

# یک روش جدید مبتنی بر معیارهای آماری توزیع برای تنظیم خودکار نرخ یادگیری آتاماتای یادگیر در محیط‌های پویا

محمدرضا ملاخلیلی میبدی، محمدرضا میبدی

چکیده - یکی از مسائل مطرح در ساخت سیستم های یادگیر نظیر شبکه های عصبی و یا آتاماتای یادگیر، تعیین نرخ یادگیری است. در اکثر موارد از یک الگوریتم کاهش یابنده در طول زمان برای تنظیم نرخ یادگیری استفاده می شود.

در این مقاله یک روش جدید برای تغییر نرخ یادگیری و انطباق سیستم یادگیرنده با وضعیت محیط، برای استفاده در آتاماتای یادگیر پیشنهاد شده است. این روش جدید از برخی معیارهای آماری مربوط به توزیع فعلی به دست آمده برای بردار احتمالات متناظر با اقدامات به منظور تعیین افزایش یا کاهش نرخ یادگیری استفاده می کند. مزیت این روش در آن است که بر خلاف روش های موجود فعلی، در طول فرایند یادگیری هم افزایش و هم کاهش مقدار نرخ یادگیری را - بسته به نتایج مقایسه معیارهای آماری - انجام می دهد و به صورت خودکار نرخ یادگیری را تنظیم می کند.

ضمن تشریح مبانی ریاضی این الگوریتم جدید، عملکرد این الگوریتم را در محیط های تصادفی نمونه بررسی کرده و با مقایسه نتایج به دست آمده نشان داده ایم، روش پیشنهادی جدید به دلیل اینکه در طول زمان یادگیری، همزمان و بر اساس معیارهای تعیین شده، افزایش و کاهش نرخ یادگیری را انجام می دهد، از انعطاف پذیری بیشتری نسبت به روش های قبلی برای انطباق با محیط های تصادفی پویا برخوردار است و مقادیر یادگرفته شده به مقادیر حقیقی نزدیکتر هستند.

کارهایی که در زمینه انطباقی کردن نرخ یادگیری در آتاماتای یادگیر انجام شده است را می توان به دو گروه تقسیم کرد: وابسته به مساله و مستقل از مساله. در روشهای پیشنهادی وابسته به مساله، بر اساس نوع مساله ای که آتاماتای یادگیر و یا اجتماعی از آتاماتاهای یادگیر درصدد حل آن هستند روشهایی برای تطبیقی کردن نرخ یادگیری پیشنهاد شده است [۴] [۵] [۶] [۷]. در تمام این نمونه ها، بر حسب نوع مساله، روشی برای به روزرسانی نرخ یادگیری پیشنهاد شده است که سرعت همگرایی و حل مساله را افزایش داده است. رده وسیعی از این روش های انطباقی را در [۴] می توانید مشاهده کنید که در آن نویسنده با توجه به مساله یافتن درخت پوشای کمینه در گراف های تصادفی، از پاره ای ویژگیهای توزیع احتمال اقدام ها به شکل مستقیم برای تطبیق کردن نرخ یادگیری و رسیدن به همگرایی سریع تر استفاده کرده است. مشابه همین ایده نیز در [۸] و [۶] مورد استفاده قرار گرفته است. اما در گروه دوم، مستقل از مساله یا ساختار اجتماع آتاماتاهای یادگیر، روشهایی برای این کار پیشنهاد می شود.

از آنجا که غالباً، استفاده از آتاماتاهای یادگیر به شکل شبکه ای ساختارمند صورت می گیرد [۹]، بنابراین روشهای پیشنهادی برای انطباقی کردن نرخ یادگیری غالباً وابسته به مساله و ساختار بوده و از مساله ای به مساله دیگر و از ساختاری به ساختار دیگر متفاوت هستند. پژوهش مستقلی که به بررسی تنظیم خودکار نرخ یادگیری در آتاماتای یادگیر (فارغ از مساله ای که آتاماتا در صدد حل آن است و یا ساختار شبکه ای از آتاماتاها که برای حل مساله مورد استفاده قرار می گیرد) پرداخته باشد، مشاهده نشده است. عموم پژوهش های این بخش مربوط به شبکه های

کلید واژه - آتاماتای یادگیر، نرخ یادگیری پویا، تنظیم نرخ یادگیری، نابرابری چیشف.

۱- مقدمه

بحث تنظیم خودکار نرخ یادگیری یکی از مسائلی است که در مباحث مرتبط با الگوریتم ها و سیستم های یادگیری و خصوصاً در مباحث مربوط به شبکه های عصبی به کرات مورد بررسی قرار گرفته است [۱] [۲] [۳].

اکثر الگوریتم های موجود، از یک نرخ یادگیری با مقدار بالا شروع کرده و در حین فرایند آموزش به کمک یک تابع کاهش یابنده با زمان (مثلاً

مقاله در تاریخ ۳۰ خرداد ماه ۱۳۹۱ دریافت شد.  
محمدرضا ملاخلیلی میبدی، دانشکده مهندسی کامپیوتر، دانشگاه آزاد اسلامی واحد علوم و تحقیقات تهران، حصارک تهران (email: m.meybodi@srbiau.ac.ir)  
محمدرضا میبدی، دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیرکبیر، تهران (email: mmeybodi@aut.ac.ir).

وضعیت بعدی

$$G \equiv \phi \rightarrow \alpha$$

تابع خروجی که وضعیت فعلی را به

خروجی بعدی نگاشت می‌کند.

$$\phi(n) \equiv \{\phi_1, \phi_2, \dots, \phi_k\} \quad \text{در} \quad \text{آتاماتا} \quad \text{زمان } n$$

مجموعه  $\alpha$  شامل خروجیهای (اقدامهای) آتاماتا است که آتاماتا در هر گام یک اقدام از  $r$  اقدام این مجموعه را برای اعمال بر محیط انتخاب می‌نماید. مجموعه ورودی‌ها ( $\beta$ ) ورودیهای آتاماتا را مشخص می‌کند. توابع  $F$  و  $G$  وضعیت فعلی ورودی را به خروجی بعدی (اقدام بعدی) آتاماتا نگاشت می‌کنند. اگر نگاشتهای  $F$  و  $G$  قطعی باشند، آتاماتا یک آتاماتای قطعی نامیده می‌شود. در چنین حالتی با فرض یک وضعیت اولیه و ورودی مشخص، حالت بعدی و خروجی بصورت یکتا مشخص شده‌اند. در حالیکه نگاشتهای  $F$  و  $G$  تصادفی باشند، آتاماتا بعنوان آتاماتای تصادفی معرفی می‌شود.

آتاماتای تصادفی را می‌توان به دو دسته آتاماتای تصادفی با ساختار ثابت و آتاماتای تصادفی با ساختار متغیر تقسیم‌بندی کرد. در آتاماتای تصادفی با ساختار ثابت احتمال اقدامهای آتاماتا ثابت هستند. در حالیکه در آتاماتای تصادفی با ساختار متغیر احتمالات اقدامهای آتاماتا در هر تکرار به‌روز می‌شوند (تغییر احتمالات اقدامها بر اساس الگوریتم یادگیری انجام می‌شود). وضعیت داخلی آتاماتا  $\phi$ ، توسط احتمالات اقدامهای آتاماتا بازنمایی می‌شوند. برای سادگی هر وضعیت داخلی آتاماتا را مطابق با یک اقدام مشخص آتاماتا در نظر می‌گیرند. بنابراین می‌توان وضعیت داخلی آتاماتا  $\phi$  را با بردار احتمال اقدامهای آتاماتا  $P$  که بصورت زیر نشان داده می‌شود، جایگزین نمود.

$$P(n) \equiv \{p_1(n), p_2(n), \dots, p_r(n)\} \quad (۲)$$

به گونه ای که:

$$\sum_{i=1}^r p_i(n) = 1, \quad \forall n, \quad p_i(n) = \text{Prob}[\alpha(n) = \alpha_i] \quad (۳)$$

در آغاز فعالیت آتاماتا، احتمال اقدامهای آن بصورت مساوی با هم برابر با  $p_i = \frac{1}{r}$  قرار داده می‌شوند. (که  $r$  تعداد اقدامهای آتاماتا می‌باشد).

