



پرديس علوم  
دانشکده ریاضی، آمار و علوم کامپیوتر

# به کارگیری یادگیری تقویتی عمیق برای آموزش عامل هوشمند جهت بازی Snake

نگارنده

رامتین راهوریان

استاد راهنما: دکتر هدیه ساجدی

پایان نامه برای دریافت درجه کارشناسی  
در رشته علوم کامپیوتر

بهمن ۱۴۰۰

## چکیده

یادگیری تقویتی شامل خانواده ای از الگوریتم ها است که برای کنترل و تصمیم گیری استفاده می شوند. این رویکردها مسائلی را حل می کنند که می توانند به عنوان فرایند تصمیم گیری مارکوف (MDP)<sup>۱</sup> بیان شوند. بدان معنی که مساله را باید با مجموعه ای از حالات  $S$  (به عنوان مثال بر اساس موقعیت مار)، مجموعه ای از اقدامات  $A$  (به عنوان مثال، بالا، پایین، راست، چپ)، یک عملکرد پاداش  $R$  (به عنوان مثال  $+10$  هنگامی که مار سیب می خورد،  $-10$  وقتی مارها به دیوار برخورد می کنند) و یک تابع گذار  $T$  که تغییرات بین حالت ها را توصیف می کند، بیان کرد. برای استفاده از یادگیری تقویتی، باید مساله را با استفاده از این چهار جزء فرموله کنیم. در این پایان نامه، این فرضیه بررسی می شود که کاربردهای یادگیری تقویتی برای انجام تجزیه و تحلیل دقیق مناسب نیستند. با ترکیب یادگیری عمیق تقویتی و شبکه عصبی، عامل بازی Snake آموزش داده می شود. وابستگی های این بازی به پایتون، کراس، تصویر کیت علمی و جریان تانسور است. برای آموزش عامل در بازی، از SARSA<sup>۲</sup> برای کنترل قابلیت های حرکتی مار طبق میل کاربر استفاده می شود. یادگیری عمیق Q در ماژول یادگیری استفاده می شود و نرخ یادگیری الگوریتم SARSA نتایج موثر را در تکرارهای مختلف عامل نشان می دهد.

---

Markov Decision Process<sup>۱</sup>  
State-action-reward-state-action<sup>۲</sup>

## پیشگفتار

با بزرگتر شدن و پیچیده تر شدن دنیای بازی، اطمینان از قابل بازی بودن و بدون اشکال بودن آنها برای توسعه دهندگان سخت تر می شود و شرکت های بازی به دنبال ابزارهای جدیدی از جمله هوش مصنوعی هستند تا به غلبه بر چالش های گسترده آزمایش محصولات خود کمک کنند. هوش مصنوعی و بازی برخلاف تصور عموم، با هم خوب کنار نمی آیند چرا که بین هوش مصنوعی و رفتار مصنوعی تفاوت وجود دارد. ما نمی خواهیم عوامل بازی های ما از بازیکنان پیشی بگیرند. ما می خواهیم آنها به همان اندازه که برای ارائه سرگرمی و تعامل لازم است باهوش باشند. ما نمی خواهیم از حداکثر محدودیت ربات ML<sup>۳</sup> خود استفاده کنیم، همانطور که معمولاً در صنایع مختلف نیز اینگونه انجام می دهیم. حریف نباید کامل باشد و از رفتاری شبیه به انسان تقلید کند. از آنجایی که با وجود چالش برانگیز بودن بازی های ویدیویی به راحتی قابل رسمی کردن هستند، برای مدت طولانی یکی از حوزه های محبوب تحقیقات هوش مصنوعی بوده است. برای دهه ها، توسعه دهندگان بازی سعی کرده اند عاملی را با هوش شبیه سازی بسازند، به ویژه برای ساخت بازیکنی با هوش مصنوعی که بتواند بازی را بر اساس تجربه بازی اش یاد بگیرد، نه اینکه صرفاً از یک استراتژی ثابت پیروی کند. برنامه نویسی پویا می تواند مشکل را با تعداد کمی از حالات و ساختار تصادفی زیرین ساده حل کند، اما برای حالت پیچیده کارآمد نمی باشد. البته بازی ها فقط سرگرمی نیستند. آموزش یک عامل مجازی برای عملکرد بهتر از بازیکنان انسانی می تواند به ما بیاموزد که چگونه فرایندهای مختلف را در زمینه های فرعی مختلف و هیجان انگیز بهینه کنیم. در این پایان نامه، یک عامل هوش مصنوعی بررسی می شود که قادر به یادگیری نحوه بازی کردن بازی محبوب Snake از ابتدا است.

# فهرست مطالب

۱	معرفي و پيش زمينه	۱
۱	۱.۱ معرفي	۱
۲	۲.۱ يادگيري ماشيني	۲
۳	۳.۱ يادگيري تقويتي	۳
۷	۴.۱ يادگيري $Q$	۷
۷	۱.۴.۱ الگوريتم يادگيري $Q$	۷
۸	۲.۴.۱ يادگيري $Q$ عميق	۸
۱۳	پيشينه پژوهش	۲
۱۵	۱.۲ بازي مار و يادگيري تقويتي	۱۵
۱۶	۲.۲ آناليز رياضي بازي مار	۱۶
۱۶	۳.۲ الگوريتم تقريبي مار	۱۶
۱۷	۴.۲ کاهش فضاى حالت براى يادگيري تقويتي	۱۷
۱۹	۳ پياده سازى بازي	۱۹
	۱.۳ الگوريتم State-Action-Reward-State-Action (SARSA) و يادگيري	
۲۱	تقويتي	۲۱
۲۳	۴ نتايج	۲۳
۲۳	۱.۴ پارامترهاي تنظيم	۲۳
۲۴	۲.۴ منحنى يادگيري	۲۴
۲۶	۴.۴ مقايسه عملکرد	۲۶

۲۶	.....	مقایسه عملکرد	۳.۴
۲۷	.....	نتیجه گیری	۵.۴

# فصل ۱

## معرفی و پیش زمینه

### ۱.۱ معرفی

انجام بازی های ویدیویی چالش برانگیز و سرگرم کننده است. بازی ها به عنوان جالب ترین حوزه هوش مصنوعی (AI)<sup>۱</sup> در نظر گرفته می شوند. با هوش شبیه سازی، توسعه دهندگان بازی در حال آموزش عوامل هستند و بازی هایی را برای کمک به یادگیری بسیاری از مفاهیم جدید ارائه می دهند. استراتژی های مختلفی دنبال می شود و از زبان های برنامه پویا در توسعه بازی ها استفاده می شود. ترکیبی از شبکه های عصبی و یادگیری عمیق برای فرمول بندی بازی ها استفاده می شود. چنین ترکیبی فرصت هایی را برای درک واضح چیزها و یادگیری موثر رفتار تقویتی فراهم می کند. این روزها ترکیبی از یادگیری ماشینی و یادگیری تقویتی در توسعه بازی ها استفاده می شود. عوامل آموزش می بینند، و تفسیرها برای آزمایش کارایی بازی گرفته می شود. قوانین خاصی در بازی طراحی شده و عوامل بازی برای کمک به حل مشکلات خاص آموزش می بینند. در یادگیری ماشینی، یادگیری تقویتی یکی از جذاب ترین فناوری هاست که در آن اقدامات بهینه در هر حالتی با آزمون و خطا آموخته می شود. این فرآیند آزمون و خطا به حل موثر مشکلات کمک می کند و دیدگاه منحصر به فردی از جنبه بررسی عملکرد یادگیری تقویتی ارائه می دهد.

در این پایان نامه بازی Snake برای بررسی حالات مختلف یادگیری تقویتی اصلاح شده است. از یادگیری عمیق برای آموزش حالت های بازی و ارزیابی عملکرد مار در بازی استفاده می شود. هدف اصلی این پایان نامه آموزش عامل بازی Snake با استفاده از یادگیری عمیق و هوش مصنوعی

---

<sup>۱</sup>Artificial Intelligence

است. اجزا جدید اضافه شده به تجزیه و تحلیل بازی، کار اصلی یادگیری تقویتی است. با تلفن های همراه نوکیا، بازی مار به یکی از محبوب ترین بازی ها تبدیل شده است. بازی Snake نوکیا به گونه ای طراحی شده است که یک بازیکن می تواند وضعیت کار مار را کنترل کند. مراحل مختلفی در بازی طراحی شده است تا به تدریج بر پیچیدگی آن افزوده شود. برخی از قوانین باید در بازی رعایت شود. به عنوان مثال، مار سعی می کند چیزی را بخورد و به دنبال غذا می دود. خوردن مواد غذایی مختلف باعث می شود اندازه مار بزرگتر شود و در نتیجه امتیاز جمع شود. با افزایش اندازه مار، پیچیدگی بازی نیز افزایش می یابد، بنابراین کسب امتیاز بالا دشوارتر می شود.

