

دومین کنگره مشترک سیستم‌های فازی و هوشمند ایران 2nd Joint Congress on Fuzzy and Intelligent Systems

۷ الی ۹ آبان ماه ۱۳۸۷ 28-30 October 2008

تحلیل یادگیری تقویتی در فرایندهای مارکوف به صورت سیستم‌های دیجیتال

محمدباقر نقیبی سیستانی
استادیار دانشکده‌ی مهندسی
دانشگاه فردوسی مشهد
mb-naghibi@ferdowsi.um.ac.ir



سیدمصطفی کلامی هریس
دانشجوی کارشناسی ارشد مهندسی کنترل
دانشگاه فردوسی مشهد
sm.kalami@gmail.com

چکیده - فرایند تصمیم‌گیری مارکوف یا MDP، یکی از مسائلی است که برای کاربردهای وسیعی در زمینه‌های مختلف علمی، مهندسی، اقتصادی و مدیریت است. بسیاری از فرایندهای تصمیم‌گیری دارای خاصیت مارکوف می‌باشند و به صورت یک مسأله‌ی تصمیم‌گیری مارکوف قابل بیان هستند. یادگیری تقویتی یکی از مسائلی است که برای حل MDP به کار می‌رود که به نوبه‌ی خود از برنامه‌ریزی پویا یا DP استفاده می‌کند. در این مقاله معادله‌ی بازگشتی مورد استفاده در بحث یادگیری تقویتی و DP برای حل MDP، به صورت یک معادله‌ی دینامیکی یک سیستم دیجیتال یا گسسته-زمان بازنویسی شده است. به این ترتیب این امکان به وجود آمده است که بتوان با بهره‌گیری از روش‌های موجود در کنترل دیجیتال، به بررسی خواص معادلات به دست آمده پرداخت و تحلیل مناسبی از رفتار عامل یادگیرنده، تحت سیاست‌های مختلف، به عمل آورد. به عنوان مثال، روش مذکور برای تحلیل یک مسأله‌ی جدولی استفاده شده است. نتایج به دست آمده، نشان می‌دهند که یک سیاست بهینه در بازجوب کنترل دیجیتال، به صورت سیستم مرده‌نوش قابل توصیف است.

کلید واژه - برنامه‌ریزی پویا، سیستم‌های کنترل دیجیتال، فرایندهای تصمیم‌گیری مارکوف، کنترل تصادفی، یادگیری تقویتی.

یادگیری تقویتی را حل کند. پسخوردهایی که محیط به عامل برمی‌گرداند، اطلاعاتی در مورد خوبی یا بدی کارهایی که عامل انجام می‌دهد، در بر دارند و عامل موظف است با استفاده از این پسخوردها، یاد بگیرد که چه عملی در چه شرایطی خوب و در چه شرایطی بد است. در یادگیری تقویتی، پاسخ‌های محیط به صورت اسکالر و با نام پاداش تعریف می‌شوند و عامل وظیفه دارد، در یک بازه‌ی زمانی، مجموع پاداش‌های دریافتی را بیشینه کند. یادگیری تقویتی در جانوران و البته انسان، سهم قابل توجهی از یادگیری را به خود اختصاص می‌دهد. هنگامی که کسی دست خود را در معرض حرارت قرار می‌دهد و دستش می‌سوزد، در اثر ناشی از این کار، یاد می‌گیرد که دوباره این کار را انجام ندهد. لذت‌ها و دردها، در واقع پاسخ‌هایی هستند که از طرف محیط به موجودات زنده داده می‌شوند، و تعیین‌کننده‌ی القوی رفتاری بسیاری از موجودات می‌باشند [1-4].

