



دانشگاه خوارزمی تهران

دانشکده فنی و مهندسی

پایان نامه کارشناسی ارشد مهندسی کامپیوتر، گرایش علوم تصمیم و مهندسی دانش

عنوان :

یادگیری تقویتی در فضای پیوسته و

تقریب تابع به وسیله شبکه عصبی خودسازمان ده **DIGNET**

دانشجو :

نجمه علی بابایی

استاد راهنما :

دکتر میر محسن پدرام

پاییز ۱۳۹۱

یادگیری تقویتی در فضای پیوسته

و تقریب تابع به وسیله شبکه عصبی خودسازمان ده

DIGNET

ای دو جهان از قلمت یک رقم

بی قلمت لوح دو عالم عدم

در کف من مشعل توفیق نه

ره به نهانخانه تحقیق ده

شمع زبانم سخن افروز ساز

شام من از صبح سخن روز ساز

تقدیم به پدر و مادر

بزرگوارانی که همواره

مدیون زحماتشان

شرمنده محبتشان

ومنتظر دعای خیرشان

. هستم.

با سپاس از سه وجود مقدس:

آنان که ناتوان شدند تا ما به توانایی برسیم...

موهایشان سپید شد تا ما روسفید شویم...

و عاشقانه سوختند تا گرمابخش وجود ما و روشنگر راهمان باشند...

پدر

مادر

استاد

چکیده

یادگیری تقویتی عبارت است از قالب بندی یک مسئله به فرم یادگیری از طریق تعامل برای رسیدن به هدف. زمانی که فضای حالات و یا کنش‌ها پیوسته و یا خیلی بزرگ شود استفاده از عناصر حافظه برای نگهداری ارزش حالات بسیار زیاد خواهد شد. این مسئله در رابطه با زمانی که در طول یادگیری ارزش کنش‌ها به دست آورده می‌شود، بحرانی تر خواهد شد. علاوه بر حافظه مصرفی مسئله، داده و زمان لازم برای پر کردن آنها نیز مهم است. بنابراین مسئله تعمیم پیش خواهد آمد. روش پیشنهاد شده برای مسائل یادگیری تقویتی با فضای پیوسته‌ای از حالات و کنش‌ها مناسب است.

در این روش از شبکه عصبی خودسازمان ده **DIGNET** برای نمایش فضای حالت و کنش استفاده شده است. استفاده از این شبکه سبب می‌شود در یک فضای پیوسته، عامل بتواند با استفاده از یک حافظه مصرفی مناسب، میزان داده و زمان قابل قبول به هدف دست یابد. در پیشنهاد این پایان نامه مفاهیم اساسی یک مسئله یادگیری تقویتی و همچنین یک شبکه عصبی خودسازمان ده **DIGNET** بررسی می‌شود، و ساختار این شبکه به عنوان یک روش خوشبندی کارا برای حل چالش پیوستگی فضای حالات و کنش‌ها، در مسئله یادگیری تقویتی پیشنهاد می‌شود و کارایی آن بررسی و مقایسه می‌شود.

کلمات کلیدی : یادگیری ماشین ، یادگیری تقویتی، مسیرهای شایستگی ، تعمیم، یادگیری Q ، شبکه عصبی خودسازمان ده **DIGNET**

فهرست مطالب

فصل ۱: مقدمه ۱
۱ ۱-۱: هدف و دستاوردها
۲ ۱-۲: مروری بر فصل‌های رساله
فصل ۲: مروری بر بادگیری تقویتی و شبکه عصبی خودسازمان ده DIGNET ۴
۴ ۴-۱: مقدمه
۵ ۴-۲: ایده کلی
۶ ۴-۳: معادله بلمن
۷ ۴-۴: بهره برداری و اکتشاف
۸ ۴-۵: مفاهیم مورد نیاز
۸ ۴-۵-۱: مدل مارکوف
۹ ۴-۵-۲: تکرار سیاست
۱۰ ۴-۵-۳: تکرار مقدار
۱۰ ۴-۵-۴: روش‌های مقدار کنش
۱۱ ۴-۶: مفاهیم مبتنی بر سیاست و مستقل از سیاست
۱۱ ۴-۷: روش‌های اختلاف زمانی
۱۲ ۴-۸: بادگیری Q
۱۲ ۴-۹: روش‌های نقاد و کنشگر
۱۳ ۴-۱۰: مسیرهای شایستگی
۱۴ ۱۱-۱: شبکه عصبی خودسازمان ده DIGNET
۱۷ ۱۱-۲: مقدمات ریاضی مورد نیاز
۱۸ ۱۲-۱: روال اصلی روش DIGNET در ساختار خوشه‌بندی الگوهای
۲۱ ۱۲-۲: پایداری، همگرایی و ظرفیت DIGNET
۲۴ ۱۳-۱: استفاده از معیارهای شباهت مختلف
۲۵ ۱۳-۲: خلاصه

