



دانشگاه صنعتی سهند

DOR:

[20.1001.1.23223146.1401.9.2.1.9](https://doi.org/10.1001.1.23223146.1401.9.2.1.9)

نشریه سالانه علمی غیرخطی در مهندسی برق

دوره ۱۰ - شماره ۲

پاییز و زمستان ۱۴۰۲

صفحات ۵۵ الی ۷۳

ISSN: 2322-3146

<http://journals.sut.ac.ir/jnsee>

تسهیم طیف برای مقابله با تداخلگر با استفاده از یادگیری عمیق تقویتی

چندعاملی

ندا کاظمی^۱ و معصومه آذغانی^۲

^۱ دانشجوی دکترا، گروه مخابرات، آزمایشگاه مخابرات بی سیم و پردازش سیگنال (WCSP)، دانشکده مهندسی برق، دانشگاه صنعتی سهند، تبریز، n_kazemi97@sut.ac.ir

^۲ نویسنده مسئول، دانشیار، گروه مخابرات، آزمایشگاه مخابرات بی سیم و پردازش سیگنال (WCSP)، دانشکده مهندسی برق، دانشگاه صنعتی سهند، تبریز، mazghani@sut.ac.ir

چکیده

در تسهیم طیف بین کاربران مختلف، مقابله با تداخلگر موجود در شبکه بسیار حیاتی است. اهمیت این کار به ویژه زمانی دوچندان می شود که تداخلگر، هوشمند بوده و قادر به یادگیری الگوهای ارتباطی بین کاربران باشد. در این مقاله، روشی برای تسهیم طیف با هدف مقابله با تداخلگر شبکه پیشنهاد می گردد. شیوه پیشنهادی بر پایه یادگیری تقویتی چندعاملی است. در مرحله اول، طیف های مورد حمله تداخلگر توسط ایستگاه پایه به دست می آیند. در مرحله بعد، یک چارچوب یادگیری تقویتی عمیق طراحی می شود که به هر کدام از کاربران طیف امن و مناسب را پیشنهاد می دهد. برای این منظور از همبستگی بین سیگنال راهنمای ارسالی و سیگنال دریافتی در هر بازه فرکانسی، میزان تاثیر تداخلگر در آن زیرباند حاصل می شود. سپس ضریب تخفیف الگوریتم یادگیری تقویتی بر پایه مقدار همبستگی به صورت وفقی تنظیم می گردد. به این ترتیب، زیرباندهایی که کمتر تحت تاثیر تداخلگر واقع شده اند، به کاربران اختصاص می یابند. روش پیشنهادی در سناریوهای مختلف مورد شبیه سازی و ارزیابی قرار گرفته و نتایج حاکی از آن است که نرخ مجموع شبکه با سرعت بالایی به نرخ مطلوب همگرا می شود و شیوه پیشنهادی مقاومت خوبی در برابر تداخلگر از خود نشان می دهد.

واژه های کلیدی

سیستم های بی سیم،

یادگیری تقویتی چند عاملی،

تسهیم طیف،

ارتباطات امن،

روش های ضد اختلال گر.



Sahand University
of Technology

DOR:

[20.1001.1.23223146.1401.9.2.1.9](https://doi.org/10.1001.1.23223146.1401.9.2.1.9)

Journal of Nonlinear
Systems in Electrical
Engineering

Vol.10, No.2

Autumn and Winter 2023

ISSN: 2322 – 3146

<http://journals.sut.ac.ir/jnsee>

Spectrum Sharing in Anti-jamming Using Multi-agent Reinforcement Learning

Neda Kazemi¹ and Masoumeh Azghani²

¹Ph.D student, Laboratory of Wireless Communication and Signal Processing (WCSP), Faculty of Electrical Engineering, Sahand University of Technology, Tabriz, n_kazemi97@sut.ac.ir

²**Corresponding Author**, Laboratory of Wireless Communication and Signal Processing (WCSP), Faculty of Electrical Engineering, Sahand University of Technology, Tabriz, mazghani@sut.ac.ir

ABSTRACT

Keywords

Wireless systems,
Multi-agent reinforcement
learning,
Spectrum sharing,
Secure communication,
Anti-jamming methods.

In spectrum sharing among different users, anti-jamming process in the network is very critical. The importance of this work doubles especially when the jammer is intelligent and able to learn patterns of communications between users. In this article, a method for spectrum sharing is proposed with the aim of dealing with the jammer. The proposed method is based on multi-agent reinforcement learning. In the first stage, the jammed slots of spectrum are acquired by the base station. In the next stage, a deep reinforcement learning framework is designed that offers each user a safe and appropriate range. For this purpose, the influence of the jammer in that sub-band is obtained from the correlation between the transmitted pilot signal and the received signal in each frequency range. Then, the discount factor of the reinforcement learning algorithm is adjusted based on the correlation value. In this way, the sub-bands that are less affected by the jammer are allocated by the users. The proposed method has been simulated and evaluated in different scenarios and the results indicate that the total rate of the network converges to the desired rate at a high speed and the proposed method shows good resistance against the jammer.

۱- مقدمه

انفجار اطلاعات و رشد ترافیک داده و تجاری سازی شبکه های ارتباطی بی سیم نسل پنجم^۱ (5th)، چالش های بزرگی در زمینه استفاده از طیف و امنیت داده ها تحمیل می کند. اقدامات ضد اختلالگر^۲ برای خنثی نمودن تهاجم اطلاعاتی، یک موضوع بسیار مهم در ارتباطات است. مخصوصا در ارتباطات بی سیم، اختلال گر خارجی، کیفیت لینک ارتباطی را به شدت تهدید می کند. در سال های اخیر، با توسعه هوش مصنوعی، مقابله هوشمند با حملات اختلال گرها بسیار مرسوم شده است [۱]، [۲]. در پژوهش های پیشین، تعاملات بین کاربران قانونی^۳ و تداخلگر^۴ها به شکل تئوری بازی ها مدل سازی شده است [۳]–[۵]. بیشتر پژوهش های ضد اختلال گر موجود، نیاز به دستیابی به اطلاعات قبلی از الگوی سیگنال ارسالی اختلال گر در فرمول بندی مربوط به استراتژی خود دارند. دو چالش مهم در این زمینه وجود دارد، چالش اول اینکه امکان دارد اطلاعات به دست آمده از تداخلگر برای تخمین الگوی آن کافی نباشد و کاربر برای اقدام ضد تداخلگر با کمبود اطلاعات مواجه شود اما چالش دوم زمانی به وجود می آید که اقدام ضد تداخلگر در یک محیط پویا اتفاق بیافتد و تداخلگر مدام الگوی خود را تغییر دهد. این درحالی است که با توسعه سریع هوش مصنوعی تداخلگرها به راحتی می توانند حملات هوشمند و پویا ایجاد کنند [۶]. برای غلبه بر این چالش ها، روش پیشنهادی در این مقاله برای عملیات ضد اختلال گری خود، نیازی به دستیابی به الگوی سیگنال اختلالگر ندارد. علاوه بر این، عملیات ضد اختلال گری در فرستنده و قبل از ارسال اطلاعات انجام می شود. بسیاری از پژوهشگران الگوریتم های یادگیری ماشین را برای تصمیم گیری های ضد اختلال گر آنلاین پیشنهاد کرده اند [۷]. یادگیری تقویتی (RL)^۵، یک روش مناسب برای تصمیم گیری در محیط های پویا است. در پژوهش های پیشین، یادگیری تقویتی به عنوان یک الگوریتم مناسب در اقدامات ضد تداخلگر به کار گرفته شده است [۸]–[۱۰]. یادگیری تقویتی شاخه ای از یادگیری ماشین^۶ است که در آن یک یا چند عامل با یک محیط معین تعامل می کنند و خط مشی^۷ بهینه را برای دستیابی به هدف از پیش تعریف شده می آموزند. تا به امروز، روش های یادگیری تقویتی بسیاری توسعه پیدا کرده اند [۱۱]–[۱۳]. این روش ها در مسایل مختلفی از جمله تخصیص منابع و طیف مورد استفاده قرار گرفته اند [۱۴]–[۱۷]. اما بسیاری از مسائل دنیای واقعی، فضای حالت بزرگ با پاداش تاخیری دارند. ساختار ابتدایی روش های یادگیری تقویتی به صورت کامل بهینه نبوده و فرآیند یادگیری آن ها بسیار کند است. یادگیری تقویتی عمیق روشی است که از شبکه های عصبی برای تخمین تابع ارزش^۸ در یادگیری تقویتی استفاده می کند [۱۸]. در یادگیری تقویتی تک عاملی^۹ (SARL)، یک عامل^{۱۰}، یک اقدام^{۱۱} مشخص را در یک وضعیت معین مطابق با خط مشی خود اتخاذ می کند و اقدام مربوطه بر روی محیط اعمال می شود. سپس عامل، حالت بعدی را بر مبنای پاداش مربوطه از محیط کسب می کند. برای محاسبه عمل بعدی، عامل مجموع