**محیط:** محیط تصادفی را به‌طور ریاضی می‌توان بصورت یک سه‌تایی  $E \equiv \{\alpha, \beta, c\}$  توصیف کرد، به‌طوری‌که:

$$\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\} \quad \text{مجموعه ورودیهای محیط}$$

$$\beta \equiv \{\beta_1, \beta_2, \dots, \beta_r\} \quad \text{مجموعه خروجیهای محیط}$$

$$c \equiv \{c_1, c_2, \dots, c_r\} \quad \text{مجموعه احتمالات جریمه}$$

ورودی محیط یکی از  $r$  اقدام انتخاب شده آتاماتا است. خروجی (پاسخ)

محیط به هر اقدام  $i$  توسط  $\beta_i$  مشخص می‌شود [۱۰]

ارتباط آتاماتای تصادفی با محیط در شکل ۱ نشان داده شده است. از این مجموعه به همراه الگوریتم یادگیری تحت عنوان آتاماتای یادگیر تصادفی نام برده می‌شود. به این ترتیب آتاماتای یادگیر تصادفی را

عصبی می‌باشند که برخی از آنها به لحاظ عمومی که در مبانی دارند- نظیر نرخ یادگیر کاهش‌یابنده با زمان- برای سایر سیستم‌های یادگیر و از جمله آتاماتای یادگیر نیز قابل استفاده هستند.

به نظر می‌رسد ایده‌های تطبیقی کردن نرخ یادگیری با هدف تسریع در همگرایی فرآیند یادگیری صورت گرفته است و مساله انطباقی کردن نرخ یادگیری در محیط‌های پویا چندان مدنظر قرار نگرفته است.

ادامه مقاله بدین صورت سازماندهی شده است. در بخش دوم ضمن بررسی مختصر آتاماتای یادگیر، نامساوی چبیشف را به عنوان یک نامساوی کاربردی در توزیع‌های تصادفی که مستقل از توزیع متغیر تصادفی رابطه‌ای میان مقدار متغیر تصادفی، مقدار میانگین و واریانس آن وضع می‌کند، را بررسی خواهیم کرد. در ادامه همین بخش ایده به کارگیری این نامساوی به عنوان یک معیار تشخیص پویایی محیط مورد بررسی قرار می‌گیرد. در بخش سوم، الگوریتم پیشنهادی مبتنی بر نامساوی چبیشف ارائه خواهد شد. بخش چهارم به بررسی تجربی این الگوریتم در تعدادی محیط پویا اختصاص داده شده است. بخش پنجم به جمع بندی مطالب ارائه شده در مقاله اختصاص یافته است.

## ۲- مبانی روش پیشنهادی

### ۲-۱- بررسی نحوه یادگیری در آتاماتای یادگیر

یک آتاماتای یادگیر به عنوان مدلی از یک سیستم یادگیر است که در محیط‌های تصادفی ناشناخته عمل می‌کند. آتاماتا در هر دور یک اقدام از میان مجموعه محدود اقدام‌های خود انتخاب کرده و با بررسی عکس العمل محیط نسبت به این اقدام، احتمال انتخاب اقدام‌های بعدی را بهبود می‌بخشد [۹].



شکل ۱- آتاماتای یادگیر و نحوه تعامل آن با محیط

یک آتاماتای تصادفی را می‌توان بصورت یک ماشین حالت متناهی در نظر گرفت. به بیان ریاضی نیز می‌توان آنرا بصورت یک پنج‌تایی مانند زیر نشان داد:

$$SA \equiv \{\alpha, \beta, F, G, \phi\} \quad (۱)$$

که

$r$  تعداد اقدامهای آتاماتا

$\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\}$  مجموعه اقدامهای آتاماتا

$\beta \equiv \{\beta_1, \beta_2, \dots, \beta_r\}$  مجموعه ورودیهای آتاماتا

$F \equiv \phi \times \beta \rightarrow \phi$  تابع نگاشت وضعیت فعلی و ورودی به

می‌توان با چهارتایی

$$SLA \equiv \{\alpha, \beta, p, T, c\} \quad (۴)$$

تعریف کرد. به طوریکه:

مجموعه اقدامهای آتاماتا/مجموعه ورودیهای محیط

$$\alpha \equiv \{\alpha_1, \alpha_2, \dots, \alpha_r\}$$

مجموعه ورودیهای آتاماتا/مجموعه خروجیهای محیط

$$\beta \equiv \{\beta_1, \beta_2, \dots, \beta_r\}$$

بردار احتمال اقدامهای آتاماتا

$$p \equiv \{p_1, p_2, \dots, p_r\}$$

الگوریتم یادگیری

$$T \equiv p(n+1) = T[\alpha(n), \beta(n), p(n)]$$

مجموعه احتمالات جریمه که معرف محیط می‌باشند.

$$c \equiv \{c_1, c_2, \dots, c_r\}$$

**الگوریتم یادگیری خطی:** الگوریتم یادگیری یک رابطه بازگشتی

است که برای انجام تغییرات و به روزرسانی در بردار احتمال اقدام های آتاماتا در یک آتاماتای یادگیر تصادفی با ساختار متغیر مورد استفاده قرار می‌گیرد. فرض کنید یک آتاماتای یادگیر تصادفی ساختار متغیر در زمان  $k$  از میان مجموعه اقدامهای  $\alpha$  عمل  $\alpha_i(k)$  را انتخاب کرده باشد. همچنین فرض کنید بردار احتمال انتخاب اقدامهای آتاماتا را با  $p(k)$  نمایش داده‌ایم. اگر  $a$  و  $b$  پارامترهایی باشند که به ترتیب میزان افزایش یا کاهش احتمالات اقدام ها را مشخص می‌کنند و  $r$  تعداد اقدامهای قابل انجام توسط آتاماتای یادگیر باشد، بردار  $p(k)$  توسط الگوریتم یادگیری خطی ارائه شده در روابط زیر به روزرسانی می‌شود. مقدار  $a$  را پارامتر پاداش و  $b$  را پارامتر جریمه می‌نامند

$$p_j(k+1) = \begin{cases} (1-a)p_j(k) + a & j = i \\ (1-a)p_j(k) & \forall j \neq i \end{cases} \quad (۵)$$

$$p_j(k+1) = \begin{cases} (1-b)p_j(k) & j = i \\ (1-b)p_j(k) + \frac{b}{r-1} & \forall j \neq i \end{cases} \quad (۶)$$

رابطه (۵) زمانی مورد استفاده قرار می‌گیرد که عمل  $\alpha_i(k)$  منجر به دریافت پاداش از محیط شده باشد و رابطه (۶) زمانی مورد استفاده قرار می‌گیرد که این عمل به دریافت جریمه از محیط منجر شده باشد.

اگر  $a=b$  باشد روابط یادگیری خطی (معادله‌های (۵) و (۶)) را الگوریتم  $L_{R-P}$  می‌نامند. اگر  $a \gg b$  باشد آن را  $L_{R-EP}$  و اگر  $b=0$  باشد آن را  $L_{R-I}$  می‌نامند [۱۱].

عامل موثر در کارایی آتاماتای ساختار متغیر، الگوریتم های یادگیری هستند که برای به روزرسانی احتمال اقدام ها استفاده می‌شود.

با این مقدمات فرض کنید یک آتاماتای یادگیر با  $r$  اقدام در یک محیط تصادفی فعالیت می‌کند. محیط تصادفی اقدام انجام شده توسط آتاماتا را ارزیابی می‌کند و آتاماتا بر اساس این ارزیابی، بردار احتمالات اقدام های خود را به روزرسانی می‌کند. فرض کنید  $P_i^t$  مقدار احتمال انتخاب اقدام  $i$  ام آتاماتا را در زمان  $t$  نشان دهد. ضمناً

$$\sum_{i=1}^r P_i^t = 1 \quad \forall t \in \{0, 1, \dots\} \quad (۷)$$

بردار احتمال انتخاب اقدام های آتاماتا در زمان  $t$  را با  $\overline{P}(t)$  نشان می‌دهیم که  $\overline{P}(t) = [P_1^t, P_2^t, \dots, P_r^t]$  علاوه بر این محیط تصادفی بر اساس بردار احتمالات  $\overline{C} = [C_1, C_2, \dots, C_r]$  اقدام های محیط را پاداش می‌دهد که در آن  $\sum_{i=1}^r C_i = 1$  است.

در این آتاماتای یادگیر هر اقدام آتاماتا دارای یک نرخ یادگیری است. این بردار را با  $\bar{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_r]$  نشان می‌دهیم. در این مقاله فرض می‌کنیم آتاماتای یادگیر برای به روزرسانی بردار احتمال انتخاب های خود از الگوریتم  $L_{R-I}$  استفاده می‌کند.