فصل ۳: مروری بر کارهای انجام شده در محیط‌های پیوسته ۲۶	
فصل ۴: روش پیشنهادی و آزمون آن ۳۱	
۳۱.....	۱- مقدمه ۴
۳۱.....	۲- روش پیشنهاد شده ۴
۳۴.....	۳- آزمون روش ۴
۳۴.....	۴- ۱- مساله کنترل بازوی دو اتصالی با هدف پایدار ۴
۳۸.....	۴- ۲- مساله کنترل بازوی دو اتصالی ناپایدار ۴
۴۶.....	۴- ۳- مساله کنترل شناور زیر دریایی ۴
۵۱.....	فصل ۵: جمع بندی و پیشنهادها ۵۱
۵۱.....	۱- جمع بندی ۵
۵۳.....	۲- پیشنهادهایی برای تحقیقات آتی ۵
۵۴.....	واژه نامه انگلیسی به فارسی ۵۴
۵۷.....	واژه نامه فارسی به انگلیسی ۵۷
۶۰.....	فهرست منابع ۶۰
۶۳.....	چکیده (به زبان انگلیسی) ۶۳

فهرست جدول‌ها، شکل‌ها و نشانه‌های اختصاصی

فصل ۲: موری بر یادگیری تقویتی و شبکه عصبی خودسازمان ده DIGNET ۴	
شکل ۲-۱: الگوریتم یادگیری Q ۱۲	
شکل ۲-۲: شماتی روش نقاد و کنشگر ۱۳	
شکل ۲-۳: همسایگی‌ها بر روی یک کره یکپارچگی ۱۸	
شکل ۲-۴: شماتی عملیات در DIGNET ۲۱	
فصل ۴: روش پیشنهادی و آزمون آن ۳۱	
شکل ۴-۱: الگوریتم پیشنهادی ۳۳	
شکل ۴-۲: دیاگرام مدل پیشنهادی ۳۵	
شکل ۴-۳: کترل بازوی دو اتصالی ۳۶	
شکل ۴-۴: یادگیری نگاشتی از فضای هدف به فضای زوایای تراوم با استفاده از سیگنانل پاداش ۳۶	
شکل ۴-۵: نتایج به دست آمده برای کترول بازوی دو اتصالی ۳۸	
شکل ۴-۶: مساله ربات بازوی دو اتصالی ناپایدار ۳۹	
شکل ۴-۷: مقایسه نتایج به دست آمده در مساله بازوی دو اتصالی ناپایدار ۴۰	
شکل ۴-۸: پراکندگی جاذب‌های شبکه کنش پس از ۵۰۰۰۰ دنباله ۴۲	
شکل ۴-۹: پراکندگی جاذب‌های شبکه کنش پس از تغییر مکان هدف ۴۲	
شکل ۴-۱۰: عمق جاذب‌های شبکه کنش پس از تغییر مکان هدف ۴۳	
شکل ۴-۱۱: عمق جاذب‌های شبکه کنش پس از تغییر مکان هدف ۴۳	
شکل ۴-۱۲: پراکندگی جاذب‌های شبکه حالت پس از ۵۰۰۰۰ دنباله (پیش از تغییر موقعیت هدف) ۴۴	
شکل ۴-۱۳: پراکندگی جاذب‌های شبکه حالت پس از تغییر مکان هدف ۴۴	
شکل ۴-۱۴: عمق جاذب‌های شبکه کنش پیش از تغییر مکان هدف ۴۵	
شکل ۴-۱۵: عمق جاذب‌های شبکه حالت پس از تغییر مکان هدف ۴۵	