¹ 5th Generation

² Anti-jamming

³ Legitimate users

⁴ jammer

⁵ Reinforcement learning

⁶ Machine learning

⁷ policy

⁸ Value function

⁹ Single agent reinforcement learning

¹⁰ agent

¹¹ action

سیگنال‌های پاداش^۱ تخفیف‌یافته را در یک افق محدود یا نامتناهی به حداکثر می‌رساند. در یادگیری تقویتی چندعاملی^۲ (MARL)، چند عامل با اشتراک‌گذاری پاداش‌های خود با یک تصمیم‌گیری جمعی، اقدام بعدی خود را تعیین می‌کنند [۱۹]. از آنجایی که پاداش آینده ناشناخته است، RL عمیق از شبکه‌های عصبی برای تخمین انتظارات به نام تابع ارزش برای آینده استفاده می‌کند. الگوریتم‌های یادگیری تقویتی مرسوم فقط می‌توان در تداخلگرهای با نظم و مدل مشخص اعمال کرد و اگر حالت تداخلگر خارجی تغییر کرد، نیاز به آموزش مجدد برای رسیدن به همگرایی خواهند داشت [۱]، [۲۰]. با توجه به اینکه تابع ارزش توسط پاداش و ضریب تخفیف محاسبه می‌شود، انتخاب صحیح آنها از اهمیت بالایی برخوردار است. در محیطی که عدم قطعیت وجود دارد، اگر تجربه به دست آوردن پاداش بدون در نظر گرفتن عدم قطعیت ارزیابی شود، می‌توان مجموعه را نسبت به عدم قطعیت هم آموزش داد. بنابراین، روش‌های یادگیری تقویتی برای تخمین توزیع پاداش تجمعی معرفی شده است. به طور خاص، تمرکز این مقاله بر روی ضریب تخفیف است که مقدار آن در پژوهش‌های موجود اغلب ثابت است. اگر ضریب تخفیف در محیطی با تغییرات سریع نسبت به زمان، مقدار بزرگی انتخاب شود، تاثیر مقادیر پاداش آینده در پاداش فعلی زیاد خواهد بود و به همین دلیل تخمین مجموع پاداش‌ها در یک دوره زمانی طولانی‌تر به دلیل عدم اطمینان زیاد منجر به بروز موقعیت‌های غیرمنتظره می‌شود. به عنوان مثال، عوامل می‌توانند در آینده دور پاداش‌ها را بیش از حد برآورد کنند، یا ممکن است همگرایی تابع ارزش با مشکل روبرو شود [۲۱]. از سوی دیگر، اگر ضریب تخفیف خیلی کوچک باشد، تاثیر پاداش آینده در تابع پاداش بسیار کم می‌شود و تصمیم‌گیری بدون آینده‌نگری کافی انجام می‌گیرد. به عبارت دیگر، با ضریب تخفیف کوچک، پاداش تنها در آینده نزدیک در ارزیابی تابع ارزش در نظر گرفته می‌شود که می‌تواند منجر به اقدامات کوتاه‌بینانه یا تهاجمی گردد. علاوه بر این، مقدار کوچک ضریب تخفیف سبب کاهش سرعت همگرایی یادگیری می‌شود [۲۲]. این مشاهدات ما را برمی‌انگیزد تا یک قانون تطبیقی برای به‌روزرسانی ضریب تخفیف پیشنهاد نماییم. برای این منظور ابتدا تعدادی سیگنال راهنما از ایستگاه پایه به کاربران ارسال می‌گردد. کاربران در مرحله بعدی سیگنال دریافتی را به ایستگاه پایه فیدبک می‌نمایند. در ایستگاه پایه همبستگی بین سیگنال‌های ارسالی به کاربران و سیگنال دریافتی از هر کاربر محاسبه می‌گردد و زیرباندی که میزان همبستگی بالاتری دارد، احتمال تخصیص بالاتری خواهد داشت. به این ترتیب، ضریب تخفیف تطبیقی‌یافته در روش پیشنهادی منجر به دسترسی به طیف فرکانسی امن‌تر برای خنثی نمودن اثر تداخلگر می‌شود. با توجه به اینکه زمانی که فرآیند تصمیم‌گیری ابعاد بالایی دارد و محیط دارای عدم قطعیت است و یافتن ضریب تخفیف ثابت بهینه، یک کار پرچالش است، وجود ضریب تخفیف تطبیقی‌یافته می‌تواند به حل مساله کمک شایانی کند. ضریب تخفیف تطبیقی‌یافته پیشنهادی، عامل را قادر می‌سازد تا برآورد درست‌تری از مبلغ پاداش آینده داشته باشد. برای ارزیابی عملکرد روش پیشنهادی، مجموع نرخ ارسال داده، توان تداخلگر و توان نویز و تعداد تکرار برای رسیدن به پاداش مطلوب مدنظر قرار گرفته است و عملکرد روش پیشنهادی بهتر از حالت با ضریب تخفیف ثابت است. نتایج حاصل نشان می‌دهد که روش پیشنهادی با ضریب تخفیف تطبیقی‌یافته در مقایسه با سایر روش‌ها در محیط‌های پویا بهتر عمل می‌کند.

¹ reward

² Multi agent reinforcement learning

۲- یادگیری تقویتی چندعاملی

یادگیری تقویتی چندعاملی نوع مهمی از یادگیری ماشین است که در آن چند عامل^۱ در همکاری یا رقابت با هم می آموزند که چگونه در محیط با انجام اقدامات و بررسی نتایج آن‌ها رفتار کنند. در هوش مصنوعی^۲، عامل هوشمند، موجودیت خودکاری است که محیط اطراف خود را می بیند و با استفاده از محرک‌ها اقداماتی را در محیط انجام می‌دهد و فعالیت‌های خود را در مسیر کسب اهداف هدایت می‌کند. عامل‌های هوشمند ممکن است از یادگیری یا دانش برای کسب اهداف خود بهره بگیرند. این عامل‌ها ممکن است بسیار ساده یا خیلی پیچیده باشند. در روش یادگیری تقویتی، یک عامل در حال یادگیری محیط از طریق پاداش اقدام‌ها است، به این صورت که عامل حالت^۳ S_0 را از محیط دریافت می‌کند، بر اساس حالت S_0 ، عمل^۴ A را انجام می‌دهد و محیط به حالت جدید S_1 انتقال پیدا می‌کند و محیط پاداش^۵ R_1 را به عامل می‌دهد. این حلقه یادگیری تقویتی دارای یک توالی از حالت، عمل و پاداش است. هدف عامل آن است که پاداش انباره‌ای (تجمعی)^۶ مورد انتظار را بیشینه کند. طرز کار یادگیری تقویتی در شکل ۱ نشان داده شده است.



