

تسريع فرايند يادگيري تقويتی با شکل دهی پاداش به کمک تحليل گراف محیط

مریم مرعشی^۱، علیرضا خلیلیان^۲ و محمد ابراهیم شیری^۳

^۱دانشگاه صنعتی امیرکبیر، marashi_maryam@aut.ac.ir

^۲دانشگاه علم و صنعت ایران، khalilian@comp.iust.ac.ir

^۳دانشگاه صنعتی امیرکبیر، shiri@aut.ac.ir

چکیده - یادگیری تقویتي، به مجموعه روش‌هایی گفته می‌شود که در آن عامل هوشمند با استفاده از تعامل پویا با محیط و دریافت سیگنال‌های تقویتي، رفتار خود را بهبود می‌بخشد. اما این فرایند اغلب بسیار زمان‌گیر، هزینه‌بر و گاهی پر خطر است. پاداش ساختگی، روشی موفق در افزایش سرعت یادگیری عامل در یادگیری تقویتي است. گرچه ایده اصلی این پاداش، دادن یک بازخورد عددی به غیر از پاداش محیط، به عامل یادگیرنده می‌باشد، اما چگونگی محاسبه این پاداش به نحو موثر در محیط‌های بزرگ و واقعی هنوز یک موضوع چالش برانگیز است. الگوریتم پیشنهادی در این مقاله، پاداش ساختگی جدیدی به عامل تزریق می‌کند تا بتواند سرعت یادگیری آن را افزایش دهد. این پاداش بر اساس ساخت گراف محیط، شناسایی اهداف میانی بر اساس معیار مرکزیت میانگی و شناسایی وضعیت‌های کم اهمیت با تحلیل خودکار گراف محیط تنظیم می‌شود. میزان موفقیت روش پیشنهادی روی محیط‌های محک مختلفی چون *maze* و برج هانوی آزمایش گردیده است. نتایج بدست آمده کارایی این روش را نشان می‌دهد.

کلید واژه- بازخورد مصنوعی، پاداش ساختگی، یادگیری تقویتي، یادگیری کبیر

بسیار زیادی به عامل داده می‌شود که در نتیجه عامل برای رسیدن به رفتار بهینه نیازمند صرف زمان بسیار است. راهکارهای مختلفی تا کنون برای غلبه بر این مشکل ارائه شده است [4-5]. یکی از راهکارها، ارائه پاداش ساختگی توسط طراح به عامل یادگیری تقویتي است [6]. ایده اصلی پاداش شکل داده شده، دادن یک بازخورد مصنوعی ایجاد شده توسط طراح به غیر از پاداش محیط، به عامل به منظور افزایش نرخ سرعت یادگیری است [7].

یکی از مشکلات این روش این است که تعیین مقادیر پاداش ساختگی عملاً برای محیط‌های واقعی و بزرگ غیر ممکن است. یکی از راه‌های غلبه بر این مشکل، ایجاد خودکار پاداش ساختگی می‌باشد [8-9-10]. در این مقاله رویکردی مبتنی بر گراف برای ایجاد خودکار پاداش ساختگی در یادگیری تقویتي ارائه شده است. در این رویکرد، عامل بر اساس تحلیل گراف حاصل از محیط، گلوگاه‌های محیط به عنوان وضعیت‌های موثر در رساندن عامل به هدف، ارزش بالایی داده شده و به وضعیت‌هایی که قرارگیری عامل در آنها، تنها منجر به افزایش زمان اکتشاف می‌گردد به عنوان وضعیت‌های کم اهمیت، ارزش کمتری داده می‌شود. در نهایت تمامی این ارزش‌های برآوردی به شکل پاداش

1- مقدمه

یادگیری تقویتي، به مجموعه روش‌هایی گفته می‌شود که در آن عامل هوشمند با استفاده از تعامل با محیط پویا، رفتار خود را بهبود می‌بخشد. در یادگیری تقویتي، با عاملی روبرو هستیم که از طریق سعی و خطا با محیط تعامل کرده و یاد می‌گیرد تا عملی بهینه را برای رسیدن به هدف انتخاب نماید.

