

استفاده از عاملهای یادگیرنده در کنترل هوشمند ترافیک شهری

منیره عبدوس^۱ و ناصر مزینی^۲

^۱دانشکده مهندسی کامپیوتر، دانشگاه علم و صنعت ایران، abdoos@iust.ac.ir

^۲دانشکده مهندسی کامپیوتر، دانشگاه علم و صنعت ایران، mozayani@iust.ac.ir

چکیده - سیستمهای چندعامله ابزار مناسبی جهت مدلسازی سیستمهای توزیع شده به شمار می رود. در اغلب این سیستمها اطلاعات دقیق و از پیش مشخص از محیط در دسترس نیست و نیاز به روشی است که عاملها بدون داشتن اطلاعات قبلی از محیط بتوانند در آن عمل نمایند. یادگیری Q این امکان را برای عاملها فراهم می آورد که بتوانند با استفاده از تجربیات خود، خط مشی مناسبی را یاد بگیرند. از طرفی دیگر، با افزایش تعداد خودروها در شبکه های شهری، روشهای کلاسیک در کنترل ترافیک شهری قابل استفاده نبوده و نیاز به روشهای هوشمند افزایش می یابد. در این مقاله، روشی برای کنترل ترافیک شهری با استفاده از سیستمهای چند عامله مبتنی بر یادگیری Q ارائه شده است. یک شبکه شامل تعداد زیادی از تقاطعها، با یک سیستم چندعامله مدل شده است که در آن عاملها از یادگیری Q برای یافتن بهترین زمانبندی چراغهای راهنمایی بهره می برند. نتایج حاصل از پیاده سازی روش بر روی شبکه گرید شامل تقاطعهای سه فاز و چهار فاز نشان می دهند که روش پیشنهادی باعث کاهش میزان تاخیر در زمان سفر شده و همچنین در حجم ترافیک سنگین زمان اشباع شبکه را به تعویق می اندازد. کلیدواژه - زمانبندی چراغهای راهنمایی، عاملهای یادگیرنده، کنترل هوشمند ترافیک، یادگیری Q.

حوزه از تکنیکهای مبتنی بر عامل استفاده می کنند.

در این مقاله، روشی جهت کنترل چراغهای راهنمایی مبتنی بر سیستمهای چندعامله ارائه می شود. در این روش عاملها از یک مدل یادگیری تقویتی، یادگیری Q، جهت یافتن خط مشی بهینه استفاده می کنند. بنابراین روش ارائه شده در این مقاله به حوزه کنترل ترافیک با استفاده از زمانبندی تقاطع ها مربوط می باشد. تحقیقات بسیاری در زمینه حمل و نقل هوشمند با استفاده از سیستمهای چندعامله ارائه شده است که جهت بررسی آنها می توان به [۱] ارجاع نمود. در این بخش به طور اجمالی به بررسی برخی از کارهای انجام شده در راستای کنترل ترافیک با استفاده از سیستمهای چندعامله و مبتنی بر یادگیری تقویتی می پردازیم.

روشهای یادگیری تقویتی به دو دسته مبتنی بر مدل و مستقل از مدل تقسیم می شوند. از روشهای مبتنی بر مدل در کنترل ترافیک می توان به روشی که در [۲] آورده شده است، اشاره نمود که در آن عاملها از تعامل مستقیم با محیط سعی می کنند جریان ترافیک را بهینه نمایند. اگرچه نشان داده شده است که این

۱- مقدمه

یکی از اساسی ترین زیرساختهای لازم برای توسعه صنایع و افزایش سطح رفاه اجتماعی در هر کشور، سیستم حمل و نقل می باشد. امروزه با افزایش روز افزون حجم ترافیک خودروها در شبکه های شهری، مدت زمان تلف شده و به همراه آن آلودگی محیط زیست و مصرف انرژی نیز افزایش یافته است. روند رشد سریع تقاضای حمل و نقل به ویژه در ساعات اوج در کلان شهرها، از جمله مشکلات اصلی در بسیاری از شهرهای دنیا محسوب می شود. از این رو تامین حمل و نقل ایمن و کارا یکی از مهم ترین مسائل مطرح در جوامع مختلف می باشد. در سالهای اخیر استفاده از روشها و تکنیکهای مختلف هوش مصنوعی در حوزه سیستم های حمل و نقل هوشمند به طور گسترده مورد توجه قرار گرفته است. از آنجا که مفهوم عاملهای هوشمند بر بخشهای مختلفی از سیستم مانند چراغ راهنمایی، خودروها و عابران پیاده انطباق دارد، لذا بسیاری از روشهای مطرح در این

