

مدل یادگیری Q سلولی و کاربردهای آن

رضا رستگار محمد رضا میبیدی

آزمایشگاه محاسبات نرم

دانشکده مهندسی کامپیوتر و فناوری اطلاعات

دانشگاه صنعتی امیرکبیر

تهران ایران

(Rastegar,meybodi)@ce.aut.ac.ir

محیطی بلادرنگ انجام میگیرد، می توان آنرا همزمان با فعالیت محیط انجام داد که در این صورت با تمام رخدادهای پیش بینی نشده بصورت یک تجربه جدید برخورد می شود و می توان از آنها برای بهبود کیفیت یادگیری استفاده کرد. مزیت عمده یادگیری تقویتی نسبت به سایر روشهای یادگیری عدم نیاز به اطلاعات بجز سیگنال تقویتی از محیط میباشد.

یکی از مدل هایی که در شبیه سازی و یا مدل کردن سیستمها مورد استفاده قرار می گیرد، اتوماتای سلولی است که اجزای آن به صورت مکانی توزیع شده اند و اطلاعات از طریق قوانین محلی حاکم بر سیستم به صورت جزئی بین اجزا رد و بدل می شوند [3]. در اتوماتای سلولی، فضا بصورت یک شبکه ای از سلولها تعریف می گردد، زمان بصورت گسسته پیش می رود و قوانین آن بصورت سرتاسری است که از طریق آن در هر مرحله هر سلول، وضعیت جدید خود را با در نظر گرفتن وضعیت همسایه های خود بدست می آورد. قانون اتوماتای سلولی، نحوه تاثیر پذیرفتن یک سلول از سلولهای همسایه خود را مشخص می کند. یک سلول، همسایه سلول دیگر گفته میشود اگر بتواند آن سلول را در یک مرحله و براساس قانون حاکم تحت تاثیر قرار دهد. ویژگی های اساسی اتوماتای سلولی، فضای گسسته، زمان گسسته، محدودیت تعداد وضعیتهای ممکن هر سلول، یکسان بودن تمام سلولها، قطعی بودن قوانین، وابستگی قانون در هر سلول به مقادیر سلولهای اطراف آن و وابستگی قانون به مقادیر تعداد محدودی از مراحل قبل همسایه ها و خود سلول می باشند. در اتوماتای سلولی همگام^۳ عمل بروز در آوردن سلولها به صورت همگام و در اتوماتای سلولی نا همگام^۴ عمل بروز در آوردن سلولها به بصورت نا همگام انجام میگیرد.

یکی از مشکلات اتوماتای سلولی تعیین فرم قطعی قوانین مورد نیاز برای یک کاربرد خاص است. اتوماتای سلولی برای مدل کردن سیستمهایی مناسب است که قطعیت در تغییر حالات سیستم وجود داشته باشد در

چکیده: اتوماتای سلولی برای مدل کردن سیستمهایی مناسب است که قطعیت در تغییر حالات سیستم وجود داشته باشد. در حالیکه اغلب سیستمهای واقعی پیچیده بوده و ویژگی نویزی بودن و عدم قطعیت و احتمالی بودن در آنها دیده می شود و به همین دلیل برای مدل کردن چنین سیستمهایی استفاده از اتوماتای سلولی با قوانین قطعی منطقی به نظر نمی رسد. در این مقاله از ترکیب یادگیری Q با اتوماتای سلولی مدل جدیدی به نام یادگیری Q سلولی^۱ معرفی می گردد. این مدل جدید با استفاده از قابلیت های یادگیری Q مشکل نبود عدم قطعیت در تغییر حالات در اتوماتای سلولی را تا حدودی مرتفع می سازد. کاربرد این مدل ترکیبی در مساله تخصیص کانال در شبکه های سلولی مخابراتی مورد بررسی قرار میگیرد.

کلمات کلیدی: اتوماتای سلولی، یادگیری Q، یادگیری Q سلولی، شبکه سیار سلولی، تخصیص کانال، یادگیری