۱- مقدمه

از دیدگاه هوش مصنوعی، هدف از یادگیری، به دست آوردن قابلیت تصمیم‌گیری برای حل مسائل مختلف است. این توانایی هنگامی محقق می‌شود که، نگاهی خاص از فضای حالات به فضای اعمال پیدا شود، که از جهاتی بهتر از سایر نگاهی‌های ممکن باشد [1,2]. یادگیری هنگامی رخ می‌دهد که عامل یادگیرنده (یا به اختصار، عامل) در اثر تجاربی که کسب می‌کند، به نحوی متفاوت و احتمالاً بهتر، عمل کند. به این ترتیب می‌توان گفت: یادگیری، تاثیری است که عامل از محیط اطراف می‌گیرد. تاثیری که محیط بر روی عامل دارد، غالباً به صورت یک پسخورد یا فیدبک حاوی اطلاعات است که تابعی از حالت عامل و همچنین عملی است که روی انجام می‌دهد [1-4].

اگر عاملی موظف باشد که صرفاً با تکیه بر پاسخ‌هایی که از محیط دریافت می‌کند، نحوه‌ی عملکرد بهینه را یاد بگیرد، در واقع از عامل خواسته شده است که یک مسأله‌ی

یکی از انواع مسائلی که به وفور در کاربردهای علمی و

بیش از ۲۰ سال (۱۳۸۷) تحقیقاتی انجام شده بر روی یادگیری تقویتی، توأم با این فرض بوده‌اند که، تعامل بین عامل و محیط اطرافش را می‌توان به صورت یک فرایند تصمیم‌گیری مارکوف یا MDP گسسته در زمان مدل‌سازی محدود و گسسته عبارتند از [2,5-7]:

مارکوف^۲ یا MDP می‌باشد [5-7]. این مسأله را می‌توان در قالب یک مسأله‌ی یادگیری تقویتی بیان نمود و یا استفاده از روش‌هایی که در یادگیری تقویتی مطرح هستند، به حل آن پرداخت [1-3]. یکی از روش‌های مهم در حل MDP، روش برنامه‌ریزی پویا^۳ یا DP است. این روش که توسط پلمن [8] در سال ۱۹۵۷ معرفی شده است، برای حل انواع مسائل برنامه‌ریزی است.

• یک زمان‌بندی سراسری به صورت $t = 0, 1, \dots, T$ برای شمارش زمان گسسته. (T می‌تواند نامحدود باشد).
• فضای حالت S که مجموعه‌ی محدودی از حالت‌ها به صورت $\{s^1, s^2, \dots, s^m\}$ است.

اما کاربردی که این روش در حل MDP^۴ دارد، در زمینه‌ی یادگیری تقویتی، بیشتر مورد توجه بوده است. ارتباطی که امروزه بین یادگیری تقویتی و روش DP وجود دارد، محصول تحقیق و مطالعه‌ی است که از سال ۱۹۶۱ توسط مینسکی شروع شده است و تا امروز ادامه دارد [1]. در این میان نگارش‌های جدیدی از روش DP^۵ با اعمال تغییرات بر روی این روش و یا ترکیب آن با سایر زمینه‌ها، به وجود آمده‌اند. برنامه‌ریزی پویای تقریبی^۶، برنامه‌ریزی پویای افزایشی^۵ و برنامه‌ریزی پویای عصبی^۶ از جمله نگارش‌های تغییر یافته‌ی DP می‌باشند [9,10].

• مجموعه‌ی محدودی از اعمال به صورت $A_s = \{a^{1,s}, a^{2,s}, \dots, a^{m,s}\}$ که در هر حالت $s \in S$ قابل انتخاب هستند. اجتماع مجموعه‌ی اعمال ممکن، در تمام حالات ممکن، مجموعه‌ی A را می‌سازد.

در این نوشتار فرایند حل مسائل MDP با استفاده از روش DP، به صورت یک دینامیک گسسته-زمان و یا سیستم دیجیتال بیان شده است. سیستم دیجیتالی که به دست می‌آید با روش حلی که برای مسأله ارائه شده است، متناظر است و می‌توان با تحلیل خصوصیات کنترلی این سیستم، خواص حل متناظر با آن را مشخص نمود و کیفیت جواب نهایی را حدس زد. استفاده از این معادل‌سازی، این امکان را به وجود می‌آورد که تحلیل عملی^۶ یادگیری تقویتی در محیط‌های مارکوف پرداخت. سایر بخش‌های اسن مقاله به صورت زیر می‌باشند. در بخش ۲، یادگیری تقویتی در محیط‌های مارکوف و روش DP به اختصار توضیح داده می‌شوند. در بخش ۳، روش DP به صورت یک سیستم دیجیتال بیان می‌شود. بخش ۴ نیز، حاوی حل و بررسی یک مسأله‌ی نمونه با استفاده از مطالب مذکور در بخش‌های قبلی می‌باشد.