مانند روشهای [۶] و [۸] شبکه های گرید منظم و یا خطی را در نظر گرفته اند. این روشها به دلیل افزایش نمایی حجم محاسبات در شبکه های بزرگتر و غیر منظم قابل پیاده سازی نیستند. از میان روشهای موجود، برخی محققان روشی ارائه کرده اند که قابل اجرا بر روی شبکه های بزرگ می باشد. یادگیری تقویتی همکار، سعی به استخراج دانش از عاملهای همسایه در راستای یادگیری زمانبندی دارند، این روش که در [۱۰] معرفی شده است در قسمتی از شهر دوبلین شامل ۶۴ تقاطع پیاده سازی شده است. در مقاله دیگری نیز ۲۹ تقاطع با استفاده از یادگیری تقویتی و اطلاعاتی که میان عاملها رد و بدل می شود، کنترل می گردند [۹].

در این مقاله، روشی برای کنترل چراغهای راهنمایی با استفاده از یادگیری تقویتی ارائه می شود. شبکه ای بزرگ متشکل از ۳۶ تقاطع که در یک شبکه گرید در کنار یکدیگر قرار گرفته اند، با یک سیستم چندعامله مدل می شود. به این صورت که هر یک از عاملها مسئول زمانبندی یک چراغ راهنمایی می باشد. عاملها از یادگیری تقویتی برای زمانبندی استفاده می کنند که وابسته به مدل از پیش تعیین شده محیط ندارند. نحوه نمایش فضای حالت و عمل عاملها در روش ارائه شده به صورت پارامتری است که با تغییر پارامترها می توان ابعاد فضا را تغییر داد.

در بخش بعدی، مقدمه ای کوتاه در مورد یادگیری تقویتی مورد استفاده در این مقاله که یادگیری Q می باشد ذکر می گردد. قسمتهای مختلف روش پیشنهادی در بخش ۳ معرفی می شوند. در بخش ۴، نتایج تجربی حاصل از اجرای الگوریتم بر روی شبکه ترافیک شهری گزارش می شود. در بخش ۵، نتیجه گیری و کارهای آینده بیان می گردند.

۲- یادگیری Q

در یادگیری تقویتی، سیستم تلاش می کند تا تعامل خود با یک محیط پویا را از طریق خطا و آزمایش بهینه نماید. یادگیری تقویتی راهی برای یادگیری رفتارها به کمک تعامل با محیط بدون داشتن ناظر می باشد. یادگیری تقویتی در واقع چگونگی نگاشت موقعیت های مختلف به اعمال برای حصول بهترین نتیجه یا بیشترین پاداش می باشد. در بسیاری موارد، اعمال نه

روش بهتر از روش زمانبندی ثابت عمل کرده است، اما ارزیابی روش تنها برای یک تقاطع صورت گرفته است و چنانچه تعداد تقاطع بیشتر شود، تغییرات محیطی بیشتر شده و تضمینی بر عملکرد این روش یادگیری وجود ندارد. در [۳] نیز یادگیری تقویتی برای کنترل ترافیک در یک شبکه گرید مورد استفاده قرار گرفته است. در این روش از نوعی یادگیری همکار استفاده شده است که به طور همزمان خط مشی جهت کنترل سیگنالهای ترافیکی و مسیریابی بهینه صورت می گیرد. یکی از مشکلات این روش مربوط به هزینه بالای ارتباطات و دانش می باشد، به خصوص زمانی که تعداد تقاطع های شبکه افزایش می یابد. روش یادگیری مشابهی که مبتنی بر مدل می باشد نیز در [۴] معرفی شده است. هدف یادگیری طوری مدل شده است که نمایش حالت، مبتنی بر مجموع زمان انتظار خودروها در اطراف تقاطع می باشد. بدیهی است هرچه اطلاعات بیشتری از خودروها دریافت شود، مدل پیچیده تر و فضای حالت بزرگتر خواهد بود، این مساله برای شبکه های بزرگ یک مشکل اساسی محسوب می شود.

از آنجا که روشهای مبتنی بر مدل در محیطهای ایستا قابل استفاده هستند، قابلیت تطبیق با تغییرات در محیطهای پویا را ندارند. در محیطهای پویا و غیر ایستا، روشهای یادگیری تقویتی مستقل از محیط قابل استفاده می باشد. در [۵] روشی برای دسته مشخصی از محیطهای غیرایستا بیان شده است که قادر است مدل ناقصی از محیط را تخمین بزند. این روش در شبکه ای متشکل از ۹ تقاطع مورد استفاده قرار گرفته است. نتایج گزارش شده نشان می دهد که با فرض برخی خصوصیات، این روش کارایی بهتری نسبت به روشهای مستقل از مدل و مبتنی بر مدل دارد، اما وجود این مفروضات باعث می شود این روش قابل تعمیم و استفاده در شبکه های بزرگتر نباشد.

