבר לגבי המצב החדש. אתגר זה משליך גם על גודל הבעיות אותן ניתן לפתור – אם אי אפשר להכליל ממצב אחד למצב אחר, ממילא זה מגביל מאוד את גודל הבעיה איתה ניתן להתמודד בעזרת אלגוריתמים אלו. במקרים בהם מרחב המצבים הוא רציף, אז מרחב המצבים הוא אינסופי ואז בכלל לא ניתן להשתמש בגישות אלו.

### 11.3.3 Function Approximation

כאמור, אלגוריתמים שמנסים לשערך את ה-table  $Q$  אינם ישימים בבעיות גדולות בעיקר בגלל חוסר היכולת שלהם להכליל ממצב אחד למצב אחר. גם אם יש בידינו אפשרות לשמור טבלאות ענקיות של ה-value  $Q$ , לא נוכל לשפר את הטבלה ולעדכן בה ערכים אופטימליים, כיוון שלא ניתן להשליך ממצבים בהם ביקרנו על מצבים חדשים. בכדי להתמודד עם בעיה זו, נרצה להחליף את ה-table  $Q$  בפונקציה שמחזירה ערך עבור כל זוג של state-action. במקום לחשב טבלה ענקית  $|Q| = |S| \times |A|$ , נרצה למצוא פונקציה עם סט פרמטרים  $\theta$  כך שיתקיים:

$$Q(S, A) \approx Q_\theta(S, A), V(S) \approx V_\theta(S)$$

היתרון של שימוש בפונקציה שמנסה לשערך את הערכים הינו כפול: א. לא צריך לשמור טבלה בגודל של מרחב המצבים. ב. כן ניתן ללמוד ממצבים שיש לנו עליהם ידע על מצבים חדשים. בכדי למצוא פרמטרים שישערכו את המודל בצורה איכותית יש לפתור בעיית אופטימיזציה, כפי שנגדיר בהמשך, אך הלמידה יכולה להיות מאוד לא יציבה. ראשית, כפי שראינו ב-TD, המטרה אותה רוצים לשערך ( $V(S_t)$  או  $Q(S_t, A_t)$ ) משתנה בכל צעד. בנוסף, בניגוד לאלגוריתמים הקודמים, כעת יש קשר בין המצבים, ולכן אם אנחנו משנים ערך מסוים, זה גם ישפיע על ערכים אחרים, מה שמקשה מאוד על יציבות הלמידה.

הדוגמה הפשוטה ביותר לפונקציה כזו הינה מודל לינארי מהצורה:

$$Q(S, A) = w^T \phi(S, A)$$

מודל זה מניח שיש סט של פרמטרים המקיים קשר לינארי בין ה-state-action ומפת פיצ'רים כלשהי  $\phi(S, A)$  לבין הערך שלהם  $Q(S, A)$ . באופן הפשטני ביותר, בשביל למצוא את הפונקציה הזו, נרצה למזער כמה שיותר את המרחק שבין  $V^\pi(S)$  לבין  $V_\theta(S)$ , ולשם כך נבנה את פונקציית המטרה הבאה:

$$L(\theta) = \frac{1}{2} \mathbb{E}_\theta \left[ (V^\pi(S) - V_\theta(S))^2 \right]$$

כמובן שחישוב זה אינו אפשרי משום שאנחנו לא יודעים מה הוא  $V^\pi(S)$  – זה בדיוק הביטוי אותו אנו רוצים לשערך! בנוסף, חישוב התוחלת יכול להיות מאוד מסובך, בטח במקרה בו אנו לא יודעים את כל הערכים האמיתיים של  $V^\pi(S)$ . בשביל להתגבר על בעיות אלו נשים לב כיצד כן ניתן לשפר את  $V_\theta(S)$  – אם נגזור את פונקציית המטרה ונבצע צעד בכיוון הגרדיאנט, נקטין את ערכה, ולכן נרצה להיות מסוגלים לחשב את  $\Delta\theta$ . נבצע זאת בדומה בעזרת דגימות סטוכסטיות, כפי שכבר ראינו בפרקים קודמים, כלומר נרצה לדגום מצבים מהאסטרוגיה בשביל לחשב את הביטוי הבא:

$$\Delta\theta = (V^\pi(S_t) - V_\theta(S_t)) \nabla_\theta V_\theta(S_t)$$

עם ביטוי זה עדיין לא ניתן להתקדם כיוון ש- $V^\pi(S)$  לא ידוע, אך גם אותו ניתן לשערך, למשל בעזרת MC או TD. בעזרת שיטת MC ניתן להריץ episode שלם, לחשב את התגמול המצטבר ולשים אותו במקום  $V^\pi(S)$ , ובעזרת TD ניתן להריץ צעד אחד ולהחליף את  $V^\pi(S)$  בתגמול המיידי והשערוך של המצב הבא:

$$\text{MC: } \Delta\theta = (G_t - V_\theta(S_t)) \nabla_\theta V_\theta(S_t)$$

$$\text{TD: } (R_{t+1} + \gamma V_\theta(S_{t+1}) - V_\theta(S_t)) \nabla_\theta V_\theta(S_t)$$

יש לשים לב שבביטוי האחרון יש שני איברים שתלויים ב- $\theta$ ,  $V_\theta(S_t)$  ו- $\gamma V_\theta(S_{t+1})$ , אך נגזור רק את הביטוי הראשון כיוון שנרצה לשנות אותו כך שיתקרב לביטוי השני ולא להיפך. עם זאת, בניגוד למה שראינו לעיל בנוגע ל-TD, שם שינוי של  $V_\theta(S_t)$  לא השפיע על  $V_\theta(S_{t+1})$  כיוון שהערכים בטבלה היו בלתי תלויים אחד בשני, כאן כן יש השפעה בין שני ערכים אלו, מה שמקשה על היציבות של הלמידה.

בסופו של דבר, אנו יכולים לחשב את הביטוי  $\Delta\theta$  ובעזרתו לנסות לאפסם את הפרמטרים באיזה שיטות אופטימיזציה סטנדרטיות (למשל – stochastic gradient descent).

בדומה לשערוך שהצענו עבור MC ו-TD, ניתן גם לבצע control ולשערך את ה-table-Q:

$$\text{SARSA: } \Delta\theta = (R_{t+1} + \gamma Q_\theta(S_{t+1}, A_{t+1}) - Q_\theta(S_t, A_t)) \nabla_\theta Q_\theta(S_t, A_t)$$

$$Q - \text{Learning: } \Delta\theta = \left( R_{t+1} + \gamma \max_A Q_\theta(S_{t+1}, A) - Q_\theta(S_t, A_t) \right) \nabla_\theta Q_\theta(S_t, A_t)$$

וגם כאן, בגזירה נרצה להתייחס ל- $Q_\theta(S_{t+1}, A_t)$  ו- $Q_\theta(S_{t+1}, A_{t+1})$  כקבועים ולגזור רק את  $Q_\theta(S_t, A_t)$ , כיוון שנרצה לשנות את  $Q_\theta(S_t, A_t)$  כך שיתקרב בערכו לשאר הביטויים המבוססים על מידע נוסף.