شکل ۱- طرز کار یادگیری تقویتی

اما چالش اصلی در مبحث پاداش این است که پاداشی که در گام‌های زمانی ابتدایی برای اقدامات در نظر گرفته شده‌اند، احتمال وقوع بیشتری دارند، زیرا از پاداش‌های بلند مدت آینده قابل پیش‌بینی تر هستند. بنابراین یک ضریب تخفیف با عنوان β تعریف می‌شود. این مقدار باید بین ۰ و ۱ باشد. هر چه β بزرگ‌تر شود، تخفیف کمتر است. این یعنی عامل یادگیرنده به پاداش‌های بلند مدت اهمیت بیشتری می‌دهد. از سوی دیگر، هر چه β کوچک‌تر باشد، تخفیف بیشتر است. این یعنی عامل توجه بیشتری به پاداش‌های کوتاه مدت می‌کند. پاداش مورد انتظار انباره‌ای تخفیف داده شده به صورت جمع پاداش‌ها در گام‌های زمانی مختلف محاسبه می‌شود. هر پاداش با β به توان گام زمانی، تخفیف داده می‌شود و به پاداش تخفیف داده شده سایر زمان‌ها افزوده می‌گردد. در سناریو پیاده‌سازی یادگیری تقویتی، یک نقطه آغازین و یک نقطه پایانی وجود دارد. این منجر به ایجاد یک اپیزود یعنی لیست حالت‌ها، اعمال، پاداش‌ها و حالت‌های جدید می‌شود. در یادگیری به روش مونت کارلو، هنگامی که اپیزود به پایان می‌رسد، عامل به پاداش تجمعی کل نگاه می‌کند تا از نحوه عملکرد خود آگاه شود. در روش مونت کارلو، پاداش‌ها تنها در پایان بازی دریافت می‌شوند. بنابراین، بازی جدید با یک

¹ Agent

² Artificial intelligence

³ State

⁴ Action

⁵ Reward

⁶ Cumulative reward

دانش افزوده آغاز می شود و عامل در هر تکرار تصمیمات بهتری اتخاذ می کند. در یادگیری تقویتی سیاست محور، قصد بهینه سازی تابع سیاست^۱ P است. سیاست چیزی است که رفتار عامل را در یک زمان داده شده، تعیین می کند. عامل یک تابع سیاست را می آموزد. این امر به او کمک می کند تا هر حالت را به بهترین عمل ممکن نگاهت کند. یادگیری تقویتی عمیق از شبکه های عصبی عمیق برای حل مسائل یادگیری تقویتی استفاده می کند، از این رو در نام آن از کلمه عمیق استفاده شده است. با در نظر گرفتن Q-Learning که یادگیری تقویتی کلاسیک محسوب می شود و Deep Q-Learning (DQN)^۲ می توان تفاوت آن ها با یکدیگر را دید. در رویکرد اول، از الگوریتم های سنتی برای ساخت جدول Q استفاده می شود تا به عامل در یافتن اقدامی که باید در هر حالت انجام شود کمک کند. در دومین رویکرد، از شبکه عصبی برای تخمین پاداش استفاده می شود. به عبارت دیگر، به جای استفاده از یک جدول اقدام-حالت که تعداد عناصر محدودی دارد، از یک شبکه عصبی آموزش دیده برای تخمین پاداش در هر حالت استفاده می گردد. نحوه عملکرد یادگیری تقویتی در مرحله آموزش در الگوریتم^۱ نشان داده شده است.

الگوریتم ۱. یادگیری تقویتی در مرحله آموزش

▪ تولید مقادیر اولیه: β

▪ برای هر اپیزود

به ازای هر زمان t

a. تعیین وضعیت S_t

b. {شیوه ϵ -greedy}

- با احتمال ϵ یک اقدام تصادفی انتخاب گردد

- در غیر این صورت (با احتمال $1 - \epsilon$) اقدام متناظر با بیشینه پاداش تجمعی انتخاب گردد:

$$A_t = \arg \max_a Q_P(S_t, a) = \arg \max_a E_P \left\{ \sum_{k=0}^{\infty} \beta^k R_{k+t+1} \mid S_t, a, P \right\}$$

c. اجرای اقدام A_t و مشاهده وضعیت بعدی S_{t+1} و محاسبه پاداش R_{t+1}

d. ذخیره $(S_t, A_t, R_{t+1}, S_{t+1})$ در حافظه D

e. {انتخاب تصادفی mini-batch} چند دسته (S_j, A_j, R_j, S_{j+1}) به عنوان نمونه از

D انتخاب می شود

$$f. \text{ محاسبه } y_j = \begin{cases} r_j & S_j = \text{terminal} \\ r_j + \beta \max_a Q_P(S_j, a) & S_j = \text{non-terminal} \end{cases}$$

g. آموزش شبکه DQN بر مبنای روش بهینه سازی RMSProp و با تابع هزینه

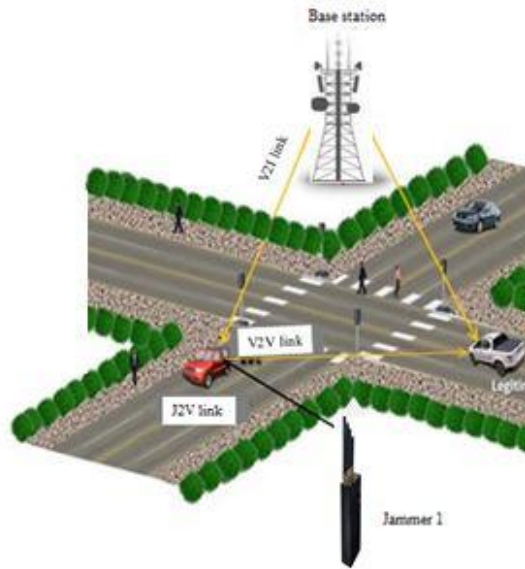
$$(y_j - Q_P(S_j, A_j))^2$$

¹ Policy

² Deep Q-learning

۳- مدل سیستم

شبکه کاربران و ایستگاه پایه در حضور تداخل گر در شکل ۲ نشان داده شده است.



شکل ۲- مدل سیستم

در این شکل، ما یک ارتباط بی سیم را در نظر می گیریم که شامل یک ایستگاه پایه، چند کاربر و یک اختلال گر است. هدف کاربران ارسال داده در محدوده فرکانسی کمتر درگیر شده با اختلال گر است. کاربرها در بازه های فرکانسی مشخص، داده خود را به ایستگاه پایه ارسال می کنند. تصمیم گیری در مورد تخصیص طیف مورد استفاده هر کاربر توسط ایستگاه پایه انجام می شود. ایستگاه پایه به عنوان عامل در یادگیری تقویتی می تواند اطلاعات مربوط به توان ارسالی و شرایط کانال کاربر را بررسی کرده و تصمیم خود را برای اقدام هر کاربر مبنی بر استفاده یا عدم استفاده از طیف مورد نظر به کاربران ارسال کند. الگوی رفتاری اختلال گر به وسیله اطلاعات به دست آمده از محیط بین جفت فرستنده و گیرنده به دست می آید. برای اعمال الگوریتم یادگیری تقویتی در این سناریو هر بازه فرکانسی یک حالت^۱ است و هر اقدام^۲ تخصیص یا عدم تخصیص فرکانس به یک کاربر به صورت باینری خواهد بود. لازم به ذکر است که هر بازه فرکانسی تنها می تواند به یک کاربر تخصیص پیدا کند. هدف این است که زیرباندهای فرکانسی را به نحوی بین کاربران توزیع کنیم که اثر تداخلگر را تا جای ممکن خنثی نماییم.

¹ state

² action

۴- روش پیشنهادی

در این بخش روشی برای جلوگیری از شنود اطلاعات توسط شنودگر پیشنهاد شده است که در آن نحوه تسهیم طیف به صورتی است که باندهای فرکانسی کمتر تخریب شده، برای ارسال داده کاربران شناسایی می‌شوند و به این ترتیب با مدیریت تسهیم فرکانس، با تخریب اطلاعات توسط تداخلگر مقابله می‌گردد. برای این منظور تکنیک انتخاب باند فرکانسی امن مبتنی بر یادگیری عمیق تقویتی پیشنهاد شده است. تابع پاداش در یادگیری تقویتی به صورت مجموع نرخ ارسال داده در یک باند فرکانسی مشخص تعریف می‌شود:

$$R(f_i) = \sum_{k=1}^K \rho_k C_{V2B}^k(f_i) \quad (1)$$

در این رابطه، ρ_k نشانگر باینری برای دسترسی یا عدم دسترسی کاربر k ام به زیر باند f_i است. اگر کاربر k ام به زیر باند f_i دسترسی داشته باشد، این ضریب برابر یک خواهد بود و در غیر این صورت مقدار آن صفر لحاظ می‌شود. ظرفیت $C_{V2B}^k(f_i)$ امین لینک $V2B$ در زیرباند f_i به صورت زیر محاسبه می‌شود:

$$C_{V2B}^k[f_i] = W \log(1 + \gamma_{V2B}^k[f_i]) \quad (2)$$

که در آن، W پهنای باند و $\gamma_{V2B}^k[f_i]$ نسبت سیگنال به تداخل و نویز^۱ ($SINR$) برای k امین لینک $V2B$ به صورت زیر نوشته می‌شود:

$$\gamma_{V2B}^k[f_i] = \frac{P_{V2B}[f_i] g_{V2B}^k}{I_k + I_J + \sigma^2} \quad (3)$$