یادگیری تقویتی تطبیقی در مقالات [۶] و [۷] جهت کنترل محیط مستقل از مدل ارائه شده است. کنترل تطبیقی در [۸] نیز معرفی شده است که از تقریب تابع به عنوان نگاشتی از حالات و زمان بندی استفاده می نماید.

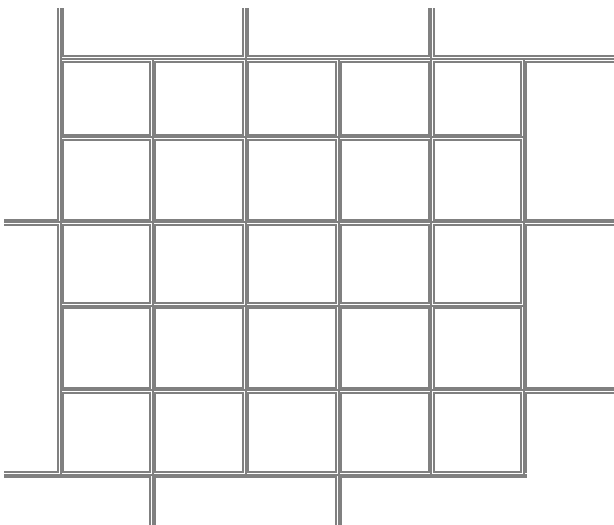
اگرچه روشهای یادگیری تقویتی مستقل از مدل در محیطهای غیر ایستا که مدلی از محیط در دسترس نیست قابل استفاده می باشد، اما برخی روشها مانند [۹] در یک تقاطع و برخی دیگر

$$\pi^*(s) = \max_a Q^*(s, a) \quad (2)$$

در این مقاله از روش یادگیری Q برای زمان بندی چراغ های راهنمایی و رانندگی توسط عامل های تعریف شده در هر تقاطع، استفاده می کنیم. جزئیات بیشتر در رابطه با نحوه تعریف فضای حالت و عمل برای کنترل ترافیک، در بخش بعدی به تفصیل شرح داده خواهد شد.

۳- یادگیری Q جهت کنترل هوشمند ترافیک

همان گونه که پیش از این اشاره شد، در راه حل پیشنهادی برای مسئله ترافیک، به هر یک از تقاطع ها یک عامل نسبت داده می شود و مجموعه ی تمامی این عامل ها، یک محیط چند عامله می سازند. این عامل ها با استفاده از حسگرهایی که در سطح خیابان های اطراف خود نصب گردیده اند، دانشی را از محیط دریافت کرده و اقدام به زمان بندی چراغ های راهنمایی و رانندگی نصب شده در تقاطع مربوط به خود می کنند. شبکه ترافیکی مورد نظر، شبکه ای متشکل از ۳۶ تقاطع و ۷۰ خیابان دوطرفه می باشد که در شکل ۱ نشان داده شده است. سیستم چندعامله متشکل از ۳۶ عامل است که هر کدام مسئول کنترل یک تقاطع می باشند. تقاطع ها به صورت چهار راه و سه راه می باشند که در یک شبکه گرید قرار گرفته اند.



شکل ۱- شبکه ترافیک شهری مورد استفاده در شبیه سازی

در این روش، به هر خیابان ورودی یک فاز اختصاص داده شده

تنها روی پاداش همان مرحله بلکه مراحل بعد هم تاثیرگذار است. این دو خصوصیت، یعنی «سعی و خطا» و «پاداش با تاخیر» مهمترین خصوصیات یادگیری تقویتی می باشند. یادگیری تقویتی از این رو مورد توجه است که راهی برای آموزش عاملها برای انجام یک عمل از طریق دادن پاداش و تنبیه است، بدون اینکه لازم باشد نحوه انجام عمل را برای عامل مشخص نمائیم.