**گزاره ۱:** رابطه به روزرسانی بردار احتمال انتخاب اقدام آتاماتای یادگیر برای زمانی که اقدام  $k$  ام توسط آتاماتا صورت گرفته و آتاماتا از الگوریتم یادگیر خطی  $L_{R-I}$  استفاده می‌کند را می‌توان به صورت رابطه (۸) نوشت:

$$P_k^{t+1} = P_k^t * (1 - \alpha_k) + \alpha_k R^{t+1} \quad (۸)$$

در رابطه (۸)،  $R$  پاسخ محیط به اقدام انجام شده توسط آتاماتا است و داریم

$$R^{t+1} = \begin{cases} 1 & \text{with probability } C_k \\ 0 & \text{with probability } 1 - C_k \end{cases} \quad (۹)$$

رابطه یادگیری آتاماتا در این حالت را به شکل برداری (۱۰) میتوان نشان داد

$$\overline{P}^{t+1} = \overline{P}^t * (1 - \alpha) + \alpha \overline{R}^{t+1} \quad (۱۰)$$

رابطه (۱۰) در حقیقت نشان می‌دهد که بردار احتمال انتخاب اقدام های آتاماتای یادگیر در طول فرآیند یادگیری، به بردار میانگین پاسخ های محیط همگرا می‌شود.

## ۲-۲ محیطهای پویا

موفقیت آتاماتای یادگیر در محیطهای پویا بستگی به تغییرات محیط و اطلاعاتی که توسط آتاماتای یادگیر از محیط قابل جمع آوری است، دارد. در [۹] نویسنده یک تقسیم بندی مفصل از محیطهای پویا ارائه داده است. اما از یک دیدگاه کلی می‌توان محیطهای پویا را به دو گروه کلی تقسیم کرد: محیطهای پویا  $MSE$  و محیطهای پویا با تابع احتمال جریمه متغیر با زمان.

در  $MSE$ ها فرض بر این است که محیط پویای واقعی خود از چند محیط ایستا<sup>۱</sup> تشکیل شده است و پویایی محیط ناشی از یک توالی از فرآیندهای جابجایی بین این محیطهای ایستا است [۱۲]. به بیان ریاضی، محیط پویای  $\mathcal{E}$  توسط یک مجموعه  $\{E_1, E_2, \dots, E_H\}$  از محیطهای ایستا و یک ماتریس جابجایی  $T$  تعریف می‌شود. در این ماتریس  $T$ ، عنصر  $T_{i,j}$  نشان دهنده احتمال آن است که اگر آتاماتای یادگیر در حال حاضر با محیط تصادفی  $E_i$  در تعامل است، در گام بعدی با محیط تصادفی  $E_j$  در تعامل باشد. در حقیقت در  $MSE$ ها محیط تصادفی پویا از

<sup>۱</sup> در این مقاله محیط ایستا را معادل با Stationary Environment و محیط پویا را برای non-Stationary Environment در نظر گرفته ایم

نابرابری چندبعدی چیشف تعمیمی از نابرابری چیشف است که به کمک آن میتوان مرزی را برای اینکه یک بردار تصادفی از بردار مقادیر میانگین اش بیش از یک مقدار معین فاصله داشته باشد را تعیین کرد. این نابرابری بدین شکل فرمول بندی می شود.

**نابرابری چیشف (تعمیم یافته):** فرض کنید  $X$  یک بردار تصادفی با میانگین  $\mu = E[X]$  و ماتریس کوواریانس  $V = E[(X - \mu)(X - \mu)^T]$  باشد. اگر  $V$  یک ماتریس تعریف شده مثبت باشد در این صورت برای هر عدد حقیقی  $t > 0$  داریم

$$\Pr((X - \mu)V^{-1}(X - \mu)^T) > qN \leq 1/q^2 \quad (۱۴)$$

که در رابطه (۱۴)  $N = \text{trace}(V^{-1}V)$  است.

## ۲-۴ استفاده از نابرابری چیشف به عنوان یک معیار در تطبیقی کردن نرخ یادگیری

در این قسمت با توجه به گزاره ۱ که نحوه تغییرات بردار احتمال انتخاب اقدام های آتاماتا را نشان میدهد و نابرابری چیشف به ارائه معیاری برای کاهش یا افزایش نرخ یادگیری آتاماتای یادگیر خواهیم پرداخت

**گزاره ۲:** فرض کنید  $X$  متغیری است که بر اساس رابطه (۱۵) به روز رسانی می شود

$$X^{t+1} = X^t * (1 - \alpha) + \alpha R^{t+1} \quad (۱۵)$$

که در آن

$$R^{t+1} = \begin{cases} 1 & \text{with probability } C^* \\ 0 & \text{with probability } 1 - C^* \end{cases} \quad (۱۶)$$

است. ثابت می شود [۱۶]:

$$\lim_{t \rightarrow \infty} E[X^t] = \mu = C^* \quad \text{الف-}$$

ب- واریانس  $X^t$  محدود و دارای مقدار حدی زیر است

$$\text{Var}(X^t) = \delta^2 = \frac{\alpha}{2-\alpha} C^*(1 - C^*)$$

ج- مطابق با نامساوی چیشف داریم:

$$\forall q > 0 \quad P(|X^t - C^*| \geq q\delta) \leq \frac{1}{q^2}$$

**اثبات:**

الف-  $\langle R^{t+1} \rangle$  یک فرآیند تصادفی برنولی با پارامتر  $C^*$  است. علاوه بر این داریم:

$$E[X^{t+1}] = (1 - \alpha)E[X^t] + \alpha E[R^{t+1}] \\ = (1 - \alpha)E[X^t] + \alpha C^* \quad (۱۷)$$

از رابطه بازگشتی (۱۷) داریم:

$$E[X^{t+1}] = (1 - \alpha)^t E[X^0] + (1 - (1 - \alpha)^t) C^* \quad (۱۸)$$

با توجه به اینکه  $0 < \alpha < 1$  رابطه (۱۸) نشان می دهد که

$$\lim_{t \rightarrow \infty} E[X^{t+1}] = C^* \quad (۱۹)$$

ب: به طریق مشابه میتوان نشان داد که

$$\text{Var}[X^{t+1}] = \frac{\alpha}{2-\alpha} C^*(1 - C^*)(1 - (1 - \alpha)^{2t+2}) \quad (۲۰)$$

با توجه به اینکه  $0 < \alpha < 1$  رابطه (۲۰) نشان میدهد که

$$\lim_{t \rightarrow \infty} \text{Var}[X^{t+1}] = \frac{\alpha}{2-\alpha} C^*(1 - C^*) \quad (۲۱)$$

ج- با فرض  $m = E[X]$  برای متغیر تصادفی  $X$  نامساوی مارکوف

مجموعه ای از محیطهای تصادفی ایستا تشکیل شده است که مجموعه حالات یک زنجیره مارکوف را تشکیل می دهند [۹] [۱۳] [۱۲].

مدل دیگری که در متون مربوطه برای محیطهای پویا ارائه شده است، مدلی است که احتمالات جریمه اقدامها را غیر ثابت و متغیر با زمان فرض کرده است [۱۴] [۱۳]. یعنی احتمال جریمه اقدام در عین اینکه به اقدام انجام شده بستگی دارد، به زمان انجام آن نیز بستگی دارد.

مقالاتی که به ارائه الگوریتم های یادگیر مبتنی بر آتاماتای یادگیر برای محیطهای MSE پرداخته اند به دلیل تنوع این محیط ها بسیار زیاد هستند. ایده غالب آنها تشکیل سلسله مراتبی از آتاماتاهای یادگیر است که برخی سطوح وظیفه تعیین محیط ایستا را بر عهده دارند و برخی سطوح دیگر فرآیند یادگیری در آن محیط ایستا را بر عهده دارند [۱۳].

همانگونه که قبلا هم متذکر شدیم، هیچ پژوهش مستقلی که فرآیند یادگیری در آتاماتای یادگیر را از طریق تطبیقی کردن نرخ یادگیری دنبال کرده باشد، در متون و مقالات تخصصی این حوزه یافت نشد و آنچه که از طریق تطبیقی کردن نرخ یادگیری صورت می گیرد، غالبا وابسته به مساله و ساختاری از آتاماتاها که برای حل مساله مورد استفاده قرار گرفته است، می باشد.