۱ - مقدمه

یادگیری می تواند به عنوان یک راه کار برای ایجاد تطبیق پذیری در اکثر سیستمهایی که دارای فرایندهای تصمیم گیری بر اساس عدم قطعیت و اطلاعات ناقص می باشند مورد استفاده گیرد. با استفاده از یادگیری در جایگاههای مناسب در سیستم، هرگز سازنده سیستم می تواند حتی با دریافت اطلاعات ناقص و غیر قطعی، به صورت تدریجی و بر اساس معیارهای تعریف شده در سیستم به استراتژی بهینه کنترلی مورد نیاز خود دست یابد. در یادگیری تقویتی که یکی از انواع مهم مدل های یادگیری میباشد، یک عامل یادگیرنده در طی فرایند یادگیری با تعاملات^۲ مکرر با محیط، به یک سیاست کنترل بهینه می رسد. کارایی این تعاملات با محیط بوسیله بیشینه (کمینه) بودن پاداشی (جریمه ای) که از محیط گرفته می شود، ارزیابی می گردد. از آنجاییکه این روش یادگیری در

³ Synchronous Cellular Automata

⁴ Asynchronous Cellular Automata

¹ Cellular Q-Learning

² Interaction

صورتیکه اغلب سیستمها نویزی و دارای عدم قطعیت می باشند و وضع قوانین برای آنها به صورت قطعی، منطقی به نظر نمی رسد. روشهای متفاوتی برای حل این مشکل پیشنهاد شده است. یکی از این روشها احتمالاتی کردن قوانین می باشد اما مشکل این رهیافت، محاسبه این احتمالات برای سیستمهای ناشناخته می باشد. با معرفی اتوماتای یادگیر سلولی گامی در حل این مساله برداشته شده است. [1][19][20][21]

در این مقاله مدل دیگری با نام "یادگیرنده Q سلولی" برای مشکل فوقالذکر پیشنهاد می شود. یادگیرنده Q سلولی یک اتوماتای سلولی است که هر سلول آن به یک یا چند یادگیرنده Q که یکی از انواع یادگیری تقویتی میباشد مجهز است که وضعیت سلول را مشخص میکند. مانند اتوماتای سلولی، قانون محلی در محیط حاکم است و این قانون تعیین می کند که آیا عمل انتخاب شده توسط یک یادگیرنده Q در سلول باید پاداش داده شود ویا اینکه جریمه شود. دادن پاداش ویا جریمه منجر بروز درآوردن ساختار یادگیرنده Q سلولی بمنظور نیل به یک هدف مشخص می گردد. برای نشان دادن کاربرد یادگیرنده Q سلولی، از آن برای حل مسئله تخصیص کانال در شبکه های سلولی مخابراتی استفاده میشود.

ادامه مقاله بدین صورت سازماندهی شده است. در ابتدا در بخش ۲ به بررسی اتوماتای سلولی می پردازیم. در بخش ۳ یادگیری Q شرح داده میشود. در بخش ۴ مدل یادگیرنده Q سلولی ارائه می شود. در بخش ۵، کاربرد یادگیرنده Q سلولی در حل مسئله تخصیص کانال در شبکه های سلولی مخابراتی مورد بررسی قرار می گیرد. بخش پایانی نتیجه گیری میباشد.

۲- اتوماتای سلولی

اتوماتای سلولی شبکه ای سلولی است که هر سلول می تواند k حالت (وضعیت) داشته باشد. در هر سلول یک اتوماتا با حالات محدود^۵ قرار دارد. در حالت یک بعدی، هر سلول دو همسایه نزدیک به خود دارد. در این حالت، وضعیت سلول i در زمان t+1 یعنی $a_i^{(t+1)}$ مطابق فرمول زیر بدست می آید.

$$a_i^{(t+1)} = \phi(a_{i-1}^{(t)}, a_i^{(t)}, a_{i+1}^{(t)}) \quad (1)$$

تابع ϕ اتون اتوماتای سلولی نامیده میشود. همسایگی در اتوماتای سلولی یک بعدی را می توان بگونه ای بسط داد که از دو همسایه بیشتر را نیز شامل شود. یعنی می توان شعاع I را برای همسایگی در نظر گرفت. البته معمولاً نزدیک ترین همسایه ها را در نظر می گیریم. همچنین سلولها در

اتوماتای سلولی می توانند در شبکه ای با هر ابعادی قرار گیرند که متناسب با بعد، تعاریف مربوط به همسایگی و قانون تغییر می یابند. متداولترین اتوماتای سلولی، اتوماتای سلولی دو بعدی است. چند نوع همسایگی مهم در این نوع اتوماتای سلولی همسایگی مور^۶، ون نیومن^۷، کول^۸ و اسمیت^۹ می باشند [۳].

