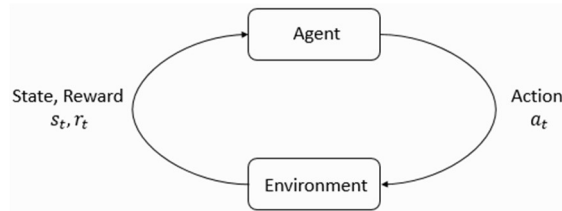


11. Reinforcement Learning (RL)

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 עבורות 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:

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

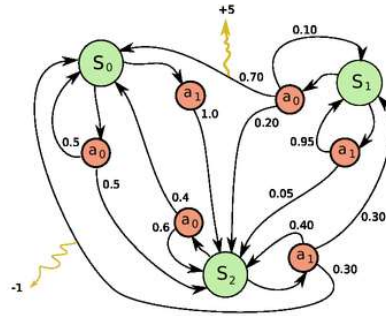
תהליך קבלת החלטות מרקובי מתואר על ידי סט הפרמטרים $\{\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}\}$:

- 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- π , 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:

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

אם $\gamma = 0$, אז מתעניינים רק בתגמול המיידי, וככל ש- γ 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:

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

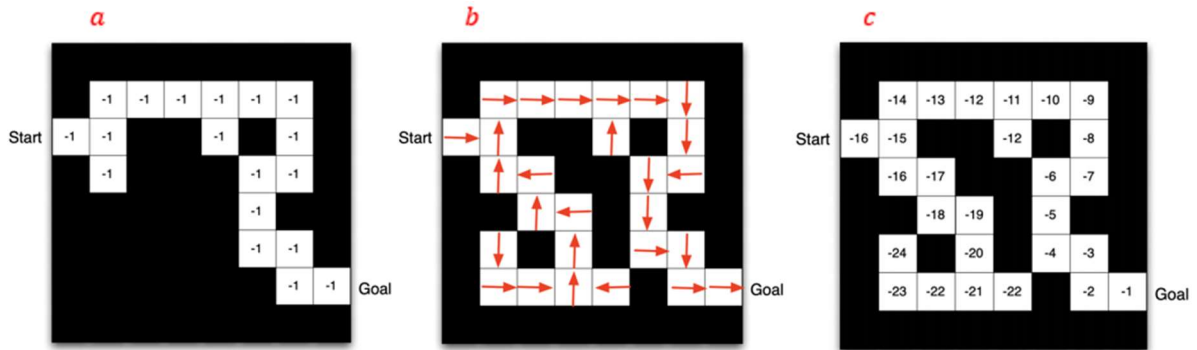
בעזרת ביטוי זה ניתן לחשב את **האסטרטגיה האופטימלית**, כאשר ניתן לנקוט בגישה ישירה ובגישה עקיפה. הגישה הישירה מנסה למצוא בכל מצב מה הפעולה הכי כדאית. בהתאם לכך, חישוב האסטרטגיה האופטימלית יעשה באופן הבא:

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

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:

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

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

- 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 π , 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 π , כאשר במצב s מבצעים את פעולה a , ולאחר מכן ממשיכים לפי האסטרטגיה π :

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

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- $Q^*(s, a)$ את הערכים של האסטרטגיה האופטימלית π^* – 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 π^* 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:

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

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

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}^{π} . כשמוצאים את V ניתן למצוא גם את Q^{π} 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:

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

לפי משפט נקודת השבת של בנך, 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}^π 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 לאסטרטגיה האופטימלית. נניח והצלחנו לחשב את $\mathcal{V}^\pi(s)$ וממנו אנו יודעים לגזור אסטרטגיה, עדיין יתכן שקיימת פעולה a שיותר משתלמת מאשר הפעולה המוצעת לפי האסטרטגיה הנגזרת מ- $\mathcal{V}^\pi(s)$. באופן פורמלי ניתן לתאר זאת בצורה פשוטה – נניח שחישבנו את $\mathcal{V}^\pi(s)$ ואת $Q^\pi(s, a)$ יתכן וקיימת פעולה עבורה:

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

אם קיימת פעולה כזו, אז ישתלם לבחור בה ורק לאחר 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 $\mathcal{V}^\pi(s)$. 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 גרדית – בכל מצב לנקוט בפעולה הכי משתלמת בטווח של צעד יחיד. השאלה העולה היא כמובן – מדוע זה בהכרח נכון? כלומר, האם הרעיון של לבחור בכל צעד את a האופטימלי בהכרח תוביל לקבלת אסטרטגיה אופטימלית עבור כל הצעדים כולם יחד? בכדי להוכיח זאת ננסח זאת כמשפט:

בהינתן 2 אסטרטגיות π, π' , כאשר π' 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- π' , 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} :

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

וזוהי 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- $Q^\pi(s, a)$.

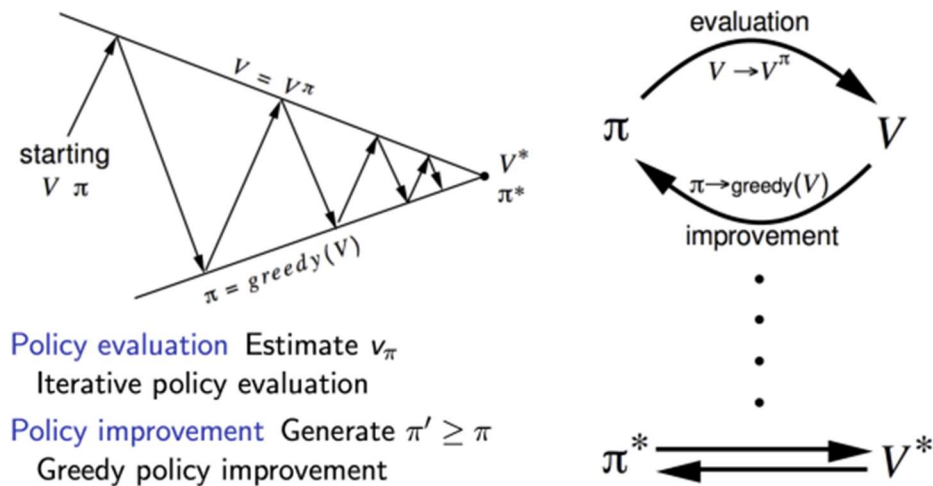
Improve (I) – בהינתן הערך של value function,

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 שבטוח 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- π^* ונקבל את הקשר הבא:

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

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:

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

ניתן להראות שאופרטור זה הינו העתקה מכווצת, וממילא לפי המשפט של בנך 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 \mathcal{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 חישובנו 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:

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

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

ובנוסף פיתחנו את המשוואות עבור הפתרון האופטימלי:

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

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

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 לסוּך 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 θ הממקסם את $\pi_\theta(a|s)$. פתרון בעיה זו יכול להיעשות באופן ישיר על ידי שיטת Gradient Ascent עבור פונקציית המחיר $J(\pi_\theta) = \mathbb{E}[R(\tau)]$, או בעזרת קירוב פונקציה זו ומציאת מקסימום עבודה.
- ב. Q-learning – שערך $Q^*(s, a)$ על ידי $Q_\theta(s, a)$. מציאת המשערך האופטימלי יכולה להתבצע על ידי חיפוש θ 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 π .
- ב. Model-based RL with a known model – אלגוריתמים המנסים למצוא את ה-Value function 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 אסימפטומטיים. במקרים בהם המודל אינו ידוע, הדגש העיקרי הוא על למידת המודל.

11. References

Reinforcement Learning Course Notes-David Silver