یادگیری تقویتي راهی است برای آموزش عامل جهت انجام یک عمل از طریق دادن پاداش و جریمه بدون آنکه لازم باشد نحوه انجام عمل را برای عامل مشخص نمود. در هر مرحله، از یادگیری، عامل کنشی را در محیط انجام می‌دهد و وضعیت او در محیط تغییر می‌کند. عامل در ازای انجام این کنش، پاداشی دریافت می‌کند که این پاداش در بهبود کارایی رفتار عامل هوشمند به کار برده می‌شود. عامل در هر مرحله، کنشی را انتخاب می‌کند که در مجموع بیشترین پاداش را از محیط دریافت کند [1-2]. می‌توان ایده‌ی یادگیری تقویتي را تعامل عامل با محیط برای رسیدن به هدف دانست. اما در بسیاری از کاربردهای واقعی همانند مسابقه فوتبال، پاداش محیط با تاخیر

شکل داده به یادگیری تقویتی داده می‌شود.

ادامه مقاله به شرح زیر است. در بخش دوم به فرآیند تصمیم‌گیری مارکوف و یادگیری تقویتی پرداخته شده است. در بخش سوم شکل‌دهی پاداش در یادگیری تقویتی معرفی می‌شود. در ادامه در بخش چهارم روش پیشنهادی ارائه شده است. در بخش پنجم آزمایشات تجربی و نتایج آن آورده شده است. تحلیل و نتیجه‌گیری هم در بخش ششم ارائه شده است.

2- معرفی

یک نمایش رسمی برای مدل نمودن مسائل یادگیری تقویتی، فرآیند تصمیم‌گیری مارکوف است که به‌طور گسترده در محیط‌های گسسته به‌کار گرفته شده است.

1-2- فرآیند تصمیم‌گیری مارکوف

فرآیندهای تصمیم‌گیری مارکوف (MDP)، که اولین بار در سال 1957 توسط بلمن معرفی گردید [11]، به‌صورت مجموعه پنج تایی (S, A, T, R, α) نشان داده می‌شود که اجزای این مجموعه به‌صورت زیر تعریف می‌شود: S مجموعه وضعیت‌های ممکن در محیط، A مجموعه کنش‌های ممکن برای یک عامل در هر مرحله از تصمیم‌گیری، T نمایانگر احتمال گذر عامل از وضعیت s به وضعیت s' است به شرط انتخاب عمل a از مجموعه اعمال ممکن، $R(s, a, s')$ پاداش دریافتی عامل بعد از انتخاب عمل a و انتقال از وضعیت s به وضعیت s' از محیط، و α یک فاکتور کاهنده.

2-2- یادگیری تقویتی

در یادگیری تقویتی، هدف یافتن یک سیاست بهینه برای مساله مدل شده به یکی از انواع فرآیند مارکوف است. یکی از راه‌های پیدا کردن سیاست بهینه، استفاده از الگوریتم‌های برنامه‌نویسی پویا است. الگوریتم‌های یادگیری کیو [12] و یادگیری سارسا [13] از معروف‌ترین روش‌های یادگیری تقویتی هستند. در الگوریتم یادگیری کیو در حالت گسسته به هر زوج وضعیت-کنش یک مقدار $Q(s, a)$ ، معرف مجموع پاداش‌های دریافتی عامل به شرط شروع از حالت s ، انجام کنش a در ادامه پیروی از سیاست موجود، نسبت داده می‌شود. تابع $Q(s, a)$ ، تا زمانی که به مقدار بهینه همگرا شود با استفاده از فرمول (1) برورسانی می‌شود:

$$Q_{t+1}(s, a) \leftarrow (1 - \alpha) Q_t(s, a) + \alpha [r(s, a) + \gamma \max_{a' \in A_s} Q_t(s', a')] \quad (1)$$