روش یادگیری Q [۱۱]، یک روش مستقل از مدل است که در آن، عامل هیچ نوع دسترسی به مدل انتقال ندارد. این روش، یکی از بهترین و پرکاربردترین روشهای یادگیری تقویتی در حل مسائل یادگیری می باشد. در این روش، عامل ارزش انتخاب عمل در حالت s که با $Q^*(s, a)$ نشان داده می شود را با استفاده از تعامل پیوسته با محیط و با سعی و خطا تخمین می زند. در این روش، عامل با مقادیر تصادفی از تخمین ها شروع کرده و بعد از هر عمل، چندتایی به صورت $\langle s, a, r, s' \rangle$ را دریافت می کند که در آن s حالت فعلی، a عمل انجام شده در حالت s ، r پاداش فعلی و s' حالت بعد از اجرای a می باشد. عامل برای هر چندتایی می تواند ارزش حالت-عمل مربوطه را به صورت زیر محاسبه کند:

$$Q(s, a) = (1 - \alpha) \times Q(s, a) + \alpha \times [r + \gamma \times \max_{a'} Q(s', a')] \quad (1)$$

که در آن، $\alpha \in [0, 1]$ نرخ یادگیری عامل است و مشخص می کند که تا چه حدی اطلاعات جدید بدست آمده، جایگزین اطلاعات قدیمی شوند. مقدار ۱ برای این نرخ سبب می شود که عامل فقط جدیدترین اطلاعات را در نظر گیرد و مقدار صفر باعث می گردد عامل یادگیری نداشته باشد؛ و $\gamma \in [0, 1]$ فاکتور کاهش نامیده می شود، برای مشخص کردن اهمیت پاداش های آینده است. مقدار صفر برای این فاکتور، عامل را فرصت طلب می کند-یعنی عامل فقط پاداش فعلی را در نظر می گیرد. از سوی دیگر، نزدیکی به مقدار ۱ سبب می شود که عامل برای یک پاداش بالا در طولانی مدت منتظر بماند.

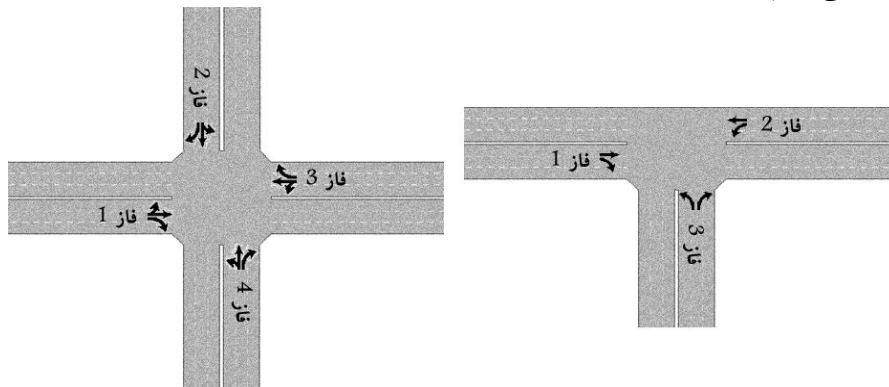
اثبات می شود که اگر تمامی زوج های حالت-عمل به صورت مکرر تجربه شوند و نرخ یادگیری در طول زمان کاهش یابد، یادگیری-

Q با احتمال ۱ به مقدار بهینه $Q^*(s, a)$ همگرا می شود. با توجه به این امر، سیاست بهینه به صورت زیر تعریف می شود:

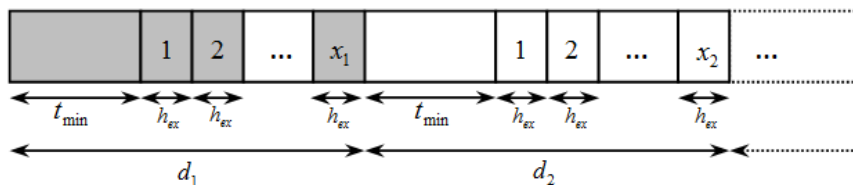
است. عاملها برای تخمین حالت از میانگین طول صف در خیابانهای ورودی استفاده می کنند، به این ترتیب که خیابانهای ورودی بر اساس این مقدار مرتب می شوند و در نتیجه تعداد حالتهای ممکن برابر با تعداد جایگشتهای ممکن خیابانهای ورودی می باشد.

فرض کنید عامل α^i دارای k خیابان ورودی است که با l_j ، $j=1, \dots, k$ نشان می دهیم. در این صورت،

فرض کنید عامل α^i دارای k خیابان ورودی است که با l_j ، $j=1, \dots, k$ نشان می دهیم. در این صورت،



شکل ۲- ترتیب فازها در تقاطعهای سه فاز و چهار فاز



شکل ۳- نمایش گرافیکی زمان اختصاص داده شده به فازها

تقاطعهایی که دارای سه فاز هستند، ۶ حالت و برای تقاطعهای چهار راه، ۲۴ حالت خواهیم داشت.