در بازی Snake که تغییر یافته است، آب نبات های تصادفی ظاهر می شوند و مار هر بار که یک آب نبات می خورد رشد می کند. تفاوت این است که یک آب نبات ها سمی است. اگر مار آب نبات سمی را بخورد، می میرد. امتیاز نهایی به اندازه مار در زمان مرگ بستگی دارد یا مانند حالت کلاسیک بازی، امتیازات با تعداد آب نبات های خورده شده توسط مار محاسبه می شود. الگوریتم Q-learning با استفاده از Keras به عنوان مدل شبکه عصبی وارد بازی شده است. استفاده از یادگیری عمیق در طراحی جدید Game Snake قابلیت های بالقوه Q-learning را افزایش می دهد. این بازی در Pygame طراحی شده است که به مدیریت موثر کد کمک می کند تا عوامل بتوانند برای ارائه کنترل کامل به کاربر آموزش ببینند. SARSA در آموزش عواملی که منجر به کنترل کاربر برای حرکت قابلیت های مار می شود استفاده شد. الگوریتم SARSA از الگوریتم on-policy و استفاده می کند و برای انجام عمل بر مبنای Q-value است. الگوریتم SARSA به جای استفاده از خط مشی حریصانه، از خط مشی متداول استفاده می کند. با استفاده از الگوریتم on-policy عملکرد بازی موثرتر و اجرای فرآیند سریعتر می شود.

## ۲.۱ یادگیری ماشینی

یادگیری ماشینی (ML) نوعی هوش مصنوعی (AI) است که به برنامه های نرم افزاری اجازه می دهد تا در پیش بینی نتایج دقیق تر شوند، بدون اینکه به صراحت برای این کار برنامه ریزی شده باشند. الگوریتم های یادگیری ماشین از داده های گذشته به عنوان ورودی برای پیش بینی مقادیر خروجی جدید استفاده می کنند. موتورهای توصیه یک مورد رایج برای یادگیری ماشین هستند (سیستم های توصیه کننده کلاس مهمی از الگوریتم های یادگیری ماشین هستند که پیشنهادات «مرتبط» را به کاربران ارائه می دهند). سایر کاربردهای محبوب عبارتند از: تشخیص تقلب، فیلتر هرزنامه

(اسپم)، شناسایی تهدید بدافزار، اتوماسیون فرآیند کسب و کار (BPA)<sup>۲</sup> و پیش‌بینی تعمیرات. یادگیری ماشینی مهم است زیرا به شرکت‌ها دیدی از روند رفتار مشتری و الگوهای عملیاتی تجاری می‌دهد و همچنین از توسعه محصولات جدید پشتیبانی می‌کند. بسیاری از شرکت‌های پیشرو امروزی، مانند فیس‌بوک، گوگل و اوپن، یادگیری ماشینی را به بخش مرکزی عملیات خود تبدیل کرده‌اند. یادگیری ماشینی به یک تمایز رقابتی مهم برای بسیاری از شرکت‌ها تبدیل شده است.

یادگیری ماشینی کلاسیک اغلب بر اساس نحوه یادگیری الگوریتم در پیش‌بینی دقیق‌تر، طبقه‌بندی می‌شود. چهار رویکرد اساسی وجود دارد: یادگیری با نظارت، یادگیری بدون نظارت، یادگیری نیمه نظارتی و یادگیری تقویتی. نوع الگوریتمی که دانشمندان داده‌های الگوریتمی انتخاب می‌کنند، بستگی به نوع داده‌هایی دارد که می‌خواهند پیش‌بینی کنند.

### ۳.۱ یادگیری تقویتی

یادگیری تقویتی (RL)<sup>۳</sup> حوزه‌ای از یادگیری ماشین است. این حوزه این‌که چگونه عوامل هوشمند باید در یک محیط اقداماتی را انجام دهند تا مفهوم پاداش تجمعی را به حداکثر برسانند، مربوط می‌شود. یادگیری تقویتی یکی از سه نمونه اصلی یادگیری ماشین در کنار یادگیری تحت نظارت و یادگیری بدون نظارت است.

یادگیری تقویتی در عدم نیاز به ارائه جفت‌های ورودی/خروجی برجسب دار و عدم نیاز به اصلاح صریح اقدامات زیر بهینه (اپتیمال) با یادگیری نظارت شده متفاوت است. در عوض تمرکز بر یافتن تعادل بین شناسایی (چیزهای ناشناخته) و بهره‌برداری (دانش فعلی) است [۱]. الگوریتم‌های RL با نظارت جزئی می‌توانند مزایای الگوریتم‌های نظارت شده و RL را ترکیب کنند [۲].

محیط معمولاً در قالب فرآیند تصمیم‌گیری مارکوف (MDP) بیان می‌شود، زیرا بسیاری از الگوریتم‌های یادگیری تقویتی برای این متن از تکنیک‌های برنامه‌نویسی پویا استفاده می‌کنند. تفاوت اصلی بین روش‌های برنامه‌نویسی پویا کلاسیک و الگوریتم‌های یادگیری تقویتی در این است که دومی دانش یک مدل ریاضی دقیق از MDP را فرض نمی‌کند و های MDP بزرگ را در جایی که روش‌های دقیق غیرممکن می‌شوند هدف قرار می‌دهند.

---

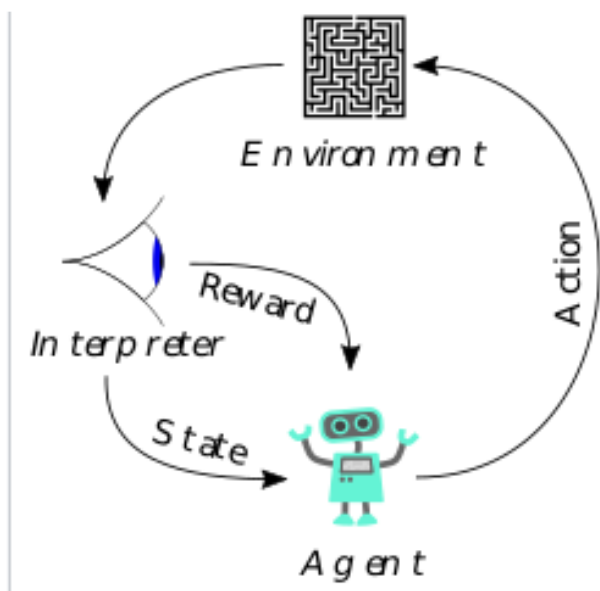
<sup>۲</sup> Business Process Automation  
<sup>۳</sup> Reinforcement Learning



به دلیل عمومیت آن، یادگیری تقویتی در بسیاری از رشته‌ها مانند نظریه بازی، نظریه کنترل، تحقیق در عملیات، نظریه اطلاعات، بهینه‌سازی مبتنی بر شبیه‌سازی، سیستم‌های چند عاملی، هوش گروه و آمار مورد مطالعه قرار می‌گیرد. در ادبیات تحقیق و کنترل عملیات، یادگیری تقویتی، برنامه‌ریزی پویا تقریبی یا برنامه‌ریزی عصبی پویا نامیده می‌شود. مسائل مورد علاقه در یادگیری تقویتی نیز در تئوری کنترل بهینه بررسی شده است. که بیشتر به وجود و توصیف راه‌حل‌های بهینه و الگوریتم‌هایی برای محاسبه دقیق آن‌ها می‌پردازد و کمتر به یادگیری یا تقریب، به‌ویژه در غیاب یک مدل ریاضی از محیط می‌پردازد. در اقتصاد و تئوری بازی‌ها، یادگیری تقویتی ممکن است برای توضیح چگونگی ایجاد تعادل در عقلانیت محدود استفاده شود.

تقویت اساسی به عنوان یک فرآیند تصمیم‌گیری مارکوف (MDP) مدل می‌شود:

- مجموعه‌ای از حالت‌های محیط و عامل،  $S$ .
  - مجموعه‌ای از اقدامات،  $A$ ، از عامل.
  - $P_a(s, s') = \Pr(s_{t+1} = s' | s_t = s, a_t = a)$  احتمال انتقال (در زمان  $t$ ) از حالت  $s$  به حالت  $s'$  تحت عمل  $a$  است.
- $R$  پاداش فوری پس از انتقال از  $s$  به  $s'$  با عمل  $a$  است.