3- شکل‌دهی پاداش در یادگیری تقویتی

توجه الگوریتم‌های موجود در بحث یادگیری تقویتی به یافتن سیاستی با بیشترین مجموع پاداش اکتسابی است. بنابراین می‌توان گفت تابع پاداش به‌طور ضمنی رفتار بهینه را برای عامل توصیف می‌کند. لازم به ذکر است که گرچه پاداش شکل داده شده ابزار قدرتمندی در افزایش سرعت یادگیری عامل است، اما به همان نسبت به کارگیری غلط آن می‌تواند در گمراه نمودن عامل تاثیر شایانی داشته باشد [14]. بررسی‌های مختلف صورت گرفته روی چگونگی تعریف و مقداردهی پاداش ساختگی از سال 1992 را می‌توان به دو دسته تعیین و مقدار دهی توسط متخصص حوزه و برآورد آن به صورت خودکار تقسیم‌بندی نمود.

3-1- استفاده از دانش متخصص حوزه

در سال 1996، Bishop از ایده شبکه عصبی چند لایه برای مقداردهی به پاداش ساختگی استفاده کرد [15]. در سال 2003، بر اساس ایده Laud & Dejong [16] عامل تنها قسمتی از محیط که از نظر فرد خبره، احتمال بیشتری جهت رسیدن به هدف و یافتن سیاست بهینه دارد را مورد بررسی قرار می‌دهد. سپس Wiewiora [17] در سال 2003 به‌صورت عملی ثابت کرد که اگر و تنها اگر پاداش شکل داده شده به‌صورت تابع پتانسیلی به صورت فرمول (2) در اختیار عامل قرار گیرد، در نهایت به سیاست بهینه همگرا خواهد شد.

$$F(s, s') = \gamma Q(s') - Q(s) \quad (2)$$

در سال 2005، Ng و Abbeel [18] روشی ارائه نمود که در آن عامل با مشاهده رفتارهای نشان داده شده توسط متخصص محیط و اجتناب از تعامل با محیط تنها با استفاده از دانش قبلی جهت تعیین چگونگی کار، نه چگونگی اکتشاف، سعی در رسیدن به هدف و یافتن راه حل بهینه می‌کند.

3-2- محاسبه پاداش شکل داده شده به‌صورت خودکار

در سال 2007، Marthi برای اولین بار بحث شکل‌دهی اتوماتیک پاداش را مطرح کرد [19]. او الگوریتمی معرفی می‌کند که مجموعه وضعیت‌های ممکن محیط را به‌عنوان ورودی دریافت کرده و وضعیت‌هایی که عملکرد مشابهی در خصوص یک وظیفه دارند را تجرید می‌نماید و بر اساس آنها تابع پتانسیل جدیدی تعریف کرده و از آن در ساختار شکل‌دهی پاداش استفاده می‌کند. در سال 2010، Kudenko و Grzes [20] به ارائه راه‌کاری برای یادگیری تابع پتانسیل به موازات فرآیند

تخصیص ارزش کمتر می پردازد. در نهایت ارزش مجازی هر وضعیت مشخص شده و بر اساس آن ماتریس پاداش ساختگی، ساخته می شود. هر چند تزریق این پاداش، به تنهایی به فرمول یادگیری کیو تاثیر چشم گیری در بالا بردن سرعت یادگیری عامل دارد، اما آزمایشات ثابت کرده اند که ترکیب این پاداش با تابع پتانسیل، تاثیر به سزایی در سرعت یادگیری عامل موجب می گردد که در قسمت نتایج نشان داده شده است.

