

מערכות דינמיות ובקרה

לביא שפיגלמן

מבוא ל Reinforcement Learning

- הקשר בין RL ו DP

הקשר בין RL ו DP

כזכור, בבעית הבקרה הדיסקרטית, משוואת הדינמיקה היא

$$x_{t+1} = f(x_t, u_t, t)$$

הו J^0 (משמעותו C^0) על המסלול האופטימי הוא cost to go

$$V^0(x_k, k) = \min_u \left\{ \phi(x_N, N) + \sum_{t=k}^{N-1} L(x_t, u_t, t) \right\}$$

או, כנוסחת נסיגה

$$V^0(x_k, k) = \min_{u_k} \left\{ L(x_k, u_k, k) + V^0(f(x_k, u_k, k), k+1) \right\}$$

ב-DP משתמשים בנוסחה הנ"ל (ובפונקציית המחיר הרגעי, L ומשוואות הדינמיקה, f) למציאת $V^0(x_k, k)$ לכל x ו k ע"י איטרציה. כפי שצויין,zman החישוב, (עבור מספר צעדים, N , הוא $O(|\mathcal{X}| |\mathcal{U}|^N)$) עלול להיות גדול ופונקליות הדינמיקה והמחיר הרגעי לא תמיד ידועות.

התחום, Reinforcement Learning (RL) בא להתמודד עם הבעיה הנ"ל. לשם כך יש צורך להתאמות מסווגים קלה. loss הרגעי נלקח כשלילי ונקרא גם reinforcement. מחיר הנזומה הסופית נלקח גם הוא כשלילי ונקרו גם כן reinforcement. RL רוצחים למקסם כשם שבבעיות בקרה האופטימלית רצינו למזער את המחיר הכללי, ב-RL רוצחים למקסם את (תוחלת) סכום reinforcements.

סימונים (ותרגום):

שם RL	סימון RL	שם בברכה	סימון בברכה
state	$s \in S$	x	state
action	$a \in A(S)$	u	control sig.
(stochastic) policy	$\pi(s, a) = Pr(a s)$	u(x, t)	control function
(stochastic) state transition	$P_{ss'}^a$	f(x, u, t, ω)	state dynamics
(stochastic) reward	r	-L	instantaneous loss
commulative expected reward	V	-J, -V	cost to go

התמונה הלאטנטית בRL היא של סוקן (agent) שחש את סביבתו באופן חלקי ומיצג את ההערכה שלו של מצבו בעולם ע"י s . בכל צעד, הסוקן מבצע פעולה באופן סטוכסטי (או דטרמיניסטי). התפלגות הפעולות מיוצגת ע"י $\pi(s, a) = Pr(a|s)$ (במונחים של ברחה, (x, u) זו פונקציה הסתובבות של המחב). לאחר ביצוע הפעולה המחב משתנה ע"פ התפלגות (ידועה או לא ידועה)

$$P_{ss'}^a = Pr\{s_{t+1} = s' | s_t = s, a_t = a\}$$

במונחים של ברחה, $P_{xx'}^u$ מקבילה להסתברות $sh(x, u, \omega) = f(x, u, \omega)$ (הסתוכריות נובעת מרווח תחילה ω). לאחר מכן מתקבל reward (loss) במנוחים של ברחה. reward נתון מוגן התפלגות ותוחלתו

$$R_{ss'}^a = E_{r_{t+1}}\{r_{t+1} | s_t = s, a_t = a, s_{t+1} = s'\}$$

מטרת הסוקן היא לבצע policy אופטימלי, כלומר לבצע פעולות כך שסכום הrewards יהיה מרבי. סכום זה נלקח בדרך כלל כסכום אינסופי מונחת (discounted). הסכום של ריצה ספציפית יסומן כך:

$$\begin{aligned} R_t &= r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \\ &= r_{t+1} + \gamma R_{t+1} \end{aligned}$$

ותוחלת הסכום (הרוחות) מסומן

$$\begin{aligned} V^\pi(s) &= E_\pi\{R_t | s_t = s\} \\ &= E_\pi\{r_{t+1} + \gamma V(s_{t+1}) | s_t = s\} \\ &= \sum_a \pi(s, a) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^\pi(s')] \end{aligned} \quad (1)$$

זה כמו cost to go (בהצגה סטוכסטית), $J(x, t)$. נשים לב כי כאן מניחים שההתקפה-גויות אינן תלויות בזמן (הנתחת סטציונריות) והמחair הוא עבר זמן ריצה אינסופי (ותוך דיעיכת rewards ע"י גורם γ) ולכן הפונקציה היא של המחב בלבד. במקרה זה פתרון DP מצריך מלאי של וקטור ערכים, v , על סמך ערכו הקודם במקומות מלאי של עמדות עוקבות במטריצה.