יש כמה אתגרים שיכולים להקשות מאוד על שיטות אלה להביא את המודל להתכנסות:

א. ראשית, הדרישה שאומרת שאנו רוצים לשנות את  $Q_\theta(S_t, A_t)$  ביחס למידע העתידי הינה בעייתית, כיוון שכעת כל הביטויים תלויים באותו פרמטר  $\theta$ .

ב. באלגוריתם שהוא off-policy אנו דוגמים מאסטרטגיה שאנו לא בהכרח פועלים בהתאם אליה (וממילא היא לא בהכרח האופטימלית עבורנו). כאשר ניסינו ללמוד tabular Q-table, זה היה יכול לגרום לכך ששכיחות המצבים שנדגמו אינה תואמת את השכיחות שלהם לפי האסטרטגיה האמיתית, מה שיכול להשפיע על קצב הלמידה, אך זה לא יגרום לשגיאות באסטרטגיה הנלמדת. כעת כאשר יש קשר בין המצבים על ידי הפרמטר  $\theta$ , אז דגימה לא לפי השכיחות האמיתית יכולה להטות את הלמידה לכיוונים שגויים.

אתגרים אלו משפיעים על יכולת המודל להתכנס לאסטרטגיה האופטימלית. בעוד שעבור יכולנו להוכיח התכנסות, עבור functional-approximation זה כלל לא מובטח, כפי שמוצג בטבלה הבאה:

Algorithm	Tabular	Linear	Non-Linear
MC	Y	Y	Y
SARSA	Y	Y	N
Q-Learning	Y	N	N

אחת השיטות פורצות הדרך התחום הייתה שימוש ברשתות נוירונים בתור ה-functional-approximation, ובשם הנפוץ שלה – Deep Q – Learning. השימוש ברשת נוירונים בפני עצמה לא הייתה מספיקה כיוון שכאמור הלמידה היא מאוד לא יציבה ויש להכניס מנגנונים נוספים שיעזרו למודל להשתפר. נציין שתי טכניקות שהוטעו בתהליך הלמידה:

1. כאשר דוגמים כל מיני מצבים, יש ביניהם קורלציה יחסית גבוהה, מה שמונע מהשונות של הגרדיאנט לקטון. כדי להתגבר על כך, במקום לעדכן את  $Q$  בכל שלב, שומרים את הרביעיות  $(S_t, A_t, R_{t+1}, S_{t+1})$  של  $N$  המצבים האחרונים שראינו (למשל:  $N = 1,000$ ), ואז מעדכנים לפי  $k$  מצבים (למשל:  $k = 10$ ) אקראיים מתוך  $N$  המצבים הקודמים שיש לנו בזיכרון (לאחר שהזיכרון מתמלא, אז כל מצב חדש שמגיע דורס את המצב הכי ישן בזיכרון). בצורה הזו השונות תקטן כי הדוגמאות שנעדכן לפיהן יהיו בעלות קורלציה נמוכה הרבה יותר. המחיר של טכניקה זו הוא השימוש בהרבה זיכרון, אבל יש לה השפעה משמעותית על הביצועים.

2. הזכרנו לעיל שאנו רוצים לשנות את  $Q_\theta(S_t, A_t)$  ביחס למידע העתידי, אך זה יכול להיות בעייתי, כיוון שכעת כל הביטויים תלויים באותו פרמטר  $\theta$ . במילים אחרות – אנו רוצים לקרב את האסטרטגיה למטרה מסוימת, אל שתיהן תלויות באותו פרמטר. בכדי לייצב את הרגרסיה, ניתן לשנות את המטרה באופן הרבה פחות תדיר. באופן פורמלי, נשמור סט פרמטרים  $\hat{\theta}$  ונשנה אותו כל מספר קבוע של צעדים לסט הפרמטרים החדש, כלומר

– כל  $k$  צעדים נעדין את  $\hat{\theta}$  להיות שווה ל- $\theta_t$ . יחד עם הפרמטרים החדשים, נשנה מעט את הביטוי אותו נרצה להביא לאופטימום:

$$\mathbb{E}_{S_t, R_{t+1}, A_t, S_{t+1}} [(R_{t+1} + \gamma \max_A Q_{\hat{\theta}}(S_{t+1}, A_t) - Q_{\theta}(S_t, A_t))$$

### 11.3.4 Policy-Based RL

יש כמה מגבלות וחסרונות בשיטות שראינו עד כה. שיטות אלו מתמקדות בלמידת ה- $Q$ -table, בין אם על ידי עדכון ישיר של ערכי הטבלה ובין אם על ידי למידת הפונקציה  $Q_{\theta}(S, A)$  התלויה בפרמטר  $\theta$ . ראינו שבבעיות גדולות הייצוג והלמידה יכולים להיות קשים. אתגר זה מתעצם ואף נהיה בלתי אפשרי עבור בעיות רציפות, בהן מרחב המצבים הוא כביכול אינסופי (למשל – תנועה של רובוט). חסרון נוסף שיש בגישות שראינו נעוץ במטרה של הלמידה – אנו מנסים ללמוד את  $Q_{\theta}(S, A)$  אך בפועל מה שמעניין אותנו זה  $\max_A Q_{\theta}(S, A)$ , מה שיכול להביא לבזבז משאבים עבור למידה של דברים לא הכרחיים. למשל – אם אנחנו צריכים להחליט בין להשקיע כסף בשוק ההון לבין השקעה בתרמית פירמידה, אז אנחנו רק צריכים להחליט איזו אופציה עדיפה לנו, אך אנו לא נדרשים לדעת מה בדיוק כוללת כל אופציה. ההחלטה בין האופציות היא פשוטה, ואילו ללמוד מה ואיך להשקיע בשוק ההון זה דבר מאוד מורכב. לכן במקום ללמוד את כל המערכת, נוכל פשוט ללמוד מה הפעולות הנדרשות עבורנו, וגישה ישירה זו יכולה לחסוך במשאבים.

כאמור, במקום ללמוד את  $Q$ , ננסה ללמוד באופן ישיר את האסטרטגיה האופטימלית –  $\pi_{\theta}(a|s)$ . באופן מעט יותר פורמלי נגדיר את פונקציית המטרה:

$$J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[ \sum_{t=0}^{\infty} \gamma^t R_{t+1} \right]$$

מציאת האסטרטגיה האופטימלית שקולה למציאת המקסימום של פונקציית המטרה  $J(\theta)$ .