شکل ۱.۱: چارچوب بندی معمول سناریوی یادگیری تقویتی (RL): یک عامل اقداماتی را در یک محیط انجام می دهد که به پاداش و نمایش وضعیت تفسیر می شود که به عامل بازخورد داده می شود.

هدف از یادگیری تقویتی این است که عامل یک خط مشی بهینه یا تقریباً بهینه را بیاموزد که ”عملکرد پاداش” یا سایر سیگنال های تقویتی ارائه شده توسط کاربر را که از پاداش های فوری جمع می شود، به حداکثر می رساند. این شبیه به فرآیندهایی است که به نظر می رسد در روانشناسی حیوانات رخ می دهد. به عنوان مثال، مغزهای بیولوژیکی برای تفسیر سیگنال هایی مانند درد و گرسنگی به عنوان تقویت کننده های منفی، و لذت و دریافت غذا به عنوان تقویت کننده های مثبت طراحی شده اند. در برخی شرایط، حیوانات می توانند یاد بگیرند که در رفتارهایی شرکت کنند که این پاداش ها را بهینه می کند. این نشان می دهد که حیوانات قادر به یادگیری تقویتی هستند.

یک عامل یادگیری تقویتی اساسی هوش مصنوعی با محیط خود در مراحل زمانی گسسته تعامل دارد. در هر زمان  $t$ ، عامل وضعیت فعلی  $s_t$  و پاداش  $r_t$  را دریافت می کند. سپس یک اقدام  $a_t$  را از مجموعه اقدامات موجود انتخاب می کند که متعاقباً به محیط ارسال می شود. محیط به حالت جدید  $s_{t+1}$  حرکت می کند و پاداش  $r_{t+1}$  مرتبط با انتقال  $(s_t, a_t, s_{t+1})$  تعیین می شود. هدف یک عامل یادگیری تقویتی یادگیری یک خط مشی است:

که  $\pi : A \times S \rightarrow [0, 1], \pi(a, s) = \Pr(a_t = a | s_t = s)$  پاداش تجمعی مورد انتظار را به حداکثر می رساند.

فرمول بندی مسئله به عنوان یک MDP فرض می کند که عامل مستقیماً وضعیت محیطی فعلی را مشاهده می کند. در این مورد گفته می شود که مشکل قابل مشاهده کامل است. اگر عامل فقط به زیرمجموعه ای از حالت ها دسترسی داشته باشد، یا اگر حالت های مشاهده شده توسط نویز خراب شده باشند، گفته می شود که عامل دارای قابلیت مشاهده جزئی است و به طور رسمی مشکل باید به عنوان فرآیند تصمیم گیری مارکوف تا حدی قابل مشاهده فرموله شود. در هر دو مورد، مجموعه اقدامات در دسترس برای عامل می تواند محدود شود. به عنوان مثال، وضعیت موجودی حساب می تواند مثبت باشد. اگر مقدار فعلی حالت ۳ باشد و انتقال حالت سعی کند مقدار را تا ۴ کاهش دهد، انتقال مجاز نخواهد بود.

وقتی عملکرد عامل با عملکردی که بهینه عمل می کند مقایسه می شود، تفاوت در عملکرد باعث ایجاد مفهوم پشیمانی می شود. به منظور عمل تقریباً بهینه، عامل باید در مورد پیامدهای بلندمدت اقدامات خود (یعنی به حداکثر رساندن درآمد آینده) استدلال کند، اگرچه پاداش فوری مرتبط با این مسئله، ممکن است منفی باشد.

بنابراین، یادگیری تقویتی به ویژه برای مشکلاتی که شامل مبادله پاداش بلندمدت در مقابل کوتاه مدت است، مناسب است. با موضوع موفقیت برای مشکلات مختلف، از جمله کنترل ربات، برنامه ریزی آسانسور، مخابرات، تخته نرد، checkers و Go (AlphaGo) اعمال شده است.

دو عنصر یادگیری تقویتی را قدرتمند می کند: استفاده از نمونه ها برای بهینه سازی عملکرد و استفاده از تقریب تابع برای مقابله با محیط های بزرگ. به لطف این دو جزء کلیدی، یادگیری تقویتی را می توان در محیط های بزرگ در شرایط زیر استفاده کرد:

- مدلی از محیط شناخته شده است، اما یک راه حل تحلیلی در دسترس نیست.
- فقط یک مدل شبیه سازی از محیط ارائه شده است (موضوع بهینه سازی مبتنی بر شبیه سازی)؛
- تنها راه جمع آوری اطلاعات در مورد محیط، تعامل با آن است.

دو مورد اول از این مسائل را می توان مسائل برنامه ریزی در نظر گرفت (از آنجایی که نوعی مدل در دسترس است)، در حالی که آخرین مورد را می توان یک مساله یادگیری واقعی در نظر گرفت. با این حال، یادگیری تقویتی هر دو مساله برنامه ریزی را به مسائل یادگیری ماشین تبدیل می کند.

## ۴.۱ یادگیری Q

یادگیری Q یک الگوریتم یادگیری تقویتی بدون مدل برای یادگیری ارزش یک عمل در یک حالت خاص است. نیازی به مدلی از محیط ندارد (بنابراین «بدون مدل»)، و می تواند مشکلات مربوط به تحولات و پاداش های تصادفی را بدون نیاز به سازگاری حل کند.

برای هر فرآیند تصمیم گیری محدود مارکوف (FMDP)<sup>۴</sup>، یادگیری Q یک خط مشی بهینه به معنای به حداکثر رساندن ارزش مورد انتظار کل پاداش در تمام مراحل متوالی، با شروع از وضعیت فعلی، پیدا می کند. یادگیری Q می تواند یک خط مشی انتخاب عمل بهینه را، برای هر FMDP مشخص، با توجه به زمان شناسایی نامحدود و یک خط مشی نیمه تصادفی شناسایی کند. "Q" به تابعی اشاره دارد که الگوریتم پاداش های مورد انتظار برای یک اقدام انجام شده در یک وضعیت خاص را محاسبه می کند.

### ۱.۴.۱ الگوریتم یادگیری Q

در اینجا مدل مسئله تشکیل شده از یک عامل، وضعیت ها  $S$  و مجموعه از اقدامات  $A$  برای هر وضعیت. با انجام یک اقدام  $a \in A$ ، عامل از یک وضعیت به وضعیت بعدی حرکت کرده و هر

<sup>۴</sup> Factored Markov Decision Processes

وضعیت پاداشی به عامل می‌دهد. هدف عامل حداکثر کردن پاداش دریافتی کل خود است. این کار با یادگیری اقدام بهینه برای هر وضعیت انجام می‌گردد. الگوریتم دارای تابعی است که ترکیب حالت اقدام را محاسبه می‌نماید:

$$Q : S \times A \rightarrow \mathbb{R}$$

قبل از شروع یادگیری،  $Q$  مقدار ثابتی را که توسط طراح انتخاب شده برمی‌گرداند. سپس هر بار که به عامل پاداش داده می‌شود، مقادیر جدیدی برای هر ترکیب وضعیت اقدام محاسبه می‌گردد. هسته الگوریتم از یک بروز رسانی تکراری ساده تشکیل شده‌است. به این ترتیب که بر اساس اطلاعات جدید مقادیر قبلی اصلاح می‌شود.

$$Q(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha_t(s_t, a_t)}_{\text{learning rate}} \times \left[ \underbrace{R(s_t)}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \underbrace{\max_{a_{t+1}} Q(s_{t+1}, a_{t+1})}_{\text{max future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right]$$

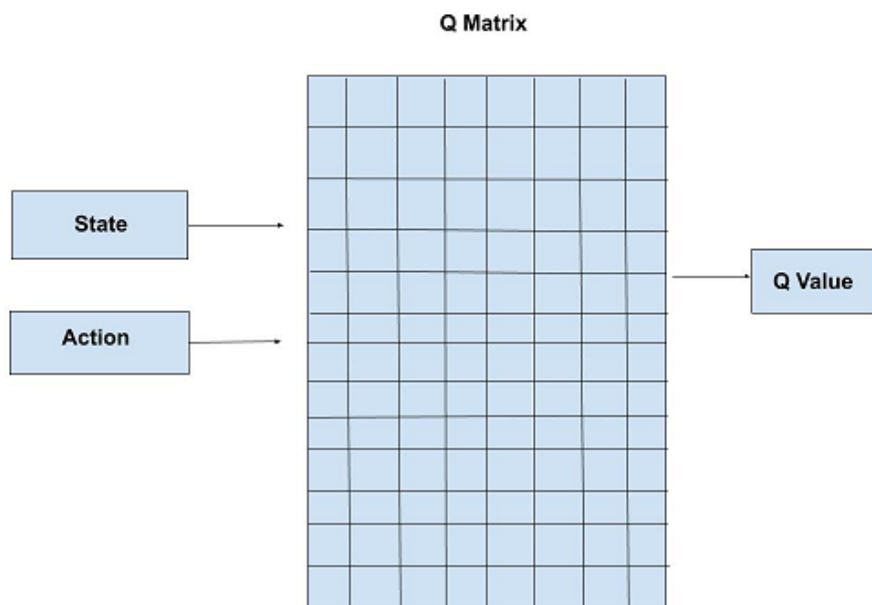
که  $R(s_t)$  پاداش  $s_t$  و  $\alpha_t(s, a)$  است. نرخ یادگیری  $(0 < \alpha \leq 1)$  ممکن است برای همه زوج‌ها یکسان باشد. مقدار عامل تخفیف  $\gamma$  بگونه است که  $0 \leq \gamma < 1$ . فرمول فوق معادل عبارت زیر است:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t)(1 - \alpha_t(s_t, a_t)) + \alpha_t(s_t, a_t)[R(s_t) + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})]$$