## ۲-۳ نابرابری چیشف

برای تنظیم پویای نرخ یادگیری، به منظور انطباق در محیطهای پویا، نیاز به معیاری داریم تا بر اساس آن معیار میزان توانمندی نرخ یادگیری فعلی سیستم یادگیر را در رصد کردن تغییرات محیط بسنجیم. یکی از ابزارهای آماری مناسب در این مورد نابرابری موسوم به نابرابری چیشف است. در نظریه احتمالات نابرابری چیشف تضمین می کند که در هر نمونه تصادفی یا در هر توزیع احتمال، "تقریبا تمامی" مقادیر در نزدیکی میانگین خواهند بود. به طور دقیق تر این قضیه بیان می کند که حداکثر مقادیری که در هر توزیع می توانند بیش از  $k$  برابر انحراف معیار با میانگین فاصله داشته باشند  $\frac{1}{k^2}$  است. [۱۵]

**قضیه (نامساوی مارکوف):** اگر  $X$  یک متغیر تصادفی و  $a$  یک عدد حقیقی مثبت باشد، در این صورت

$$\Pr(|X| > a) \leq \frac{E(X^2)}{a^2} \quad (۱۱)$$

برای اثبات این قضیه میتوانید به مرجع [۱۶] مراجعه کنید.

با استفاده از نامساوی مارکوف میتوان به نامساوی چیشف رسید. اگر  $m$  نشان دهنده میانگین متغیر تصادفی  $X$  باشد با جایگذاری  $X - m$  در رابطه بالا به نامساوی، موسوم به نامساوی چیشف خواهیم رسید.

$$\Pr(|X - m| \geq a) \leq \frac{\text{Var}(X)}{a^2} \quad (۱۲)$$

اگر  $\text{Var}(X)$  را با  $\delta$  نشان دهیم تفسیر دیگری از نامساوی چیشف به دست می آید

$$\Pr(|X - m| \geq a\delta) \leq \frac{1}{a^2} \quad (۱۳)$$

رابطه (۱۳) در حقیقت بیان می کند که احتمال اینکه یک متغیر تصادفی در خارج از بازه ای حول میانگین به شعاع  $a$  برابر واریانس باشد از  $\frac{1}{a^2}$  کمتر است. این رابطه برای مقادیر  $a > 1$  واجد اطلاعات مفید است.

محلی وجود داشته باشد و یا محیط تصادفی، تغییر کند. در الگوریتم پیشنهادی جدید از ویژگی سوم مطرح در گزاره ۲ به عنوان عامل تشخیص دهنده اختلاف میان تخمین جدید و میانگین تخمین های قبلی استفاده می کنیم و به کمک آن در موارد لازم، نرخ یادگیری را افزایش می دهیم.

در اکثر سیستم های یادگیر و از جمله در آتوماتای یادگیر، نرخ یادگیری در مدت آموزش، به شکل پویا - کاهش یابنده با زمان - تغییر می کند. دلیل این امر هم واضح است. در ابتدای یادگیری، مقادیر بزرگتر نرخ یادگیری، باعث آموزش سریع تر می گردند. به تدریج و با افزایش تعداد نمونه ها، سیستم یادگیرنده، سعی می کند بیشتر متکی بر تجربیات آموزشی گذشته باقی بماند تا اینکه سعی کند از نمونه های جدید برای یادگیری استفاده کند. این منطق باعث می شود، در ابتدای فرآیند یادگیری، سیستم یادگیرنده جسورانه تر و با گذر زمان محافظه کارانه تر عمل کند.

اکثر سیستم های یادگیری، از یک نرخ یادگیر پویای کاهش یابنده در حین فرآیند یادگیری استفاده می کنند. نرخ یادگیری با این مفهوم پارامتری است که میزان فراموش کاری سیستم را تعریف می کند. مقادیر کوچکتر این پارامتر، یعنی اتکای بیشتر سیستم به تجربیات گذشته و مقادیر بزرگتر به معنای نادیده گرفتن تجربیات گذشته است.

در بیشتر سیستم ها، منطق بالا پاسخگوی نیازها می باشد. مساله اینجاست که کاهش نرخ یادگیری در طول زمان باعث می شود به مرور زمان، انطباق پذیری سیستم یادگیر کاهش یابد و نسبت به تغییرات در محیط عکس العمل مناسب نشان ندهد. بنابراین برای تنظیم اتوماتیک نرخ یادگیری در طول مدت آموزش بایستی معیاری داشته باشیم تا بر اساس آن نسبت به افزایش (در صورت بروز تغییرات جدی در محیط) یا کاهش (در صورت یکنواخت بودن پاسخ محیط و نزدیک شدن به همگرایی) نرخ یادگیری اقدام کنیم. این معیار بایستی قادر به تعیین میزان اهمیت تغییرات در محیط باشد.

### ۳- الگوریتم پیشنهادی

با توجه به رابطه یادگیری مورد استفاده توسط آتوماتای یادگیر برای هر اقدام و تفسیر برداری آن، می توان از نابرابری چیشف به شکل ساده یا برداری آن برای تنظیم اتوماتیک نرخ یادگیری استفاده کرد. منطق کار بدین صورت است که:

الف- در حالت ساده: در این حالت برای هر یک از اقدام های آتوماتای یک نرخ یادگیری در نظر می گیریم. بر اساس نامساوی چیشف (یا همان مارکوف رابطه (۱۳))، چنانچه اختلاف مقدار احتمال انتخاب یک اقدام آتوماتا با میانگین مقادیر قبلی آن از یک مقدار آستانه بیشتر باشد (نامساوی چیشف)، به منزله تغییرات زیاد در محیط بوده و نرخ یادگیری را افزایش می دهد. همزمان با افزایش نرخ یادگیری در هر بار، مقدار میانگین با صفر مقداردهی شده تا الگوریتم قابلیت گریز از بیشینه های محلی را داشته باشد (گرچه همان طور که توضیح دادیم، این فرار از بیشینه های محلی تضمین شده نیست)

$$P(|X| \geq a) \leq \frac{E(X^2)}{a^2}$$

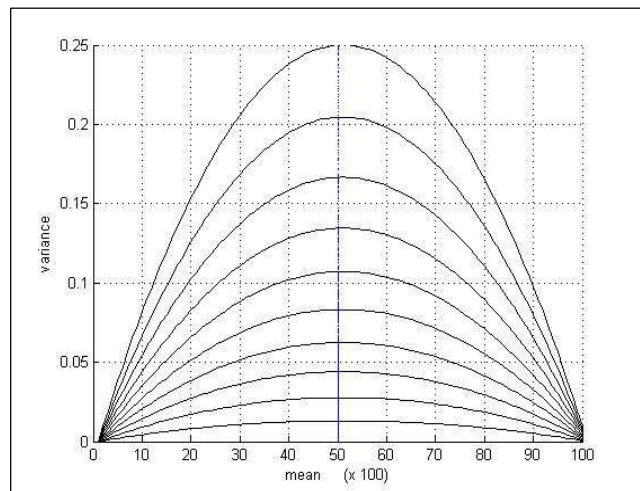
به صورت رابطه (۲۲) تبدیل خواهد شد

$$P(|X - m| \geq a) \leq \frac{\text{Var}(X)}{a^2} \quad (22)$$

با جایگذاری  $m = E[X^t] = C^*$  و  $a = q \cdot \sqrt{\text{Var}(X^t)} = q \cdot \delta$  در مورد متغیر تصادفی  $X = X^t$  در رابطه (۲۲) نتیجه ج حاصل خواهد شد ■

گزاره ۲ چند نکته را نشان می دهد:

اولا مقدار  $\alpha$  کنترل کننده نرخ همگرایی است. روابط بازگشتی مربوط به میانگین و واریانس نشان می دهند که اگر  $\alpha$  برابر ۱ یا ۱ باشد سرعت همگرایی بیشینه است. اگرچه مقدار  $\alpha$  برابر ۱ با سرعت همگرایی بالایی دارد اما منجر به واریانس بزرگتری می شود (شکل ۲). مقادیر  $\alpha$  کوچک تر گرچه همگرایی کندتری را سبب می شوند، اما جواب نزدیکتری به مقدار احتمال واقعی دارند. نکته حائز اهمیت در آن است که واریانس به صفر نمی رسد.