נרצה להשתמש באסטרטגיה סטוכסטית מכמה סיבות. ראשית, הרבה פעמים יותר קל למצוא אופטימום עבור בעיה רציפה מאשר בעיה דיסקרטית, ואסטרטגיה סטוכסטית יכולה להפוך בעיה דיסקרטית לרציפה. שנית, יהיו מקרים בהם נתעניין דווקא באסטרטגיה הסטוכסטית ולא בזו הדטרמיניסטית. בנוסף, הסטוכסטיות מאפשרת exploration, וכפי שראינו זה מאפיין הכרחי בשביל להגיע לכל המצבים. במבט יותר רחב, אסטרטגיה סטוכסטית היא הכללה של אסטרטגיה דטרמיניסטית – אם כל ההסתברויות של המודל הן 0 או 1. במובן הזה, השימוש באסטרטגיה סטוכסטית לא ממעיט מהיכולת של אסטרטגיה דטרמיניסטית אלא הוא רק מכליל אותה.

באופן הכי פשוט, עבור בעיה בה מרחב הפעולות הינו דיסקרטי, ניתן להפוך אסטרטגיה שמספקת ערך לכל פעולה  $\phi(a; x, \theta)$  להתפלגות על ידי SoftMax:

$$\pi_{\theta}(a|s) = \frac{\exp(\phi(a; x, \theta))}{\sum_{a'} \exp(\phi(a'; x, \theta))}$$

אם מרחב הפעולות הוא רציף, ניתן להשתמש בגאוסיאן:

$$\pi_{\theta}(a|s) \sim \mathcal{N}(\mu(x; \theta), \Sigma(x; \theta))$$

ואם רוצים לאפשר יותר חופש פעולה, ניתן להשתמש ב-Gaussian mixture model. הפונקציה שבאמצעותה רוצים לייצג את האסטרטגיה יכולה להיות כרצוננו – פונקציה ליניארית, רשת נוירונים וכו'.

רוב האלגוריתמים שעושים אופטימיזציה ל- $J(\theta)$  מבוססי גרדיאנט, אך ישנם גם אלגוריתמים אחרים, כמו למשל אלגוריתמים גנטיים, אך עליהם לא נדבר. הדרך הסטנדרטית לבצע אופטימיזציה היא להשתמש ב-gradient descent/ascent ובפיתוחים שלו (כמו למשל stochastic gradient descent/ascent), והמוקד שלנו יהיה בשאלה כיצד לשערך את הגרדיאנט ופחות נתעסק בשיטות האופטימיזציה השונות. האלגוריתם הבסיסי מבוסס על עדכון מהצורה הבאה:

$$\theta_{t+1} = \theta_t + \alpha \nabla f(\theta)$$

נניח ויש לנו episodic MDP, ומהאסטרטגיה הנתונה נוכל לדגום  $N$  מסלולים  $\tau_1, \dots, \tau_N \sim \pi_{\theta}$ , כאשר כל מסלול הינו  $G_i = \sum_{t=0}^{T(i)} \gamma^t R_{t+1}^{(i)}$ . או בקצרה:  $\tau_i = (S_0^i, A_0^i, R_1^i, S_1^i, \dots, R_T^i, S_T^i)$ . שערך מונטה-קרלו הרגיל ישערך את המטרה באופן הבא:

$$\mathcal{J}(\theta) = \mathbb{E}_{\pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{R}_{t+1} \right] \approx \frac{1}{N} \sum_i G_i$$

הבעיה היא שהאפקט של הפרמטר  $\theta$  על המסלול  $G_i$  היא לא ישירה – אנחנו דוגמים מסלול מאסטרטגיה  $\pi_\theta$  התלויה בפרמטר  $\theta$ , ואז מחשבים את התגמול של המסלול וממנו מרכיבים את  $G$ . לכן אי אפשר פשוט לגזור את  $G_i$  לפי  $\theta$ , שהרי הקשר ביניהם הוא עקיף ותלוי ב- $f(\theta)$ . בכדי להתמודד עם בעיה זו נשתמש ב-log-derivative trick (לעיתים התהליך הזה נקרא Reinforce). נגדיר:

$$g(\theta) = \mathbb{E}_{\pi_\theta} [f(x)]$$

ואז הנגזרת הינה:

$$\nabla_\theta g = \mathbb{E}_{\pi_\theta} [f(x) \nabla \log(\pi_\theta(x))]$$

הוכחה:

$$\nabla_\theta g = \nabla_\theta \mathbb{E}_{\pi_\theta} [f(x)] = \nabla_\theta \int \pi_\theta(x) f(x) dx$$

נחליף את הסדר של הגרדיאנט והאינטגרל:

$$= \int \nabla_\theta \pi_\theta(x) f(x) dx$$

כעת נשתמש בקשר הבא (המתקיים עבור כל פונקציה):  $\frac{\nabla_\theta h(x)}{h(x)} = \nabla_\theta \log(h(x))$ , ונקבל:

$$= \int \nabla_\theta \log(\pi_\theta(x)) \pi_\theta(x) f(x) dx$$

וזוהו בדיוק שווה לתוחלת הבאה:

$$= \mathbb{E}_{\pi_\theta} [f(x) \nabla \log(\pi_\theta(x))]$$

כעת, במקום לחשב את הגרדיאנט של  $G_i$ , נוכל לחשב את התוחלת  $\mathbb{E}_{\pi_\theta} [f(x) \nabla \log(\pi_\theta(x))]$  על ידי דגימות ושערוך מונטה קרלו, רק כעת הנגזרת אותה אנו צריכים לחשב אינה תלויה בפרמטר  $\theta$  ולכן ניתן לחשב אותה בצורה ישירה. בצורה מפורשת, פונקציית המטרה שלנו הינה:

$$\mathcal{J}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau)], G(\tau) = \sum_{t=0}^{T(\tau)} \gamma^t R_{t+1}^{(\tau)}$$

ושימוש ב-log-derivative trick יביא אותנו לביטוי:

$$\nabla \mathcal{J}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau) \nabla \log(\pi_\theta(\tau))]$$

נשים לב שהאסטרטגיה  $\pi_\theta(\tau)$  מורכבת משני אלמנטים – הפעולה שאני בוחר לבצע, והמצב אליו אני מגיע יחד עם התגמול המתקבל מצעד זה. החלק הראשון הוא לגמרי בשליטה שלנו ואנו יודעים אותו, אך החלק השני הוא דינמי ותלוי בסטוכסטיות של המערכת, ולכן לכאורה הוא לא ידוע ולא יאפשר לבצע את הגזירה. אך אם נפתח את הביטוי שקיבלנו נראה שכל החלקים התלויים בדינמיות של המערכת אינם תלויים ב- $\theta$  ויתאפסו בגזירה. ניתן לרשום את האסטרטגיה כמכפלת ההסתברויות של כל הצעדים באופן הבא:

$$\pi_\theta(\tau_i) = P(S_0^i) \cdot \pi_\theta(A_0^i | S_0^i) \cdot P(S_1^i, R_1^i | S_0^i, A_0^i) \cdots \pi_\theta(A_{T^i-1}^i | S_{T^i-1}^i) \cdot P(S_{T^i}^i, R_{T^i}^i | S_{T^i-1}^i, A_{T^i-1}^i)$$

כעת אם נוציא לוג, המכפלה תהפוך לסכום:

$$\log(\pi_\theta(\tau_i)) = \log(P(S_0^i)) + \sum_{t=0}^{T^i-1} \log(\pi_\theta(A_t^i | S_t^i)) + \sum_{t=0}^{T^i-1} \log(P(R_{t+1}^i, S_{t+1}^i | S_t^i, A_t^i))$$

הביטוי האמצעי מבטא את בחירות הצעדים של המשתמש, והוא היחיד שתלוי בפרמטר  $\theta$ . שני הביטויים הנוספים מבטאים את הדינמיות של המערכת שאין לנו מידע כלפיהם, והם אינם תלויים ב- $\theta$ , לכן בנגזרת לפי  $\theta$  הם מתאפסים.