در رابطه فوق، $P_{V2B}[f_i]$ توان یکسان مورد استفاده کاربران به عنوان فرستنده در لینک $V2B$ ، g_{V2B}^k ضریب کانال لینک $V2B$ مورد نظر در فرکانس f_i ، σ^2 توان نویز در زیر باند f_i می‌باشند، I_k تداخل ناشی از ارتباطات مربوط به لینک‌های دیگر کاربران با ایستگاه پایه و I_J تداخل ناشی از ارتباطات مربوط به لینک تداخلگر با کاربر در همین فرکانس است که به صورت زیر تعریف می‌شوند:

^۱ Signal to interference and noise ratio

$$I_k = \sum_{k' \neq k} \rho_{k'} P_{V2B}[f_i] g_{V2B}^k[f_i] \quad (4)$$

$$I_J = \int_{f_i - \frac{b_u}{2}}^{f_i + \frac{b_u}{2}} g_{J2V}[f_i] J(f - f_i) df \quad (5)$$

$J(f)$ چگالی طیفی توان تداخلگر و b_u پهنای باند سیگنال باند پایه تداخلگر را نشان می‌دهد. سه نوع تداخلگر نقطه‌ای^۱، تداخلگر جاروب^۲ و تداخلگر رگباری^۳ در این پژوهش در نظر گرفته شده‌است [۲۳]. تداخلگر نقطه‌ای شکلی از تداخلگر است که در آن یک تداخلگر تمام توان خود را بر روی یک فرکانس متمرکز می‌کند. تداخلگر جاروب توانایی انتقال قدرت کامل از یک فرکانس به فرکانس دیگر را دارد. این حرکت فراگیر چندین فرکانس را به سرعت و نه در یک زمان مسدود می‌کند. تداخلگر رگباری، چند فرکانس را به طور همزمان توسط یک تداخلگر واحد مورد حمله قرار می‌دهد. داده ارسالی از تداخلگر نقطه‌ای، جاروب و رگباری به ترتیب در روابط (۶-۸) تعریف شده است.

$$J(f) = a(f_0) \delta(f - f_0) \quad (6)$$

$$J(f) = \sum_{m=1}^M a(f_m) \delta(f - f_m) e^{-j2\pi f_m t_m} \quad (7)$$

$$J(f) = \sum_{n=1}^N \sum_{m=1}^M a(f_n) \delta(f - f_n) e^{-j2\pi f_n t_m} \quad (8)$$

که در آن $a(f_n)$ دامنه داده ارسالی توسط تداخلگر در فرکانس f_n و $\delta(f_n)$ تابع ضربه نشان‌دهنده فرکانس مورد حمله f_n است. $e^{-j2\pi f_n t_m}$ نشانگر تاخیر زمانی داده ارسالی توسط تداخلگر در فرکانس مورد نظر است. ضرایب کانال بین کاربرها و ایستگاه پایه بر مبنای محوشدگی مقیاس بزرگ (شامل تلفات مسیر^۴ و محوشدگی سایه^۵) و محوشدگی مقیاس کوچک برحسب زمان باید به‌روز رسانی شوند. محوشدگی مقیاس کوچک به صورت یک پارامتر وابسته به زمان با توزیع نمایی در نظر گرفته شده است. ضرایب کانال لینک $V2B$ به طریق زیر محاسبه می‌شود:

$$g_{V2B}[f_i] = \alpha_{V2B} h_{V2B}[f_i] \quad (9)$$

پارامتر $h_{V2B}[f_i]$ محوشدگی مقیاس بزرگ کانال در فرکانس f_i و کانال لینک $V2B$ است و پارامتر α_{V2B} محوشدگی مقیاس کوچک کانال در لینک مورد نظر و مستقل از فرکانس است.

در سیستم مورد نظر، برخی از زیر باندهای طیف فرکانسی موجود هستند که کاربرها خواستار ارسال داده در آن‌ها هستند. این باندها کم و بیش مورد حمله اختلال‌گر هستند. در روش پیشنهادی ابتدا یک داده راهنما از ایستگاه پایه به کاربرها ارسال می‌شود. داده راهنمای ارسالی در فرکانسهای مختلف و در گام زمانی t برابر است با:

$$X_t = (X_t^{f_1}, X_t^{f_2}, \dots, X_t^{f_T}) \quad (10)$$

¹ Spot jammer
² Sweep jammer
³ Barrage jammer
⁴ Path loss fading
⁵ Shadow fading

داده ارسالی توسط تداخلگر به صورت زیر تعریف می‌شود:

$$J_t = (J_t^{f_1}, J_t^{f_2}, \dots, J_t^{f_T}) \quad (11)$$

در این صورت داده دریافتی توسط کاربر برابر خواهد بود با:

$$Z_t = h_{B-U} X_t + h_{j-U} J_t + n_t \quad (12)$$

که در آن n_t نویز سفید گوسی با میانگین صفر است. کاربر مورد نظر پس از دریافت داده راهنما همان داده را به ایستگاه پایه برمی‌گرداند. داده‌های دریافتی توسط ایستگاه پایه $Y_t = (Y_t^{f_1}, Y_t^{f_2}, \dots, Y_t^{f_T})$ است که در آن علاوه بر داده ارسالی از BS (X_t)، عامل دیگر مربوط به تداخلگر نیز حضور دارد:

$$Y_t = h_{B-U}^T Z_t = h_{B-U}^T (h_{B-U} X_t + h_{j-U} J_t + n_t) = X_t + h_{B-U}^T h_{j-U} J_t + h_{B-U}^T n_t \quad (13)$$

در ایستگاه پایه، همبستگی میان X_t و Y_t برای باندهای مختلف فرکانسی برای هر کاربر محاسبه می‌شود:

$$\beta^k = [\rho_{X_t^{f_1}, Y_t^{f_1}}, \rho_{X_t^{f_2}, Y_t^{f_2}}, \dots, \rho_{X_t^{f_T}, Y_t^{f_T}}] \quad (14)$$

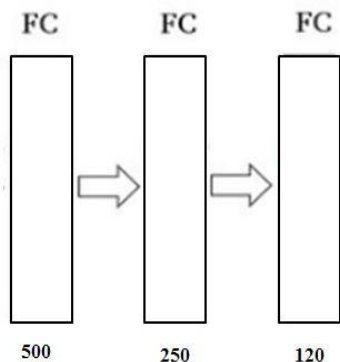
در این رابطه $\rho_{X_t^{f_i}, Y_t^{f_i}}$ نشان‌دهنده همبستگی میان $X_t^{f_i}$ و $Y_t^{f_i}$ در فرکانس f_i است. واضح است که هرچه میزان همبستگی نزدیک به عدد ۱ باشد نشان می‌دهد که داده راهنما کمتر تحت تاثیر عوامل اختلال‌گر قرار گرفته است و فرکانس مورد نظر صلاحیت بیشتری برای انتخاب توسط کاربر مورد نظر دارد و باید تاثیر بیشتری در پاداش در نظر گرفته شده در الگوریتم یادگیری تقویتی عمیق داشته باشد. ضرایب همبستگی محاسبه شده به عنوان ضریب تخفیف در الگوریتم یادگیری تقویتی مورد بهره‌برداری قرار می‌گیرند. بنابراین، مقدار پاداش تجمعی به ازای شروع اقدام از حالت S ، به ازای اقدام مشترک اولیه A تحت خط مشی P برابر خواهد بود با:

$$Q(t) = E\left\{\sum_{k=0}^{\infty} \beta^k(i) \rho_k R(f_i) \mid S_0 = S, A_0 = A, p\right\} \quad (15)$$