یک اپیزود الگوریتم وقتی  $s_{t+1}$  به وضعیت نهایی می‌سد پایان می‌یابد. توجه کنید که برای همه وضعیت‌های نهایی  $s_f$  و  $Q(s_f, a)$  مربوطه هیچگاه بروز نمی‌شود و مقدار اولیه خود را حفظ می‌کند.

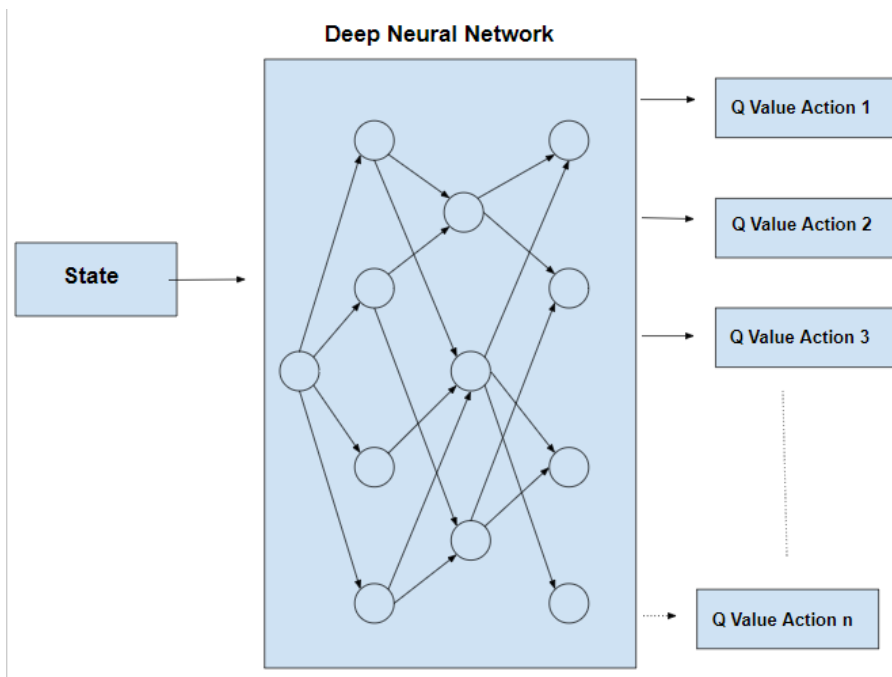
## ۲.۴.۱ یادگیری Q عمیق

فرآیند یادگیری Q یک ماتریس دقیق برای عامل کار ایجاد می‌کند که می‌تواند به آن رجوع کند تا پاداش خود را در دراز مدت به حداکثر برساند. اگرچه این رویکرد به خودی خود اشتباه نیست، اما فقط برای محیط‌های بسیار کوچک عملی است و با افزایش تعداد حالت‌ها و اقدامات در



شکل ۲.۱: در این شکل State و Action به Q-matrix داده می شود و Q-value بدست می آید

محیط، به سرعت ممکن بودن خود را از دست می دهد. راه حل مشکل فوق از درک این موضوع ناشی می شود که مقادیر موجود در ماتریس فقط اهمیت نسبی دارند، یعنی مقادیر فقط با توجه به سایر مقادیر اهمیت دارند. بنابراین، این تفکر ما را به یادگیری Q عمیق هدایت می کند که از یک شبکه عصبی عمیق برای تقریب مقادیر استفاده می کند. تا زمانی که اهمیت نسبی حفظ شود، این تقریب مقادیر ضرری ندارد. مرحله کار اساسی برای یادگیری Q عمیق این است که حالت اولیه به شبکه عصبی وارد می شود و Q-value تمام اقدامات ممکن را در خروجی برمی گرداند. تفاوت بین یادگیری Q و یادگیری Q عمیق را می توان به صورت زیر نشان داد:



شکل ۳.۱: در این شکل State به Deep Neural Network داده می شود و تعدادی Action Q-value بدست می آید

کد شبه برنامه:

```

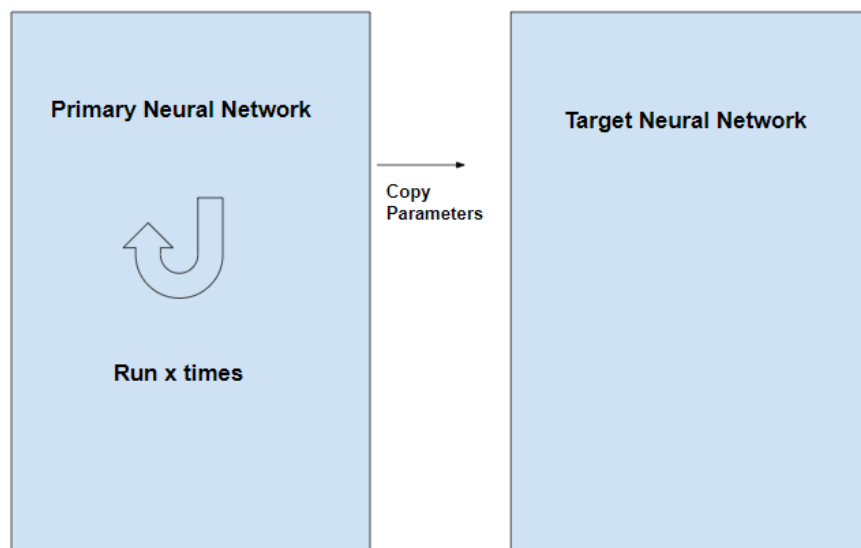
Initialize  $Q_0(s, a)$  for all pairs  $(s, a)$ 
s = initial state
k = 0
while(convergence is not achieved)
{
    simulate action a and reach state s'
    if(s' is a terminal state)
    {
        target = R(s, a, s')
    }
    else
    {
        target = R(s, a, s') +  $\gamma \max_{a'} Q_k(s', a')$ 
    }
     $\theta_{k+1} = \theta_k - \alpha \Delta_{\theta} E_{s' \sim P(s'|s, a)} [(Q_{\theta}(s, a) - \text{target}(s'))^2] |_{\theta=\theta_k}$ 
    s = s'
}

```

#### شکل ۴.۱: شبه کد برای Q-Learning Deep

توجه داشته باشید که در معادله هدف  $R(s, a, s') + \gamma \max_{a'} Q_k(s', a')$  عبارت  $\gamma \max_{a'} Q_k(s', a')$  یک عبارت متغیر است. بنابراین در این فرآیند، هدف برای شبکه عصبی بر خلاف سایر فرآیندهای یادگیری عمیق معمولی که در آن هدف ثابت است، متغیر است. این مساله با داشتن دو شبکه عصبی به جای یکی برطرف می شود. یک شبکه عصبی برای تنظیم پارامترهای شبکه و دیگری برای محاسبه هدف استفاده می شود که دارای معماری مشابه با شبکه اول است اما دارای پارامترهای ثابت است.





شکل ۵.۱: پس از تعداد  $x$  تکرار در شبکه اصلی، پارامترها در شبکه هدف کپی می شوند.

## فصل ۲

### پیشینه پژوهش

کارهای قابل توجهی در زمینه هوش مصنوعی و یادگیری ماشین انجام شده است. تکنیک ها و روش های متعددی برای کمک به حل مسائل محاسباتی پیشنهاد شده است. با استفاده از یادگیری تقویتی، بازی های مختلفی برای حل مسائل چالش برانگیز و ارائه درک عمیق از مدل های عصبی مختلف اجرا می شوند. با استفاده از زبان های مختلف، بازی ها با استفاده از مفهوم یادگیری تقویتی توسعه می یابند. مدل های ترکیبی با استفاده از یادگیری عمیق و یادگیری تقویتی در بازی ها تهیه شده اند.