عملی که برای عامل تعریف شده است، نحوه تقسیم زمان سبز مربوط به هر فاز را مشخص می کند. در این روش یک طول ثابت زمانی را در نظر گرفته و آن را طوری میان فازهای مختلف تقسیم می کنیم که به هر فاز، زمانی به عنوان حداقل زمان سبز اختصاص یابد و علاوه بر آن تعدادی بازه سبز برای تمدید در نظر گرفته شده است که هر عمل عامل، نحوه تخصیص تمدیدهای زمانی به فازهای مختلف را مشخص می نماید. فضای عمل با چهارتایی $\langle n_{ph}, t_{min}, n_{ex}, h_{ex} \rangle$ مشخص می گردد که در آن:

n_{ph} : تعداد فازها

t_{min} : حداقل زمان سبزی است که باید به هر سیگنال تخصیص یابد

n_{ex} : تعداد تمدیدهای زمانی

h_{ex} : طول هر تمدید زمانی بر حسب ثانیه

زمان سبز تخصیص داده شده به فازهای مختلف که با δ نشان داده می شود به صورت زیر قابل محاسبه است.

$$\delta = h_{ex} \times n_{ex} + n_{ph} \times t_{min} \quad (3)$$

این رابطه به صورت گرافیکی در شکل ۳ آورده شده است. زمان سبز اختصاص داده شده به فاز i را d_i مشخص می نماید.

تخصیص تمدیدهای زمانی از نظر تئوری مشابه با حل مسئله زیر می باشد:

تقاطعهایی که دارای سه فاز هستند، ۶ حالت و برای تقاطعهای چهار راه، ۲۴ حالت خواهیم داشت.

عملی که برای عامل تعریف شده است، نحوه تقسیم زمان سبز مربوط به هر فاز را مشخص می کند. در این روش یک طول ثابت زمانی را در نظر گرفته و آن را طوری میان فازهای مختلف تقسیم می کنیم که به هر فاز، زمانی به عنوان حداقل زمان سبز اختصاص یابد و علاوه بر آن تعدادی بازه سبز برای تمدید در نظر گرفته شده است که هر عمل عامل، نحوه تخصیص تمدیدهای زمانی به فازهای مختلف را مشخص می نماید. فضای عمل با چهارتایی $\langle n_{ph}, t_{min}, n_{ex}, h_{ex} \rangle$ مشخص می گردد که در آن:

n_{ph} : تعداد فازها

جدول ۱- تاثیر ρ بر روی تعداد جوابهای معادله (۷)

ρ	تعداد جوابهای معادله (سه فاز)	تعداد جوابهای معادله (چهار فاز)
۱	غیرقابل حل	۱
۲	۶	۱۹
۳	۱۲	۳۱
۴	۱۵	۳۵

در این مقاله برای تقاطعهای سه فاز $\rho=4$ و تقاطعهای چهارفاز $\rho=2$ در نظر گرفته شده است. بنابراین تعداد حالتیهای مختلفی که برای تخصیص زمان سبز می توان در نظر گرفت، برای تقاطعهای مختلف با روابط زیر به دست می آید:

$$x_1 + x_2 + x_3 = 4, \quad x_i \in 0,1,2,3,4 \quad (8)$$

$$x_1 + x_2 + x_3 + x_4 = 4, \quad x_i \in 0,1,2 \quad (9)$$

زمان سبز که به هر فاز اختصاص می یابد برابر است با مجموع حداقل زمان سبز و میزان تمدید های زمان سبز که با رابطه زیر حاصل می گردد:

$$d_i = t_{\min} + x_i \times h_{ex} \quad (10)$$

به طور خلاصه، مقادیر پارامترهایی که در یادگیری سطح اول مورد استفاده قرار گرفته است، در جدول ۲ نشان داده شده است. با توجه به این مقادیر، مجموعه اعمال برای عملهای مربوط به تقاطعهای سه فاز شامل ۱۵ عمل و برای عملهای مربوط به تقاطعهای چهار فاز ۱۹ عمل می باشد که چند نمونه از آنها در جدول ۳ آورده شده است. توجه داشته باشید که ۲ ثانیه برای زمان زرد پس از هر فاز در نظر گرفته شده است.

جدول ۲- پارامترهای مورد استفاده در یادگیری عملهای سطح اول

پارامتر	توضیح	مقدار (۳ فاز)	مقدار (۴ فاز)
δ	طول دوره زمانی موثر	۹۴	۹۲
n_{ph}	تعداد فازها	۳	۴
t_{\min}	حداقل زمان سبز برای هر فاز	۱۸	۱۳
n_{ex}	تعداد بازه های زمانی تمدید فاز	۴	۴