• یک تابع تغییر حالت به صورت $P_{ss'}^a$ که احتمال تغییر حالت از s به s' را، در اثر انتخاب عمل a در حالت s مشخص می‌کند.

برای استفاده از یادگیری تقویتی در چنین محیطی، جزء دیگری نیز به این صورت افزوده می‌شود:

• تابع پاداش که به ازای هر سه‌تایی s, a و s' ، یک کمیت تصادفی با توزیع ثابت را ایجاد می‌کند. امید ریاضی این کمیت تصادفی $R_{ss'}^a$ است.

فرض کنید عامل در زمان t در حالت s_t قرار دارد و عمل $a_t \in A_{s_t} \subseteq A$ را انجام می‌دهد. عامل با احتمال $P_{s_t, s_{t+1}}^{a_t}$ در زمان $t+1$ به حالت s_{t+1} می‌رود و خواهیم داشت: $s_{t+1} = s'$. ضمناً عامل در نتیجه‌ی این عمل، پاداشی اسکالر به اندازه‌ی r_t خواهد گرفت که در حالت کلی، کمیتی تصادفی و با امید ریاضی $R_{s_t, s_{t+1}}^{a_t}$ می‌باشد.

علاوه بر خواص فوق، اصلی‌ترین خاصیت هر فرایند مارکوف به این شکل قابل بیان است: در هر زمان t ، احتمال رسیدن به حالت s_{t+1} ، صرفاً به s_t و a_t وابسته است و مستقل از حالات و اعمال قبل از زمان t می‌باشد. به عبارت دیگر می‌توان نوشت [2,5-7]:

$$\Pr(s_{t+1} | s_t, a_t, s_{t-1}, a_{t-1}, \dots) = \Pr(s_{t+1} | s_t, a_t) \quad (1)$$

² Markov Decision Process

³ Dynamic Programming

⁴ Approximate Dynamic Programming

⁵ Incremental Dynamic Programming

⁶ Neural Dynamic Programming

حالت‌ها، بیشینه شوند [1-3]. روش‌های متفاوتی برای

تعریف خروجی وجود دارند. روش در اکثر کاربردها معمول است و در این مقاله نیز مورد توجه است، تعریف خروجی به صورت تنزیلی⁹ می‌باشد. اگر ضریب تنزیل به صورت $\gamma \in [0,1]$ باشد، خروجی تنزیلی به صورت زیر خواهد بود:

$$z_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$$

در این حالت ارزش حالت s به صورت زیر قابل بیان است:

$$V^\pi(s) = E_\pi(z_t | s_t = s) = \sum_{a \in A} \pi(s, a) \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V^\pi(s')) \quad (\Delta)$$

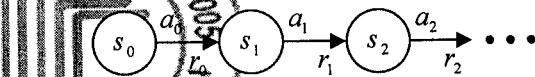
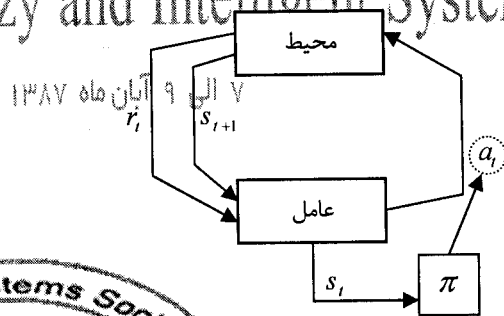
که در آن، منظور از E ، عملگر امید ریاضی است. اندیس‌های π نیز، صرفاً برای تاکید بر این که عامل از سیاست π پیروی می‌کند، نوشته شده‌اند. رابطه‌ی (Δ)، به معادله‌ی (بهینگی) بلمن معروف است [1-3, 5-10]. یکی از روش‌هایی که برای حل این معادله و یافتن مقدار ارزش تمام حالات استفاده می‌شود، بازگشتی کردن این معادله است. این روش که مبتنی بر قضیه‌ی نقطه‌ی ثابت¹¹ [11] است، پیشنهاد می‌کند که معادله‌ی (Δ) به صورت زیر بازنویسی شود:

$$V_{k+1}^\pi(s) = \sum_{a \in A} \pi(s, a) \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V_k^\pi(s')) \quad (\Xi)$$