به این ترتیب، پاداش‌های آینده با ضریب تخفیف متغیر محاسبه شده به هم مربوط می‌شوند. روند انتخاب اقدام‌ها با طرح ϵ -greedy انجام شده است [۲۴]. طرح ϵ -greedy برای بهینه‌سازی تابع هزینه در یک اکتشاف اتخاذ می‌شود. به عبارت دیگر، با استفاده از ϵ -greedy، اقدام با احتمال $1-\epsilon$ با حداکثر پاداش انتخاب شده است، در حالی که یک اقدام تصادفی با احتمال ϵ انتخاب می‌شود. در پایان مرحله آموزش، جدولی از مقدار پاداش‌های تجمعی به ازای سیاست‌های اتخاذ شده (حالت‌های و اقدام‌های موجود در مسیر) تشکیل می‌شود تا در مرحله تست، به ازای سیاست‌های در نظر گرفته شده در محیط جدید، بهترین سیاست برای رسیدن به بیشترین پاداش تجمعی انتخاب شود.

$$P = \arg \max_p Q \quad (16)$$

تعداد اقدام‌ها و حالت‌ها در روش یافتن بهترین سیاست طبق رابطه (۱۶) می‌تواند هزاران نفر باشد، که مدیریت مقادیر Q را در جدول بسیار ناکارآمد می‌کند. اینجاست که می‌توانیم به جای استفاده از جدول، از شبکه‌های عصبی (DQN) برای پیش‌بینی مقادیر Q برای اقدامات در یک حالت معین استفاده کنیم. به عبارت دیگر، به جای نگاشت یک جفت حالت-اقدام به یک مقدار Q ، یک شبکه عصبی حالت‌های ورودی را به جفت عمل و Q نگاشت می‌کند. شبکه استفاده شده در این مقاله در شکل ۳ نشان داده شده است.



شکل ۳- شبکه عمیق استفاده شده به عنوان DON

همان‌طور که در شکل ۳ نشان داده شده است، شبکه DQN استفاده شده در این مقاله از سه لایه تمام متصل کننده^۲ با تعداد ۵۰۰، ۲۵۰ و ۱۲۰ نورون تشکیل شده است که از واحد خطی اصلاح شده^۳ (Relu) به عنوان تابع فعال سازی^۴ در رابطه (۱۷) و ریشه میانگین مربعات^۵ (RMSProp) استفاده می‌کند.

$$f(x) = \max(0, x) \quad (17)$$

در مرحله نهایی آموزش، سیاست‌های در نظر گرفته و Q متناظر آن‌ها به عنوان ورودی و خروجی به یک شبکه متشکل از سه لایه تمام جمع کننده با ۵۰۰، ۲۵۰ و ۱۲۰ نورون اعمال می‌شوند. شبکه آموزش یافته در مرحله تست، بهترین سیاست را برای داشتن بیشترین پاداش تعیین می‌کند. مراحل روش پیشنهادی مقابله با تداخلگر برای تسهیم طیف در الگوریتم ۲ نشان داده شده است. در این الگوریتم ابتدا محیط آموزش با تولید موقعیت، سرعت و ضرایب کانال تشکیل می‌شود. ضریب تخفیف با استفاده از دنباله راهنما و بازخورد آن برای هر کاربر محاسبه می‌شود. با اعمال الگوریتم یادگیری تقویتی و محاسبه پاداش جمعی به ازای تخصیص یا عدم تخصیص بازه‌های فرکانسی به کاربرها، شبکه DQN با تابع بهینه‌سازی انتشار ریشه میانگین مربعات آموزش داده می‌شود.

شبکه DQN آموزش یافته قادر خواهد بود بهترین سیاست از نظر تخصیص یا عدم تخصیص بازه‌های فرکانسی به کاربران را به ازای بیشینه کردن پاداش جمعی تعیین کند.

¹ Deep Q-network

² Fully-connected

³ Rectified linear unit

⁴ Activation function

⁵ Root Mean Squared Propagation (RMSProp)

الگوریتم ۲- روش پیشنهادی مقابله با تداخلگر برای تسهیم طیف

- تولید مقادیر اولیه: موقعیت، سرعت کاربرها، ضرایب کانال
 - ارسال دنباله راهنما: X_t
 - دریافت داده بازخورد شده از کاربران در ایستگاه پایه: Y_t
 - تعیین مقدار ضریب تخفیف با محاسبه همبستگی بین داده راهنما و داده بازخورد داده شده:
- $$\beta^k = [\rho_{X_t^{f1}, Y_t^{f1}}, \rho_{X_t^{f2}, Y_t^{f2}}, \dots, \rho_{X_t^{fT}, Y_t^{fT}}]$$
- مرحله آموزش:
- به ازای هر اپیزود
- به روزرسانی مکان وسایل نقلیه
 - به ازای هر زمان t
 - برای هر کاربر k
 - تعیین وضعیت $S_t^{(k)}$: انتخاب یک بازه فرکانسی برای کاربر
 - انتخاب اقدام $A_t^{(k)}$: (تخصیص بازه فرکانسی به کاربر) از بین اقدامها با ε -greedy
 - با احتمال ε یک اقدام تصادفی انتخاب گردد
 - در غیر این صورت (با احتمال $1 - \varepsilon$) اقدام متناظر با بیشینه پاداش تجمعی انتخاب گردد:
- $$A_t^{(k)} = \arg \max_a Q_p^{(k)}(S_t^{(k)}, a) = \arg \max_a E_P \left\{ \sum_{m=0}^{\infty} \beta^m R_{m+t+1}^{(k)} \mid S_t^{(k)}, a, P \right\}$$
- اجرای اقدام تعیین شده $A_t^{(k)}$ توسط هر وسیله نقلیه (تخصیص بازه فرکانسی به کاربر) و مشاهده وضعیت بعدی $S_{t+1}^{(k)}$ و محاسبه پاداش $R_{t+1}^{(k)}$ (ظرفیت کانال)
 - ذخیره $(S_t^{(k)}, A_t^{(k)}, R_t^{(k)}, S_{t+1}^{(k)})$ در حافظه D_k
 - {انتخاب تصادفی mini-batch} چند دسته $(S_j^{(k)}, A_j^{(k)}, R_j^{(k)}, S_{j+1}^{(k)})$ به عنوان نمونه از D انتخاب می شود
 - محاسبه
$$y_j^{(k)} = \begin{cases} r_j^{(k)} & S_j^{(k)} = \text{terminal} \\ r_j^{(k)} + \beta \max_a Q_p(S_j^{(k)}, a) & S_j^{(k)} = \text{non-terminal} \end{cases}$$
 - آموزش شبکه DQN بر مبنای روش بهینه سازی RMSProp و تابع هزینه $(y_j^{(k)} - Q_p(S_j^{(k)}, A_j^{(k)}))^2$
- مرحله تست:
- تولید محیط تست: موقعیت، سرعت کاربرها، ضرایب کانال
 - به دست آوردن سیاست بهینه P (تخصیص یا عدم تخصیص بازه های فرکانسی به کاربران) و Q بیشینه با استفاده از شبکه DQN

در الگوریتم ε -greedy یک مقدار برای ε تعیین می شود که در این مقاله ۰/۰۲ در نظر گرفته شده است.