۳- یادگیری Q

یک سیستم با مجموعه حالات محدود و قابل شمارش، S، را در نظر بگیرید. یک کنترل کننده در هر حالت $s \in S$ ، یک عمل a را از میان مجموعه اعمال مجاز، A(s)، انتخاب و انجام می دهد. پس از آن سیستم با احتمال $p(s, a, s')$ ، از حالت s به s' می رود و پاداش $r(s, a, s')$ را به کنترل کننده می دهد. هدف کنترل کننده ماکزیمم کردن تابع زیر در تمامی حالات است:

$$J(s) = E \left\{ \int_0^\infty e^{-\beta t} r(t) dt \mid s \right\} \quad (3)$$

که $E\{ \cdot \mid s \}$ امید ریاضی پاداش در یافت شده برای تمام مسیرهای ممکن آغاز شوند از s، $r(t)$ ، نرخ عایدی کل در زمان t، و β نرخ تخفیف^۹ می باشند. در حالتی که زمان پیوسته نباشد می توانیم فرمول (۳) را به صورت زیر بنویسیم:

$$J(s) = E \left\{ \sum_{t=0}^T \gamma^t r(t) \mid s \right\} \quad (4)$$

به طوری که γ نرخ تخفیف گسسته است. مقدار بهینه تابع J، J^* ، طبقه معادله بلمن^{۱۰} به صورت زیر تعریف می شود:

$$J^*(s) = \max_{a \in A(s)} [E_{\Delta, s'} \{ \alpha(s, a, s') + \gamma(\Delta t) J^*(s') \}] \quad (5)$$

که Δt زمان تصادفی تا رویداد بعدی، و $\gamma(\Delta t)$ نرخ تخفیف موثر برای حالت بعدی s' ، $\gamma(\Delta t) = e^{-\beta \Delta t}$ می باشند.

تخمین J^* می تواند از طریق یادگیری Q انجام گیرد. معادله بلمن را می توان با استفاده از مقادیر Q به صورت زیر بازنویسی کرد:

$$J^*(s) = \max_{a \in A(s)} Q^*(s, a) \quad (6)$$

در آغاز مقادیر به صورت تصادفی مقدار دهی می شوند و در ادامه با هر انتقال حالت در سیستم مقادیر Q بروز رسانی می شوند. و در حالت S عملی که بیشترین مقدار $Q(s, \cdot)$ را داشته باشد انتخاب می گردد. اگر عامل یادگیرنده با انتخاب عمل a باعث انتقال سیستم به s' شود و

⁶ Moore

⁷ Cole

⁸ Smith

⁹ Discount Factor

¹⁰ Bellman Equation

⁵ Finite State Automata

پاداش $r(s, a, s')$ را دریافت کند، مقدار $Q(s, a)$ به صورت زیر تنظیم می شود:

$$Q(s, a) = (1 - \alpha)Q(s, a) + \alpha(r(s, a, s') + \gamma \max_{b \in A(s')} Q^*(s', b)) \quad (7)$$

به طوریکه $0 \leq \alpha \leq 1$ نرخ یادگیری می باشد. برای اطلاعات بیشتر در باره یادگیری Q میتوان به [۵ و ۱۶] مراجعه کرد.

۴- یادگیرنده Q سلولی

"یادگیرنده Q سلولی" که از این پس آن را با نام CQL می شناسیم مدلی است که از ترکیب اتوماتای سلولی و یادگیری Q حاصل میشود. این مدل برای سیستم های طراحی شده که اجزای آنها از طریق تعامل با یکدیگر از تجربیات گذشته همدیگر اطلاع پیدا میکنند و از این طریق میتوانند رفتار خود را اصلاح کنند. یک CQL را می توان به صورت یک شش تایی $\langle E, A, N, \Phi, Q, C \rangle$ تعریف کرد که $E = \{e_1, \dots, e_n\}$ مجموعه مکانهای تعریف شده در اتوماتای سلولی هستند که می توانند در بشکلهای مختلفی مانند خطی، دو بعدی و سه بعدی درکنار هم قرار گیرند. $A = \{a_1, \dots, a_k\}$ مجموعه مقادیر مجازی است که یک سلول میتواند اختیار کنند. $A^t(e_i)$ نشان دهنده مقدار سلول e_i در زمان t است. Φ قوانین حاکم بر جامعه سلولی است که پاداشها و جریمه ها براساس آن تعیین میشود. $N(e_i)$ مجموعه همسایه های سلول e_i را تعریف می کند که این مجموع دارای این ویژگیهاست:

$$e_i \notin N(e_i) \quad \forall e_i \in E \quad (8)$$

$$e_i \in N(e_j) \leftrightarrow e_j \in N(e_i) \quad \forall e_i, e_j \in E \quad (9)$$