עובדה זה בעצם מספקת יתרון נוסף לשימוש בטריק של הלוג – מלבד האפשרות לחשב את הנגזרת באופן ישיר, הלוג גם מפריד בין הצעדים של המשתמש לבין הדינמיות הלא ידועה של המערכת. לכן בסך הכול נקבל:

$$\nabla_{\theta} \log(\pi_{\theta}(\tau_i)) = \sum_{t=0}^{T^i-1} \nabla_{\theta} \log(\pi_{\theta}(A_t^i | S_t^i))$$

כעת נציב את זה בגרדיאנט של פונקציית המטרה ונשתמש בדימויות מונטה קרלו:

$$\begin{aligned} \nabla J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}} [G(\tau) \nabla \log(\pi_{\theta}(\tau))] \approx \frac{1}{N} \sum_i G_i \nabla_{\theta} \log(\pi_{\theta}(\tau_i)) \\ &= \frac{1}{N} \sum_i \left( \sum_{t=0}^{T^i-1} \gamma^t \mathcal{R}_{t+1}^i \right) \left( \sum_{t=0}^{T^i-1} \nabla_{\theta} \log(\pi_{\theta}(A_t^i | S_t^i)) \right) \end{aligned}$$

חישוב הגרדיאנט בצורה זו נקרא בשם הפופולרי **Policy Gradient**. כעת משחישבנו את הגרדיאנט, נוכל להשתמש בשיטות אופטימיזציה סטנדרטיות על מנת למצוא את  $\theta$  האופטימלי. בקצרה נוכל לסכם את כל התהליך ה-Reinforce בשלושה שלבים:

1. Run the policy and sample  $\{\tau^i\}$  from  $\pi_{\theta}(A_t | S_t)$ .
2. Estimate gradients:  $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_i \left( \sum_{t=0}^{T^i-1} \gamma^t \mathcal{R}_{t+1}^i \right) \left( \sum_{t=0}^{T^i-1} \nabla_{\theta} \log(\pi_{\theta}(A_t^i | S_t^i)) \right)$ .
3. Improve the policy using gradient descent/ascent:  $\theta \leftarrow \theta_{\alpha} \nabla_{\theta} J(\theta)$ .

נתבונן על הביטוי של הנגזרת שקיבלנו ונסה לתנח אותו. נניח ויש מומחה שיודע לייצר trajectories טובים (-imitation learning) – במקרה כזה נרצה לקחת את מה שהוא מייצר ולגרום לכך שהאסטרטגיה שלנו תייצר בהסתברות יותר גבוהה דגימות כמו שקיבלנו מומחה. במילים אחרות, כאשר נתונים לנו מסלולים מוצלחים, נרצה לגרום לאסטרטגיה שלנו "לחקות" את אותם מסלולים. הדרך הסטנדרטית לבצע כזו משימה היא בעזרת maximum likelihood – לקחת את המסלולים האלה ולסכום אותם, ואז לחפש את המקסימום באמצעות גזירה של הביטוי המתקבל. באופן שקול ניתן גם לחפש מקסימום עבור הסכום של **לוג** המסלולים, ובסך הכל נרצה לחשב את הביטוי:

$$\sum_{t=0}^{T^i-1} \nabla_{\theta} \log(\pi_{\theta}(\tau_i))$$

וזה בדיוק הביטוי השני בגרדיאנט שראינו קודם. בנוסף אליו יש גם את הגורם הראשון שממשקל כל מסלול בהתאם ל-reward, כך שמסלולים בעלי reward גבוה יקבלו משקל גבוה ואילו מסלולים בעלי reward נמוך יקבלו משקל נמוך. יוצא מכך, שהגרדיאנט מנסה להביא לכך שמסלולים "מוצלחים" יופיעו בהסתברות יותר גבוהה.

השערוך של הגרדיאנט באמצעות הביטוי שראינו הוא אמנם חסר הטיה, אך יש לו שונות גבוהה. בנוסף, אם כל הדגימות הם בעלי reward חיובי (או כולן בעלי reward שלילי), אז יהיה קשה להבחין איזה צעדים יותר טובים ומה לא כדאי לעשות. נראה שתי דרכים לשיפור היציבות של השימוש בגרדיאנט:

בסיבתיות (casualty) – כיוון שהצעד  $a_t$  שהתרחש בצעד  $t$  לא יכול להשפיע על כל הצעדים שקרו לפני כן  $\{a_0, \dots, a_{t-1}\}$ , ניתן להחליף את הביטוי  $\sum_{t=0}^{T^i-1} \gamma^t \mathcal{R}_{t+1}$  בביטוי שמתחשב רק בצעדים שקרו החל מצעד  $t$ , כלומר:  $\sum_{k=t}^{T^i-1} \gamma^k \mathcal{R}_{k+1}^i$  – ביטוי זה נקרא גם **Reward to go**. החלפה זו מורידה את השונות כיוון שהיא גורמת לכך שכל צעד יהיה תלוי בפחות איבריים אקראיים. הביטוי שהתעלמנו ממנו הוא  $\sum_{k=0}^{t-1} \gamma^k \mathcal{R}_{k+1}^i$ , וניתן לראות שאכן אין לו השפעה על העתיד.

Baseline – ניתן להראות שהוספה של קבוע ל-reward משפיעה על התוצאה של הגרדיאנט. קבוע מסוים יכול לגרום לכך שנרצה לשפר צעדים מסוימים ולהימנע מצעדים אחרים, ואילו קבוע אחר יכול לגרום לכך שנרצה לשפר דווקא צעדים אחרים. קבוע זה הוא למעשה פרמטר נוסף שניתן לשלוט בו, ועל ידי בחירה מושכלת שלו ניתן להפחית את השונות. באופן פורמלי, את הביטוי  $\sum_{t=0}^{T^i-1} \gamma^t \mathcal{R}_{t+1}^i$  נחליף בביטוי זהה עם תוספת קבוע:

$$\sum_{t=0}^{T^i-1} \gamma^t \mathcal{R}_{t+1}^i - b$$

השאלה היא כיצד לבחור את הקבוע  $b$  בצורה טובה? באופן פשוט ניתן לבחור את הקבוע בתור הממוצע של ה-rewards. בחירה כזו תביא לכך שמה שמתחת לממוצע נרצה להוריד את ההסתברות שלו ומה שמעל הממוצע נרצה לחזק אותו. בחירה כזו אכן מורידה את השונות, אך היא אינה אופטימלית. ניתן לחשב את הקבוע האופטימלי בצורה מדויקת. הנגזרת של פונקציית המטרה יחד עם האיבר הנוסף הינה:

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [(G(\tau) - b) \nabla \log(\pi_\theta(\tau))]$$

הממוצע אינו תלוי ב- $b$ , לכן נוכל לרשום את השונות כך:

$$\begin{aligned} Var &= \mathbb{E}_{\tau \sim \pi_\theta} [(G(\tau) - b)^2 \nabla \log(\pi_\theta(\tau))^2] + C \\ &= \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau)^2 \nabla \log(\pi_\theta(\tau))^2] + \mathbb{E}_{\tau \sim \pi_\theta} [b^2 \nabla \log(\pi_\theta(\tau))^2] - 2 \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau) \cdot b \nabla \log(\pi_\theta(\tau))^2] + C \end{aligned}$$

כעת נגזור את השונות לפי  $b$ :

$$\begin{aligned} \frac{dVar}{db} &= 2 \mathbb{E}_{\tau \sim \pi_\theta} [b \nabla \log(\pi_\theta(\tau))^2] - 2 \mathbb{E}_{\tau \sim \pi_\theta} [G(\tau) \cdot b \nabla \log(\pi_\theta(\tau))^2] = 0 \\ \rightarrow b_{opt} &= \frac{\mathbb{E}_{\tau \sim \pi_\theta} [G(\tau) \nabla \log(\pi_\theta(\tau))^2]}{\mathbb{E}_{\tau \sim \pi_\theta} [b \nabla \log(\pi_\theta(\tau))^2]} \end{aligned}$$

הביטוי המתקבל הוא אכן הממוצע של ה-reward, אך מוכפל במשקל התלוי בנגזרת של לוג האסטרטגיה. השערוך של הביטוי המדויק יכול להיות מסובך ורועש, ולכן לרוב כן משתמשים פשוט ב-reward הממוצע.

אלגוריתם policy gradient הינו on-policy, כלומר האסטרטגיה שבעזרתה אנו דוגמים מסלולים הינה אותה אסטרטגיה שאנו מנסים להביא לאופטימום. כאמור לעיל, באלגוריתמים שהם on-policy, אי אפשר ל"מחזר" דגימות, כלומר בכל נקודת זמן ניתן להשתמש רק בדגימות שנלקחו מהאסטרטגיה הנתונה באותה נקודה. עבור מצבים בהם הדגימות יקרות, נרצה להשתמש באותן דגימות עבור נקודות זמן שונות, כלומר להפוך את האלגוריתם להיות off-policy. בכדי לעשות זאת יש לתקן את השגיאה שנוצרת מכך שבנקודות זמן מסוימת אנו משתמשים בדגימות שנלקחו מאסטרטגיה שכבר השתנתה וכעת ההתפלגות של האסטרטגיה שונה. במילים אחרות – יש לנו דגימות שנלקחו מתוחלת מסוימת, ובעזרתן אנו רוצים לשערך תוחלת אחרת. זו בעיה מתחום הסטטיסטיקה, והפתרון שלה נקרא **Importance sampling**:

נניח ואנו רוצים לשערך את  $\mathbb{E}_p[f(x)]$ , אך אנו לא יכולים לדגום מהתפלגות  $p$  אלא רק מהתפלגות  $q$ , כאשר ההנחה היחידה שלנו היא שאם  $p > 0$  אז גם  $q > 0$ . באמצעות מעבר מתמטי די פשוט ניתן לייצג את התוחלת של  $p$  באמצעות התוחלת של  $q$ :

$$\mathbb{E}_p[f(x)] = \int f(x)p(x)dx = \int f(x) \frac{p(x)}{q(x)} q(x)dx = \mathbb{E}_q \left[ f(x) \frac{p(x)}{q(x)} \right]$$

כעת נדגום  $x_1, \dots, x_N \sim q$  ונשערך את התוחלת באמצעות מוטנה קרלו:

$$\mathbb{E}_p[f(x)] = \mathbb{E}_q \left[ f(x) \frac{p(x)}{q(x)} \right] \approx \frac{1}{N} \sum_i f(x_i) \frac{p(x_i)}{q(x_i)}$$

אם ההתפלגויות  $p(x), q(x)$  הן יחסית דומות, אז התוצאה המתקבלת משערכת בצורה יחסית טובה את התוחלת הרצויה. אם זה לא המצב, השונות תהיה גדולה והתוצאה תהיה לא יציבה. מכל מקום, נוכל להשתמש ברעיון זה עבור שערך הגרדיאנט:

- For trajectory probability the dynamics does not depend on  $\theta$  so

$$\frac{\pi_{\theta'}}{\pi_{\theta}} = \frac{\prod \pi_{\theta'}(a_t | s_t)}{\prod \pi_{\theta}(a_t | s_t)}$$

- Can now Compute  $\nabla_{\theta'} J(\theta')$  with samples from  $\pi_{\theta}$

$$\bullet \sum_{i=1}^N \left( \sum_{t=0}^{T^{(i)}-1} \gamma^t r_{t+1}^{(i)} \right) \left( \sum_{t=0}^{T^{(i)}-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \frac{\prod \pi_{\theta'}(a_t | s_t)}{\prod \pi_{\theta}(a_t | s_t)}$$

- Can improve with causality

$$\bullet \sum_{i=1}^N \sum_{t=0}^{T^{(i)}-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \frac{\prod_{k=0}^t \pi_{\theta'}(a_k | s_k)}{\prod_{k=0}^t \pi_{\theta}(a_k | s_k)} \left( \sum_{k=t}^{T^{(i)}-1} \gamma^k r_{k+1}^{(i)} \right)$$

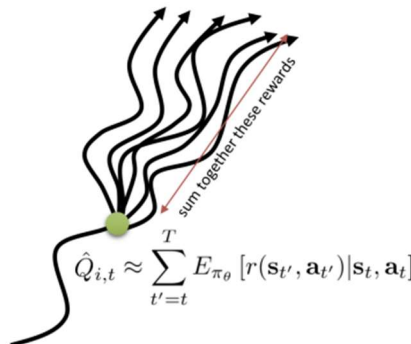
כאמור, שימוש ברעיון זה יכול להיות טוב אם ההתפלגויות  $\theta, \theta'$  יחסית דומות. אם הן רחוקות, התוצאה המתקבלת לא בהכרח מהווה שערך טוב, ולכן הלמידה לא תהיה מספיק טובה. לסיכום – אפשר להפוך את האלגוריתם של policy gradient להיות off-policy באמצעות importance sampling, אך גם שיטה יכולה לסבול משונות גבוהה.

