

۲-۲ بازی های تصادفی

بازی های تصادفی تعمیم فرآیندهای تصادفی مارکوف به حالت چندعامله بوده و می توانند به عنوان چارچوبی مناسب برای بررسی و تحقیقات یادگیری چند عامله استفاده شوند. این بازی ها توسعه ای از بازی های ماتریسی با چندین حالت بوده و به آنها بازی های مارکوفی نیز گویند. دو مثال از بازی های ماتریسی در شکل (۱) آورده شده اند. در هر بازی دو بازیکن وجود دارند که بصورت سطری و ستونی اعمالشان نشان داده شده است. درایه های ماتریس همان پاداش بازیکنان است. در بازی اول که یک بازی با جمع صفر است پاداش بازیکن دوم منفی پاداش اولی است. در حالت کلی ماتریس پاداش برای بازی های دونفره با دو ماتریس نشان داده می شود. هر حالت در بازی اتفاقی می تواند بصورت یک بازی ماتریسی بیان شود. در هر حالت از بازی بازیکنان پس از بازی و دریافت پاداش با توجه به عمل گروهی انجام شده به حالت دیگری از بازی می روند. در یک بازی تصادفی، اعمال بازیکنان همزمان صورت می گیرد. بطور کلی تعریف رسمی بازی های تصادفی بصورت زیر است [18]:

تعریف ۲. یک بازی تصادفی به صورت چندتایی $(n, S, A_{1..n}, T, R_{1..n})$ بیان میشود که در آن n تعداد عامل ها، S مجموعه حالات، A_i مجموعه اعمال هر عامل i (در فضای اعمال گروهی $A_1 \times A_2 \times \dots \times A_n$) تابع انتقال $T: S \times A \times S \rightarrow [0, 1]$ و r تابع پاداش برای عامل i ام $S \times A \rightarrow \mathcal{R}$ است.

در یک بازی تصادفی کاهش پذیر، هدف یک بازیکن بیشینه کردن مجموع پاداش یافته با پارامتر کاهش $\gamma \in [0, 1]$ است. اگر استراتژی بازیکن i باشد بازای یک وضعیت اولیه s ، بازیکن i سعی در بیشینه نمودن رابطه زیر دارد:

$$V(s, \pi^1, \pi^2, \dots, \pi^n) \equiv \sum_{t=0}^{\infty} \gamma^t E(r_t | \pi^1, \pi^2, \dots, \pi^n, s_0 = s) \quad (3)$$

	T	L		L	R
T	1,-1	-1,1	L	6,6	2,7
L	-1,1	1,-1	R	7,2	0,0

شکل ۱. دو نوع بازی ماتریسی به نامهای تطابق سکه ها (سمت چپ) و بازی Chicken (سمت راست)

با نگاه به (۳) و (۱) تشابه با مدل MDP دیده می شود منتها در اینجا چندین عامل وجود دارند و رفتن به وضعیت بعدی و دریافت پاداش به عمل گروهی عاملها بستگی دارد. در حالت یادگیری تقویتی چند عامله، بیشینه نمودن سودمندی (پاداش) مورد انتظار هر عامل در بازی های مجموع کلی ممکن نبوده و هدف پیدا کردن سیاست متعادل در بازی های مارکوف است. از جمله این سیاست ها می توان سیاست تعادل نش را نام برد. به عبارت دیگر پیدا کردن یک سیاست متعادل به عنوان یک راه حل برای بازی های اتفاقی محسوب می شود.

۲-۲ استراتژی های تعادل

در بازی های با جمع کلی، راه حل مبنا تعادل نش است [18]. نظریه‌ی تعادل‌های نش می‌گوید که در هر بازی به فرض آنکه بازی کنان معقولانه استراتژی های خود را انتخاب کرده و به دنبال بدست آوردن حداکثر بهره (سود) از بازی باشند، حداقل یک استراتژی برای بدست آوردن بهترین نتیجه برای هر بازیکن قابل انتخاب است که اگر بازیکن راهکار دیگری به غیر از آن را انتخاب کند، نتیجه‌ی بهتری بدست نخواهد آورد [19].

تعریف ۳. در یک بازی تصادفی Γ ، یک تعادل نش بصورت چندتایی استراتژی های (π_1, \dots, π_n) تعریف می‌شود بطوریکه برای هر $s \in S$ و $i = 1, \dots, n$ داریم:

$$V^i(s, \pi_*^1, \pi_*^2, \dots, \pi_*^n) \geq V^i(s, \pi_*^1, \dots, \pi_*^{i-1}, \pi_*^i, \pi_*^{i+1}, \dots, \pi_*^n) \quad (4)$$

for all $\pi^i \in \Pi^i$

که در آن Π^i مجموعه استراتژی های قابل دسترس برای عامل i است. استراتژی هایی که یک تعادل نش را می‌سازند می‌توانند استراتژی های ایستا باشند. قضیه زیر نشان می‌دهد که حداقل یک تعادل در استراتژی های ایستا وجود دارد: قضیه ۱. (Fink 1964) هر بازی اتفاقی کاهش پذیر n نفره حداقل یک نقطه تعادل نش در استراتژی های ایستا را داراست [20].