که در آن $V_k^\pi(s)$ تخمین k ام از مقدار واقعی $V^\pi(s)$ است. با توجه به این که $|\gamma| < 1$ است، می‌توان استدلال کرد که $V_{k+1}^\pi(s)$ با یک نگاه تکراری¹² [11] به $V_k^\pi(s)$ مرتبط است. طبق قضیه‌ی نقطه‌ی ثابت [11]، این نگاه دارای نقطه‌ی ثابت منحصر به فردی است که جواب معادله‌ی (Ξ) نیز می‌باشد. با توجه به معادله‌ی (Ξ)، جواب معادله‌ی (Δ) به صورت زیر خواهد بود:

$$V^\pi(s) = \lim_{k \rightarrow \infty} V_k^\pi(s) \quad (\Upsilon)$$

فابند محاسبه‌ی $V^\pi(s)$ برای تمام حالت‌ها، در بحث یادگیری تقویتی به نام ارزیابی سیاست¹² [1-3, 5] معروف است.



شکل ۱ - نحوه‌ی ارتباط عامل با محیط اطراف

در مورد یادگیری تقویتی، فرض بر این است که تابع پاداش نیز، دارای خاصیت مارکوف است [2].

$$\Pr(r_t | s_t, a_t, s_{t-1}, a_{t-1}, \dots) = \Pr(r_t | s_t, a_t) \quad (2)$$

احتمال انتخاب عمل a از طرف عامل، هنگامی که در حالت s قرار دارد، با نگاهی به صورت $\pi: S \times A \rightarrow [0,1]$ تعریف می‌شود و می‌توان نوشت:

$$\Pr(a_t = a | s_t = s) = \pi(s, a) \quad (3)$$

نگاشت π معمولاً با نام سیاست⁷ شناخته می‌شود و معمولاً اصلی یک مسأله‌ی یادگیری تقویتی و هر مسأله‌ی تصمیم‌گیری می‌باشد [1-3, 5]. نحوه‌ی عملکرد یک عامل با محیط اطرافش، در شکل ۱ به وضوح نشان داده شده است.

برای مقایسه‌ی سیاست‌های مختلف با یکدیگر، می‌توان معیاری را برای سنجش آن‌ها تعریف نمود. این معیار مقداری است که سیاست در هر حالت از فرایند برمی‌گرداند و به عنوان خروجی⁸ سیاست در حالت مذکور از آن یاد می‌شود. خروجی یک سیاست، میزانی از پاداش که در اثر اتخاذ تصمیمات متوالی با تبعیت از آن سیاست به دست آمده است. برای هر کدام از حالت‌ها در MDP، ارزشی در نظر گرفته می‌شود. ارزش هر حالت برابر با امید ریاضی خروجی که با شروع کردن از همان حالت و تبعیت از یک سیاست خاص است. در حالت کلی، منظور از حل یک مسأله‌ی یادگیری تقویتی، پیدا کردن سیاست π^* است به نحوی که مقدار خروجی سیاست و یا ارزش هر کدام از

⁹ Discounted
¹⁰ Fixed Point Theorem
¹¹ Contraction Map
¹² Policy Evaluation

$$\gamma < \frac{\max_{1 \leq i \leq n} |\lambda_i(P)|}{\rho(P)} = \frac{\rho(P)}{\rho(P)} \quad (13)$$