با استفاده از نمایش نمادین، بازی های مختلفی از جمله بازی های کارتی و بازی های رومیزی تهیه می شود. بازی های ورق و بازی های رومیزی پیچیدگی خاصی ندارند. تکنیک های هوش محاسباتی (CI)<sup>۱</sup> برای کمک به حل مشکلات چالش برانگیز در این بازی ها استفاده شده است، اما روش های اختیار داشتن با استفاده از تکنیک CI شکست خوردند [۱].

تخته نرد شناخته شده ترین بازی است که با استفاده از یادگیری تقویتی ساخته شده است. در تخته نرد از الگوریتم مدل - آزاد، که یک الگوریتم تقویتی است، استفاده شد. این الگوریتم مبتنی بر لایه های چندگانه و پرسپترون چند لایه ای است که در این بازی برای ارزیابی موثر مقدار عملکرد، استفاده شده است. الگوریتم مدل - آزاد، همان عملکردهایی را فراهم می کند که یادگیری Q ارائه می کند.

مشارکت دیگری که توسط Tsitisklis نشان داده شد، ترکیب الگوریتم مدل - آزاد بود که باعث واگرایی در Qnetworking می شود. در یادگیری تقویتی، اکثر کارها بر اساس تقریبگرهای

---

<sup>۱</sup>Continuous Integratio

تابع خطی است. برای تخمین محیط می توان از شبکه عصبی عمیق برای ارائه تخمین بهتر نتایج مطابق با یادگیری تقویتی استفاده کرد [۲]. برای تخمین مقدار تابع، می توان از ماشین های بولترمن برای ارائه نتایج دقیق در ترکیب با یادگیری تقویتی استفاده کرد.

تکامل عصبی برای کمک به تولید پارامتر شبکه عصبی مورد استفاده در بازی ها پیشنهاد شد. تکامل عصبی بخشی از الگوریتم یادگیری تقویتی است که به الگوریتم های تکاملی کمک می کند. تکامل عصبی بر اساس توپولوژی است و به عنوان یکی از بهترین روش های مورد استفاده در یادگیری تقویتی در نظر گرفته می شود. این کمک می کند تا وظایف یادگیری تقویتی را کارآمدتر و سریعتر انجام دهید. در مقایسه با روش تکامل عصبی قراردادی (معمولی) (CNE)<sup>۲</sup>، تکامل عصبی بهترین گزینه است، زیرا CNE را فقط می توان بر روی توپولوژی های ثابت اعمال کرد و این روش در ارائه دسترسی سریع برای انجام وظایف یادگیری تقویتی شکست خورده است [۳].

با استفاده از یادگیری، بازی های مختلفی وجود دارند که مفهومی دقیق برای حل مسائل مختلف ارائه می دهند. برای پیاده سازی الگوریتم یادگیری تقویتی باید با چالش های متعددی مواجه شد که پیچیده ترین آنها یافتن محیط مناسبی است که در آن عامل به گونه ای آموزش ببیند که نتایج را به طور موثر ارزیابی کند [۴]. با تعمیم رفتار مناسب عامل در بازی، الگوریتم های مختلفی برای کمک به حل پیچیدگی محیط، با توجه به رفتار عامل پیشنهاد شده است.

روش ها و رویکردهای یادگیری تقویتی مختلف برای کمک به حل مسائل پیچیده محیطی پیشنهاد شده است. بیشتر رویکردها در بازی های آتاری استفاده شده است. رویکرد برنامه پویا برای کمک به حل استراتژی بهینه استفاده شده است [۵]. روش دیگر، شاخص تخصیص است. روش شاخص تخصیص به یافتن انتخاب عمل بهینه کمک می کند. با استفاده از این روش می توان اقدامات مقایسه ای انجام داد و می توان از آن برای دستکاری رباتیک استفاده کرد. با استفاده از روش شاخص تخصیص، مشکلات تقویت تاخیری یافت می شود که مقدار شاخص آنالوگ ها نمی تواند آنها را پیدا کند [۶].

در یادگیری تقویتی از تکنیک ها و استراتژی های Ad-Hoc مختلف از جمله راهبردهای حریصانه استفاده شده است. تکنیک های مختلفی در الگوریتم های یادگیری تقویتی استفاده می شود. استراتژی تصادفی یک استراتژی اکتشافی است که به انجام عمل کمک می کند [۷]. استراتژی دیگری که در الگوریتم یادگیری تقویتی استفاده می شود، اکتشاف بولترمن است.

مدل های تقویتی مختلفی پیشنهاد شده اند که در آنها عامل با استفاده از اعمال و ادراک به محیط متصل می شود. کل اصطلاحات این بازی ها بر اساس محیط و عامل است، اما این امر منجر

به کندی عملکرد می شود زیرا انجام اقدامات و ارزیابی نتایج زمان زیادی می برد. رویکردهای مختلفی برای انجام عملکردهای مختلف و ساده کردن پیچیدگی یک مساله خاص استفاده می شود. از رویکرد پویا می توان برای حل یک برنامه ساده استفاده کرد و این شامل استفاده از استدلال بیزی می شود. استدلال بیزی به حل پیچیدگی یک مسئله ساده کمک می کند [۸]. علاوه بر این، برای بدست آوردن مقدار تقریبی، از روش های مدل - آزاد استفاده می شود که در آن محاسبات مطابق با تکرار خط مشی انجام می شود.

روش دیگری که اخیرا مورد استفاده قرار گرفته است، Critic Actor نام دارد. Actor یک خروجی دو لایه برای استفاده از شبکه تکی در اجرای وظایف مجزا تولید می کند. برای انجام این کار، هم مقدار تابع کامل و هم مقدار گسسته را ارائه می کند. الگوریتم دیگری که استفاده می شود، الگوریتم تکرار سیاست است که در آن معادله خطی و پیچیدگی را حل می کند. الگوریتم تکرار سیاست ( Policy Iteration ) نتایج ارزیابی را ارائه می دهد و هر بار که یک عمل انجام می شود خط مشی را در بازی تغییر می دهد. با استفاده از این الگوریتم، محاسبه تعداد تکرارهای انجام شده در موقعیت های پیچیده تر دشوار است. در بدترین شرایط، الگوریتم به اندازه کافی برای محاسبه تکرارهایی که نتایج ارزیابی را تولید می کند، معتبر نیست.

## ۱.۲ بازی مار و یادگیری تقویتی

با استفاده از بازی مار مفهوم یادگیری تقویتی آشکار می شود. بازی قدیمی مار با استفاده از رویکردها و الگوریتم های تقویتی اصلاح شده است. علاوه بر این، در بازی قدیمی مار، تولید ایده های جدیدی تقویت شدند که در آن قوانین و مقررات به طور مؤثر تعریف می شدند. بازی مار به گونه ای بود که یک بازیکن طبق قوانین از پیش تعیین شده بازی می کرد. آب نبات در بازی ظاهر می شود و اگر مار آن را بخورد، به اندازه ۱ واحد رشد می کند.

گسترش بازی مار، با معرفی آب نبات های سمی به بازی داده می شود. اگر مار آب نبات سمی را بخورد، فوراً می میرد. نشانه ها به گونه ای در بازی طراحی شده اند که قبل از تمام شدن ساعت باید تعداد معینی آب نبات خورده شود. امتیاز نهایی بازی به تعداد کل آب نبات های خورده شده توسط مار بستگی دارد. قوانین بازی شامل ۲ روش دیگر برای مردن مار می شود: (۱) مار می تواند از دم خود عبور کند اما نه بیش از سه بار. در گذر چهارم خواهد مرد. (۲) مار وقتی سرش به ناحیه مرزی برخورد کند می میرد. علاوه بر این، از تکنیک های جدید نیز استفاده می شود و با استفاده از الگوریتم SARSA، فرآیند اقدام انجام شده، قابل اعتمادتر می شود. قوانینی برای

افزایش پیچیدگی و افزایش علاقه بازیکن تعریف شده است.