شکل ۴-۱: مدل شماتیک شناور زیردریایی	۴۶
شکل ۴-۲: نتایج حاصل از شبکه DIGNET در مساله یادگیری تقویتی کنترل شناور	۴۸
شکل ۴-۳: چند مسیر طی شده توسط عامل در طول زمان یادگیری در مسئله کنترل شناور	۵۰
جدول ۴-۱: مجموعه پارامترهای استفاده شده برای مدل پیشنهادی در مسئله کنترل بازوی دو اتصالی	۳۷
جدول ۴-۲: مراکز جاذب‌های شبکه حالت و کنش پیش از تغییر هدف و پس از تغییر هدف	۴۱

فصل ۱ : مقدمه

چالش اساسی یک تیم روبوکاپ بر سر توسعه فوتbalیستهای مستقل و خودمختار است . گاهی هدف اصلی یک مساله مشخص است اما راه رسیدن به آن دشوار و یا معلوم نیست در این موقع یادگیری تقویتی به عنوان یک ابزار مناسب می تواند با اینگونه مسائل روبرو شود. ایده اصلی این تئوری برپایه جعبه پریچ و خم^۱ است.

هدف اساسی یک تیم روبوکاپ توسعه رباتهایی است که بتواند در مسابقه فوتbal روباتها برنده جام جهانی شوند. در این زمینه روش‌های زیادی استفاده شده است و یکی از آنها یادگیری تقویتی است. ربات به وسیله این روش می تواند خود مختاری^۲ و استقلال را از طریق تجربیات خود بیاموزد. این نوع یادگیری بر پایه پاداش است همانند روشی که برای تعلیم سگها و آموزش ترفندها به آنان به کار می برد.

از جمله روش‌های حل یک مسئله یادگیری تقویتی می توان به برنامه‌ریزی‌پویا^۳، روش‌های مونت‌کارلو^۴ و یادگیری اختلاف زمانی^۵ اشاره کرد. هر یک از این روش‌ها مزایا و معایب خاص خود را دارند. مثلاً نقطه قوت روش‌های برنامه‌ریزی‌پویا خود را در توسعه آن به صورت ریاضی نشان می دهد درحالی که نقطه ضعف آن ریشه در نیاز آن به یک مدل کامل و دقیق از محیط دارد.

۱-۱: هدف و دستاوردها

زمانی که فضای حالات و یا عمل‌ها پیوسته و یا خیلی بزرگ شود استفاده از عناصر حافظه برای نگه داری ارزش حالات بسیار زیاد خواهد شد و این مسئله در رابطه با زمانی که ارزش کنش‌ها به دست می آید، بحرانی‌تر خواهد شد. علاوه بر حافظه مصرفی مسئله ، داده و زمان لازم برای پر کردن آنها نیز مهم است.

¹ Maze

² Autonomous

³ Dynamical programming

⁴ Monte Carlo

⁵ Temporal Difference Learning

بنابراین مسئله تعمیم پیش خواهد آمد. آیا می‌توان با تعداد محدودی حالت یک مسئله را تجربه کرد و نتایج را به مسئله بزرگتر تعمیم داد و تقریبی برای مقادیر به دست آورده تعمیم در این بحث همان ترکیب روش‌های تعمیم و روش‌های یادگیری تقویتی است.

در مبحث یادگیری تقویتی ، تعمیم به عنوان مسئله ای سخت و شاق مطرح می‌شود. به این ترتیب زمانی که یادگیری تقویتی در بسیاری از زمینه‌ها به کار می‌رود، حالاتی پیش می‌آید که پیش از این با آنها برخوردي صورت نگرفته است. این حالت زمانی رخ می‌دهد که فضاهای کنش یا حالت دارای متغیرهای پیوسته و ظواهر پیچیده باشند و تنها راه برای یادگیری این است که عامل بتواند اطلاعاتی که پیش از این آموخته است را برای حالاتی که تا به حال ندیده است ، تعمیم دهد.

در مساله پیش رو پویایی محیطی موجود معلوم نیست. بنابراین از الگوریتم‌های برنامه ریزی پویا نمی‌توان استفاده کرد. دوم اینکه تاحد امکان هم متغیرهای حالت پیوسته و هم متغیرهای کنش پیوسته در نظر گرفته شوند. سوم اینکه با توجه به اینکه در الگوریتم‌هایی مثل مونت کارلو و... که پویایی محیط در آنها مشخص نیست ممکن است عامل با حالاتی مواجه شود که پیش از این هرگز با آنها مواجه نشده است. پس باید به گونه ای عمل کرد که در هر لحظه با مواجه شدن با موقعیت جدید ، عامل بتواند نتایج حاصل از تجربه این موقعیت را به دانش خود بیفزاید. روش پیشنهاد شده برای مسائل یادگیری تقویتی با فضای پیوسته‌ای از حالت‌ها و کنش‌ها مناسب است. در این روش از شبکه عصبی خودسازمان ده *DIGNET* برای نمایش فضای حالت و کنش استفاده شده است. استفاده از این شبکه سبب می‌شود در یک فضای پیوسته، عامل بتواند با استفاده از یک حافظه مصرفی مناسب ، میزان داده و زمان قابل قبول به هدف دست یابد.