המחיר תחת ברחה אופטימלית, $J^0(x, t)$, יסומן כ

$$V^*(s) = \max_\pi V^\pi(s)$$

הpolicy האופטימלי הוא $\pi^*(s)$. כדי שהגדכנו את $J^0(x, t)$ להיות מחיר ביצוע אותן בברחה כלשהו, u ולאחריו ביצוע בברחה אופטימלית, כך בRL מוגדרת פונקציית

הן optimal action-value (באון רקורסיבי):

$$\begin{aligned} Q^*(s, a) &= E \{r_{t+1} + \gamma V^*(s_{t+1})\} \\ &= E \left\{ r_{t+1} + \gamma \max_{a'} Q^*(s', a') \right\} \\ &= \sum_{s'} P_{ss'}^a \left[R_{ss'}^a + \gamma \max_{a'} Q^*(s', a') \right] \end{aligned}$$

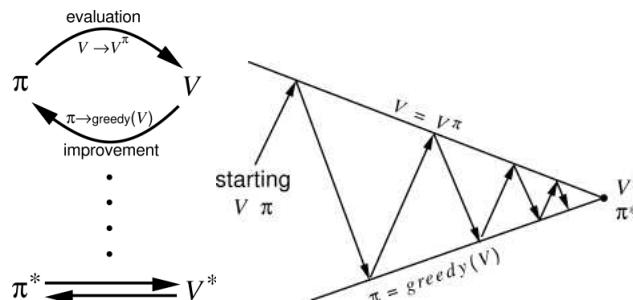
זה האנלוג לתנאי האופטימליות של Bellman מבקרה אופטימלית. policy האופטימלי במנחים של Q הוא

$$\pi^*(s) = \arg \max_{a \in A(s)} Q^*(s, a)$$

כש שמעצנו את $J(x, t)$ (ואת $J^1(x, u, t)$) בשיטות DP, ניתן למצוא (בהינתן פונקציות L ו- f , סטוכסיות, ככלומר בהינתן $(P_{ss'}^a, R_{ss'}^a)$) policy אופטימלי באוטה סיבוכיות זמן ריצה.

כשהדינמייקה לא ידועה במלואה או המחיר הרגעי אינו ידוע במלואו או זמן הרצה אינו ארוך דיו (כפי שקרה למשל בש בש בו יש כ- 10^{20} מצבים) RL מציב אלגוריתמים שמתכנסים (בד"כ, תחת תנאים מסוימים אך ריאליים) לפיתרון האופטימלי (כלומר מותכנים לקיום תנאי האופטימליות של Bellman). ישנה גם הרחבת RL לזמן רציפים.

כל שיטות RL וכן DP (וגם EM) מבצעות איטרציות בהן בצד אחד ישנו הערכה של הערך עבור policy הנתונה וצד שני בו יש שיפור של policy על סמך המעודכנות value function.



שיטות מעין זו מכונות (GPI) Generalized Policy Iteration (GPI) והקלاسي שהוצע כאן מבצע סריקה של כל המצבים ועדכן של $V(s)$ לכל s לפני π מעודכן.¹ דבר זה אינו הכרחי. Asynchronous DP מעדכן את $V(s)$ עבור קבוצה חלקית של S (אפילו s בודד) ולאחר מכן מעדכן את policy, במקומות לעדכן את וקטור הערכים במלואו על סמר וקטור קודם, מעכינים מצבים בודדים על סמך הוקטור הקודם או, אפילו ערך בודד על סמך שאר הערכים (לא זיכרו זמני נספחים). באופן זה ניתן לרכז את מאמץ החישוב במצבים שסבירות ההגעה אליהם גבוהה.

¹ בהצגת הנוסחאות שלנו הדבר נבע מכך שוחשב לכל i ע"י בחירת a שמייצר את V על סמך $V_{i,k+1}$