

در یادگیری تقویتی (Reinforcement Learning)، مجموعه‌ای از نمادها و تعاریف پایه‌ای وجود دارد که برای درک الگوریتم‌ها و معادلاتی مانند معادله بلمن باید با آن‌ها آشنا باشیم. در ادامه مهمترین نمادها را همراه با توضیح و مثال‌های ساده به زبان فارسی مرور می‌کنیم.

۱. محیط و عامل (Agent & Environment)

- حالت (s – State)

نشان‌دهنده وضعیت فعلی محیط است.

مثال: در بازی شطرنج، یک «حالت» ممکن است آرایش مهره‌ها روی صفحه باشد.

- اقدام (a – Action)

عملی که عامل می‌تواند در حالت s انجام دهد.

مثال: در همان بازی، حرکت قلعه از خانه e1 به g1 یک اقدام است.

- سیاست ($\pi(a | s)$ – Policy)

تابعی که احتمال انجام اقدام a را در حالت s مشخص می‌کند.

- سیاست قطعی (Deterministic)

$$[s = tS \mid a = tP[A = (s \mid \pi(a : Stochastic))]$$

۲. پویایی‌های محیط (Environment Dynamics)

- تابع گذار احتمال (Transition Probability)

$$(s, a) \mid P(s'$$

احتمال این که با انجام اقدام a در حالت s . به حالت بعدی s' منتقل شویم.

- گاهی به اختصار P_{ss}^a یا P_{ss}^a , هم نوشته می‌شود.

- تابع پاداش (Reward Function)

$$R(s, a, s')$$

پاداشی که محیط پس از انتقال از s به s' به واسطه اقدام a میدهد.

- گاهی به اختصار R_{ss}^a یا R_{ss}^a نیز می‌آید.

- پاداش فوری (Immediate Reward)

مقدار عددی که بلا فاصله پس از انجام اقدام tA در حالت tS و رسیدن به $t+1S$ دریافت می‌شود

$$R(S_t, A_t, S_{t+1}) = r_{t+1}$$

۳. بازگشت (Return)

- بازگشت (Return)

مجموع پاداشهای آینده از گام t به بعد با اعمال ضریب تنزیل γ :

$$\dots + t+3\gamma^2 R + t+2\gamma R + t+1R = \text{t}G$$

که در عمل معمولاً تا یک پایان (Terminal State) یا تعداد محدود گام محاسبه می‌شود.

- ضریب تنزیل (Discount Factor) $0 < \gamma \leq 1$

تعیین می‌کند که پاداشهای دوردست تا چه حد ارزش دارند.

- اگر $\gamma = 0$. فقط پاداش فوری اهمیت دارد

- اگر $\gamma \approx 1$. پاداشهای آینده تقریباً به اندازه پاداشهای فوری اهمیت دارند.

۴. ارزش‌ها (Value Functions)

- تابع ارزش حالت (State-Value Function)

$$[s = {}_t S \mid {}_t \mathbb{E}_\pi [G] = v_\pi(s)]$$

یعنی انتظار بازگشت G وقتی عامل از سیاست π پیروی کند و در حالت s باشد.

- تابع ارزش جفت حالت-اقدام (Action-Value Function)

$$[a = {}_t s, A = {}_t S \mid {}_t \mathbb{E}_\pi [G] = q_\pi(s, a)]$$

یعنی انتظار بازگشت وقتی در s اقدام a اجرا شود و سپس از π پیروی شود.

۰. معادله بلمن (Bellman Equation)

معادله بلمن رابطه بازگشت کوتاه‌مدت و بلندمدت ارزش را برقرار می‌کند.

۰.۱. بلمن برای v_π

$$[\gamma v_\pi(s') + s, a) [r \mid s) \sum_{s', r} P(s', r \mid \pi(a) \sum_a = v_\pi(s)$$

$P(s', r \mid s, a)$ احتمال مشترک رسیدن به حالت s' و دریافت پاداش r .

این معادله می‌گوید: ارزش یک حالت برابر است با میانگین پاداش فوری به اضافه ارزش حالت بعدی (تنزیل شده)، که بر اساس سیاست π و دینامیک محیط وزن‌دهی شده‌اند.

۰.۲. بلمن بهینه (Optimal Bellman Equation)

برای سیاست بهینه π^* و ارزش بهینه v^* :

$$[\gamma v^*(s') + s, a) [r \mid P(s', r \sum_{s', r} \max_a = v^*(s)$$

و به طور مشابه برای q^* :

$$\left[\gamma \max_{a'} q^*(s', a') + s, a) [r \mid P(s', r \sum_{s', r} = q^*(s, a) \right]$$

۶. مثال ساده

فرض کنید یک عامل در خانه‌های شماره ۱ تا ۳ حرکت می‌کند، و اگر به خانه ۳ برسد، بازی تمام می‌شود (بایان) با پاداش ۰.۹ هر حرکت • هزینه دارد. ضریب تنزیل $\gamma = 0.9$.

- حالت‌ها: $s \in \{3, 2, 1\}$

- اقدام‌ها: $a \in \{\text{چپ}, \text{راست}\}$ (به ترتیب حرکت به خانه قبلی یا بعدی)
- دینامیک:

- از ۱ با «چپ» درجا می‌ماند.
- از ۲ با «راست» به ۳ می‌رود.
- پاداش:

- هر حرکت: $r = 0$

- رسیدن به ۳: $r = 1$ و سپس بایان.

محاسبه بازگشت و ارزش‌ها

- اگر از ۲ با «راست» برویم:
 $R_{t+1} = 1$, بازی تمام، پس

$$G = 1$$

- تابع ارزش بهینه (بدون جزئیات):

$$v_*(3) = 1 \times 0.9 + 0 = v_*(2)$$

چون $v_*(3) = 1$ (حالت بایان با ۱ پاداش).

معنی

حالت فعلی و حالت بعدی

نماد

 $'s, s$

اقدام

 a سیاست (احتمال یا قاعده انتخاب a در s) $(s \mid \pi(a$ احتمال انتقال به s' با اقدام a در s $(s, a \mid 'P(s$ پاداش دریافت شده در انتقال $s \rightarrow s'$ با a $R(s, a, s')$ پاداش فوری در گام t $t+1r$ بازگشت (مجموع پادash‌های تنزیل شده پس از گام t) tG

ضریب تنزیل

 γ ارزش حالت تحت سیاست π $v_\pi(s)$ ارزش جفت حالت-اقدام تحت سیاست π $q_\pi(s, a)$

ارزش‌های بهینه (حداکثر کردن بازگشت)

 $v_*(s), q_*(s, a)$