۱-۲: مروری بر فصل‌های رساله

محورهای اصلی مطرح شده در این پایان نامه بر سه قسمت کلی استوار است. در قسمت اول بررسی روش‌های موجود در یادگیری تقویتی و جزئیات مربوط به آن است، تئوری یادگیری تقویتی مطرح شده در این مکتب براساس کتاب "Reinforcement Learning: An introduction" (۱) است. قسمت دوم شبکه عصبی

خودسازمان ده *DIGNET* را معرفی و نکات موجود در آن را مطرح می کند و در نهایت قسمت سوم ترکیب دو قسمت پیش را برای به دست آوردن روش پیشنهای ما و پیاده سازی و نتایج کار را نشان می دهد.

در فصل دوم پس از معرفی مفاهیم اولیه سعی می شود به جزئیات روش های برنامه ریزی پویا برای حل یک مسئله یادگیری تقویتی پرداخته شود. پیش از آن خاصیت مهم مارکوف^۶ به عنوان یک پیش نیاز برای تئوری های بعدی معرفی می شود. به علاوه المان های کلیدی در ساختار ریاضی مسئله نظریه توابع مقدار^۷ و معادله بلمن^۸ نیز توضیح داده خواهد شد. به علاوه روش تکرار سیاست^۹ برای به دست آوردن توابع ارزش و اساس روش یادگیری از اختلاف زمانی عنوان خواهد شد.

در فصل سوم به معرفی شبکه عصبی خودسازمان ده *DIGNET* و چند اصل ریاضی به کار رفته در آن ، تشریح روال اصلی این روش در ساختار خوشه بندی الگوها ، پایداری ، همگرایی ، ظرفیت و معیارهای شباخت مختلفی که می تواند در آن به کار رود پرداخته می شود. در فصل چهارم نیز ساختار شبکه عصبی *DIGNET* را درون چند مسئله یادگیری تقویتی به کار برده و ضمن ارائه روش پیشنهادی، نتایج به دست آمده ارزیابی می شود. در نهایت در فصل پنجم به جمع بندی و پیشنهادها پرداخته می شود.

⁶ Markov

⁷ Value Functions

⁸ Bellman Equation

⁹ Policy Iteration

فصل ۲: مرواری بر یادگیری تقویتی و شبکه عصبی خود سازمان ده DIGNET

۱-۲: مقدمه

یادگیری ماشین^{۱۰} علمی است که در آن سعی بر طراحی و گسترش الگوریتم هایی است که از طریق آنها یک ماشین بتواند یاد بگیرد.^(۲) انواع روش های یادگیری در یادگیری ماشین عبارتند از :

۱. یادگیری ناظارت شده^{۱۱}
۲. یادگیری بدون ناظارت^{۱۲}
۳. یادگیری تقویتی^{۱۳}

۱. یادگیری ناظارت شده :

در این نوع یادگیری ، تعدادی ورودی و خروجی معلوم وجود دارد و در طی این یادگیری نگاشتی از ورودی ها به خروجی انجام می شود سپس رفتار مناسب ثبت می شود و عناصر یادگیری سعی می کنند که این رفتار را برای موارد مشاهده نشده کپی کنند^(۲). یک مثال برای این روش الگوریتم پس انتشار خطای^{۱۴} می باشد که در شبکه های پرسپترون^{۱۵} برای آموزش ، مورد استفاده قرار می گیرد .

۲. یادگیری بدون ناظارت:

در این نوع یادگیری ، سعی بر یافتن ارتباطات موجود بین ورودی ها است و تنها از ورودی ها برای ارائه رفتار خروجی استفاده می کنند. نمونه هایی از این نوع یادگیری در مساله خوشه بندی^{۱۶} و طبقه بندی و تخمین ابعاد و مواردی همچون تشخیص چهره^{۱۷} مشاهده می شود. این روش برای تفسیر و پردازش

¹⁰ Machine learning

¹¹ Supervised learning

¹² Unsupervised learning

¹³ Reinforcement Learning

¹⁴ Error Backpropagation

¹⁵ Perceptron Network

¹⁶ Clustering

¹⁷ Face recognition

مقدایر زیادی داده به کار رفته و عنوان می‌شود که هنگامی که بتوان حواس پنجگانه دریافتی را تفسیر و پردازش کرد به نوعی فرمی از هوشمندی وجود دارد^(۲). یک مثال برای این روش الگوریتم کوهن^{۱۸} می‌باشد.