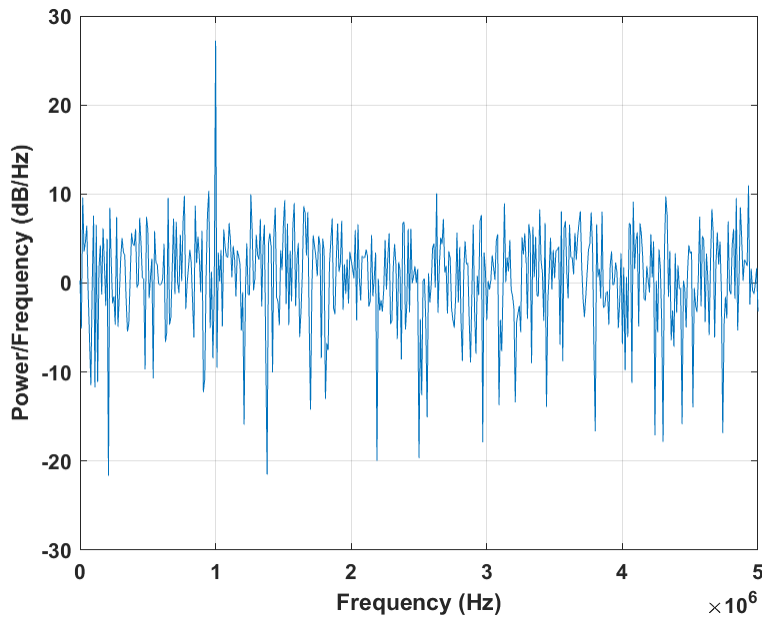
۵- نتایج شبیه‌سازی

در این بخش، نتایج شبیه‌سازی روش پیشنهادی ارائه شده است. شبکه ارتباطی وسایل نقلیه مطابق [۲۵]، مدل‌سازی شده است. پارامترهای به کار رفته در شبیه‌سازی در جدول ۱ آورده شده است.

جدول ۱- پارامترهای استفاده شده در شبیه‌سازی

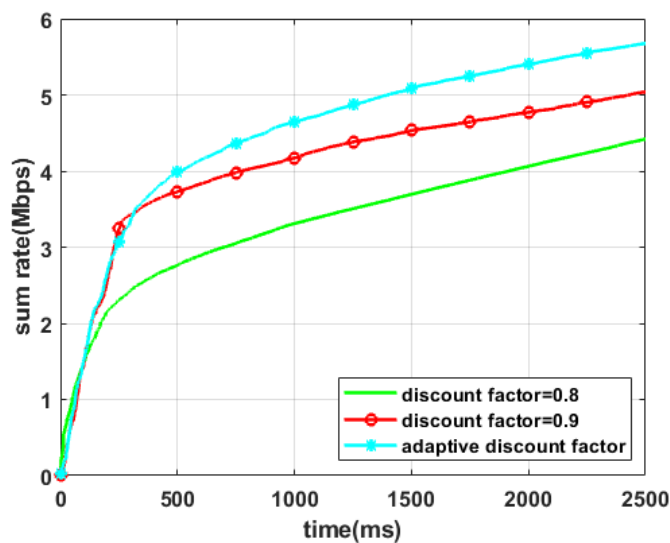
پارامتر	مقدار
تعداد کاربران	۴
پهنای باند	۴ مگاهرتز
توان نویز کانال	λ dBm
محدوده زمانی انتقال	۱۰۰ ms
توزیع محوشدگی سایه	توزیع لوگ نرمال
فاصله همبستگی	۵۰ m
توزیع محوشدگی سریع	توزیع رایلی
توان تداخلگر	۲۲ dBm
توان کاربر	۲۳ dBm
توان نویز	λ dBm
زمان به روز رسانی محوشدگی سریع	۱ ms
زمان به روز رسانی تلفات مسیر و محوشدگی سایه	۱۰۰ ms

در این پژوهش، تداخلگر نقطه‌ای با تولید تابع ضربه در یک فرکانس تصادفی و نیز تداخلگر جاروب به صورت اندازه ضرایب تبدیل فوریه مجموع تابع $\cos(2\pi ft)$ و شیفت یافته‌های آن پیاده‌سازی شده است. تداخلگر رگباری حالت کلی‌تر تداخلگرهای نقطه‌ای و جاروب است. چگالی طیفی توان در نظر گرفته شده برای تداخلگر رگباری مورد استفاده در این پژوهش، با استفاده از مدل موجود در **Array System Toolbox Phased** در برنامه **Matlab** پیاده‌سازی شده که در شکل ۴ نشان داده شده است.



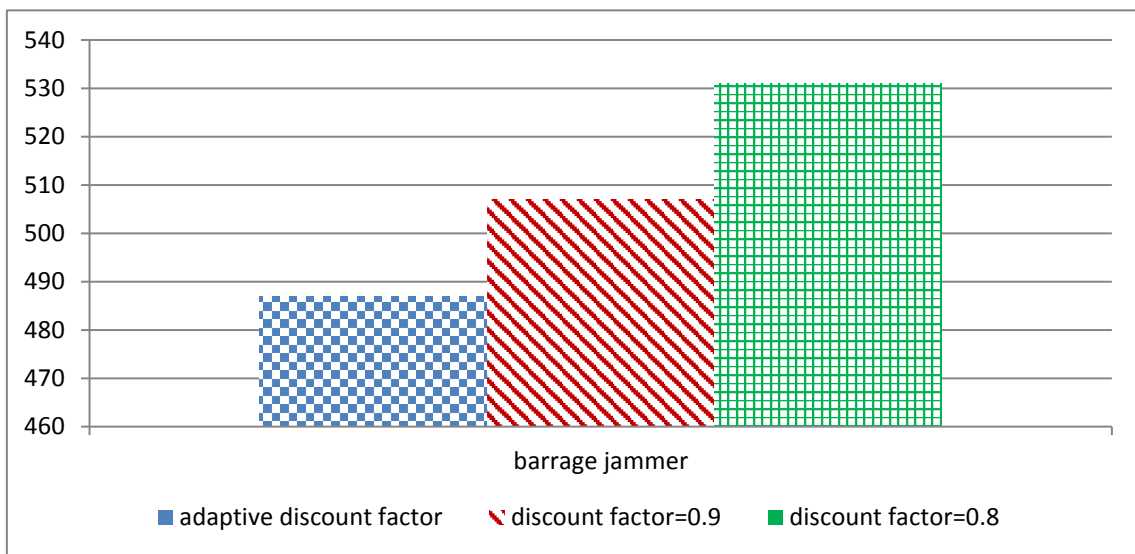
شکل ۴- چگالی طیفی توان در نظر گرفته شده برای تداخلگر رگباری

میانگین چگالی طیفی توان نشان داده شده در این نمودار ۲۲ dBm است. گفتنی است برای مقابله با تداخلگر در تسهیم طیف روشی در مقالات ارائه نشده است. از این رو عملکرد خود روش پیشنهادی را در سناریوهای مختلف مورد ارزیابی قرار می دهیم و نیز برای اینکه تکنیک پیشنهادی برای تنظیم ضریب تخفیف یافته را مورد ارزیابی قرار دهیم، از ضریب تخفیف های ثابت در روش پیشنهادی استفاده کرده ایم. نرخ مجموع بدست آمده در مرحله آموزش توسط روش پیشنهادی در حضور تداخلگر رگباری در شکل ۵ نشان داده شده است. توان نویز ۲۲ dBm و توان تداخلگر ۹ dBm است. در این شکل، روش پیشنهادی با ضریب تخفیف تطبیق یافته با روش پیشنهادی به ازای ضرایب تخفیف ۰/۸ و ۰/۹ مقایسه شده اند.



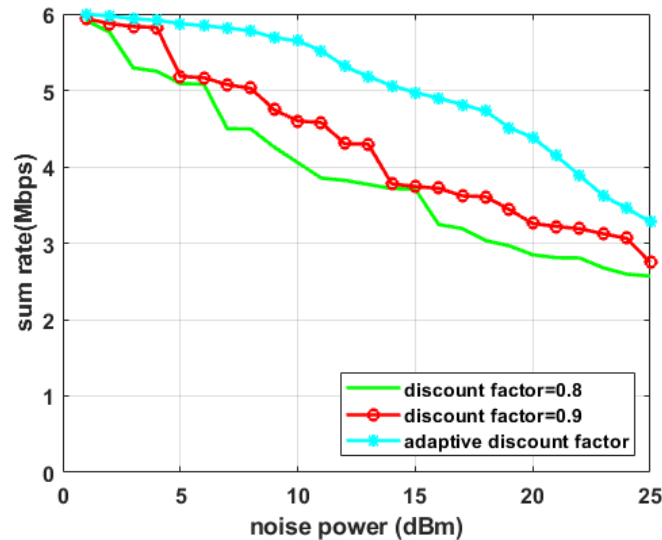
شکل ۵- نرخ مجموع بر حسب زمان به ازای ضریب تخفیف های متفاوت در حضور تداخلگر رگباری

مطابق شکل پاداش در شرایط ضریب تخفیف تطبیق یافته نسبت به ضریب تخفیف های ثابت همگرایی سریع تری دارد. دلیل این امر انتخاب طیف های کمتر تخریب شده توسط روش پیشنهادی است. در این شرایط کاربرانی که در طیف مورد نظر در معرض تداخلگر هستند، همبستگی کمتری بین سیگنال ارسالی و سیگنال بازخورد داده شده تجربه می کنند و این مقدار همبستگی به عنوان ضریب تخفیف تاثیر مشارکت این کاربران برای انتخاب به عنوان اقدام آینده در پاداش را کاهش می دهد. بنابراین همگرایی سریع تری با استفاده از ضریب تخفیف تطبیق یافته رخ می دهد. تعداد تکرار برای رسیدن به همگرایی در صورت وجود تداخلگر رگباری برای ضریب تخفیف ثابت 0.8 ، 0.9 و ضریب تخفیف تطبیق یافته در شکل ۶ مقایسه شده است.