28-30 October 2008

۷ الی ۹ آبان ماه ۱۳۸۷

فرض کنید برداری به صورت

که در آن، $\lambda_i(P)$ نشان دهنده مقدار ویژه i ام، و $\rho(P)$ نیز شعاع طیفی^{۱۳} ماتریس P می‌باشد. لذا مشاهده

$$v^\pi = [v^\pi(s^1) \ v^\pi(s^2) \ \dots \ v^\pi(s^n)]^T \quad (8)$$

می‌شود که شرط $\gamma \leq 1$ ، الزاما تضمین کننده همگرایی

تعریف شده باشد. این بردار حاوی ارزش تمام حالات یک

مدل است. در این صورت می‌توان رابطه بازگسی

به شکل زیر برای تمام حالت بازنویسی و آن را به

برای (۱۳) رابطه می‌باشد. رابطه (۱۳)، شرط دقیق تری برای γ

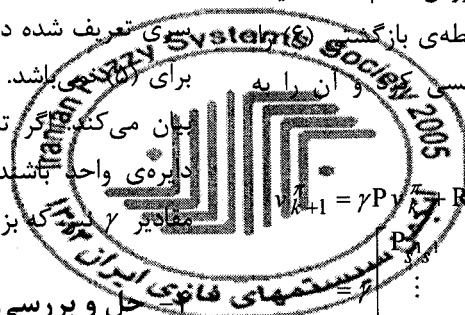
صورت یک معادله تبدیل نمود:

نشان می‌کند که اگر تمام مقادیر ویژه ماتریس P ، داخل

دایره واحد باشند، آن گاه سری (۴)، به ازای برخی از

مقادیر γ بزرگتر از یک هستند، همگرا خواهد بود.

که بزرگتر از یک هستند، همگرا خواهد بود.



حل و بررسی یک مسأله نمونه

جدولی به صورت شکل نشان داده شده در شکل ۲ را در نظر بگیرید. عاملی (مثلا یک روبات) در یکی از خانه‌های

که در آن دایره‌های ماتریس‌های P و R عبارتند از:

این جدول قرار دارد. این عامل می‌بایست با حرکت در یکی

$$P_{ss'} = \sum_{a \in A} \pi(s, a) P_{ss'}^a \quad (10)$$

از چهار جهت بالا، پایین، چپ و راست، خود را به یکی از دو

خانه‌ی هدف، که با رنگ خاکستری مشخص شده‌اند،

برساند. حرکت در هر جهت، پاداشی به اندازه‌ی ۱- در پی

$$R_s = \sum_{a \in A} \sum_{s' \in S} \pi(s, a) P_{ss'}^a R_{ss'}^a = \sum_{a \in A} \pi(s, a) R_s^a \quad (11)$$

دارد، که این پاداش منفی، نشان دهنده‌ی هزینه‌ای است که

می‌توان معادله‌ی (۹) را به صورت زیر بازنویسی کرد:

عامل برای هر حرکت می‌پردازد. حرکت‌هایی که باعث خارج

شدن عامل از صفحه می‌شوند، بر موقعیت عامل تاثیری

ندارند و محل قرار گرفتن عامل در جدول را تغییر

نمی‌دهند. عامل بایستی یاد بگیرد که با دریافت بیشترین

پاداش و یا پرداخت کمترین جریمه، خود را به یکی از

خانه‌های هدف برساند. اگر چنین کاری محقق شود، عامل

توانسته‌است با کمترین تعداد حرکت، به هدف برسد. این

مسأله به نام Grid-world شناخته می‌شود و دارای خواص

مربوط به فرایند تصمیم‌گیری مارکوف می‌باشد [1-3]. لذا

می‌توان برای حل این مسأله، از روش DP استفاده نمود.