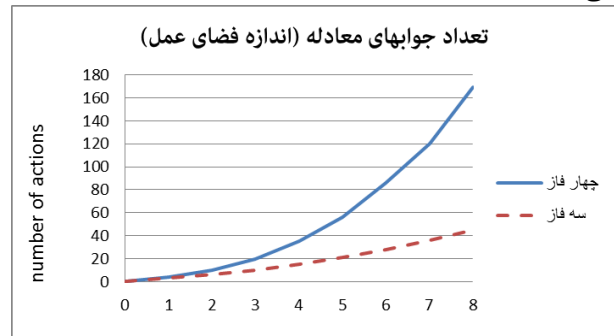
$$\sum_{i=1}^{n_{ph}} x_i = n_{ex}, \quad x_i \in \mathbb{N} \quad (4)$$

که در آن n_{ex} ، تعداد تمدیدهای زمانی و n_{ph} تعداد فازها می باشد. واضح است که x_i باید اعداد حسابی باشند. از آنجا که تعداد فازها در شبکه مورد بررسی برای تقاطعهای سه راه و چهارراه به ترتیب ۳ و ۴ می باشد، لذا رابطه (۴) به صورت زیر نوشته می شود:

$$x_1 + x_2 + x_3 = n_{ex} \quad (5)$$

$$x_1 + x_2 + x_3 + x_4 = n_{ex} \quad (6)$$

تعداد جوابهای معادله فوق، تعداد اعضای مجموعه اعمال را مشخص می نماید و این تعداد به n_{ex} ، تعداد تمدیدها، بستگی دارد. شکل ۴ تعداد اعضای مجموعه اعمال را بر حسب n_{ex} نشان می دهد. همانطور که در شکل ۴ دیده می شود، روند افزایش تعداد اعمال نسبت به n_{ex} نمای می باشد.



شکل ۴- نمودار تعداد عمل بر حسب تعداد تمدیدهای زمانی در تقاطعهای سه فاز و چهار فاز

در یادگیری تقویتی چنانچه تعداد اعمال ممکن برای عامل زیاد باشد، همگرایی بسیار کند و در برخی موارد نیز حتی عدم همگرایی رخ می دهد. در این روش، تعداد تمدیدها را برابر با ۴ در نظر گرفته ایم. به این ترتیب، برای سه فاز تعداد جوابهای معادله برابر با ۱۵ و برای چهار فاز این تعداد برابر با ۳۵ می شود. از محدودیت زیر برای کاهش تعداد جوابهای معادله و محدود کردن تغییرات زمان فازهای مختلف می توان استفاده نمود.

$$\sum_{i=1}^{n_{ph}} x_i = n_{ex}, \quad x_i \in \mathbb{N}, \quad x_i \leq \rho, \quad 1 \leq \rho \leq n_{ex} \quad (7)$$

حداکثر تعداد تمدیدهایی که به هر فاز می توان اختصاص داد با پارامتر ρ کنترل می شود. چنانچه ρ مقدار کوچکی در نظر گرفته شود، از تعداد عمل های ممکن کاسته می گردد. برای انتخاب ρ ، حالتیهای مختلفی را می توان در نظر گرفت که در جدول ۱ آورده شده است:

جدول ۴- پارامترهای مورد استفاده در شبیه سازی شبکه ترافیک شهری

مقدار	خصوصیات
۳۶	تعداد تقاطع ها
۷۰	تعداد خیابانها
۶۱۲/۸۵	میانگین طول خیابانها (متر)
۳	تعداد خطوط عبوری خیابانها
۵۰	حداکثر سرعت (کیلومتر بر ساعت)
۱۰	تعداد مراکز ورودی/خروجی
نرمال	توزیع ورودی خودروها
۱۰	مدت زمان شبیه سازی(ساعت)
۱۰۰۰۰	حجم ترافیک ۱ (خودرو در ساعت)
۲۰۰۰۰	حجم ترافیک ۲ (خودرو در ساعت)

نتایج حاصل از اجرای روشهای مذکور برای حجم ترافیک ۲ در شکل ۶ نشان داده شده است. برای حجم ترافیک بالا اگرچه با استفاده از هر دو روش، شبکه بعد از مدتی اشباع شده، اما میانگین تاخیر قبل از اشباع و همچنین زمان رسیدن به اشباع متفاوت می باشد.