هر قانون $\phi \in \Phi$ را می توان با توجه به مفاهیم قوانین عمومی و جمعی و یا وزن دار به یکی از فرمهای زیر تعریف کرد:

قوانین عمومی:

$$\langle a_1, \dots, a_h \rangle \rightarrow r \quad (10)$$

برای مثال قانون $-1 \rightarrow \langle 1, 0, 0, 0, 0 \rangle$ با فرض همسایگی ون نیومن به معنای آن است که در صورت ۱ بودن سلول مرکزی و صفر بودن بقیه سلولهای همسایه، این خانواده جریمه شود. که این جریمه بر اساس تابع تخصیص پاداش C بین آنها تقسیم می شود.

قوانین جمعی:

$$\langle a_1 \times n_1, \dots, a_k \times n_k \rangle \rightarrow r \quad (11)$$

$$\sum_{i=1}^k n_i = |N| + 1 \quad (12)$$

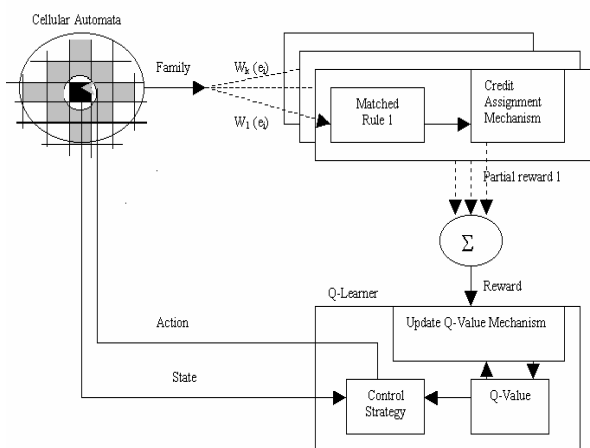
قوانین وزن دار:

$$\langle a_1 \times n_1, \dots, a_k \times n_k \rangle \rightarrow \sum_{i=1}^k n_i \times w_i \quad (13)$$

که $a_i \in A$ ، r مقادیر مجاز برای پاداش و تنبیه و w_i وزنی است که به هر وضعیت داده می شود. اگر $W(e_i) = N(e_i) \cup \{e_i\}$ ، در این صورت $\phi = \phi(W(e_j))$ و $C = C(W(e_i))$ وظیفه تقسیم پاداش بدست آمده از ϕ مابین اعضای $W(e_i)$ را دارد. بر طبق تابع تخصیص پاداش C به هر سلول براساس کارایی آن پاداش داده می شود. حال که با این مدل آشنا شدیم به شرح عملکرد آن می پردازیم. در نتیجه در یافتی یک درخواست توسط سلول e_i ، این سلول فعال می شود (در مدل آسنکرون تمام سلولها به طور همزمان فعال نمی شوند) پس از فعال شدن، سلول، e_i ، از میان از مجموعه A یک مقدار انتخاب می کند و سپس اعمالی بدین شرح را انجام می دهد. تابع ϕ برای سلول ارزیابی می شود و سپس طبق تابع C، مقدار پاداش بدست آمده برای مجموعه $W(e_i)$ را بین اعضای آن تقسیم می کنیم. با چنین مکانیسمی در صورت سنکرون بودن اتوماتای سلولی، هر سلول e_i در گام t ام $|W(e_i)|$ خرده پاداش دریافت می کند که با جمع این خرده پاداشها، پاداش کل سلول طبق (۱۴) بدست می آید که به یادگیرنده Q داده می شود و بر اساس آن مقادیر Q بهنگام می شوند.

$$sum = \sum_{\substack{\forall W(e_j) \\ \exists e_i \in W(e_j)}} C(W(e_j); e_i) \quad (14)$$

اگر اتوماتای سلولی به صورت آسنکرون عمل کند پس از دریافت سیگنال تقویتی، تنها سلول فعال شده و همسایه های آن بر طبق تابع تخصیص پاداش C پاداش را دریافت می کنند. سپس مقادیر پاداشها محاسبه شدند و به یادگیرنده های هر سلول داده می شود. نهایتا سلول عمل بروز رسانی را انجام می دهد. مراحل فوق تا رسیدن سیستم به تعادل که در انجا مقدار هر سلول بهینه میباشد ادامه پیدا میکند.