شکل ۲- رابطه میان واریانس با میانگین و نرخ یادگیری

دوم اینکه گرچه وجود واریانس ناخوشایند به نظر می رسد، اما یک مقدار محدود و کوچک واریانس باعث می شود با مشاهده داده های جدید، الگوریتم از گیرافتادن در نقاط بیشینه محلی نجات یابد. علاوه بر اینکه باعث می شود قابلیت انطباق با محیط های متغیر را نیز پیدا کند. سرعت این انطباق پذیری به کمک  $\alpha$  یا همان نرخ یادگیری قابل کنترل است. ضمناً گیر نیفتادن در نقاط بیشینه محلی، تضمین نشده است و به مقدار  $\alpha$  و شکل تابع توزیع بستگی دارد.

سوم اینکه نامساوی چیشف یک بازه اطمینان از مقدار احتمالی تخمین زده شده توسط آتوماتا نسبت به واریانس تخمین به دست می دهد. این ویژگی می تواند به انتخاب یک  $\alpha$  مناسب کمک کند.

تمام الگوریتم هایی که سعی می کنند به گونه ای پویا نرخ یادگیری را تنظیم کنند، از این نکته استفاده می کنند که در زمانی که همگرایی رخ می دهد، نرخ یادگیری را کاهش می دهند. به بیانی دیگر، چنانچه اختلاف میان احتمال تخمین جدید و میانگین تخمین های قبلی بزرگ باشد، نرخ یادگیری افزایش می یابد. این حالت زمانی رخ می دهد که یک بیشینه

بزرگتر از ۱ ضرب (برای افزایش) یا تقسیم (برای کاهش) می‌کنیم.

با این اصلاحات الگوریتم نهایی پیشنهادی، در حالت ساده آن به شکل ۳ خواهد بود.

ب- در حالت برداری: در این حالت، بردار احتمال انتخاب اقدام‌های آتاماتای یادگیر را در نظر گرفته و به این صورت بر خلاف حالت قبل، آتاماتای یادگیر یک نرخ یادگیری دارد که بر اساس نامساوی چیشف در حالت برداری، مقدار آن تنظیم می‌شود. هر زمان که بردار جدید یادگرفته شده توسط آتاماتای یادگیر، بیش از یک میزان مشخص از بردار میانگین مقادیر قبلی فراگرفته شده توسط آتاماتای، فاصله داشته باشد، به معنای وجود تغییرات گسترده در محیط است. برای اینکه بتوان اثر این تغییرات را منعکس کرد آتاماتای بایستی نرخ یادگیری را افزایش دهد. ملاک سنجش میزان تغییرات نابرابری چیشف برداری است. همانند حالت ساده، در اینجا نیز با هر بار افزایش مقدار نرخ یادگیری برای گریز از بیشینه محلی، مقدار میانگین با صفر مقداردهی می‌شود.

به این ترتیب الگوریتم پیشنهادی جدید در حالت برداری آن به شکل ۴ خواهد بود

#### Proposed Algorithm L<sub>R-1</sub> (2)

```

1: Parameters: Real  $q > 1$ ,  $\alpha < 1$  Learning Rate
2: Initial:  $p_j \leftarrow \frac{1}{K}$ ,  $\mu_j \leftarrow p_j$ ,  $dt \leftarrow 0$ ,  $t \leftarrow 0$  for  $j \leftarrow 1$  to  $K$ 
3: loop
4:   Draw randomly an action  $i$  according to probabilities  $p_0, \dots, p_K$ 
5:   Receive either reward or penalty
6:   if reward then
7:     for  $j \leftarrow 1$  to  $K$  do
8:       if  $j \neq i$  then
9:          $p_j \leftarrow (1 - \alpha)p_j // \text{penalt}$ 
10:      else
11:         $p_i \leftarrow p_i + \alpha(1 - p_i) // \text{reward}$ 
12:      end if
13:    end for
14:     $t \leftarrow t + 1$ ;
15:     $dt \leftarrow dt + 1$ ;
16:     $\vec{\mu} = \text{mean}(\vec{P}(i - dt - 1), \dots, \vec{P}(i))$ 
17:     $\vec{V} = \text{Cov}(\vec{\mu}, \vec{P})$ 
18:     $N = \text{trace}(\vec{V} \vec{V}^{-1})$ 
19:    if  $(\vec{P} - \vec{V})^T \vec{V}^{-1} (\vec{P} - \vec{V}) > qN$ 
20:      increment  $(\alpha)$ ;  $t \leftarrow 0$ ;  $dt \leftarrow 0$ ;
21:    else if  $(1 - \alpha)^t < \text{TSH}$ 
22:      decrement  $(\alpha)$ ;  $t \leftarrow 0$ ;
23:    end if
24:  end if
25: end loop

```

شکل ۴: الگوریتم پیشنهادی جدید مبتنی بر نامساوی چیشف به شکل برداری

در هر دو الگوریتم پیشنهادی مبتنی بر نابرابری چیشف، اگر با نرخ یادگیری  $\alpha$  پس از  $t$  بار، مقدار  $(1 - \alpha)^t$  از یک مقدار آستانه ثابت کمتر باشد، بیانگر آن است که واریانس به مقدار حدی خود- به ازای آن نرخ یادگیری خاص- نزدیک شده است. این نشان‌دهنده ثبات در محیط است و میتوان با کاهش نرخ یادگیری، اثر پذیری آتاماتای از نمونه‌های ورودی جدید را کاهش داد. دلیل انتخاب این ملاک برای کاهش نرخ یادگیری به قسمت ب اثبات گزاره ۲ برمی‌گردد که واریانس را به شکل تابعی از میانگین واقعی نشان میدهد. بر اساس آنچه که در قسمت ب گزاره ۲

در این الگوریتم برای اجتناب از سربار محاسباتی، از یک کران بالا برای مقدار واریانس استفاده می‌کنیم. برای رسیدن به این کران بالا قسمت ب گزاره ۲ و اثبات آن را در نظر بگیرید. نشان داده شد که مقدار واریانس در زمان  $t+1$  برابر است با:

$$\text{Var}[X^{t+1}] = \frac{\alpha}{2-\alpha} C^* (1 - C^*) (1 - (1 - \alpha)^{2t+2}) \quad (23)$$

رابطه (۲۳) نشان میدهد که واریانس تابعی از میانگین  $(C^*)$  و نرخ یادگیری  $(\alpha)$  است. شکل ۲ رابطه میان واریانس و میانگین را به ازای نرخ‌های مختلف یادگیری نشان می‌دهد. همانگونه که این شکل نیز نشان می‌دهد، هرچقدر میانگین به صفر یا یک نزدیک‌تر باشد، واریانس کوچک‌تر است. برعکس در میانگین برابر با ۰.۵ واریانس، بیشینه مقدار را دارد. علاوه بر این، می‌توان دید که هرچقدر نرخ یادگیری به ۱ نزدیک‌تر باشد (نمودار بالایی) واریانس مقدار بیشتری دارد و بالعکس، در مقادیر کوچک‌تر نرخ یادگیری (پایین‌ترین نمودار)، واریانس کمتری داریم. برای سهولت در محاسبه واریانس، از یک کران بالا برای آن استفاده می‌کنیم. بدین صورت که مقدار واریانس به ازای میانگین ۰.۵ را در محاسبات در نظر می‌گیریم. بدین ترتیب محاسبه واریانس در هر دور تنها تابعی از نرخ یادگیری خواهد بود. این کران بالا با مقدار  $q$  کوچک‌تر در الگوریتم جبران می‌شود.