G	s^1	s^2	s^3
s^4	s^5	s^6	s^7
s^8	s^9	s^{10}	s^{11}
s^{12}	s^{13}	s^{14}	G

شکل ۲- شبکه‌ی مربوط به مسأله‌ی مورد بررسی

$$v_{k+1} = \gamma P v_k + R \alpha_k \quad (12)$$

که در آن فرض شده است که u_k به ازای تمام مقادیر

$k \geq 0$ برابر با واحد باشد، که همان تعریف تابع پله‌ی واحد

[12] می‌باشد. معادله‌ی (۱۲) نشان دهنده‌ی یک سیستم

گسسته-زمان یا دیجیتال [12] می‌باشد که متغیرهای حالت

آن، ارزش‌های مربوط به حالات مختلف می‌باشند. ماتریس

حالت این سیستم، از ترکیب اطلاعات مربوط به محیط در

قالب $P_{ss'}^a$ ، اطلاعات مربوط به سیاست در قالب $\pi(s, a)$ و

ضریب تنزیل به دست آمده است. بردار وروی این سیستم

R می‌باشد که اطلاعات مربوط به محیط، سیاست و

پاداش‌ها را در بر دارد. طبق قرارداد، ورودی این سیستم،

همواره برابر با پله‌ی واحد در نظر گرفته می‌شود. شرط

پایداری سیستم فوق، عبارت است از این که، همه‌ی مقادیر

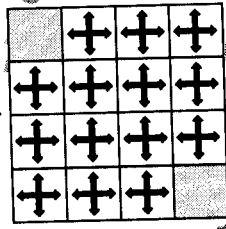
ویژه‌ی ماتریس γP ، که قطب‌های سیستم (۹) هستند، در

داخل دایره‌ی واحد قرار بگیرند [12]. برای تحقق این شرط،

می‌بایست ضریب تنزیل γ در نامساوی ذیل صدق کند:

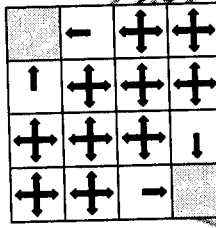
$$k=0$$

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0



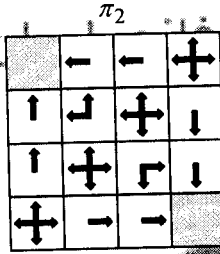
$$k=1$$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0



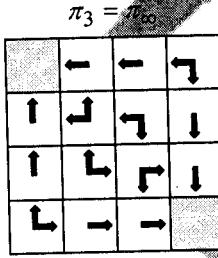
$$k=2$$

0	-1.75	-2	-2
-1.75	-2	-2	-2
-2	-2	-2	-1.75
-2	-2	-1.75	0



$$k=3$$

0	-2.44	-2.94	-3
-2.44	-2.88	-3	-2.94
-2.94	-3	-2.88	-2.44
-3	-2.94	-2.44	0



برای استفاده از روش DP، ارزش اولیه‌ی هر کدام از خانه‌ها برابر با صفر در نظر گرفته می‌شود [1-3,5]. طبق قضیه نقطه‌ی ثابت، نتیجه‌ی نهایی، مستقل از ارزش اولیه‌ی خانه‌ها می‌باشد و الگوریتم همواره به یک نقطه در فضای جستجو، همگرا می‌شود [11]. سیاستی که تا پایان حل مسأله مورد استفاده قرار گرفته است، سیاست تصادفی است. به این معنی که، در تمام خانه‌های جدول، حرکت به تمام جهات مساوی، و همگی برابر با یک چهارم باشد می‌باشد. در شکل ۳، چند مرحله از حل تکراری معادله‌ی بلمن، توسط معادله‌ی (۶)، نشان داده شده است. با استفاده از نتایج مربوط به هر مرحله، می‌توان سیاست‌ها را پاره‌ریزی کرد. به این ترتیب که، عامل می‌بایست در هر خانه از

جدول، به سمت خانه‌هایی حرکت کند که بیشترین ارزش را دارند. سیاستی که با استفاده از ارزش‌های به دست آمده در مرحله‌ی k ام به دست می‌آید، به صورت π_k نشان داده شده است. π_0 همان سیاست تصادفی است. π_∞ نیز سیاستی است که با استفاده از ارزش‌های نهایی به دست آمده است و π_∞ و تمام سیاست‌های بعد از آن، همگی معادل هستند.

فرض کنید با استفاده از هر کدام از سیاست‌های به دست آمده، معادله‌ی سیستم معرفی شده در معادله‌ی (۹) محاسبه شوند، و ماتریس P در معادله‌ی (۹) برای سیاست π_i به صورت P_i باشد. استفاده از اندیس i ، صرفاً به دلیل جلوگیری از تداخل اندیس‌ها در معادله‌ی (۹) می‌باشد.