۵- نتیجه گیری

سیستمهای چندعامله اغلب برای مدل کردن سیستمهای توزیع شده در مقیاس بالا مورد استفاده قرار می گیرند. در اغلب این سیستمها اطلاعات دقیق و از پیش مشخص از محیط در دسترس نیست و نیاز به روشی است که عاملها بدون داشتن اطلاعات قبلی از محیط بتوانند در آن عمل نمایند که یادگیری Q این امر را ممکن می سازد. در این مقاله یک شبکه ترافیکی شامل ۳۶ تقاطع با استفاده از سیستمهای چندعامله مدل گردید که در آن عاملها با به کارگیری یادگیری Q به کنترل و زمانبندی چراغهای راهنمایی در هر تقاطع می پردازند.

h_{ex}	طول زمان تمدید	۱۰	۱۰
ρ	حداکثر تعداد تمدید برای هر فاز	۲	۴
$ S $	تعداد حالات	۲۴	۶
$ A $	تعداد اعمال	۱۹	۱۵

یک عمل برابر با یک دوره کامل از هر فاز خواهد بود که در یک ترتیب ثابت اجرا می گردد. برای کلیه عاملها طول یک دوره زمانی با احتساب زمان زرد برابر با ۱۰۰ ثانیه می شود. در روش یادگیری تقویتی، برای محاسبه پاداش، مقدار نرمال شده میانگین طول صف خیابانها مورد استفاده قرار گرفته است.

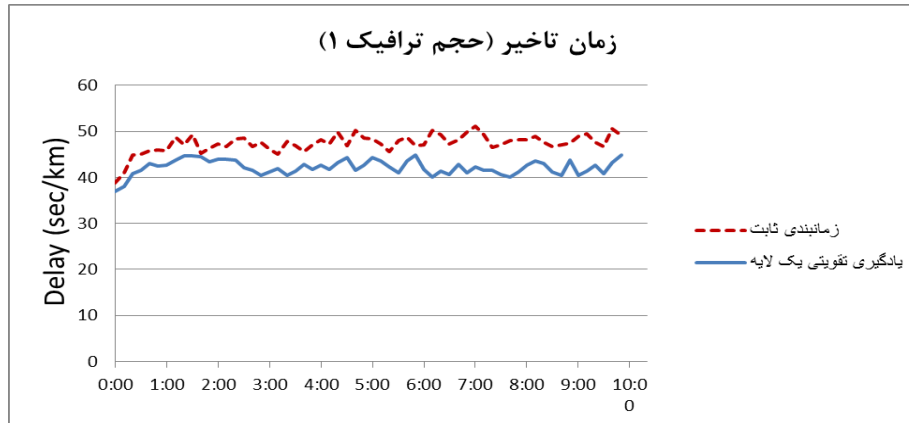
جدول ۳- مثالی از مقادیر زمان سبز اختصاص داده شده به فازهای مختلف

عملهای دارای سه فاز		عملهای دارای چهار فاز			
شماره عمل	شماره عمل	فاز ۱	فاز ۲	فاز ۳	فاز ۴
۱	۱	۳۳	۳۳	۱۳	۱۳
۲	۲	۳۳	۱۳	۳۳	۱۳
۳	۳	۳۳	۱۳	۱۳	۳۳
:	:	:	:	:	:

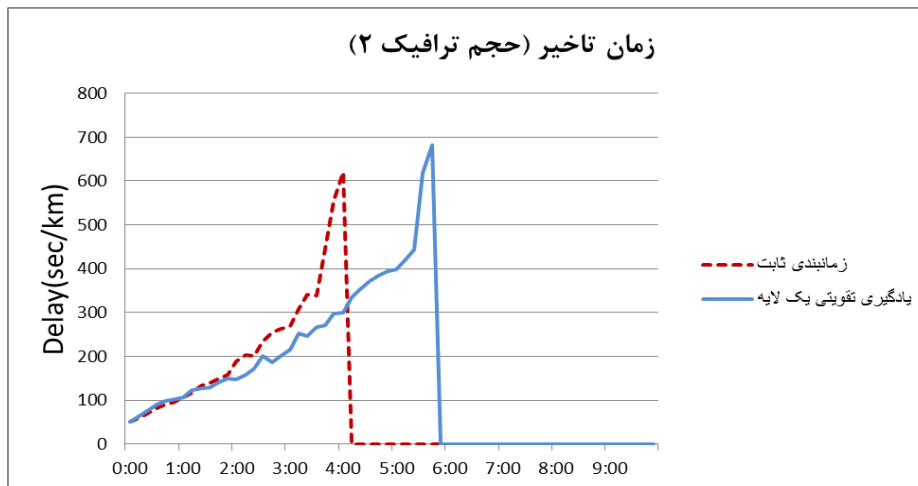
۴- نتایج تجربی

روش ارائه شده در این مقاله با استفاده از نرم افزار Aimsun بر روی شبکه شکل ۱ پیاده سازی شده است. نتایج تجربی حاصل از به کارگیری این روش در این بخش آورده شده است. پارامترهایی که در شبیه سازی ترافیک شهری در نظر گرفته شده اند، در جدول ۴ آورده شده است.