شکل ۱: معماری یک سلول در مدل CQL

همانطوریکه که گفته شد در مدل ارائه شده، پس از مشخص پاداش مربوط به یک $W(e_i)$ ، سیگنال پاداش بین سلولهای آن تقسیم می‌شود که این روش تقسیم توسط، $C(W(e_i))$ تعریف می‌شود و مقدار بازگشتی آن سهم پاداش هر سلول $e_j \in W(e_i)$ از مقدار $\phi(W(e_i))$ را نشان می‌دهد. روشهای متفاوتی برای پیاده‌سازی تابع تخصیص پاداش C وجود دارد که از مهمترین آنها می‌توان به تقسیم مساوی، تقسیم تصادفی، تقسیم براساس ویژگیهای سلول از جمله خبرگی [۲] سلول یا موقعیت جغرافیایی اشاره کرد.

۵- حل مساله تخصیص کانال در شبکه های سیار سلولی توسط یادگیرنده Q سلولی

مسئله تخصیص کانال در شبکه های سیار سلولی [۶] یک مسئله NP complete میباشد و بهمین دلیل الگوریتمهای تقریبی متعددی برای حل این مساله از جمله الگوریتمهای مبتنی بر الگوریتمهای ژنتیکی [۸]، تابکاری فلزات [۸]، جستجوی TABU [۸]، شبکه‌های عصبی [۹]، برنامه سازی پویای نرونی [۴] و یادگیری Q [۱۰] گزارش شده اند.

در تمامی استراتژی‌های تخصیص کانال به صورت پویا، دانش و تجربه بدست آمده در طول کار سیستم به دست فراموشی سپرده می‌شود. اگر چه استراتژی‌های مبتنی بر شبکه‌های عصبی از آموزش بهره می‌برند ولی در همه آنها داشتن یک ناظر خوب (یک استراتژی شناخته شده تخصیص پویای کانال) ضروری است. به نظر می‌رسد که استفاده از روشهای یادگیری بدون نیاز به داشتن استراتژی معین و شناخته شده بتواند گره‌گشای این مشکل باشد [۱۰]. در [۴ و ۱۰] دو استراتژی تخصیص کانال مبتنی بر الگوریتم‌های یادگیری تقویتی که در آنها نیازی به داشتن استراتژی معین و شناخته شده نیست ارائه گردیده است. در [۴] از یک معماری خطی تخمینی به همراه $TD(0)$ استفاده شده است. در این استراتژی هدف ماکزیمم کردن $E[\int_0^\infty e^{-\beta t} n(t) dt]$ می‌باشد که $n(t)$ تعداد مکالمات در جریان در زمان t و β نرخ تخفیف است. قابلیت توزیع شدگی و کارایی بالا از ویژگی‌های این الگوریتم است. در الگوریتم ارائه شده در [۱۱] از یادگیری Q برای تخصیص کانال به صورت تمرکز یافته^{۱۴} استفاده شده است. آنچه الگوریتم پیشنهادی در مقاله را از

الگوریتم متمایز متمایز می‌سازد، مستقل شدن فرایند یادگیری در هرکدام از سلولهای شبکه می‌باشد که سبب می‌شود هر سلول بتواند بر اساس موقعیت جغرافیایی خود (به دلیل شرایط مرزی پوچ^{۱۵}، سیستم متقارن نمی‌باشد). استراتژی مناسب خود را یاد بگیرد. همچنین الگوریتم پیشنهادی کاملاً توزیع شده میباشد که این خود از حجم پیغامهای کنترلی بر روی شبکه میکاهد. الگوریتم پیشنهادی با دو الگوریتم فوق‌الذکر مقایسه خواهد شد.

در ادامه این قسمت به توصیف یک الگوریتم پویای توزیع شده تخصیص کانال مبتنی بر یادگیرنده Q سلولی اسنکرون می‌پردازیم. در الگوریتم پیشنهادی هر سلول در شبکه سلولی به یک سلول در یادگیرنده Q سلولی نگاشت می‌شود. زمانیکه یک سلول فعال شود (زمانیکه یک درخواست مکالمه به سلول میرسد)، یادگیرنده Q این سلول یک عمل انتخاب و سپس قانون حاکم بر یادگیرنده Q سلولی ارزیابی میشود. براساس نتیجه ارزیابی قانون، فرایند یادگیری انجام می‌شود. در ادامه این بخش پارامترهای یادگیرنده Q سلولی تعریف و سپس الگوریتم تخصیص کانال پویای مبتنی بر یادگیرنده Q سلولی ارائه میشود.