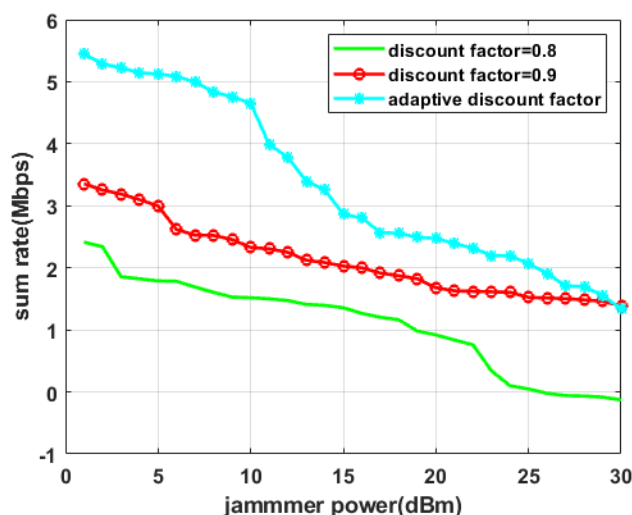
شکل ۶- تعداد تکرار برای همگرایی پاداش تجمعی در سه نوع تداخلگر

همانطور که در شکل ۶ مشخص شده است، روش تسهیم طیف امن با ضریب تخفیف تطبیق یافته تعداد تکرار کمتری را برای رسیدن به همگرایی به نرخ مطمئن مطلوب تجربه می کند. از داده های موجود در جدول ۲ می توان نتیجه گرفت که روش پیشنهادی با ضریب تخفیف تطبیق یافته نسبت به روش پیشنهادی با ضریب تخفیف ثابت، سرعت آموزش الگوریتم ی دارد. روند نرخ مجموع بر حسب توان نویز در حضور تداخلگر رگباری در شکل ۷ نشان داده شده است. توان تداخلگر برای رسم این نمودار ۹ dBm است.



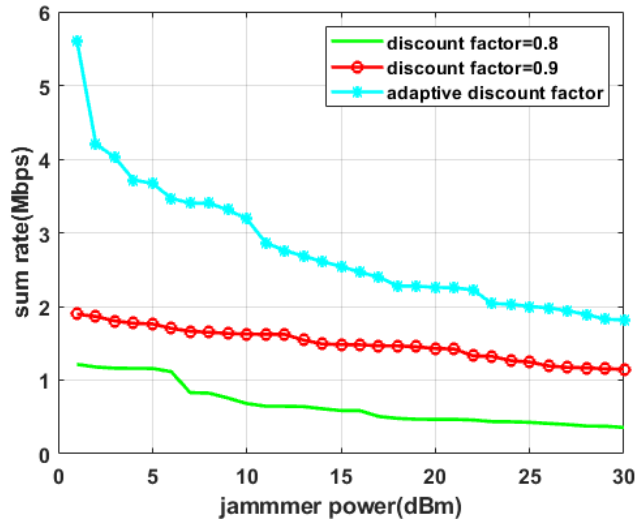
شکل ۷- نرخ مجموع بر حسب توان نویز در تداخلگر رگباری

همان طور که در شکل مشخص است، با افزایش توان نویز، نرخ مجموع به دلیل رابطه عکس با توان نویز کاهش پیدا می کند، اما در این شرایط، تسهیم طیف با ضریب تخفیف تطبیق یافته عملکرد بهتری نسبت به ضریب تخفیف ثابت دارد و نرخ مجموع ارائه شده توسط ضرایب تطبیق یافته بهتر از ضرایب ثابت است. با افزایش توان نویز سفید، مقدار همبستگی برای فرکانس هایی که نویز بیشتری دارند افزایش پیدا می کند. در این بخش می خواهیم نمودارهای نرخ مجموع بر حسب توان تداخلگر را برای انواع مختلف تداخلگر ترسیم کنیم. برای این منظور، ۱۰۰ الگوی تصادفی مختلف برای چگالی طیفی توان تولید شده است، سپس با تغییر توان تداخلگر، نرخ مجموع برای هر کدام از ۱۰۰ الگو محاسبه شده و متوسط نتایج در منحنی ها گزارش شده است. توان نویز در حضور هر سه تداخلگر ۲۲ dBm در نظر گرفته شده است. متوسط نرخ مجموع بر حسب توان های مختلف در حضور تداخلگر نقطه ای در شکل ۸ نمایش داده شده است.



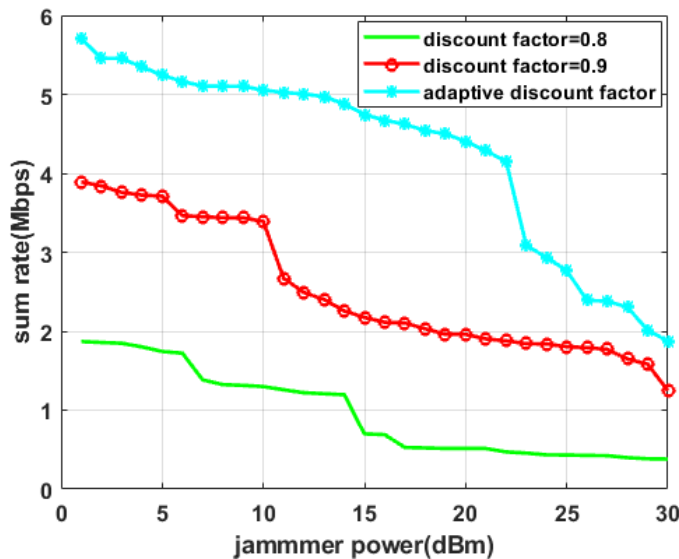
شکل ۸- نرخ مجموع به ازای توان تداخلگر نقطه ای

همان‌طور که در شکل ۸ نشان داده شده است، با افزایش توان تداخلگر نقطه‌ای، مقدار نرخ مجموع روند کاهشی دارد چرا که با افزایش توان تداخلگر نرخ مجموع کاهش می‌یابد. نرخ مجموع بر حسب توان‌های مختلف در حضور تداخلگر جاروب در شکل ۹ نمایش داده شده است.



شکل ۹- نرخ مجموع به ازای توان تداخلگر جاروب

با توجه به نمودار شکل ۹، هر سه نمودار روند نزولی دارند و روش تطبیق یافته نسبت به دو روش دیگر نرخ مجموع بالاتری در تمام توان‌های تداخلگر تولید می‌کند. نرخ مجموع بر حسب توان‌های مختلف در حضور تداخلگر رگباری در شکل ۱۰ نمایش داده شده است.



شکل ۱۰- نرخ مجموع به ازای توان تداخلگر رگباری

همان‌طور که در شکل ۱۰ نشان داده شده است، نرخ مجموع با افزایش توان تداخلگر در حضور تداخلگر رگباری، روند کاهشی دارد که این امر به دلیل تاثیر معکوس توان تداخلگر بر نرخ ارسال داده است. همچنین ملاحظه می‌کنیم که روش پیشنهادی با ضریب تطبیق یافته به نرخ مجموع بالاتری نسبت به ضریب تخفیف ثابت دست پیدا می‌کند. نکته قابل توجه در این نتایج آن است که نرخ مجموع ارائه شده توسط روش پیشنهادی با افزایش توان تداخلگر نیز قابل قبول است و به ویژه تکنیک ضریب تطبیق یافته سبب مقاومت مناسب روش پیشنهادی در برابر تداخلگر در توان‌های مختلف می‌شود.