نتایج حاصل از اجرای روش پیشنهادی برای حجم ترافیک ۱ در شکل ۵ با روش زمانبندی ثابت مقایسه شده است. این نتایج میانگین حاصل از ۱۰ اجرای مستقل روشها بر روی شبکه مذکور می باشد. نتایج نشان می دهد که میانگین تاخیر در کل شبکه کاهش یافته است و در نتیجه کارایی بهبود یافته است.



شکل ۵- میانگین تاخیر در ۱۰ ساعت شبیه سازی برای حجم ترافیک ۱



شکل ۶- میانگین تاخیر در ۱۰ ساعت شبیه سازی برای حجم ترافیک ۲

سیاسگزاری

این پژوهش با حمایت مالی مرکز تحقیقات مخابرات ایران صورت گرفته است.

مراجع

- [1] B. Chen, H. H. Cheng and J. Palen, "Integrating mobile agent technology with multi-agent systems for distributed traffic detection and management systems", Transportation Research Part C: Emerging Technologies, No. 17, Issue 1, pp. 1-10, 2009.
- [2] P. Grégoire, C. Desjardins, J. Laumônier and B. Chaib-draa, "Urban traffic control based on learning agents", Intelligent Transportation Systems Conference, Laval University, Quebec, pp. 916-921, 2007.
- [3] M. Weiring, "Multi-agent reinforcement learning for traffic light control, 7th International Conference on Machine Learning, Stanford University, US, pp. 1151-1158, 2000.
- [4] M. Steingröver, R. Schouten, S. Peelen, E. Nijhuis, B. Bakker, "Reinforcement learning of traffic light controllers adapting to

مدل یادگیری ارائه شده به صورت پارامتری بوده و لذا امکان افزایش و یا کاهش اعمال ممکن در مجموعه عمل عامل وجود دارد. از طرفی دیگر مجموعه حالت و عمل طوری تعریف شده است که به راحتی می تواند برای تقاطع هایی با تعداد فازهای مختلف تعمیم یابد. نتایج تجربی حاصل از اجرای روش بر روی شبکه گرید، پیاده سازی و با روش زمانبندی ثابت مقایسه گردیده است. میانگین تاخیر با استفاده از روش پیشنهادی بهتر از زمانبندی ثابت بوده است. در حجم ترافیک بالا، اگرچه هر سه روش در نهایت منجر به اشباع شبکه می شوند اما زمان رسیدن به اشباع فرق می کند و روش پیشنهادی باعث تعویق زمان اشباع شبکه شده است. کارهای آینده شامل به کارگیری روش پیشنهادی در شبکه های واقعی و تعمیم روش برای تقاطع هایی با تعداد فازهای مختلف می باشد.



یازدهمین کنفرانس سیستم های هوشمند ایران

۹ و ۱۰ اسفند ۱۳۹۱

(ICIS2013)

11th Iranian Conference on Intelligent Systems
February 27th & 28th, 2013



- [8] L.A. Prashanth and S. Bhatnagar, "Reinforcement learning with function approximation for traffic signal control", IEEE Transactions on Intelligent Transportation Systems, No. 99, pp. 1–10, 2011.
- [9] P.G. Balaji, X. German and D. Srinivasan, "Urban traffic signal control using reinforcement learning agents", IEEE Transactions on Intelligent Transportation Systems, No. 4, Issue 3, pp. 177–188, 2010.
- [10] A. Salkham, R. Cunningham, A. Garg and V. Cahill, "A collaborative reinforcement learning approach to urban traffic control optimization", International Conference on Web Intelligence and Intelligent Agent Technology, Sydney, Australia, pp. 560–566, 2008.
- [11] C. Watkins and P. Dayan, "Q-learning", Machine learning, Vol. 8, pp. 279-292, 1992.
- traffic congestion", 17th Belgium-Netherlands Conference on Artificial Intelligence (BNAIC 2005), Brussels, Belgium , pp. 216–223, 2005.
- [5] B.C.d. Silva, E.W. Basso, A.L.C. Bazzan and P.M. Engel, "Improving reinforcement learning with context detection", 5th International Joint Conference on Autonomous Agents and Multi-agent Systems (AAMAS), Hakodate, Japan, pp. 811–812, 2006.
- [6] Y. Dai, D. Zhao and J.a Yi, "A comparative study of urban traffic signal control with reinforcement learning and adaptive dynamic programming", International Joint Conference on Neural Networks. Barcelona, pp. 1–7, 2010.
- [7] K. Wen, S. Qu and Y. Zhang, "A stochastic adaptive control model for isolated intersections", IEEE International Conference on Robotics and Biomimetics, Sanya, China, pp. 2256–2260, 2008.