مدل همسایگی: مدل همسایگی، مدل همسایگی مور توسعه یافته^{۱۶} میباشد. مجموعه همسایه های هر سلول به دو دسته تقسیم می‌شوند: همسایه‌های ثابت و همسایه‌های متغیر. با فرض اینکه فاصله استفاده مجدد هم کانال R باشد، R حلقه حول سلول، سلولهای غیر همکانال یا همسایه‌های ثابت سلول می‌باشند که با سلول مرکزی تعامل محلی انجام می‌دهند. هر سلول در حالت کلی دارای $6(R+1)$ همسایه ثابت است که از این تعداد، حداکثر ۶ سلول همسایه مجاور می‌باشند و بقیه همسایگان غیر مجاور این سلول را تشکیل می‌دهند. هر سلول دارای یک مجموعه سلول هم کانال نیز می‌باشد که حالت آنها نیز در حالت فعلی سلول تاثیر دارد. مجموعه سلولهای همکانال یک سلول، سلولهایی هستند که در فاصله $R+1$ از سلول واقع هستند و می‌توانند از یک مجموعه کانال استفاده کنند. تعداد همسایه‌های متغیر یک سلول حداکثر ۶ تا می‌باشد. نکته‌ای که باید به آن توجه داشته باشیم آن است که با توجه به ویژگیهای فیزیکی شبکه مخابراتی سلولی اتوماتای سلولی استفاده شده در این الگوریتم دارای شرایط مرزی پوچ می‌باشد.

حالات سلول: در مدل استفاده شده هر سلول دارای متغیرهای حالت U_i (متفاوتی است که تعداد آنها برابر تعداد کانالهای شبکه مخابراتی موبایل است و هر متغیر میتواند یکی از مقادیر *free* و *used* را اختیار

کند.

¹¹ Simulated Annealing

¹² TABU Search

¹³ Neuro-Dynamic Programming

¹⁴ Centralized

¹⁵ Null Boundary Conditions

¹⁶ Extended Moore

$$e^{Q(s,k)/T} / \sum_{l=1}^M e^{Q(s,l)/T} \quad (20)$$

که T پارامتر دما است که با گذشت زمان کاهش می یابد. مقدار مینیمم T برابر ۱ می باشد.

- **نرخ یادگیری:** در الگوریتم پیاده سازی شده از نرخ یادگیری ثابت^{۱۹} و متغیر کاهش یابنده^{۲۰} استفاده شده است. نرخ یادگیری کاهش یابنده به صورت زیر تعریف می شود:

$$\alpha(s, a) = 1 / \text{visit}(s, a) \quad (21)$$

که $\text{visit}(s, a)$ تعداد دفعاتی است که عامل یادگیرنده، عمل a در حالت s را انتخاب می کند.

راههای متفاوتی برای ذخیره سازی مقادیر Q (شبکه عصبی، درخت تصمیم گیری^{۲۱} و جدول) وجود دارد. در این تحقیق از روش ذخیره سازی به شکل جدول استفاده شده است. از آنجا که حافظه مورد نیاز برای جدول، نمایی نسبت به تعداد کانالها در شبکه افزایش میابد و با توجه به اینکه در عمل تعداد کانالها زیاد است، حافظه مورد نیاز برای ذخیره سازی مقادیر Q بسیار بالا خواهد بود. برای حل این مشکل فضای حالت را به چند مجموعه افراز کرده و برای هر مجموعه یک سطر در جدول Q ذخیره شده است.

حال به شرح الگوریتم تخصیص کانال پیشنهادی میپردازیم. هر گاه یک درخواست مکالمه به سلول i میرسد مراحل زیر انجام میگیرد.

- ۱- برای تمام سلولهای در همسایگی ثابت، پیام `give_used_channels` فرستاده می شود. سلولهای همسایه پس از دریافت این پیام با ارسال پیام `get_used_channels` که شامل لیست کانالهای اشغال شده خود می باشد به سلول i پاسخ می دهند.
- ۲- سلول i براساس پیامهای دریافتی، حالت s را محاسبه می کند و از یادگیرنده Q خود بهترین کانال را درخواست می کند. یادگیرنده نیز براساس استراتژی کنترل Soft-Max، کانال k را انتخاب می کند. در صورتی که کانالی موجود نباشد، درخواست جدید مسدود می شود.
- ۳- پس از انتخاب کانال و لتساب آن، سلول i با ارسال پیام `lock_channel(k)` به تمام سلولهای در همسایگی ثابت خود از آنها می خواهد کانال k را در خود قفل کنند.
- ۴- پیام `give_channel_status(k)` به تمام همسایگان سلول متغیر ارسال می شود. هر همسایه پس از دریافت پیام، وضعیت کانال k را در پیام `get_channel_status(k)` به سلول i ارسال می کند.