4-2- معرفی الگوریتم

با تحلیل گراف حاصل از محیط، به صورت خودکار ارزش مجازی برای رسیدن به هر وضعیت محاسبه می گردد. این پاداش مجازی به دست آمده که به عنوان $Manual Reward(s, a)$ تعریف شده است، همانطور که در الگوریتم شکل 1 دیده می شود، با یک ضریب کاهنده در فرمول یادگیری کیو قرار می گیرد که باعث می گردد به مرور زمان و کسب دانش عامل از محیط، تأثیر خود را از دست بدهد؛ تا چنانچه خطایی در ارزش دهی خودکار وضعیت ها رخ دهد منجر به گمراه نمودن و به هدف نرسیدن عامل نشود و الگوریتم همچنان همگرا باقی بماند. فرمول جدید تعریف شده برای به روزرسانی مقدار $Q(s, a)$ عبارتست از:

$$Q(s, a) = Q(s, a) + \alpha * [r + \max_{a'} Q(s', a') - Q(s, a) + \beta * manual Reward(s, a) + \gamma * \max_{a'} Q(s', a') - Q(s, a)] \quad (3)$$

$$\beta = e^{-\frac{t^2}{\sigma^2}}$$

همانطور که در فرمول (3) دیده می شود، با توجه به اعمال ضریب کاهنده بر پاداش ساختگی، عامل در اپیزودهای اولیه کاملاً تصادفی حرکت کرده؛ اما به مرور زمان و با کسب دانشی قابل قبول از محیط، حرکات خود را به صورت حریصانه بر مبنای بیشترین پاداش اکتسابی انتخاب خواهد کرد.

عامل در اپیزودهای اولیه پاداشی از طرف محیط دریافت نمی کند. لذا الگوریتم پیشنهادی مقدار پاداش ساختگی را به صورت خودکار از تحلیل گراف به دست می آورد و به عامل تزریق می نماید. عامل از این پاداش مصنوعی برای بهبود رفتار خود استفاده می کند تا عملیات اکتشاف مورد نظر برای شناسایی و ارزش دهی به وضعیت ها کاهش یابد. با توجه به اعمال تابع پتانسیل و عدم نیاز به منتظر ماندن عامل تا لحظه رسیدن به هدف نهایی و نیز دریافت بازخورد از محیط، این امر می تواند سرعت یادگیری را افزایش دهد. با توجه به ضریب تعریف شده برای این پاداش، علاوه بر افزایش سرعت یادگیری عامل، نقش

یادگیری تقویتی بر اساس الگوریتم R-max پرداخته اند. الگوریتم ارائه شده به ساختن و تعریف مدلی پویا از محیط برای یادگیری تابع پتانسیل به صورت بر خط می پردازد. اما با توجه به وابستگی این تابع در تمامی الگوریتم های پیشنهادی، به دانش متخصص حوزه و یا ارزش وضعیت ها بر اساس دانش یادگرفته شده ی عامل از اکتشاف محیط، هیچ کدام نتوانستند تأثیر به سزایی در افزایش سرعت یادگیری به خصوص در محیط های بزرگ و واقعی داشته باشند.

4- روش پیشنهادی برای تخمین پاداش ساختگی

در این بخش الگوریتم جدید پیشنهادی برای تخمین تابع پاداش ساختگی، مستقل از دانش متخصص حوزه و اکتشاف عامل در محیط ارائه می شود. در قسمت اول، نحوه نگاشت محیط به گراف و تحلیل آن را بیان کرده، در قسمت دوم الگوریتم پیشنهادی به همراه توضیح آن آورده شده است. این الگوریتم مبتنی بر گراف است و در آن تابع پاداش ساختگی جدید بر اساس ارزش بدست آمده از بررسی و تحلیل گراف، یافتن زیر گراف های مستقل، شناسایی گره های گلوگاه (اهداف میانی) و نیز گره های کم اهمیت، تعریف می شود.