### 11.3.5 Actor Critic

בפרקים הקודמים דיברנו על שתי גישות בכדי לשערך את האסטרטגיה בלי לדעת המודל. גישה אחת התמקדה בשערך ה-Action-Value function (SARSA/Q-Learning), בעוד הגישה השנייה ניסתה לשערך את האסטרטגיה באופן ישיר (Policy gradient). בפרק זה נציג אלגוריתם המנסה לשלב את שתי הגישות יחד. נתבונן שוב בביטוי של הגרדיאנט, כאשר בהתייחסות ל-rewards נתייחס רק ל-go-to-reward, כפי שהוסבר לעיל:

$$\frac{1}{N} \sum_i \left( \sum_{k=t}^{T^{(i)}-1} \gamma^k \mathcal{R}_{k+1}^i \right) \left( \sum_{t=0}^{T^{(i)}-1} \nabla_{\theta} \log (\pi_{\theta}(A_t^i | S_t^i)) \right)$$

ניתן להבחין כי הביטוי  $\sum_{k=t}^{T^{(i)}-1} \gamma^k \mathcal{R}_{k+1}^i$  הוא למעשה דגימה של  $Q^{\pi_0}(s_t, a_t)$ , כלומר עבור נקודת זמן, ניתן לדגום מסלול באמצעות ה-Action-Value function ולקבל סכום של rewards. החיסרון בדגימה זו נעוץ בכך שיש לה שונות גבוהה והיא יכולה להיות מאוד רועשת. כפי שניתן לראות באיור, יתכנו הרבה מסלולים שונים היוצאים מאותה נקודה, ולכן הביטוי של ה-go-to-reward הוא מאוד רועש ביחד ל- $Q^{\pi_0}(s_t, a_t)$ :



איור 11.9 מסלולים שונים אפשריים היוצאים מאותה נקודה. המיצוע/התוחלת של הביטוי אמנם מוריד את השונות ל-0, אך כל דגימה בפני עצמה היא בעלת שונות גבוהה.



הבחנה זו של הקשר בין ה-reward to go לבין ה- $Q^{\pi_0}(s_t, a_t)$  מעלה את השאלה האם אפשר להחליף את בין הביטוי הרועש  $\sum_{k=t}^{T^i-1} \gamma^k \mathcal{R}_{k+1}^i$  לבין הפונקציה עצמה  $Q^{\pi_0}(s_t, a_t)$ . בשביל לבצע החלפה זו, יש להוכיח שהיא אכן חוקית ושומרת על תוצאה נכונה, וכדי לעשות זאת נפתח את פונקציית המטרה כך שתופיע בצורה יותר נוחה. כזכור, פונקציית המטרה אותה אנו מעוניינים להביא לאופטימום הינה התוחלת של ה-rewards (מוכפלים ב-discount factor):

$$\mathcal{J}(\theta) = \mathbb{E}_{\pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{R}_{t+1} \right]$$

המשמעות של התוחלת היא לבצע אינטגרל על כל צמד של action-state  $(s, a)$ , כפול המשקל של אותו צמד. כיוון שצמד יכול לחזור על עצמו בצעדים שונים (כלומר בשתי נקודות זמן  $t_1, t_2$  המערכת יכולה להיות באותו מצב ולבצע את אותו פעולה), אז למעשה באינטגרל ישנם ביטויים שחוזרים על עצמם. בכדי להימנע מכך, ניתן לחשב כל ביטוי באינטגרל פעם אחת, ולהכפיל אותו במשקל המתאים למספר הביקורים באותו צמד action-state. עם זאת, עדיין צריך לזכור שהמופעים של אותו צמד נבדלים זה מזה ב-discount factor, כיוון שהמשקל של reward בזמן  $t_1$  שונה מהמשקל שלו בזמן  $t_2$ , ויש לקחת את ההבדל הזה בחשבון. כעת נתבונן על הביטוי המפורש של האינטגרל המתקבל, כאשר ראשית נחליף את סדר התוחלת והסכימה:

$$\mathcal{J}(\theta) = \mathbb{E}_{\pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{R}_{t+1} \right] = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi_\theta} [\mathcal{R}_{t+1}]$$

התוחלת של כל ה-rewards מורכבת ממיצוע שלהם על פני: (1) ההסתברות להתחיל ממצב  $s_0$  (כפול 2) ההסתברות להגיע למצב  $s$  בזמן  $t$  בהנתן שהתחלנו מ- $s_0$  (כפול 3) הסיכוי לבצע פעולה  $a$  במצב  $s$  שהתקבל בזמן  $t$ . באופן פורמלי הביטוי המתקבל הינו

$$\begin{aligned} &= \sum_{t=0}^{\infty} \gamma^t \int_S \int_S \int_{\mathcal{A}} p(s_0) p_t^\pi(s|s_0) \pi(a|s) \mathcal{R}(s, a) ds_0 ds da \\ &= \int_S \left( \sum_{t=0}^{\infty} \gamma^t \int_S p(s_0) p_t^\pi(s|s_0) ds_0 \right) \int_{\mathcal{A}} \pi(a|s) \mathcal{R}(s, a) ds da \end{aligned}$$

נסמן את הביטוי שבסוגריים כ- $\rho^\pi(s)$ , והמשמעות שלו היא מה הסיכוי להגיע ב- $t$  צעדים ממצב  $s_0$  למצב  $s$ , כאשר עבור כל  $t$  יש להכפיל ב-discount factor המתאים. ביטוי זה נקרא Discount state visitation measure, ולמעשה הוא כולל בתוכו את כל המופעים של צמד  $(s, a)$ , ובכך במקום לעשות אינטגרל על אותו צמד  $(s, a)$  מספר פעמים, עושים זאת פעם אחת עבור כל נקודות הזמן השונות ומכפילים ב-discount factor המתאים. באופן מפורש:

$$\begin{aligned} &= \int_S \rho^\pi(s) \int_{\mathcal{A}} \pi(a|s) \mathcal{R}(s, a) ds da, \\ \rho^\pi(s) &= \sum_{t=0}^{\infty} \gamma^t \int_S p(s_0) p_t^\pi(s|s_0) ds_0 \end{aligned}$$

יש לשים לב שהביטוי הכולל הוא לא בצורה הפשוטה של תוחלת – אינטגרל של הסתברות כפול פונקציית צפיפות, אך הוא עדיין משקף תוחלת ולכן ניתן לשערך אותו באמצעות מונטה-קרלו. כעת נתבונן על הגרדיאנט של פונקציית המטרה. אנחנו רוצים לגזור את הביטוי הבא לפי  $\theta$ :

$$\mathcal{J}(\theta) = \int_S \rho^\pi(s) \int_{\mathcal{A}} \pi(a|s) \mathcal{R}(s, a) ds da$$

כל האיברים תלויים ב- $\theta$  אך ניתן להראות (שקופיות 10-11 במצגת 7) שהנגזרת תלויה רק ב- $\pi(a|s)$ , כלומר:

$$\nabla_\theta \mathcal{J} = \int_S \rho^\pi(s) \int_{\mathcal{A}} \nabla_\theta (\pi(a|s)) \mathcal{Q}^{\pi_\theta}(s, a) ds da$$