قانون: قانون از نوع قانون وزن دار می باشد و تابعی از همسایه های متغیر و ثابت سلول می باشد. پاداش $r(s, a, s')$ پاسخ آنی سیستم در برابر تخصیص کانال a در حالت s می باشد که به صورت زیر تعریف میشود:

$$r(s, a, s') = n_1(k)w_1 + n_2(k)w_2 \quad (15)$$

که $n_1(k)$ تعداد سلولهای همکانال سلول iام می باشد که در آنها از کانال K استفاده شده است و $n_2(k)$ تعداد سلولهای همکانال می باشد که در لایه سوم همسایگی (با شرط این که فاصله استفاده مجدد برابر دو باشد) قرار دارند و کانال K در آنها موجود است. ضرایب $w_1 = -1, w_2 = +1$ مقادیر ثابتی هستند.

تابع تخصیص پاداش: تابع تخصیص پاداش C به صورت زیر تعریف می شود:

$$C(W(e_i), e_j) = \begin{cases} 0 & \text{otherwise} \\ \phi(W(e_i)) = r(s, a, s') & e_j \text{ is central cell of } W(e_i) \end{cases} \quad (16)$$

یادگیرنده Q: در زیر پارامترهای یادگیرنده Q هر سلول e_i تعریف میشود.

- **حالت:** با فرض وجود N سلول و M کانال، حالت S به صورت زیر تعریف می شود (این تعریف مربوط به ساختار درونی یادگیرنده بوده و با تعریف حالات اتوماتای سلولی متفاوت است):

$$s = \sum_{i=1}^M H(i) 2^{i-1} \quad (17)$$

به طوری که:

$$H(i) = \begin{cases} 1 & \text{Channel } i \text{ is not available} \\ 0 & \text{Otherwise} \end{cases} \quad (18)$$

- **عمل:** تخصیص یک کانال a از میان کانالهای $A(i)$ به درخواست رسیده به سلول i ام:

$$a = k, \quad k \in A(i) \quad (19)$$

که $A(i)$ ، مجموعه کانالهای موجود^{۱۷} در سلول i است.

- **حالت بعدی:** با توجه به تعریف حالت که قبلاً به آن اشاره شد، حالت بعدی از طریق واکنش الگوریتم به رویداد رسیده و حالت فعلی قابل محاسبه است.
- **استراتژی کنترل:** در اینجا از استراتژی کنترلی انتخاب Soft-Max^{۱۸} استفاده شده است که در آن احتمال انتخاب کانال k در حالت S بصورت زیر تعریف می شود:

¹⁹ Constant Learning Rate

²⁰ Constant Learning Rate

²¹ Decision Tree

¹⁷ Available

¹⁸ Soft-Max Selection Strategy Control

- [8] Chen, J., Seah, D. and Xu, W., "Channel Allocation for Cellular Networks Using Heuristic Methods", unpublished report, 1999.
- [9] Funabiki, N., "A Neural Network Parallel Algorithm for Channel Assignment Problems in Cellular Radio Networks", IEEE Transactions on Vehicular Technology, Vol 41, No. 4, 1992.
- [10] Hykin, S. and Nie, J., "A Dynamic Channel Assignment Policy through Q-learning", IEEE Transactions on Neural Networks, Vol. 10, No. 6, 1999.
- [11] Rastegar, R. and Meybodi, M. R., "CQL and Its Applications", Technical Report, Computer Engineering Department, Amirkair University, 2005.
- [12] Krumke, S., Marathe, M. and Ravi, S., "Approximation Algorithms Assignment in Radio Networks", Dallas, International Workshop on Discrete Algorithms and Methods for Mobile Computing and Communications, 1996.
- [13] Lawrence, K. and Peter, Yum T., "Phantom Cell Analysis of Dynamic Channel Assignment in Cellular Mobile Systems", IEEE Transactions on Vehicular Technology, Vol. 47, No. 1, 1998.
- [14] Lawrence Young, K. and Yum, K., "Compact Pattern Based Dynamic Channel Assignment for Cellular Mobile Systems", IEEE Transactions on Vehicular Technology, Vol. 43, No. 4, 1994.
- [15] Smith K., "A Genetic Algorithm for the Channel Assignment Problem", unpublished report, 1998.
- [16] Sutton, R., and Barto, A., Reinforcement Learning: An Introduction, MIT Press, 1998.
- [17] Tong, H. and Brown, T. X., "Reinforcement Learning for Call Admission Control and Routing under Quality of Service Constraints on Multimedia Networks", Accepted in Machine Learning Journal, 2000.
- [18] Meybodi, M. R. and Kharazmi, M. R., "Cellular Learning Automata and Its Application to Image Processing", Journal of Amirkabir, Vol. 14, No. 56A, pp. 1101-1126, 2004.
- [19] Meybodi, M. R., Beigy, H. and Taherkhani, M., "Cellular Learning Automata and Its Applications", Journal of Science and Technology, University of Sharif, No. 25, pp.54-77, Autumn/Winter 2003-2004.
- [20] Beigy, H. and Meybodi, M. R., "A Mathematical Framework for Cellular Learning Automata", Advances on Complex Systems, Vol. 7, Nos. 3-4, pp. 295-320, September/December 2004.
- [21] Beigy, H. and Meybodi, M. R., "Open Synchronous Cellular Learning Automata", Journal of Computer Science and Engineering, ۲005, to appear.