## ۲.۲ آنالیز ریاضی بازی مار

- به طور انتزاعی، هر مرحله از بازی مار را می توان به عنوان یافتن، راه رفتن خود اجتنابی (بدون اجتناب) (SAW) از یک شبکه در  $\mathbb{R}$  در نظر گرفت. تحقق بازی مار را می توان به صورت ریاضی نشان داد:  
- صفحه  $m \times n$  به صورت  $G = (V, E)$  معرفی می شود که در آن  $V = \{v_i\}$  مجموعه ای از رئوس است، که در آن هر رأس مربوط به یک مربع در تخته است، و  $E = \{e_{ij} | v_i \text{ در مجاورت } v_j\}$ .  
- مار به صورت یک مسیر  $\{u_1, \dots, u_k\}$  معرفی می شود که در آن  $u_1$  و  $u_k$  ابتدا و انتها می باشند.
- *NP – Hard* (Viglietta: ۲۰۱۳) [۶] ثابت کرد که هر بازی شامل آیتم کلکسیونی، شامل پیمایش مکان و مسیر یک طرفه NP-hard است. مار یک بازی است که شامل عبور از یک مسیر یک طرفه است، یعنی مار نمی تواند از خود عبور کند. بنابراین، اگر اطلاعات قبلی در مورد توالی ظهور غذا نداشته باشیم، انتخاب کوتاهترین SAW برای هر مرحله در یک قسمت NP-hard است.

## ۳.۲ الگوریتم تقریبی مار

با توجه به سختی NP (NP-hardness) در یافتن راه حل بهینه برای مار، یک الگوریتم اکتشافی توسعه داده شده است که به طور قابل توجهی برای مسئله مار، خوب عمل می کند و از آن به عنوان معیار برای الگوریتم های یادگیری تقویتی که بعداً اشاره شده می شود استفاده شده است. این الگوریتم، در الگوریتم ۱ نشان داده شده است.

---

Simple Additive Weighting<sup>۳</sup>

الگوریتم ۱: الگوریتم ابتکاری قطعی برای مار:

---

```

1 Build the graph  $G = (V, E)$ , each square  $v_{ij} \in V$ 
2 while Snake  $S = \{s_1, \dots, s_k\} \in V$  and Food  $F = \{f\} \in V$  do
3   build  $\{v_{ij}, v_{kl}\} \in E$ , where  $v_{ij}, v_{kl} \notin S$  with
    $(i = \pm k, j = l)$  or  $(i = k, j = \pm l)$ 
4   initialize  $m \times n$  array  $D$ 
5   assign  $D[i][j] =$  length of path between  $v_{ij}$  and
    $f$  or  $MAX\_NUM$  if no path exists
6   find the  $\min\{D[i][j]\}$  with  $\{v_{ij}, s_1\} \in E$ 
7   if  $D[i][j] \neq MAX\_NUM$  then
8     find Path  $P = \{v_{ij}, p_2, \dots, p_{n-1}, f\}$ 
9     assign  $\bar{s}_1 = f, \bar{s}_2 = p_{n-2}, \dots, \bar{s}_k = p_{n-k}$ 
10    if  $\bar{s}_1$  and  $\bar{s}_k$  is connected then
11      assign  $s_1 = v_{ij}, s_2 = s_1, \dots$ 
12      continue to step 2
13    end
14    else
15      Find the longest Path between  $s_1$  and  $s_k$ ,
16       $P = \{s_1, p_2, p_3, \dots, p_{n-1}, s_k\}$ 
17      assign  $s_1 = p_2, s_2 = s_1, \dots$ 
18      continue to step 2
19    end
20  end
21 end

```

---

## ۴.۲ کاهش فضای حالت برای یادگیری تقویتی

برای انجام و پیاده سازی الگوریتم های یادگیری تقویتی موفق برای انجام بازی، یکی از موانع اساسی، اندازه عظیم فضای حالت است. به عنوان مثال، اندازه صفحه بازی را به صورت  $n \times n$  مشخص کنید. خیلی ساده، با استفاده از موقعیت دقیق مار و غذا، هر سلول را می توان با یکی از چهار شرط پارامتر کرد: { حاوی غذا، حاوی سر مار، حاوی بدن مار، خالی }. با اصل شمارش

ساده ، اندازه فضا ي حالت  $|S|$  در  $n^8 > |S|$  صدق مي کند، به طوري که وقتي  $n$  بسيار بزرگ است و يادگيري آن در فضاي حالت از چنين سايزي مشکل است، بسيار بزرگ است. از اين رو، براي تسريع سرعت يادگيري، يک روش کاهش ساده اين است که فقط موقعيت نسبي مار و غذا را با وضعيت فعلي مار ثبت کنيد. دقيقاً، هر حالت در فضاي کاهش يافته به شکل

$$\{w_s, w_t, w_r, q_f, q_t\}$$

است. در عبارت فوق،  $w_s$ ،  $w_t$  و  $w_r$  توابع نشانگر هستند که آيا ديوارى در مجاورت سر به ترتيب در جهت جلو، چپ و راست وجود دارد.  $q_t$  و  $q_f$  موقعيت نسبي غذا و دم نسبت به سر است. با استفاده از چنين نقشه‌اي، اندازه کل فضاي حالت تنها  $2^3 \cdot 4^2 = 128$  است، و عملکرد به طور برجسته در فرآيند يادگيري بهبود مي يابد.

## فصل ۳

# پیاده سازی بازی

بازی مار با استفاده از پایتون ۶.۳ با برخی وابستگی های دیگر پیاده سازی شده است. از الگوریتم یادگیری Q با کراس به عنوان مدل شبکه عصبی استفاده شده است. وابستگی های این بازی به پایتون، کراس، تصویر کیت علمی و جریان تانسور است. بخشی از کد در این بازی با <sup>۱</sup>CNN ساخته می شود. الگوریتم یادگیری Q تا جایی ادامه دارد که در آن پارامتر گاما حل می شود و ماتریس اولیه Q به صفر می رسد.

الگوریتم یادگیری تقویتی عمیق با استفاده از Keras انجام می شود. بازی مار در پایتون با استفاده از Pygame نوشته شده است. یادگیری D-Q به جای استفاده از رویکرد سنتی یادگیری ماشین انجام می شود. از معادله بلمن استفاده می شود که در آن مقادیر Q به روز می شوند. از ۱۲۰ نوروں پنهان استفاده می شود و یک لایه حذفی برای بهینه سازی استفاده می شود. برای آموزش بازی، موقعیت - عمل - جایزه - موقعیت، اعمال می شود. پس از آموزش داده ها، تست بازی اعمال می شود. در تست بازی، قابلیت حرکت مار تا پایان بازی تست و آزمایش می شود.

یادگیری Q عمیق برای کمک به افزایش کارایی یادگیری Q و ارائه یک جریان ثابت برای کنترل موقعیت مار در بازی استفاده می شود [۹]. چندین حالت در بازی استفاده می شود که در آن مار حرکت می کند و طبق میل کاربر قدم بعدی را پیش بینی می کند. برای پیش بینی حالت ها، جدول Q به طور مداوم به روز می شود و مقادیر به روز رسانی ها با استفاده از معادله بلمن ساخته

---

<sup>۱</sup>Convolutional Neural Network



می شوند.

$$NewQ(s, a) = Q(s, a) + \alpha [R(s, a) + \gamma \max_{a'} Q'(s', a') - Q(s, a)]$$

New Q-Value      Current Q-Value      Learning rate      Reward      Discount rate      Maximum predicted reward, given new state and all possible actions

این الگوریتم طوری عمل می کند که بازی با "s" شروع می شود و مقدار "Q" مقداردهی اولیه می شود. بر اساس وضعیت سیستم، "s" یک عمل را اجرا می کند و انجام اقدامات تصادفی را انتخاب می کند. همانطور که Q مقدار به روز می شود و اقدامات بیشتری در بازی انجام می شود. اقدامات لازم برای جمع آوری پاداش و ارزیابی نهایی پس از خاتمه حالات انجام می شود. این نوع عملیات به عنوان حافظه رله برچسب گذاری شده است.

عامل به گونه ای آموزش می بیند که اقدام بعدی انجام شده در بازی بر اساس الگوریتم SARSA [۱۰] باشد. این الگوریتم یک الگوریتم مبتنی بر خط مشی است که به تخمین ارزش و امتیازات پیش بینی شده بر اساس آن کمک می کند. عامل در بازی با الگوریتم SARSA [۱۱] آموزش دیده است. در بازی، الگوریتم آموزشی بر اساس الگوریتم SARSA با استفاده از کد زیر می باشد:

```
def train_short_memory(self, state, action, reward, next_state, done):
    target = reward
    if not done:
        target = reward + self.gamma *
np.amax(self.model.predict(next_state.reshape((1, state_dim))))[0])
    target_f = self.model.predict(state.reshape((1, state_dim)))
    target_f[0][np.argmax(action)] = target
    self.model.fit(state.reshape((1, state_dim)), target_f, epochs=1, verbose=0)
```

با استفاده از کد پایتون در بازی، الگوریتم SARSA برای آموزش عاملی که در آن اقدام بعدی انجام می شود، به کار می رود. الگوریتم SARSA طوری کار می کند که وقتی عامل در حالت "S" است و اقدامات انجام شده توسط عامل منجر به پاداش "R" می شود. عامل به حالت "S" برمی گردد تا اقدام بعدی "A" را انجام دهد.