۳-۲ حل بازی های تصادفی

در بازی های اتفاقی بر خلاف MDP ها یک راه حل بهینه ای که مستقل از عاملهای دیگر باشد محتمل نیست. لذا نسبت به حالت تک عامله مساله مشکل تر است. همچنین ممکن است که چندین نقطه تعادل وجود داشته و هماهنگی عاملها برای توافق کار مشکلی است. مفهوم راه حل در بازی های اتفاقی به معنای پیدا کردن نقطه تعادل منحصر به فرد می‌باشد. الگوریتمهای مختلفی برای حل این بازی های پیشنهاد شده اند که در یک رده بندی می‌توان به روشهایی که از تئوری بازی ها استفاده می‌کنند و روشهای مبتنی بر یادگیری تقویتی اشاره نمود. در روشهای مبتنی بر تئوری بازی ها، نیاز به مدل محیط داشته و توابع T و R نیز بایستی مشخص باشند. هدف این الگوریتمها محاسبه مقدار تعادل بازی (پاداش مورد انتظار کاهش یافته برای هر عامل) است، لذا نیاز به فرضیات قوی تری از رفتار عاملها دارند. در مقابل، الگوریتمهای مبتنی بر یادگیری تقویتی چند عامله فرض می‌کنند محیط ناشناخته بوده و مشاهدات T و R از طریق تجربی قابل مشاهده هستند. هدف این الگوریتمها پیدا کردن سیاست تعادل بوده و معمولا نیازهای کمتری را درباره رفتار عاملهای دیگر دارند. یک چارچوب کلی برای یادگیری Q چند عامله در شکل ۲ نشان داده شده است. در قالب کلی، الگوریتم های مختلف یک ورودی تابع انتخاب تعادل را دریافت می‌کنند تا بتوانند با توجه به بردار ماتریسی نظیر $Q = (Q_1, \dots, Q_n)$ تابع ارزش V را محاسبه نمایند. Hu و Wellman از تابع با نام $Nash-Q$ برای محاسبه مقدار V طبق رابطه (۵) استفاده نمودند [8]. Qio و همکاران نیز از مکانیزم چانه زنی^۱ (۶) [21] Meiping Song نیز از الگوریتم Pareto-Q استفاده کردند [11].

$$V_i(s) \in \text{Nash}_i(Q_1(s), \dots, Q_n(s)) \quad (5)$$

$$\text{Nash}_i(Q_1(s), \dots, Q_n(s)) = \pi^1(s) \dots \pi^n(s) \cdot Q_i(s)$$

$$V_i(s) \in \text{NBS}_i(Q_1(s), \dots, Q_n(s)) \quad (6)$$

$$\text{NBS}_i(Q_i(s)) = \text{Max}_{\bar{a}}(Q_i(s, \bar{a}) \times \dots \times Q_n(s, \bar{a}))$$

Multi-agent Q-learning (Stochastic Game, α, γ, M)

Inputs: discount factor γ , learning rate α , M total training time

Output : state-value functions, action-value V_i^* function Q_i^*

Initialize: $s, a_1 \dots a_N, Q_1 \dots Q_N$

1. for $k=1$ to M
 2. simulate actions a_1, \dots, a_n in state s
 3. observe rewards R_1, \dots, R_n and next state s'
 4. for $i = 1$ to N
 - (a) compute $V_i(S')$
 - (b) $Q_i(s, \bar{a}) \leftarrow (1-\alpha)Q_i(s, \bar{a}) + \alpha[(1-\gamma)R_i + \gamma V_i(s')]$
 5. agents choose action a'_1, \dots, a'_N
 6. $s = s', a_1 = a'_1, \dots, a_N = a'_N$
 7. decay α
-

شکل ۲. الگوریتم یادگیری Q چند عامله

با توجه به پیچیدگی بالای محاسبات در هر یک از الگوریتمهای محاسبه تعادل نش در بازی های اتفاقی، به نظر می رسد استفاده از الگوریتمهای مبتنی بر اتوماتای یادگیر با توجه به سادگی و پیچیدگی محاسباتی کم این الگوریتمها بتوانند کارایی مطلوبی را از خود ارائه دهند. با توجه به این موضوع در بخش بعدی مدلی مبتنی بر اتوماتای یادگیر ارائه و نتایج ارائه گردیده اند.

matlab.ir

ⁱ Nash Bargaining

matlab | .ir

ایران منتلب