4-1- نگاشت محیط به گراف و تحلیل آن

برای ساخت گراف غیر جهت دار و بدون وزن $G(V, E)$ بر مبنای فرآیند تصمیم گیری مارکوف، هر وضعیت $s \in S$ به عنوان یک گره v در گراف و هر انتقال $s' \rightarrow s$ در تابع انتقال فرآیند، به عنوان یک یال از $v' \rightarrow v$ در نظر گرفته می شود. سپس با تحلیل خودکار گراف، ارزش تقریبی اکتشاف هر وضعیت برآورد شده و تابع پاداش شکل داده شده بر اساس آن بدست می آید. این امر باعث می گردد تا مقیاس پاداشی که شکل داده شده نه تنها وابسته به دانش متخصص حوزه که قطعاً ضریب خطای بالایی دارد نباشد، بلکه مستقل از ارزش به دست آمده از اکتشاف محیط توسط عامل بوده و عامل نیازی به یادگیری و برآورد آن نیز نخواهد داشت.

قدم اول این تحلیل، شناسایی اهداف میانی به عنوان وضعیت هایی از محیط است که رسیدن به آنها برای عامل مفید خواهد بود. برای این منظور، روش مبتنی بر گراف با استفاده از خاصیت مرکزیت میانگی گراف بکار گرفته شده است [21]. با شناسایی این نقاط و تعریف پاداش مجازی بیشتر برای آنها، در قدم دوم، به شناسایی نقاط کم اهمیت (کمترین همسایگی به شرط هدف میانی نبودن) بر اساس ماتریس مجاورت گراف و

دانش یادگرفته شده عامل از اکتشاف محیط و همگرایی الگوریتم یادگیری تقویتی حفظ می شود.

5- آزمایش ها و ارزیابی

مؤثر بودن الگوریتم پیشنهادی در افزایش سرعت یادگیری و نرخ همگرایی، با انجام آزمایش روی محیط محک maze با 625 حالت و برج هانوی با 27 حالت نشان داده شده است. برج هانوی یک محیط دارای ویژگی غیر ناوبری است. در حالیکه محیط maze یک محیط ناوبری شمرده می شود. محیط ناوبری محیطی است که عامل یک شی را از نقطه الف به نقطه ب انتقال می دهد. در آزمایش های انجام شده مقادیر زیر مورد استفاده قرار گرفته است: $\alpha=0.5$, $\gamma=0.9$, $\sigma=0.9$. همچنین مقدار پاداش و جریمه برای اهداف میانی و حالت های کم اهمیت 1 درصد پاداش هدف در نظر گرفته شده است. علامت مقدار پاداش اهداف میانی مثبت و علامت مقدار جریمه حالت کم اهمیت منفی است. آزمایش های هر دو محیط برای 2000 اپیزود انجام شده اند. همچنین، در 50 اپیزود اول شبیه سازی از روش اکتشاف تصادفی [1-21] و در سایر اپیزودها از اکتشاف با روش ϵ -greedy [1-21] استفاده شده است. از آنجائیکه نمودار روش پیشنهادی در هر دو آزمایش حدوداً از اپیزود 1000 به بعد همگرا شده است، برای صرفه جویی در فضا و نمایش بهتر، نمودار

دو محیط فقط تا اپیزود 1000 رسم شده اند. شکل 2 نتایج شبیه سازی ها را به صورت میانگین مجموع پاداش دریافتی و همگرایی به سیاست بهینه بر حسب تعداد اپیزودها به ترتیب روی محیط های آزمایشی برج هانوی (قسمت الف) و maze (قسمت ب) نشان داده اند. در این شکل ها نتایج روش پیشنهادی با نتایج یادگیری کیو مقایسه شده است. طبق آنچه در شکل ها مشاهده می شود، سرعت همگرایی توسط روش پیشنهادی به طور قابل ملاحظه ای بهبود یافته است. برای تعیین اینکه آیا بهبود به دست آمده توسط روش پیشنهادی از لحاظ آماری معنی دار است، از یک روش آماری به نام آزمون فرض برای تفاضل دو میانگین استفاده شده است [22]. با استفاده از فرمول (4)، مقدار z را محاسبه کرده و در جدولی از مقادیر بحرانی جستجو می کنیم. با محاسبه روی داده های آزمایش، مقدارهای 12/08 و 7/75 به ترتیب در محیط های maze و برج هانوی برای z محاسبه شد. با جستجو در جدول مقادیر بحرانی [22]، فرض صفر با احتمال بیشتر از 99/99 درصد برای هر دو محیط رد می شود. رد فرض صفر دلالت می کند بر اینکه بهبود به دست آمده توسط روش پیشنهادی از نظر آماری معنی دار است.