۳. یادگیری تقویتی

قالب بنده یک مسئله به فرم یادگیری از طریق تعامل برای رسیدن به هدف^(۳). آنچه که واقعاً در یک مسئله یادگیری تقویتی روی می‌دهد به شرح ذیل است: یک یادگیرنده و تصمیم گیرنده با نام عامل^{۱۹} شناخته می‌شود. این عامل مرتب در مقاطع زمانی با محیط تعامل دارد. به این ترتیب که عامل در ابتدا در یک حالت^{۲۰} از محیط قرار دارد. سپس با توجه به حالتی که در آن قرار دارد یک کنش^{۲۱} از میان کنش‌های ممکن در آن حالت را انجام می‌دهد. سپس محیط، یک حالت جدید را به عامل نشان می‌دهد. به علاوه یک مقدار پاداش^{۲۲} هم به ازای کنش مورد نظر به عامل اختصاص می‌دهد. هدف یک عامل بیشینه ساختن این پاداشها است که در طی یک دوره بلند مدت به دست می‌آید.

در هر گام ، عامل یک نگاشت از حالات به احتمال انتخاب کنش‌ها انجام می‌دهد که این نگاشت با نام سیاست^{۲۳} شناخته می‌شود و با π_t نشان داده می‌شود. به عبارت دیگر $\pi_t(s, a)$ یعنی اینکه اگر حالت جاری s باشد ($s_t = s$) احتمال اینکه کنش a باشد برابر است با $\pi_t(s, a)$.

۲-۲: ایده کلی

در یک مسئله یادگیری تقویتی عامل سعی می‌کند بر اساس تجربیات خود یاد بگیرد و از طرف دیگر ناظری نیز وجود ندارد. یعنی عاملی وجود دارد که به طور همزمان هم برای به دست آوردن دانش جدید تلاش می‌کند و هم تصمیمات خود را بر اساس دانش موجود بهینه می‌کند .

¹⁸ Kohonen

¹⁹ Agent

²⁰ State

²¹ Action

²² Reward

²³ Policy

در همه این مسائل عملی یک موضوع اساسی وجود دارد و آن هم به وجود آوردن نوعی تعادل بین بهدست آوردن حداکثر پاداش بر اساس دانشی که قبل بهدست آمده است و همینطور تلاش برای انجام یک کنش جدید است به طوری که بتواند دانش بعدی را افزایش دهد که این مسئله با نام اکتساف در مقابل بهره‌برداری^{۲۴} در یادگیری تقویتی شناخته می‌شود. به عبارت دیگر در این مسئله نیاز به جست و جو در فضای کنش‌ها می‌باشد. لذا گفته می‌شود که بازخورد در این مسئله صرفاً ارزیابی^{۲۵} است و این برخلاف حالتی است که در مسئله یادگیری نظارت شده وجود دارد.

۲-۳: معادله بلمن^{۲۶}

تقریباً تمامی الگوریتم‌های یادگیری تقویتی بر پایه تخمین توابع ارزش هستند و آنها توابعی از حالات و یا حالت - کنش‌ها هستند. منظور این است که برآورده شود که چقدر خوب است که یک عامل در یک حالت معلوم باشد. یا اینکه برآورده شود که چقدر خوب است که عامل در یک حالت معلوم یک کنش مشخص را انتخاب کند(۴).

ارزش اتخاذ کنش ، a ، در یک حالت s و تحت سیاست ، π ، با $Q^\pi(s, a)$ نشان داده می‌شود. به صورت زیر برای هر π و هر حالت s :

$$\sum_a \pi(a) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^\pi(s')] \quad 1-2$$

معادله بالا معادله بلمن نامیده می‌شود. معادله بلمن ارتباط بین ارزش یک حالت و ارزش حالت‌های جایگزینش را نشان می‌دهد. معادله بلمن میانگین وزن دار همه حالات ممکن را براساس احتمال رخداد هریک محاسبه می‌کند.