## ۱.۳ الگوریتم State - Action - Reward - State - Action (SARSA)

### و یادگیری تقویتی

در یادگیری تقویتی، الگوریتم SARSA به تخمین خط مشی ارزشی که دنبال می شود کمک می کند. این یک الگوریتم روی خط مشی است که در آن عامل اقدام را انجام می دهد و برای انجام اقدام بعدی به حالت نهایی باز می گردد. الگوریتم SARSA به گونه ای عمل می کند که عامل حالت اول را انجام می دهد، اولین اقدام را برای دریافت پاداش انجام می دهد و برای تصمیم گیری در مورد اقدام بعدی به حالت نهایی باز می گردد.

در مورد اقدام بعدی به حالت نهایی باز می گردد.

$Q(s, a)$  را به صورت دلخواه انتخاب کنید.

تکرار (برای هر دور)

$s$  را مقداردهی کنید.

با استفاده از خط مشی مشتق شده از  $Q$ ،  $a$  را از  $s$  انتخاب کنید

(برای مثال  $\epsilon$  حریص)

تکرار (برای هر مرحله از قسمت)

اقدام  $a$  را انجام دهید،  $r$ ،  $s$  را مشاهده کنید.

با استفاده از خط مشی مشتق شده از  $Q$ ،  $a'$  را از  $s'$  انتخاب کنید.

(برای مثال  $\epsilon$  حریص)

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$$

$$s \leftarrow s'; a \leftarrow a'$$

با استفاده از الگوریتم SARSA، عامل اقدام می کند و با محیط تعامل می کند. این الگوریتم روی خط مشی روی فرآیندهای حالت و عمل کار می کند و با تغییر عملکرد به روز می شود. حالت نشان دهنده موقعیتی است که عامل در آن قرار می گیرد. عامل در بازی موقعیت خاص حالت را بر اساس اقدامات کاربر دریافت می کند. مقدار  $Q$  در الگوریتم SARSA مطابق با وضعیت فعلی عامل به روز می شود.

$$Q(S, A) \leftarrow Q(S, A) + \alpha(R + \gamma Q(S', A') - Q(S, A))$$

معادله SARSA بر اساس تغییرات عامل در عمل به روز می شود. ، "S" "A" وضعیت اولیه عامل را نشان می دهد و پاداش "R" مطابق با عمل داده می شود. عامل به حالت نهایی می رسد و مطابق با محیط عمل می شود.

SARSA الگوریتم روی خط مشی است که به تصمیم گیری موثر کمک می کند. با توجه به رفتار عامل، الگوریتم SARSA مقدار جایی که عامل عمل خاصی را در حالت S انجام داده است، تخمین می زند. عملکرد الگوریتم SARSA ساده است و با هر تکرار عملکرد عامل،  $a$  نشان دهنده نرخ یادگیری است.  $\gamma$  پاداش با تخفیف را نشان می دهد.  $Q(S, A)$  زمانی که عمل توسط عامل از هر حالت خاص S انجام می شود مقدار را ذخیره می کند.

## فصل ۴

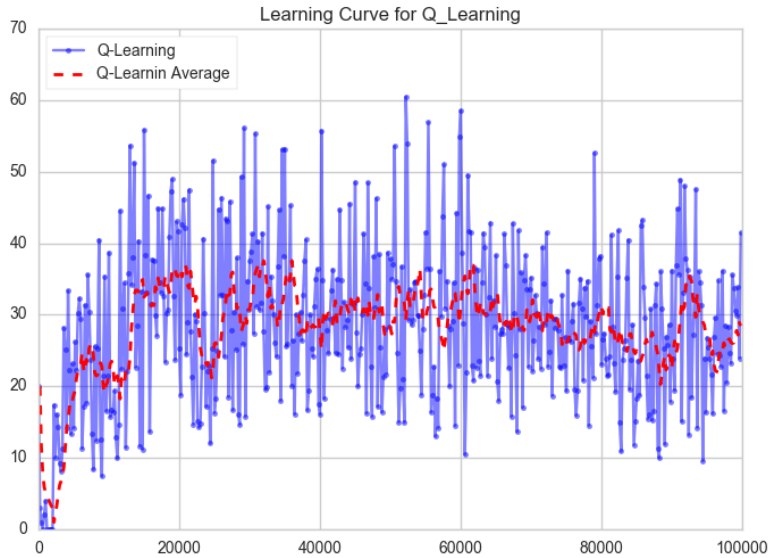
### نتایج

#### ۱.۴ پارامترهای تنظیم

ضریب تخفیف 0/95 تنظیم شد، نرخ یادگیری نسبتاً کاهش می یابد که از  $\alpha = 0/1$  شروع می شود و پاداش ها همانطور که در جدول ۱ نشان داده شده است تعیین شد. ما می خواهیم که عامل کنترل، مار را به سرعت به سمت غذا ببرد و آخرین ستون در جدول ۱، برای انجام یک حرکتی است که در آن عامل را تشویق می کند تا مسیر کوتاه تری را به سمت غذا طی کند. این پاداش های منفی مانند عملکردی مشابه ضریب تخفیف  $\gamma$  عمل می کند. در این پایان نامه آزمون و خطا را روی ترکیب های مختلف پاداش برای موارد مختلف انجام شده است تا ترکیب فعلی مقادیر پاداش به عنوان یک ترکیب بهینه در بین همه آزمایش ها به دست آورده شود.

حالت	خوردن غذا	برخورد با دیوار	برخورد با خود مار	موارد دیگر
پاداش	+500	-100	-100	-10

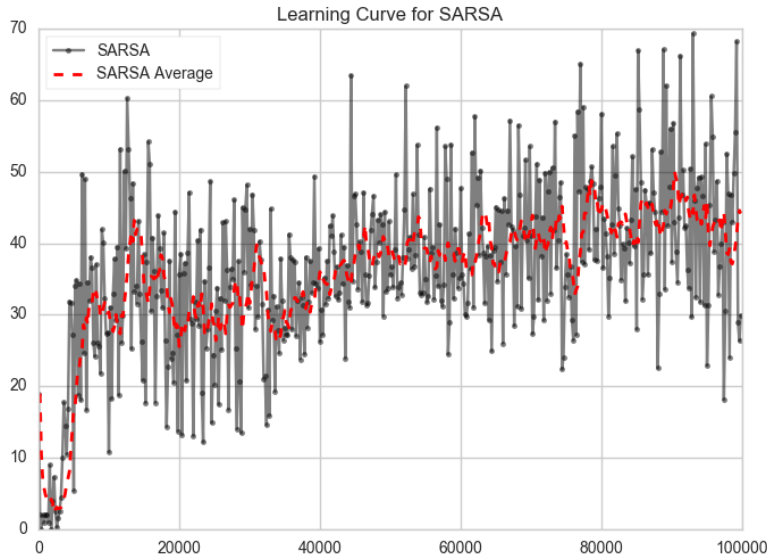
جدول ۱.۴: جدول پاداش



شکل ۱.۴: منحنی های یادگیری از یادگیری Q

## ۲.۴ منحنی یادگیری

منحنی های یادگیری برای یادگیری Q (شکل ۱.۵) و SARSA (شکل ۲.۵) به ترتیب در زیر نشان داده شده است. خطوط خط قرمز نشان دهنده میانگین منحنی های یادگیری است.



شکل ۲.۴: منحنی های یادگیری از SARSA

ما به راحتی می توانیم مشاهده کنیم که در همان ابتدا عملکرد عامل با یادگیری  $Q$  سریعتر از عامل با SARSA بهبود می یابد، یعنی در کوتاه مدت، عامل با الگوریتم یادگیری  $Q$  از عامل با SARSA بهتر عمل می کند. اما با افزایش تعداد تکرارهای آموزشی، عملکرد عامل با یادگیری  $Q$  بهبود چندانی پیدا نمی کند، در حالی که عملکرد عامل با SARSA هنوز به طور قابل توجهی بهبود می یابد. یادگیری  $Q$  زمانی که دوره آموزشی کوتاه است به خوبی انجام می شود (در مقایسه با SARSA). اما در دوره طولانی، SARSA بهتر عمل می کند.