۶- نتیجه‌گیری

این مقاله یک روش ضد تداخلگر پنهان بر اساس تسهیم فرکانس معرفی می‌کند که مکانیزم آن کاهش احتمال تخصیص زیرباند فرکانسی مورد حمله تداخلگر به کاربرها است. یک چارچوب یادگیری تقویتی چندعاملی طراحی شده است که با تجزیه و تحلیل همبستگی داده انتقالی بین ایستگاه پایه و کاربر، زیرباندهای کمتر تخریب شده توسط تداخلگر را به طور غیر مستقیم به دست می‌آورد. روش پیشنهادی بر اساس تکنیک‌های یادگیری عمیق تقویتی چند عاملی و وابسته ساختن مقدار ضریب تخفیف به میزان تخریب زیر باند فرکانسی، امکان دسترسی سریع‌تر کاربران به طیف‌های مناسب و در نتیجه مقابله با تداخلگر را فراهم می‌کند.

مراجع:

- [1] Y. Li *et al.*, "On the performance of deep reinforcement learning-based anti-jamming method confronting intelligent jammer," *Appl. Sci.*, vol. 9, no. 7, pp. 1–15, 2019, doi: 10.3390/app9071361.
- [2] Naparstek, Oshri, and Kobi Cohen. "Deep multi-user reinforcement learning for distributed dynamic spectrum access." *IEEE transactions on wireless communications* 18, no. 1 (2018): 310-323.
- [3] L. Jia, F. Yao, Y. Sun, Y. Niu, and Y. Zhu, "Bayesian Stackelberg Game for Antijamming Transmission With Incomplete Information," *IEEE Commun. Lett.*, vol. 20, no. 10, pp. 1991–1994, 2016, doi: 10.1109/LCOMM.2016.2598808.
- [4] M. K. Hanawal, M. J. Abdel-Rahman, and M. Krunz, "Joint Adaptation of Frequency Hopping and Transmission Rate for Anti-Jamming Wireless Systems," *IEEE Trans. Mob. Comput.*, vol. 15, no. 9, pp. 2247–2259, 2016, doi: 10.1109/TMC.2015.2492556.
- [5] L. Xiao, T. Chen, J. Liu, and H. Dai, "Anti-jamming transmission stackelberg game with observation errors," *IEEE Commun. Lett.*, vol. 19, no. 6, pp. 949–952, 2015, doi: 10.1109/LCOMM.2015.2418776.
- [6] H. Zhu, C. Fang, Y. Liu, C. Chen, M. Li, and X. S. Shen, "You Can Jam but You Cannot Hide: Defending Against Jamming Attacks for Geo-Location Database Driven Spectrum Sharing," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 10, pp. 2723–2737, 2016, doi: 10.1109/JSAC.2016.2605799.
- [7] Liu, Xin, Yuhua Xu, Luliang Jia, Qihui Wu, and Alagan Anpalagan. "Anti-jamming communications using spectrum waterfall: A deep reinforcement learning approach." *IEEE Communications Letters* 22, no. 5 (2018): 998-1001.
- [8] B. Wang, Y. Wu, K. J. R. Liu, and T. C. Clancy, "An anti-jamming stochastic game for cognitive radio networks," *IEEE J. Sel. Areas Commun.*, vol. 29, no. 4, pp. 877–889, 2011, doi: 10.1109/JSAC.2011.110418.
- [9] Y. Wu, B. Wang, K. J. R. Liu, and T. C. Clancy, "Anti-jamming games in multi-channel cognitive radio networks," *IEEE J. Sel. Areas Commun.*, vol. 30, no. 1, pp. 4–15, 2012, doi: 10.1109/JSAC.2012.120102.

- [10] S. Machuzak and S. K. Jayaweera, "Reinforcement learning based anti-jamming with wideband autonomous cognitive radios," *2016 IEEE/CIC Int. Conf. Commun. China, ICC 2016*, pp. 1–5, 2016, doi: 10.1109/ICCChina.2016.7636793.
- [11] A. P. Badia *et al.*, "Never Give Up: Learning Directed Exploration Strategies," pp. 1–28, 2020.
- [12] Y. Xiao, J. Hoffman, T. Xia, and C. Amato, "Learning Multi-Robot Decentralized Macro-Action-Based Policies via a Centralized Q-Net," *Proc. - IEEE Int. Conf. Robot. Autom.*, pp. 10695–10701, 2020, doi: 10.1109/ICRA40945.2020.9196684.
- [13] K. K. Nguyen, T. Q. Duong, N. A. Vien, N.-A. Le-Khac, and M.-N. Nguyen, "Non-Cooperative Energy Efficient Power Allocation Game in D2D Communication: A Multi-Agent Deep Reinforcement Learning Approach," *IEEE Access*, vol. 7, pp. 100480–100490, 2019, doi: 10.1109/access.2019.2930115.
- [14] H. Li, "Multiagent Q-learning for aloha-like spectrum access in cognitive radio systems," *Eurasip J. Wirel. Commun. Netw.*, vol. 2010, 2010, doi: 10.1155/2010/876216.
- [15] H. H. Chang, H. Song, Y. Zhang, H. He, and L. Liu, "Distributive dynamic spectrum access through deep reinforcement learning: A reservoir computing-based approach," *IEEE Internet Things J.*, vol. 6, no. 2, pp. 1938–1948, 2019, doi: 10.1109/JIOT.2018.2872441.
- [16] X. Chen, Z. Zhao, and H. Zhang, "Stochastic power adaptation with multiagent reinforcement learning for cognitive wireless mesh networks," *IEEE Trans. Mob. Comput.*, vol. 12, no. 11, pp. 2155–2166, 2013, doi: 10.1109/TMC.2012.178.
- [17] S. B. Janiar and V. Pourahmadi, "Deep-reinforcement learning for fair distributed dynamic spectrum access in wireless networks," *2021 IEEE 18th Annu. Consum. Commun. Netw. Conf. CCNC 2021*, 2021, doi: 10.1109/CCNC49032.2021.9369536.
- [18] V. François-Lavet, P. Henderson, R. Islam, M. G. Bellemare, and J. Pineau, *An introduction to deep reinforcement learning*, vol. 11, no. 3–4. 2018.
- [19] X. Zhu *et al.*, "Dynamic Spectrum Anti-Jamming with Reinforcement Learning Based on Value Function Approximation," *IEEE Wirel. Commun. Lett.*, pp. 1–5, 2022, doi: 10.1109/LWC.2022.3228045.
- [20] X. Liu, Y. Xu, L. Jia, Q. Wu, and A. Anpalagan, "Anti-jamming communications using spectrum waterfall: A deep reinforcement learning approach," *IEEE Commun. Lett.*, vol. 22, no. 5, pp. 998–1001, 2018, doi: 10.1109/LCOMM.2018.2815018.
- [21] X. Chen, C. Wang, Z. Zhou, and K. Ross, "Randomized Ensembled Double Q-Learning: Learning Fast Without a Model," pp. 1–25, 2021.
- [22] H. van Seijen, M. Fatemi, and A. Tavakoli, "Using a logarithmic mapping to enable lower discount factors in reinforcement learning," *Adv. Neural Inf. Process. Syst.*, vol. 32, no. NeurIPS, pp. 1–11, 2019.
- [23] X. Song, P. Willett, S. Zhou, and P. B. Luh, "The MIMO radar and Jammer games," *IEEE Trans. Signal Process.*, vol. 60, no. 2, pp. 687–699, 2012, doi: 10.1109/TSP.2011.2169251.
- [24] "Black, Paul E. 'greedy algorithm, Dictionary of Algorithms and Data Structures.' US Nat. Inst. Std. & Tech Report 88 (2012): 95.," vol. 88, p. 2012, 2012.
- [25] conformance specification Radio, User Equipment UE. "3rd Generation Partnership Project; Technical Specification Group Radio Access Network; Evolved Universal Terrestrial Radio Access (E-UTRA); User Equipment (UE) conformance specification Radio transmission and reception." (2011).