۵- سلول i با آگاهی از وضعیت کانال k در تمام سلولهای همسایه مقدار پاداش خود را محاسبه و به یادگیرنده می‌دهد و یادگیرنده خود را به روز می‌کند.

۶- پس از پایان مکالمه سلول i به تمام همسایگان ثابت خود پیغام `unlock_channel(k)` ارسال می‌کند و از آنها می‌خواهد کانال k را آزاد کنند.

برای ارزیابی کارایی الگوریتم پیشنهادی (CQL-CA)، یک شبکه سلولی ۷ در ۷ با تعداد کانال ۱۲ و شعاع استفاده مجدد ۲ (۱۲ حالت) در نظر گرفته شده است. متوسط زمان مکالمه و متوسط زمان تحویل کانال به ترتیب ۳ و ۲ دقیقه فرض شده است. در شبیه‌سازیها سه الگوریتم FA، Bert [۴] و Hykin [۱۱] پیاده‌سازی و با نتایج الگوریتم پیشنهادی مقایسه شده‌اند. همچنین کارایی الگوریتم پیشنهادی با نرخ یادگیری ثابت و متغیر کاهش یابنده مورد بررسی قرار دادیم. در تمام آزمایشها تحویل کانال در نظر گرفته شده است. نتایج آزمایشها که به تفصیل در [11] آمده است نشان دهنده برتری الگوریتم پیشنهادی در مقایسه با سه الگوریتم Bert، FA و Hykin می‌باشد.

۶- نتیجه گیری

در این مقاله یک مدل یادگیری جدید به نام یادگیرنده Q سلولی پیشنهاد و کاربرد آن. برای حل مساله تخصیص کانال در شبکه‌های سلولی مخابراتی ارائه گردید. از طریق شبیه‌سازی نشان داده شد که الگوریتم پیشنهادی برای تخصیص کانال مبتنی بر این مدل ترکیبی از کارایی خوبی بر خوردار است.

مراجع

- [1] Meybodi, M. R., Beigy, H. and Taherkhani, M., "Cellular Learning Automata", Proceedings of 6th Annual CSI Computer Conference, Computer Engineering Department, University of Isfahan, pp. 153 –163, 20-22 Feb. 2001.
- [2] Ahmatabadi, M. N. and Asadpour, M., "Expertness Based Cooperative Q-Learning", IEEE Transactions on Systems, Man, and Cybernetics, Vol. 32, No. 1, 2002.
- [3] Wolfram, S., Cellular Automata and Complexity, Perseus Books Group, 1994.
- [4] Bertsekas, D. P. and Singh, S., "Reinforcement Learning for Dynamic Channel Allocation in Cellular Telephone Systems", NIPS96 Proceeding, 1996.
- [5] Bertsekas, D. P. and Tsitsiklis, J. N., Neuro-Dynamic Programming, Athena Scientific, Belmont, Massachusetts, 1996.
- [6] Brown, T. and Tong, H., "Adaptive Resource Allocation in Telecommunications", Denver, Proceeding of the SPIE, 1999.
- [7] Boukerche, A. and Jacob, T., "A Distributed Algorithm for Dynamic Channel Allocation", Kluwer Academic Publishers, Netherlands, Mobile Networks Journal, Vol. 7, PP. 115-126, 2002.