$$(4) \quad \frac{\bar{X}_1 - \bar{X}_2 - \delta}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

input: The matrices Q , $reward$, and $next$ for a given environment

output: The final r value for each episode of the simulation

declare:

Beta($episode$): A function to reduce the effect of the domain expert reward values, it is defined as follows:
 $\exp[-(episode * episode) / (\sigma * \sigma)]$ in which σ is a constant value between 0 and 1.

Modify($manualReward$): A function that finds sub goals and bad states automatically, and then modifies the $manualReward$ matrix according to these special states.

algorithm QLearningWithShapingAndManualReward **begin**

Initialize $manualReward(s, a)$ with $reward(s, a)$

Modify($manualReward$)

for each episode **do begin**

Initialize s , e.g. randomly

$r := 0$

repeat (for each step, action, of episode)

Choose a from s using policy derived from Q (e.g., ϵ -greedy)

$s' := next(s, a)$

$r := r + reward(s, a)$

$Q(s, a) := Q(s, a) + \alpha * [r + \max_{a'} Q(s', a') - Q(s, a) + \text{Beta}(episode) * manualReward(s, a) + \gamma * \max_{a'} Q(s', a') - Q(s, a)]$

$s := s'$

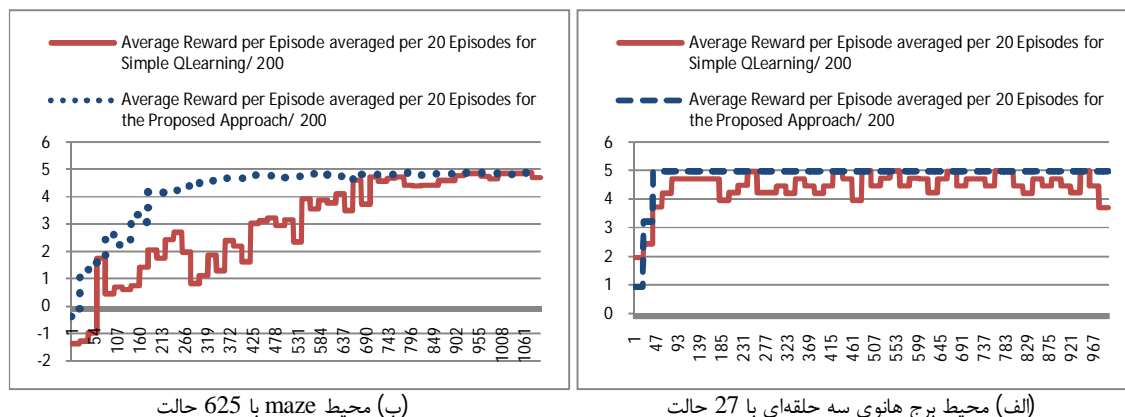
until s is terminal **or** the number of actions exceeds a threshold

Store the number of current episode and the corresponding final r value.

endfor

end QLearningWithShapingAndManualReward

شکل 1: شبه کد الگوریتم پیشنهادی



شکل 2: مقایسه میانگین پاداش دریافتی در هر 20 اپیزود و تعداد 1000 اپیزود شبیه‌سازی برای دو روش یادگیری کیو (خط صاف) و روش پیشنهادی (خط نقطه چین)

6- نتیجه‌گیری

در این مقاله، روشی جدید با استفاده از تحلیل گراف محیط و به‌دست آوردن ماتریس پاداش مجازی به‌عنوان ساختگی پیشنهاد شد. آزمایش‌های تجربی انجام شده افزایش سرعت یادگیری و نرخ همگرایی را توسط روش پیشنهادی نشان می‌دهد. با توجه به آنکه مدل‌سازی هر محیط به گراف امکان‌پذیر بوده و زمانی چندجمله‌ای را می‌طلبید، به نظر می‌رسد این الگوریتم روی اغلب محیط‌های بزرگ و واقعی قابل پیاده‌سازی خواهد بود.