#### Proposed Algorithm L<sub>R-1</sub> (1)

```

1: Parameters: Real  $\text{TSH} < 1$ ,  $q > 1$ ,  $\vec{\alpha}$  Learning Rate Vector,
2: Initialization:  $p_j \leftarrow \frac{1}{K}$ ,  $\mu_j \leftarrow \frac{1}{K}$ ,  $\delta_j \leftarrow 0$ ,  $t_j \leftarrow 0$  for  $j \leftarrow 1$  to  $K$ 
3: loop
4:   Draw randomly an action  $i$  according to probabilities  $p_0, \dots, p_K$ 
5:   Receive either reward or penalty
6:   if reward then
7:     for  $j \leftarrow 1$  to  $K$  do
8:       if  $j \neq i$  then
9:          $p_j \leftarrow (1 - \alpha_j)p_j // \text{penalt}$ 
10:      else
11:         $p_i \leftarrow p_i + \alpha_i(1 - p_i) // \text{reward}$ 
12:      end if
13:    end for
14:    for  $j \leftarrow 1$  to  $K$  do
15:      update  $(\mu_j)$ 
16:      update  $(\delta_j)$ 
17:      if  $|p_j - \mu_j| > q\delta_j$ 
18:        increment  $(\alpha_j)$ ;  $t_j \leftarrow 0$ ; reset  $\mu_j$ 
19:      else if  $(1 - \alpha_j)^{t_j} < \text{TSH}$ 
20:        decrement  $(\alpha_j)$ ;  $t_j \leftarrow 0$ 
21:      else  $t_j \leftarrow t_j + 1$ 
22:    end if
23:  end for
24: end if
25: end loop

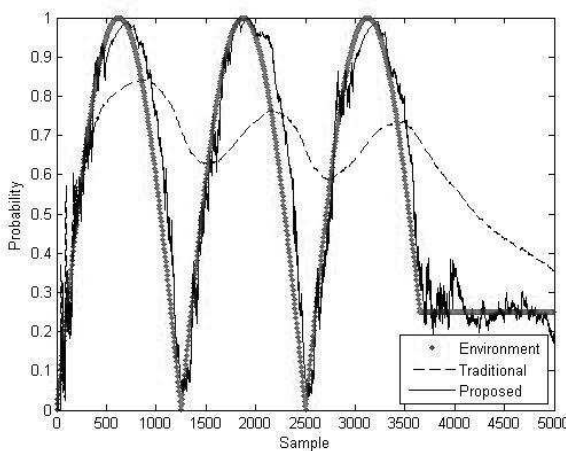
```

شکل ۳: الگوریتم شماره ۱ مبتنی بر نامساوی چیشف در حالت ساده

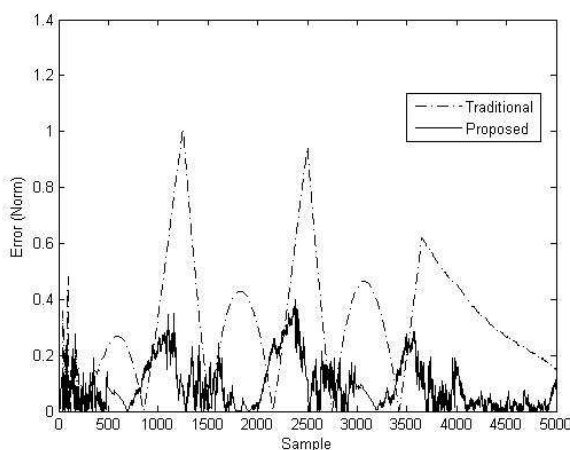
برای محاسبه میانگین نیز از یک میانگین گیری روی مقادیر مربوط به احتمالات استفاده می‌کنیم. هر زمان که نرخ یادگیری افزایش می‌یابد، میانگین‌های قبلی را در نظر نگرفته و میانگین گیری را روی مقادیر جدید آغاز می‌کنیم. این کار باعث گریز از بیشینه‌های محلی می‌شود. برای افزایش و کاهش مقدار نرخ یادگیری آنها را در یک مقدار ثابت

است. برای بررسی عملکرد روش جدید و مقایسه آن، همزمان آتاماتا را با الگوریتم ۱ و نیز الگوریتم  $L_{RI}$  با نرخ یادگیری کوچک  $\alpha = 0.002$  و تکنیک سرد کردن تدریجی با ضریب کاهش 0.001 در هر دور، آموزش داده ایم ( $\alpha^{t+1} = 0.999 * \alpha^t$ ). نتیجه مقایسه این دو در شکل ۵ آورده شده است. شکل ۶ مقدار نرم (فاصله) هر یک از دو بردار آموزش داده شده به روش معمول و روش پیشنهادی با بردار واقعی را نشان می دهد. شکل ۷ نحوه تغییر (افزایش یا کاهش) نرخ یادگیری را در الگوریتم پیشنهادی نشان می دهد.

همچنان که شکل ۷ نشان می دهد، افزایش مقدار نرخ یادگیری در موقعیت هایی رخ داده که تغییرات بزرگی در محیط ایجاد شده است و در نتیجه، بردار احتمال انتخاب اقدام های آتاماتا به بردار واقعی محیط نزدیک تر است



شکل ۵- خط تیره درشت تابع ارزیابی اقدام آتاماتا توسط محیط را نشان می دهد. خط تیره نازک نحوه تغییر بردار احتمال مربوط به اقدام آتاماتای آموزش داده شده به شیوه جدید و خطوط نقطه چین، همان بردار، آموزش داده شده به شیوه معمول را نشان می دهد. (آزمایش ۱)



شکل ۶- فاصله بردارهای آموزش داده شده به هریک از دو شیوه با بردار واقعی در آزمایش شماره ۱

دیدیم مقدار واریانس در زمان تابعی از  $(1 - \alpha)^{2t}$  است که در آن  $\alpha$  نرخ یادگیری و  $t$  تعداد دفعاتی است که نرخ یادگیری بدون تغییر برای آموزش مورد استفاده قرار گرفته است. بدین ترتیب الگوریتم می تواند از  $(1 - \alpha)^t$  به عنوان معیاری استفاده کند که میزان ثبات در محیط را نشان می دهد و در حقیقت به جای آن که منتظر بماند تا بر اثر افزایش مقدار  $t$  مقدار  $Var[X^{t+1}]$  به مقدار حدی  $\frac{\alpha}{2-\alpha} C^* (1 - C^*)$  نزدیک شود، با کاهش مقدار نرخ یادگیری  $\alpha$  این فرآیند تسریع شود

#### ۴- بررسی نتایج شبیه سازی:

برای بررسی نحوه عملکرد روش پیشنهادی جدید، یک آتاماتای تصادفی با ۲ اقدام در نظر گرفته ایم که توسط یک محیط تصادفی مورد ارزیابی قرار می گیرد. محیط تصادفی، یک محیط پویا در نظر گرفته شده است. بدین صورت که تابع ارزیابی اقدام آتاماتا توسط محیط یک مقدار ثابت فرض نشده و در طول زمان آموزش تغییر می کند. الگوریتم مورد استفاده توسط آتاماتای یادگیر الگوریتم  $L_{R-I}$  و  $L_{R-I}$  های پیشنهادی است. برای انجام شبیه سازی، توابع مختلفی بر حسب زمان به عنوان تابع پاداش یا جریمه محیط در نظر گرفته ایم. برای هر محیط، یک آتاماتا را هم زمان به دو شیوه آموزش داده ایم. روش اول، همان روش معمول مورد استفاده در آتاماتای یادگیر است به این صورت که با توجه به پویایی محیط، آتاماتای یادگیر از یک نرخ یادگیری کوچک بین 0.1 و 0.2 - آغاز کرده و با کاهش آن به اندازه 0.01 در هر دور اجرای الگوریتم (تا رسیدن به آستانه 0.001) فرآیند یادگیری را انجام داده است. روش دوم مبتنی بر الگوریتم های پیشنهادی جدید است. الگوریتم های پیشنهادی جدید نیز از یک نرخ یادگیری دلخواه (و غالباً بزرگ نزدیک به ۱) آغاز می کنند. نتایج هر دو الگوریتم پیشنهادی در مقایسه با روش معمول، در ادامه آورده شده و نتایج آن شرح داده شده است. معیار مقایسه، میزان انطباق بردار آموزش داده شده در روش جدید پیشنهادی و روش معمول با بردار ارزیابی محیط یا همان احتمال جریمه اقدام انجام شده توسط آتاماتا است (که قاعدتاً مقدار ثابتی نداشته و بر حسب زمان در تغییر است). از نرم (norm) بردار احتمال انتخاب اقدام های یاد گرفته شده و بردار احتمالی واقعی به عنوان معیاری جهت مقایسه میزان انطباق پذیری آتاماتا با محیط استفاده کرده ایم. نتایج را برای محیط های پویای مثال، در ادامه مشاهده می کنید.

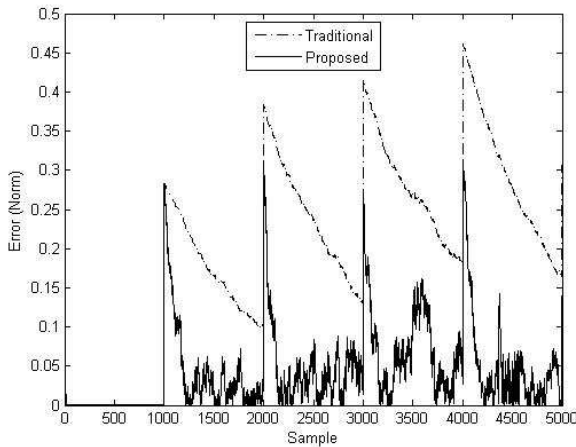
#### ۴-۱ بررسی نتایج الگوریتم پیشنهادی ۱

در این گروه از آزمایش ها از الگوریتم شماره ۱ (شکل ۳) استفاده شده است.