²⁴ Exploration vs exploitation

²⁵ Evaluation

²⁶ Bellman equation

همیشه حداقل یک سیاست وجود دارد که بهتر یا مساوی همه سیاست‌ها باشد اما ممکن است چندین سیاست برتر وجود داشته باشد به هر حال این سیاست‌های بهینه با π^* نمایش داده می‌شود و برای آنها یک تابع ارزش حالت بهینه مشترک به نام V^* استفاده می‌شود:

$$V^*(s) = \max_{\pi} V^\pi(s) \quad \forall s \in S \quad 2-2$$

و همچنین یک تابع ارزش کنش بهینه مشترک به نام Q^* به صورت زیر

$$Q^*(s, a) = \max_{\pi} Q^\pi(s, a) \quad \forall s \in S, a \in A(s) \quad 3-2$$

معادله بهینگی بلمن برای Q^* به صورت زیر است:

$$\begin{aligned} Q^*(s, a) &= E \left\{ r_{t+1} + \gamma \max_{a'} Q^*(s_{t+1}, a') \middle| s_t = s, a_t = a \right\} \\ &= \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma \max_{a'} Q^*(s', a')] \end{aligned} \quad 4-2$$

اگر V^* موجود باشد می‌توان سیاست بهینه را مشخص کرد، به این صورت که برای هر حالت s ، یک یا چند کنش وجود دارد که در معادله بهینگی بلمن مقدار حداکثر را به دست آورده است.

$$V^*(s) = \max_{a \in A(s)} \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^*(s')] \quad 5-2$$

هر سیاست که به این کنش‌های ماکسیمم احتمال غیر صفری را نسبت دهد یک سیاست بهینه است.

۴-۲: بهره برداری و اکتشاف

اغلب برای انجام یک کنش یا انتخاب یک گزینه نمی‌توان هم اکتشاف^{۲۷} و هم بهره برداری^{۲۸} انجام داد و این موضوع با عنوان جنگ میان شناسایی و اکتشاف شناخته می‌شود. اینکه اکتشاف در یک حالت خاص بهتر است

²⁷ Exploration

²⁸ Exploitation

یا بهره برداری به روش پیچیده ای براساس مقادیر غیر مبهم از مقادیر برآورده شده ، عدم قطعیت‌ها و تعداد مراحل باقیمانده بستگی دارد^(۵) . در زمینه تعادل و توازن بهره برداری و اکتشاف‌ها روشهای توانمندی وجود دارد که بیشتر آنها از این فرض استفاده می‌کنند که اطلاعات قبلی و ثابت معلومی در این کاربردها و مسائل یادگیری تقویتی وجود ندارد.

۲-۵: مفاهیم مورد نیاز

۱-۵-۲: مدل مارکوف

در چهار چوب یک مسئله یادگیری تقویتی ، یک عامل تصمیمات خود را به صورت تابعی از یک سیگنال محیط یعنی حالت می‌گیرد. در یک پرتاب گلوله مکان آن و سرعت و شتاب آن مهم است ولی این موضوع اهمیتی ندارد که این سرعت و شتاب چگونه و از کجا آمده است. لذا به این موضوع با نام "استقلال از مسیر"^{۲۹} اشاره می‌شود، چرا که تمامی چیزهایی که اهمیت دارند در سیگنال حالت فعلی هست و با این حال از مسیر یا تاریخچه مستقل است. اما توضیح ریاضی خاصیت مارکوف :

یک کنش در زمان t اتخاذ می‌شود و محیط در زمان $t+1$ به آن پاسخ می‌دهد. اما شاید پاسخی که در گام زمانی فعلی صادر شده است وابسته به مواردی باشد که قبلاً انجام شده باشد به عبارتی این مطالب حاکی از درنظرگرفتن پویایی مسئله به وسیله توزیع احتمال کل است.

$$Pr \{s_{t+1} = s', r_{t+1} = r | s_t, a_t, r_t, s_{t-1}, a_{t-1}, \dots, r_1, s_0, a_0\}$$

۶-۲

. $s_t, a_t, r_t, s_{t-1}, a_{t-1}, \dots, r_1, s_0, a_0$ برای همه 's' و 'r' ها و همه مقادیر ممکن برای رویدادهای گذشته

اما اگر سیگنال حالت ما خاصیت مارکوف داشته باشد، پاسخی که محیط در زمان $t+1$ می‌دهد فقط وابسته به حالت و کنش زمان t است ، لذا پویایی محیط به وسیله احتمال زیر نشان داده می‌شود.

$$Pr \{s_{t+1} = s', r_{t+1} = r | s_t, a_t\}$$

۷-۲

²⁹ Independence of Path

برای همه s' و r و s_t و a_t ها.