כעת ניתן להיעזר ב-log-derivative trick ולקבל:

$$\nabla_{\theta} J = \int_{\mathcal{S}} \rho^{\pi}(s) \int_{\mathcal{A}} \pi(a|s) \nabla_{\pi} \log(\pi(a|s)) Q^{\pi_{\theta}}(s, a) ds da$$

ואז לשערך את הביטוי באמצעות דגימות מונטה-קרלו:

$$\nabla_{\theta} J \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T^{(i)}-1} \nabla_{\theta} \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) Q^{\pi_{\theta}}(s_t^{(i)}, a_t^{(i)})$$

הביטוי המתקבל טוב מאוד מבחינת זה שהוא מצליח להפחית את השונות, אך הבעיה בו נעוצה בכך שהאיבר  $Q^{\pi_{\theta}}(s_t^{(i)}, a_t^{(i)})$  אינו ידוע! אלגוריתם Actor Critic בא להתמודד עם בעיה זו, והרעיון הוא לנסות לשערך את  $Q^{\pi_{\theta}}$  באמצעות  $Q^w$  במקביל לשערך של הגרדיאנט. כלומר, בתהליך הלמידה אנו מנסים לשערך שני דברים במקביל:

$\pi_{\theta}$  (Actor) – הגורם המחליט איזה צעד לנקוט בכל מצב.

$Q^w$  (Critic) – הגורם המבצע אבולוציה של  $Q$  על מנת לשפר את הבחירה של ה-Actor.

הערת אגב: לעיל ראינו שהשימוש ב-Learning –  $Q$  יכול להיות בעייתי בבעיות רציפות בגלל שמרחב המצבים הוא אינסופי ומאוד קשה לפתור את בעיית האופטימיזציה של בחירת action בעזרת הביטוי  $\arg \max_A Q_{\theta}(s, a)$ . כאן בעיה זאת איננה, כיוון שאנחנו לא נדרשים לבחור את ה-action באמצעות  $Q$  אלא רק לשפר באמצעותו את האסטרטגיה (או במילים אחרות – אנחנו לא צריכים לפתור בעיית אופטימיזציה על  $Q$  עבור מרחב מצבים רציף).

האלגוריתם פועל בהתאם לשלבים הבאים. ראשית נאתחל את  $w, \theta$ . כעת נבצע  $T$  איטרציות באופן הבא:

- נדגום מצב התחלתי  $s_0$  מתוך מרחב המצבים  $\mathcal{S}$  ופעולה מתוך האסטרטגיה  $a_0 | s_0 \sim \pi_{\theta}(\cdot | s_0)$ .
- כעת כל עוד לא הגענו למצב הסופי:
- נחשב את המצב החדש  $s$  ואת התגמול המתקבל  $r$ .
- בנוסף, נדגום פעולה חדשה  $a \sim \pi_{\theta}(\cdot | s)$ .
- נחשב את ה-TD error:  $\delta = r + \gamma Q^w(new\ s, new\ a) - Q^w(s, a)$
- נעדכן את הפרמטרים:
 
$$\begin{aligned} \theta &\leftarrow \theta + \alpha Q^w(s, a) \nabla_{\theta} \log \pi_{\theta}(a|s) \\ w &\leftarrow w + \beta \delta \nabla_w Q^w(s, a) \end{aligned}$$
- נעדכן את המצב והפעולה:  $a = new\ a, s = new\ s$

פסאודו קוד עבור האלגוריתם:

```

Initialize parameters  $s, \theta, w$  and learning rates  $\alpha_{\theta}, \alpha_w$ ; sample  $a \sim \pi_{\theta}(a|s)$ .
for  $t = 1 \dots T$ : do
  Sample reward  $r_t \sim R(s, a)$  and next state  $s' \sim P(s'|s, a)$ 
  Then sample the next action  $a' \sim \pi_{\theta}(a'|s')$ 
  Update the policy parameters:  $\theta \leftarrow \theta + \alpha_{\theta} Q_w(s, a) \nabla_{\theta} \log \pi_{\theta}(a|s)$ ; Compute
  the correction (TD error) for action-value at time t:
     $\delta_t = r_t + \gamma Q_w(s', a') - Q_w(s, a)$ 
  and use it to update the parameters of Q function:
     $w \leftarrow w + \alpha_w \delta_t \nabla_w Q_w(s, a)$ 
  Move to a  $\leftarrow a'$  and  $s \leftarrow s'$ 
end for
  
```

חשוב להדגיש שהעדכון של הגרדיאנט מתבצע בדומה ל-SARSA ולא כמו  $Q$ -Learning. כאמור לעיל, ההבדל ביניהם הוא באופי האלגוריתם – SARSA הוא אלגוריתם on-policy, כלומר הוא מנסה לאפסם את האסטרטגיה ממנו הוא דוגם את המסלולים, בעוד ש- $Q$ -Learning מנסה ללמוד את האסטרטגיה האופטימלית בלי תלות בדוגמאות

אותן הוא רואה. אלגוריתם Actor Critic רוצה לשפר את ה-Policy gradient, כלומר המטרה היא לשפר את האסטרטגיה הנתונה, ולכן יש לבצע למידה שהיא On policy.

בדומה ל-Policy gradient, גם ל-Actor Critic ניתן להוסיף Baseline על מנת לשפר את תהליך הלמידה. בניגוד ל-Policy gradient הוא התוספת היתה של קבוע, כאן נהוג להוסיף את  $V^\pi(s)$ , שהוא תלוי ב-s. אם ה-baseline תלוי רק במצב s, אז זה בסדר כפי שנוכח מיד (שקף 16), אבל אם הוא תלוי גם ב-a action, אז זה מוסיף bias והלמידה תפגע. הביטוי:  $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$  נקרא ה-advantage function, והסיבה לבחירה שלו כ-baseline היא פשוטה – אם הערך של  $A^\pi(s, a)$  חיובי, אז זה סימן לכך שהבחירה בצעד a היא טובה, ולכן נרצה לחזק אותה. אם לעומת זאת הערך שלילי, זה אומר שהפעולה לא טובה וממילא נרצה להימנע ממנה. במילים אחרות:  $V^\pi(s)$  הינו הערך הממוצע של המצב s, והביטוי  $Q^\pi(s, a)$  הינו הערך של פעולה ספציפית עבור מצב s, ולכן  $A^\pi(s, a)$  אומר לנו עד כמה הפעולה a טובה או גרועה ביחס לממוצע של המצב בו אנו נמצאים, וממילא נוכל לדעת האם נרצה לחזק את הפעולה הזו או לא.