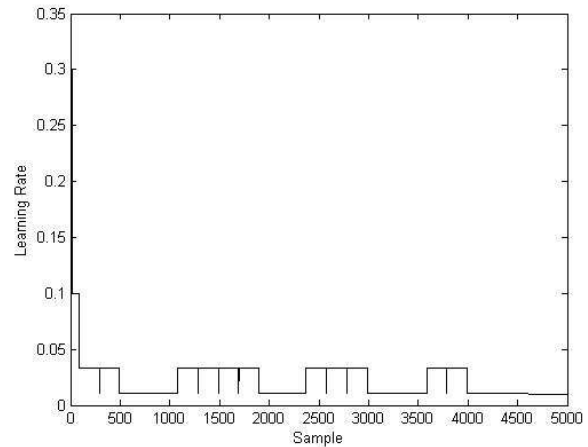
**آزمایش ۱:** در اولین نمونه، از تابع رابطه (۲۴) به عنوان تابع ارزیابی اقدام شماره ۱ آتاماتا استفاده کرده ایم:

$$f(i) = \begin{cases} \left| \sin\left(\frac{4i\pi}{n}\right) \right| & i < 0.73n \\ \left| \sin\left(\frac{4t\pi}{n}\right) \right| & t = 0.73, i \geq 0.73n \end{cases} \quad (24)$$

در رابطه (۲۴) معرف  $i$  امین نمونه و  $n$  نشان دهنده تعداد کل نمونه ها



شکل ۹-فاصله بردارهای آموزش داده شده به هریک از دو شیوه با بردار واقعی در آزمایش شماره ۲

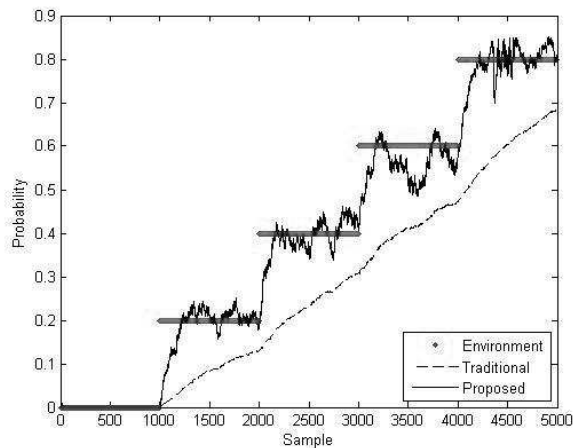


شکل ۷-تغییرات نرخ یادگیری در آزمایش ۱

## آزمایش ۲:

در این آزمایش از یک محیط تصادفی استفاده کرده ایم که تابع ارزیابی آن در طول زمان به شکل پله‌ای تغییر می‌کند. این تابع  $f(i) = 0.2 * \left\lfloor \frac{5i}{n} \right\rfloor$  است که  $i$  معرف  $i$  امین نمونه و  $n$  نشان‌دهنده تعداد کل نمونه‌ها می‌باشد.

نتایج بررسی عملکرد مقایسه ای الگوریتم پیشنهادی و الگوریتم معمول (بهینه شده)، در تنظیم نرخ یادگیری را در شکل ۸ و شکل ۹ مشاهده می‌کنید. شکل ۸ نشان می‌دهد که الگوریتم پیشنهادی انطباق بیشتری با رفتار پویای محیط داشته و شکل ۹ نیز موید میزان خطای کمتر روش جدید در مقایسه با روش معمول است.



شکل ۸-مقایسه روش معمول (خطوط نازک بریده) در مقایسه با روش پیشنهادی (خطوط تیره پیوسته) در یک محیط پویا (خطوط تیره پلکانی) در آزمایش ۲

## ۴-۲ بررسی نتایج الگوریتم پیشنهادی ۲

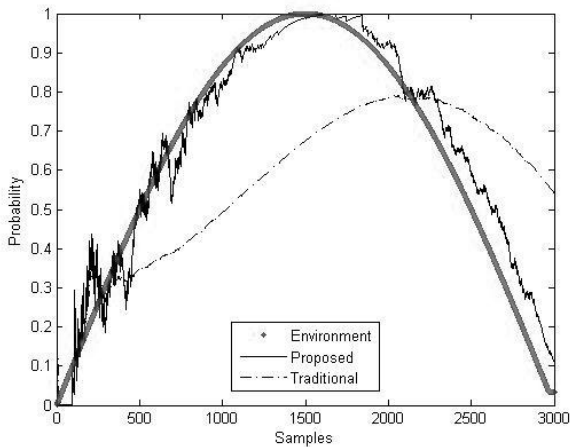
در این گروه از آزمایش‌ها، از الگوریتم ۲ (شکل ۴) استفاده کرده‌ایم. الگوریتم جدید از یک نرخ یادگیری دلخواه بزرگ آغاز کرده. در هر دور اجرا بر اساس میزان فاصله بردار جدید احتمال انتخاب به دست آمده برای اقدام‌های آتوماتا از بردار میانگین، اقدام به کاهش یا افزایش نرخ یادگیری نموده است.

در این الگوریتم، آتوماتای یادگیر یک نرخ یادگیری دارد و نامساوی چپیشف در حالت برداری آن به عنوان معیار تشخیصی برای تعیین آنکه آیا تغییر بزرگی در محیط رخ داده است یا نه مورد استفاده قرار می‌گیرد. نتیجه را در مورد نمونه‌هایی از محیط‌های پویا در ادامه بررسی کرده‌ایم.

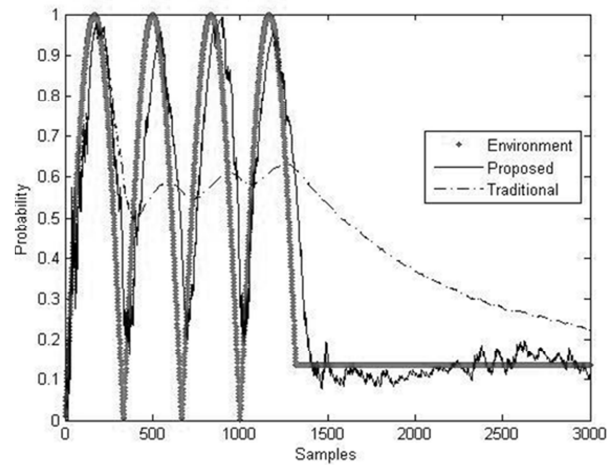
**آزمایش ۳:** محیط پویایی که برای این بررسی انتخاب شده است. تا حدود ۴۰٪ زمان آموزش از تابع پویای ارزیابی اقدام آتوماتا به صورت  $|\sin(\frac{9i\pi}{n})|$  برای ارزیابی اقدام آتوماتا استفاده می‌کند (در  $i$  امین نمونه که  $n$  تعداد کل نمونه‌ها است). اما از این مرحله به بعد، تابع به شکل خطی ثابت و بدون تغییر، اقدام آتوماتا را مورد ارزیابی قرار می‌دهد. نتایج مقایسه ای را در شکل ۱۰ و شکل ۱۱ مشاهده می‌کنید. در شکل ۱۲ نیز نحوه تغییرات نرخ یادگیری نشان داده شده است. تلاش آتوماتا برای انطباق با محیط در بخش اول نمونه‌ها که محیط از پویایی بالایی برخوردار است، از طریق افزایش نرخ یادگیری قابل مشاهده است.