یک سیگنال حالت خاصیت مارکوف دارد یا به عبارتی یک حالت مارکوف است هرگاه، به ازای تمام s و r به علاوه $s_0, a_0, r_1, s_1, a_1, r_2, \dots, s_t, a_t$ دو معادله اخیر برابر باشد. در چنین صورتی وظیفه^{۳۰} و محیط به همراه هم خاصیت مارکوف دارند. نتیجه سودمند خاصیت مارکوف برای یک محیط می‌تواند در موارد زیر دیده شود:

۱. به سادگی می‌توان حالت بعدی را با توجه به اطلاعات موجود حالت فعلی و در نتیجه خاصیت

مارکوف پیش‌بینی کرد

۲. خاصیت مارکوف برای یک حالت می‌تواند تشکیل دهنده پایه مناسبی برای انتخاب کنش باشد ،

به این معنا که بهترین سیاست برای انتخاب یک کنش به صورت تابعی از یک حالت مارکوف می‌تواند

به خوبی یک سیاست باشد که تابعی از تاریخچه کامل است(۶).

فرایند تصمیم مارکوف : هر وظیفه یادگیری تقویتی که در آن خاصیت مارکوف وجود داشته باشد با نام فرایند

تصمیم مارکوف شناخته می‌شود. اگر فضای کنش‌ها و حالت‌ها محدود و متناهی باشد آن گاه فرایند مذکور با

نام فرایند تصمیم مارکوف محدود شناخته می‌شود(۷).

۲-۵-۲: تکرار سیاست

هنگامی که با استفاده از π^* ، سیاست ، π به عنوان سیاست بهتر به دست می‌آید می‌توان π^* را محاسبه کرد و

مجددأً از آن استفاده کرد و باز سیاست بهتر π^* را به دست آورد و به همین ترتیب یک توالی از سیاست‌هایی را

به دست آورد که مرتب بهبود یافته‌اند یعنی :



³⁰ Task

³¹ Policy Evaluation

³² Policy Improvement

به این روش یافتن سیاست بهینه، تکرار سیاست گفته می‌شود که از تابع مقدار مربوط به سیاست قبلی شروع می‌شود و نهایت^{۳۳} همگرایی سریعتری از ارزیابی سیاست به دنبال خواهد داشت.

۲-۵-۳: تکرار مقدار

براساس آنچه که در تکرار سیاست گفته شد پس از اینکه π تعیین می‌شود روند ارزیابی سیاست اجرا می‌شود و این روند پس از چند دنباله اجرا به π^I همگرا می‌شود. پس از آن در فاز \rightarrow آن بهبود می‌یابد و مجدداً \rightarrow برروی π^E اجرا می‌شود. همین محاسبات و دوره اجرای مربوط به \rightarrow می‌تواند کار را طولانی کند و این موضوع اشکالی است که بر تکرار سیاست وارد می‌شود.

در حقیقت کوتاه کردن گام \rightarrow در تکرار سیاست به روش‌های مختلفی امکان پذیر است، بدون اینکه در این میان خللی به تضمین همگرایی آن وارد شود. در یکی از مهمترین آنها \rightarrow تنها پس از یک جریان متوقف می‌شود که به این الگوریتم تکرار مقدار گفته می‌شود.

۲-۵-۴: روش‌های مقدار کنش

در واقع بحث بر سر نحوه برآورد مقداربرای کنش‌ها است به عبارت دیگر منظور برآورد میانگین برای توزیع مربوط به هر گزینه است. یک روش ساده برای این منظور استفاده از میانگین پاداش‌هایی است که پیش از این به دست آمده است. به این روش میانگین نمونه‌ها^{۳۴} گفته می‌شود. این روش یکی از روش‌های ساده است اما لزوماً بهترین نیست. روش‌های جایگزین دیگری وجود دارد که در آنها در اکثر موقع اقدام حریصانه انجام داده می‌شود، اما هر از گاهی با احتمال ∞ به صورت یکنواخت و تصادفی از بین کنش‌های موجود بدون در نظر

^{۳۳} Sample-Average Method