לאחר שראינו שניתן להשתמש ב-Baseline עבור Actor Critic ויש היגיון להשתמש ב- $V^\pi(s)$  בתור ה-baseline, נראה כיצד ניתן ליישם זאת בפועל. אם מרחב המצבים הוא דיסקרטי, ניתן ללמוד את  $Q^w(s)$  בשיטות הסטנדרטיות שראינו, ואז לחשב באופן מדויק את  $V^\pi(s)$  בעזרת הקשר  $V^\pi(s) = \sum_a \pi(a|s) Q^w(s, a)$ . במצבים מסוימים ניתן להשתמש באותו רעיון גם עבור מרחב מצבים רציף – נניח ו- $Q^w(s)$  הוא לינארי ביחס למרחב הפעולות, כלומר מתקיים:  $Q^w = \phi_w(s)^T \cdot a$ , והאסטרטגיה היא גאוסיאנית:  $\pi_\theta(a|s) \sim \mathcal{N}(\mu(s), \sigma^2(s))$  אז  $V^\pi(s)$  הוא התוחלת של  $Q^w(s)$  על פני כל הפעולות:  $V^\pi(s) = \mathbb{E}_{a \sim \pi_\theta(\cdot|s)}[Q^w(s)]$ . בעזרת הנוסחה הסגורה הזו ניתן להשתמש בשערוך  $Q^w(s)$  ולחשב את  $V^\pi(s)$  ואז להשתמש בו כ-Baseline.

כאשר לא ניתן להשתמש ב- $Q^w(s)$  על מנת לחשב את  $V^\pi(s)$ , ניתן לנסות לשערך את שני הביטויים בנפרד, אך זה יכול להיות בזבזני. לחילופין, ניתן לבנות מודל (למשל רשת נוירונים) שתנסה לשערך במקביל את שני הביטויים, כלומר הפלט שלה יהיה גם  $Q^w(s)$  וגם  $V^\pi(s)$ . גישה אחרת מציעה להשתמש במשוואות בלמן. לפי משוואות אלו מתקיים בקירוב  $Q(s_t, a_t) = r_{t+1} + \gamma V(s_t)$ , ולפי זה ה-advantage function יהיה:

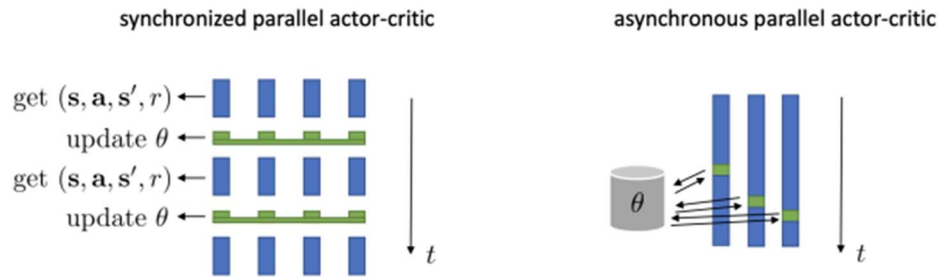
$$A(s_t, a_t) = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

את הביטוי הזה נכניס לגרדיאנט ובכך נשערך את האסטרטגיה. נשים לב שביטוי זה תלוי במצב בלבד ולא בפעולה, מה שמקל על הלמידה, כיוון שיותר קל ללמוד את  $V(s)$  השייך למרחב בגודל S מאשר את  $Q(s, a)$  השייך למרחב בגודל  $S \times \mathcal{A}$ . המשמעות הפרקטית של הביטוי היא שמסתכלים על שני מצבים, ובוחנים את הפרש הערכים שלהם ועוד ה-reward. אם הערך חיובי – אז נרצה לחזק את המעבר מהמצב הראשון למצב השני, ואם הערך שלילי, אז נרצה למנוע כזה מעבר.

חסר – שקופית 19.

אחת הבעיות שיש גם ב-Q-Learning וגם ב-Actor Critic היא הקורלציה שיש בין הדגימות – עבור מסלול שנדגם, יש קשר בין המצבים השונים, ולכן השונות עדיין גבוהה. לעיל ראינו שבאלגוריתם Q-Learning ניתן להתגבר על כך בעזרת שמירה של N מצבים ועדכון לפי k מצבים אקראיים מתוך אלה שנשמרו (למשל:  $N = 1000, k = 10$ ). כאן לא ניתן להשתמש בפתרון זה, שהרי הלמידה היא on-policy ואי אפשר להשתמש במצבים מדגימות קודמות, כיוון שהן נדגמו מאסטרטגיה שכרגע היא לא רלוונטית. ניתן בכל אופן להשתמש במצבים שנדגמו מה-epochs האחרונים, תחת הנחה שהאסטרטגיה לא השתנה הרבה (ובנוסף ניתן לתקן באמצעות Importance sampling).

פתרון יותר מוצלח הוא לאמן מספר סוכנים במקביל ולהשתמש במידע שמגיע מכולם, כאשר היתרון הגדול של שיטה זו הוא שהדגימות של הסוכנים השונים הן חסרות קורלציה. שיטה זו הציגה תוצאות מעולות ובזמן הרבה יותר מהיר מהשיטות הקודמות שראינו. יש שתי אפשרויות להריץ במקביל – ריצה סינכרונית (הסוכנים רצים במקביל, ולאחר שכולם מסיימים epoch מחשבים את הגרדיאנט ומעדכנים את האסטרטגיה) וריצה אסינכרונית (כל פעם שסוכן מסיים epoch הוא מחשב גרדיאנט ומעדכן האסטרטגיה). הריצה הסינכרונית יכולה להיות קצת יותר איטית כיוון שבכל עדכון יש לחכות שכל הסוכנים יסיימו. מצד שני הריצה האסינכרונית יכולה להיות פחות מדויקת כיוון שיתכנו מצבים שלסוכנים שונים תהיה אסטרטגיה מעט שונה (אם הם עדכנו את האסטרטגיה שלהם בזמנים שונים ובאמצע היה עדכון של הגרדיאנט). עם זאת ניתן להניח שההבדלים יחסית לא גדולים, כיוון שהסוכנים מעדכנים את האסטרטגיה שלהם בערך באותן נקודות זמן.



איור 11.10 אימון מספר סוכנים במקביל בשיטת Actor-Critic. ניתן לבצע את העדכון בצורה סינכרונית – כאשר כל הסוכנים מסיימים epoch מחשבים את הגרדיאנט ומעדכנים את האסטרטגיה, או בצורה אסינכרונית – כל סוכן שמסיים epoch מעדכן את האסטרטגיה בהתאם.

לסיכום, שיטת Actor-Critic מנסה לשלב עקרונות מתוך  $Q$  – Learning יחד עם Policy gradient. הרעיון הבסיסי הוא להחליף את הביטוי של ה-Policy gradient reward-to-go בשערוך של  $Q$ , שהוא הרבה פחות רועש. את  $Q$  ניתן לשערך במקביל לגרדיאנט עצמו, אך יש לעשות זאת on-policy. אמנם השונות יורדת אך עדיין יש הטיה, כיוון שיש צורך לשערך גם את  $Q$ . בדומה ל-Policy gradient, גם כאן ניתן להוסיף Baseline, וראינו שמקובל להוסיף את  $V$ .