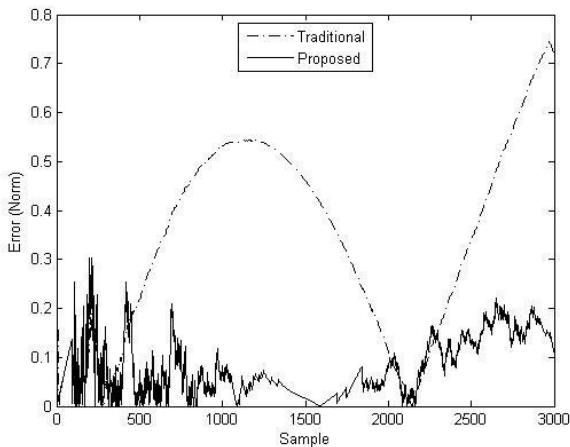
برخوردار است. نتیجه مقایسه ای عملکرد الگوریتم یادگیری  $L_R-I$  و الگوریتم جدید را در شکل ۱۳ و شکل ۱۴ مشاهده می کنید



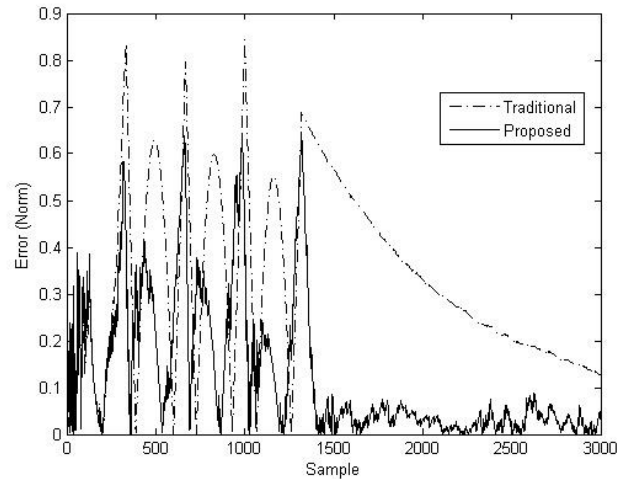
شکل ۱۳- مقایسه روش معمول (خطوط نازک بریده) در مقایسه با روش پیشنهادی (خطوط تیره پیوسته) در یک محیط پویا (شبه سینوسی) در آزمایش ۴



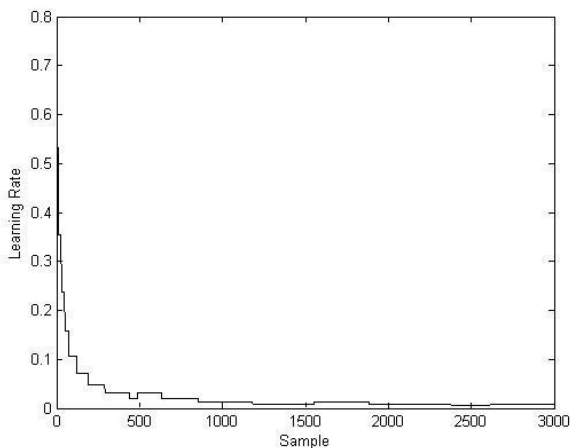
شکل ۱۰- خط تیره رنگ درشت تابع ارزیابی اقدام آتاماتا توسط محیط را نشان میدهد. خط تیره رنگ نازک نحوه تغییر بردار احتمال مربوط به اقدام آتاماتای آموزش داده شده به شیوه جدید و خطوط نقطه چین، همان بردار، آموزش داده شده به شیوه معمول را نشان میدهد



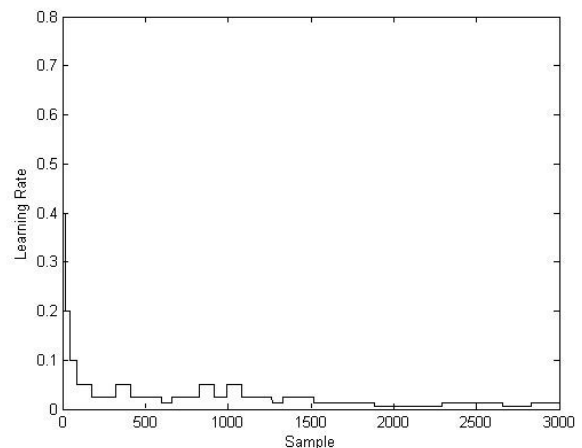
شکل ۱۴- فاصله بردارهای آموزش داده شده به هریک از دو شیوه با بردار واقعی در آزمایش شماره ۴



شکل ۱۱- فاصله بردارهای آموزش داده شده به هریک از دو شیوه با بردار واقعی در آزمایش شماره ۳



شکل ۱۵- تغییرات نرخ یادگیری در آزمایش ۴



شکل ۱۲- تغییرات نرخ یادگیری در آزمایش ۳

آزمایش ۴: برای این نمونه، محیط دیگری در نظر گرفته ایم که تابع ارزیابی محیط (احتمال پاداش به اقدام) از یک رفتار شبه سینوسی

## ۵- نتیجه گیری:

در این مقاله، به کمک شاخص های آماری حاصل از توزیع احتمال بردار انتخاب اقدام های آتاماتا، و با کمک نابرابری چیبیشف روش های جدیدی پیشنهاد شد که به کمک آن آتاماتای یادگیر ضمن تنظیم خودکار نرخ یادگیری، قادر به یادگیری در محیط های با پویایی بالا است. این الگوریتم های جدید قادر به تنظیم اتوماتیک نرخ یادگیری بوده و میتوانند بر حسب میزان تغییرات در محیط و به منظور انطباق با آن، اقدام به کاهش یا افزایش نرخ یادگیری نماید. استفاده از این الگوریتم جدید در محیط های تصادفی پویا که پاسخ محیط به اقدام انجام شده توسط آتاماتا غیر ابت و تابعی از زمان است مورد بررسی قرار گرفت و نشان داده شد که الگوریتم های جدید مبتنی بر نامساوی چیبیشف از عملکرد بهتری نسبت به روش های یادگیری معمول برخوردارند.

## ۶- منابع

- [1] H.Beigy, M. R. Meybodi, and M. B. Menhaj, "Adaptation of Learning Rate in Back Propagation Algorithm using Fixed Structure Learning Automata.," in *6th Iranian Conference on Electrical Engineering*, 1998.
- [2] S. Shah-Hosseini, "Automatic adjustment of learning rates of the self-organizing feature map," *Scientia Iranica*, vol. 8, 2001.
- [3] carlo H. sequin. chedsada chinrungrueng, "Optimal Adaptive K-Means Algorithm with Dynamic Adjustment of Learnig Rate." 1991.
- [4] J. Akbari Torkestani and M. R. Meybodi, "Learning automata-based algorithms for solving stochastic minimum spanning tree problem," *Applied Soft Computing*, vol. 11, no. 6, pp. 4064-4077, Sep. 2011.
- [5] J. Akbari Torkestani and M. R. Meybodi, "Learning Automata-Based Algorithms for Finding Minimum Weakly Connected Dominating Set in Stochastic Graphs," *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 18, no. 06, pp. 721-758, Dec. 2010.
- [6] H.Beigy and M.R.Meybodi, "Utilizing distributed learning automata to solve stochastic shortest path problems," *International Journal of Uncertainty, ...*, vol. 14, no. 5, pp. 591-615, 2006.
- [7] M. R. MollakhaliliMeybodi and M.R.Meybodi, "A New Distributed Learning Automata Based Algorithm for Solving Stochastic Shortest Path," in *6th Conference on Intelligent Systems*, 2004.
- [8] H. Beigy and M. R. Meybodi, "Solving stochastic shortest path problem using Distributed Learning Automata," in *6th Annual CSI Computer Conference (CSICC 2001)*.
- [9] H. Beigy and M. . Meybodi, "Intelligent Channel Assignment in Cellular Networks: A Learning Automata Approach," Amirkabir University of Technology, 2004.
- [10] M. L. Thathachar and P. S. Sastry, "Varieties of learning automata: an overview.," *IEEE transactions on systems, man, and cybernetics. Part B, Cybernetics : a publication of the IEEE Systems, Man, and Cybernetics Society*, vol. 32, no. 6, pp. 711-22, Jan. 2002.
- [11] K. S. Narendra and M. A. L. Thatacher, *Learning Automata*. Prentice-Hall, 1989.
- [12] M. L. Tsetlin, "on the behaviour of finite automata in random media," *Automata., Telemekh.*, vol. 22, pp. 1345-1354, 1961.
- [13] B. J. Oomen and H. MAsum, "Switching Models for Nonstationary Random Environments," *IEEE transactions on systems, man, and cybernetics.*, vol. 25, no. 9, pp. 1334-1347.
- [14] K. S. Narendra and M. A. L. Thathachar, "On the Behavior of a Learning Automaton in a Changing Environment with Application to Telephone Traffic Routing," *Systems, Man and Cybernetics, IEEE Transactions on*, vol. 10, no. 5, pp. 262-269.
- [15] S. M. Ross, *Introduction to Probability and Statistics dor Engineers and Scientists*, Third. Elsevier Academic Press, 2004.
- [16] A. Papoulis, *Probability, Random Variables, and Stochastic Processes*, 3rd ed. New York, USA: McGrawHill, 1991.