مراجع

- [8] B.Marathi, Automatic shaping and Decomposition of Reward function, Proceedings of the 24th international conference on Machine learning (ICML), Pages 601 – 608 , 2007
- [9] J.Asmuth And M.L.Littman And R.Zinkov, Potential-based shaping in model based reinforcement learning, Proceeding AAAI'08 Proceedings of the 23rd national conference on Artificial intelligence - Volume 2 , Pages 604-609, 2008
- [10] Marek Grzes, Daniel Kudenko, Learning Shaping Rewards in Model-based Reinforcement Learning, Proc. of AAMAS 2009 ,Workshop on Adaptive Learning Agents, ALA 2009
- [11] L.P. Kaelbling, et al., Reinforcement Learning : A Survey ,Journal Of Artificial Intelligence Research, vol.4, pp.237-285, 1996
- [12] Watkin, Watkin Proof of Q-learning Convergence, 1992
- [13] G. A. Rummery and M. Niranjan, "On-Line Q-Learning Using Connectionist Systems," Cambridge University Engineering Department, 1994.
- [14] Randløv and p.Alstrom, Learning to drive a bicycle using reinforcement learning and shaping, J. (Ed.), In Proceedings of the 15th international conference on machine learning, pages 463-471, Morgan Kaufmann, CA., 1998
- [15] C. M. Bishop, Neural networks for pattern recognition. Oxford University Press., 1996
- [16] A.Laud, and G.Dejong, The influence of reward on the speed of reinforcement learning: An analysis of shaping. Proceedings of the 20th International Conference on Machine Learning (ICML), pages 440-447, 2003
- [17] Wiewiora, E., Potential-based shaping and Q-value initialization are equivalent. Journal of Artificial Intelligence Research., page 205-208, 2003
- [18] Pieter Abbeel and Andrew Y.Ng, Exploration and apprenticeship learning in reinforcement learning, Appearing in Proceedings of the 22th International Conference on Machine Learning (ICML), pages 1-8, 2005
- [19] B.Marathi, S.Russell, Automatic shaping and Decomposition of Reward function., In Proceedings of the 24th International Conference on Machine Learning (ICML), pages 601-608, 2007.
- [20] Marek Grze, Daniel Kudenko., Online learning of shaping rewards in reinforcement learning., Department of Computer Science, University of York, York, YO10 5DD, UK., 2010
- [21] P.Moradi, PhD Thesis, Amirkabir University Iran, 2011
- [22] J. E. Freund, *Mathematical statistics*, 5th ed., Prentice-Hall, 1992
- [1] S.Sutton & A.G.Barto, Reinforcement Learning : An Introduction, 1998
- [2] L.P. Kaelbling, et al. Reinforcement Learning : A Survey ,Journal Of Artificial Intelligence Research, vol., page 237-285, 1996
- [3] M. J. Mataric. Reward functions for accelerated learning. In Proceedings of the 11th International Conference on Machine Learning (ICML), pages 181-189, 1994.
- [4] A.Epshteyn and G.Dejong, Qualitative Reinforcement Learning, Appearing in Proceedings of the 23rd International Conference on Machine Learning, Pittsburgh, PA, 2006.
- [5] Andrew Y.Ng, Shaping and police search in reinforcement learning, PhD Thesis, University of California, Berkeley, 2003
- [6] A.LAUD, Theory And Application Of Reward Shaping In Reinforcement Learning, PhD Thesis, University of Illinois at Urbana-Champaign, 2004
- [7] P.Alstrom and J.Randlov, learning to drive a bicycle using reinforcement learning and shaping. In Proceedings of the 15th international conference on machine learning (ICML), Pages 463 – 471, 1998