شعاع طیفی هر کدام از ماتریس‌های مذکور محاسبه شده‌اند و عبارتند از:

$$\rho(P_0) \square 0.9468, \quad \rho(P_1) \square 0.8431, \quad (14)$$

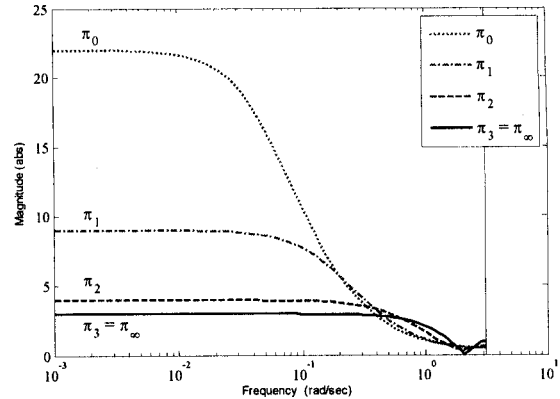
$$\rho(P_2) = 0.5, \quad \rho(P_3) = \rho(P_\infty) \square 0$$

به وضوح دیده می‌شود که هر چه قدر سیاست به کار رفته در ایجاد مدل، بهینه‌تر باشد، اندازه‌ی بزرگترین مقدار ویژه‌ی ماتریس حالت نیز کوچک‌تر می‌شود. به خصوص به ازای سیاست π_3 یا همان π_∞ ، تمامی مقادیر ویژه برابر با صفر است. به این ترتیب، حد بالای ضریب تنزیل γ ، برای همگرایی سری (۴)، به ازای سیاست‌های π_0 تا π_3 به ترتیب عبارت است از: $1/0.562$ ، $1/1.1861$ ، 2 و ∞ . به عبارت دیگر، سری تعریف شده با معادله‌ی (۴)، به شرط پیروی از سیاست π_∞ ، به ازای تمام مقادیر γ همگرا خواهد بود.

شکل ۳ - چند مرحله از برنامه‌ریزی پویا به ازای سیاست تصادفی از طرفی، قطب‌های سیستم (۹)، با مقادیر ویژه‌ی ماتریس P برابر هستند. مشاهده می‌شود که قطب‌های سیستم معادل با سیاست π_∞ ، همگی در مبدا قرار دارند. چنین سیستمی در بحث کنترل دیجیتال، به نام مرده نَوش^{۱۴} شناخته می‌شود [12]. یک سیستم مرده نَوش، سریع‌ترین پاسخ ممکن را در بین سیستم‌های هم‌مرتبه‌اش دارد [12]. اگر درجه‌ی چنین سیستمی برابر با n باشد، دقیقاً n واحد زمانی (گسسته) طول می‌کشد که پاسخ پله‌ی سیستم به مقدار نهایی برسد و نشست کند. هدف از حل مسأله‌ی فوق نیز، رسیدن به یکی از خانه‌های هدف در کمترین تعداد حرکت می‌باشد. لذا کاملاً طبیعی است که پاسخ بهینه، متناظر با یک سیستم مرده نَوش باشد.

¹⁴ Dead Beat

مقاله ۳۰ اکتبر ۲۰۰۸
 در مورد یادگیری تقویتی، فرایندهای مارکوف و برنامه‌ریزی پویا، معادلات مربوط به حل یک فرایند مارکوف با استفاده از برنامه‌ریزی پویا، به صورت یک دینامیک گسسته-زمان جمع‌بندی و حل یک مسئله یادگیری تقویتی در محیط مارکوف را، در قالب یک سیستم دیجیتال فراهم می‌آورد. به این ترتیب، می‌توان شیوه‌های مرسوم در کنترل دیجیتال را برای تحلیل یک فرایند یادگیری استفاده نمود. نتایج حاکی از آن هستند که یک سیستم بهینه، در قالب کنترل دیجیتال به صورت یک سیستم مرده نوس قابل توصیف می‌باشد. تعریف یک سیستم فازی این فضای تصمیم‌گیری و فضای سیستم‌های کنترل دیجیتال، از مطالعات و تحقیقات تکمیلی است که می‌توان در ادامه‌ی این مقاله، متصور شد.



شکل ۴ - پاسخ فرکانسی سیستم‌های تبعیت کننده از سیاست‌های π_0 تا π_∞ در خانه‌ی s^3 از جدول نشان داده شده در شکل ۳ همان طور که در شکل ۳ مشاهده می‌شود، عملکرد سیاست‌های قبل از π_∞ ، در موزه حالت‌های s^3 و s^{12} ، مانند π_∞ نمی‌باشد و احتمالاً این سیاست‌ها نتایج ضعیفی را برای این حالات در پی دارند. اگر ارزش‌های تمام حالات، به صورت خروجی‌های سیستم دینامیکی (۹) تعریف شوند، این سیستم، یک سیستم یک-ورودی و چند خروجی خواهد بود. تابع تبدیل متناظر با این سیستم، یک بردار با ۱۴ مولفه است. پاسخ فرکانسی این سیستم، حاوی اطلاعات مهمی در مورد محیط و سیاست اعمال به‌کار رفته از طرف عامل، می‌باشد. به عنوان مثال، پاسخ فرکانسی مولفه‌ی سوم تابع تبدیل سیستم به ازای سیاست‌های π_0 تا π_∞ ، در شکل ۴ ترسیم شده است. عملکرد حالت ماندگار برای این مولفه از سیستم، متناظر با مقدار پاسخ فرکانسی در فرکانس صفر است. دیده می‌شود که دامنه‌ی پاسخ فرکانسی در فرکانس صفر، به ازای سیاست‌های π_0 تا π_∞ ، به ترتیب عبارت است از: ۲۲، ۹، ۴ و ۳. این مقادیر نشان دهنده‌ی متوسط تعداد حرکت‌هایی هستند که با تبعیت از هر سیاست، عامل از s^3 به یکی از خانه‌های هدف می‌رسد. مشاهده می‌شود که به ازای π_∞ کمترین تعداد ممکن از حرکت‌ها به دست آمده است. از روی جدول نشان داده شده در شکل ۲، به وضوح معلوم است که کمترین تعداد حرکت برای رسیدن از s^3 به هر کدام از اهداف، برابر با ۲ حرکت می‌باشد. همچنین در شکل ۴، مشاهده می‌شود که سیستم‌های متناظر با سیاست‌های بهتر، پهنای باند بیشتری دارند. با توجه به این که پهنای باند، معیاری از سرعت پاسخ‌دهی هر سیستم می‌باشد، می‌توان استدلال کرد که برای سیاست‌های بهتر، سرعت پاسخ‌دهی بیشتر می‌باشد.

مراجع

- [1] Richard S. Sutton and Andrew G. Barto, *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, 2002.
- [2] S. I. Reynolds, "Reinforcement Learning with Exploration," *Ph.D. Thesis*, School of Computer Science, The University of Birmingham, UK, 2002.
- [3] Leslie Pack Kaelbling, Michael L. Littman and Andrew W. Moore, "Reinforcement Learning: A Survey," in *Journal of Artificial Intelligence Research*, Vol. 4, pp. 237-285, 1996.
- [4] Tom M. Mitchell, *Machine Learning*. McGraw Hill, 1997.
- [5] Martin L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*, John Wiley & Sons, Inc. 2005.
- [6] Qiyong Hu and Wuyi Yue, *Markov Decision Processes with Their Applications*, Springer Science+Business Media, LLC, 2008.
- [7] Heyeong Soo Chang et al., *Simulation-based Algorithms for Markov Decision Processes*, Springer-verlag, London, 2007.
- [8] R. E. Bellman, *Dynamic Programming*, Princeton University Press, Princeton, 1957.
- [9] D. P. Bertsekas, *Dynamic Programming and Optimal Control*, Athena Scientific, Belmont, 1995.
- [10] D. P. Bertsekas and J. N. Tsitsiklis, *Neural Dynamic Programming*, Athena Scientific, Belmont, 1996.
- [11] H. Royden, *Real Analysis (3rd Edition)*, Prentice Hall, 1988.
- [12] K. Ogata, *Discrete-Time Control Systems (2nd Edition)*, Prentice Hall, 1994.