دلیل اینکه عامل یادگیری  $Q$  در برخی موارد در دوره طولانی آموزش خوب عمل نمی کند این است که الگوریتم یادگیری  $Q$  ارزش  $Q$  خود درست بینی را حتی با نرخ یادگیری بسیار کمی تقویت می کند. این منجر به عملکرد قابل توجهی نوسان می شود. اگرچه به نظر می رسد عامل با SARSA در یک دوره آموزشی طولانی از عامل با الگوریتم یادگیری  $Q$  بهتر عمل می کند، اما در موارد بررسی در این جا، منحنی یادگیری نسبتاً کندی دارد.

یادگیری Q	روش		
	SARSA	بهینه	تکرارها
18/023	36/858	77/504	$10^4$
25/789	51/994	77/504	$10^5$
36/567	61/830	77/504	$10^6$

## ۴.۴ مقایسه عملکرد

### ۳.۴ مقایسه عملکرد

عملکرد الگوریتم حل بهینه تقریبی به عنوان معیار در نظر گرفته شده است. همانطور که قبلاً اشاره کردیم، با توجه به اینکه مار یک مساله NP-hard است، بهترین معیاری که می‌توانیم در اینجا بدست آوریم، راه‌حل بهینه تقریبی است. و بعداً می‌توان نشان داد که عوامل ما با الگوریتم‌های یادگیری تقویتی نمی‌توانند این راه‌حل تقریبی را حتی با یک دوره آموزشی طولانی مدت شکست دهند.

سپس عملکرد عامل‌ها بر اساس دو الگوریتم یادگیری تقویتی با معیار - عملکرد راه حل بهینه تقریبی مان مقایسه شده است. در عین حال، تأثیرات دوره آموزشی بر عملکرد عوامل مار در نظر گرفتیم.  $10^4$ ،  $10^5$  و  $10^6$  تکرار آموزشی بر اساس الگوریتم‌های یادگیری تقویتی مختلف بر روی عوامل انجام داده شده است و به ترتیب عملکرد آنها ارزیابی شده است.

محدودیت زمانی برای بازی ۱ دقیقه تعیین شده است، ۱۰۰۰ تست برای عوامل با الگوریتم‌های مختلف اجرا می‌شود و میانگین امتیازی که عوامل مختلف می‌توانند به دست آورند محاسبه شده است. دلیل انتخاب ۱ دقیقه به عنوان محدودیت زمانی این است که در عرض ۱ دقیقه، عواملی با آن الگوریتم‌های یادگیری تقویتی می‌توانند زنده ماندن مار را کنترل کنند، بنابراین می‌توان مطمئن شد که تفاوت فقط در این است که آیا عوامل مختلف می‌توانند مسیر کوتاه‌تری برای غذا پیدا کنند. ۱ دقیقه یک دوره زمانی نسبتاً طولانی برای نشان دادن تفاوت عملکرد قابل توجه بین الگوریتم‌های مختلف است در حالی که یک دوره زمانی نسبتاً کوتاه است که ۱۰۰۰ تست را می‌توان با رایانه شخصی در زمان اجرای معقول انجام داد. نتایج در شکل ۲.۵ و جدول ۲ نشان داده شده است. در اینجا شایان ذکر است که عملکرد الگوریتم حل بهینه تقریبی به تعداد تکرارهای آموزشی مربوط نمی‌شود. فقط عملکرد عوامل با الگوریتم SARSA و یادگیری Q ارتباط مستقیمی با تعداد

تکرارهای آموزشی دارد. عملکرد راه حل بهینه تقریبی به عنوان یک معیار در همان جدول و شکل عملکرد عوامل با الگوریتم های یادگیری تقویتی مختلف نشان داده شده است. نتایج نشان می دهد که در محدوده  $10^4$  تا  $10^6$  تکرار آموزشی، تعداد بیشتری از تکرارهای آموزشی به میانگین نمرات بهتری برای الگوریتم یادگیری Q و SARSA منجر می شود. با توجه به تعداد تکرارهای آموزشی یکسان، عامل آموزش دیده در SARSA عملکرد بهتری نسبت به آموزش یادگیری Q دارد.

## ۵.۴ نتیجه گیری

در این پروژه بازی مار با قوانین جدید تعریف شده پیاده سازی شده است. گسترش بازی مار با معرفی آب نبات های سمی در بازی ایجاد می شود. با استفاده از الگوریتم SARSA، بازی در پایتون نوشته شده است. عملکرد الگوریتم SARSA با یادگیری تقویتی عمیق مشاهده می شود. با استفاده از الگوریتم SARSA کارایی بازی به دلیل عملکرد سریع کار افزایش می یابد. نتایج نسبتاً بهتری را در مقایسه با تکنیک های موجود ارائه می دهد. این روش که در آن فرآیند اولیه سازی کنش عامل شکل می گیرد و سپس اقدامات انجام می شود تا به مرحله بعد می رسند، کار می کند. عملکرد بازی با توجه به عملکرد عامل و امتیاز نهایی مشاهده می شود و سپس ارزیابی می شود. الگوریتم یادگیری تقویتی عمیق با استفاده از Keras انجام می شود. میانگین امتیاز ۵۰ و امتیاز ثبت شده بازی مار ۱۵۰ است.



## کتاب نامه

- [1] S. Mukhopadhyay, O. Tilak and S. Chakrabarti, "Reinforcement Learning Algorithms for Uncertain, Dynamic, Zero-Sum Games," 2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA), vol. 5, no. 3, pp. 66-70, 2018 .
- [2] B. Tsitsiklis, "An Analysis of Temporal-Difference Learning with Function Approximation," IEEE TRANSACTIONS ON AUTOMATIC CONTROL, vol. 42, no. 5, pp. 5-6, 1997.
- [3] A. Notsu, K. Honda, H. Ichihashi and Y. Komori, "Simple Reinforcement Learning for Small-Memory Agent," 2011 10th International Conference on Machine Learning and Applications and Workshops, vol. 1, pp. 78-80, 2011.
- [4] A. Jeerige, D. Bein and A. Verma, "Comparison of Deep Reinforcement Learning Approaches for Intelligent Game Playing," 2019 IEEE 9th Annual Computing and Communication Workshop and Conference (CCWC), vol. 6, no. 2, pp. 55-60, 2019.
- [5] s. Çalışır and m. K. Pehlivanoglu, "Model-Free Reinforcement Learning Algorithms: A Survey," 2019 27th Signal Processing And Communications Applications Conference (siu), vol. 4, no. 2, pp. 67-70, 2019

- [6] G. Fenza, F. Orciuoli and D. G. Sampson, "Building Adaptive Tutoring Model Using Artificial Neural Networks and Reinforcement Learning," 2017 IEEE 17th International Conference on Advanced Learning Technologies (ICALT), vol. 8, no. 6, pp. 56-67, 2017 .
- [7] M. Wang, L. Wang and T. Yue, "An Application of Continuous Deep Reinforcement Learning Approach to Pursuit-Evasion Differential Game," 2019 IEEE 3rd Information Technology, Networking, Electronic and Automation Control Conference (ITNEC), vol. 3, no. 2, pp. 4-7, 2019.
- [8] B. Gökçe and H. L. Akin, "Implementation of Reinforcement Learning by transferring sub-goal policies in robot navigation," 2013 21st Signal Processing and Communications Applications Conference (SIU), vol. 6, no. 4, pp. 7-9, 2013.
- [9] Z. Wu, N. M. Khan, L. Gao and L. Guan, "Deep Reinforcement Learning with Parameterized Action Space for Object Detection," 2018 IEEE International Symposium on Multimedia (ISM), 2018.
- [10] X. Wang, Y. Gu, Y. Cheng, A. Liu and C. L. P. Chen, "Approximate Policy-Based Accelerated Deep Reinforcement Learning," IEEE Transactions on Neural Networks and Learning Systems, 2019.
- [11] G. Surma, "Towards data science," 10 2018. [Online]. Available: <https://towardsdatascience.com/slitherin-solving-the-classicgame-of-snake-with-ai-part-3-genetic-evolution-33186e6be110>.
- [12] Risto Miikkulainen, Bobby Bryant, Ryan Cornelius, Igor Karpov, Kenneth Stanley, and Chern Han Yong. Computational Intelligence in Games. URL:<ftp://ftp.cs.utexas.edu/pub/neural-nets/papers/miikkulainen.wcci06.pdf>

- [13] Gerald Tesauro. Temporal difference learning and td-gammon. *Communications of the ACM*, 38(3):5868, 1995.
- [14] John N Tsitsiklis and Benjamin Van Roy. analysis of temporal difference learning with function approximation. *Automatic Control, IEEE Transactions on*, 42(5):674690, 1997.
- [15] Brian Sallans and Geoffrey E. Hinton. Reinforcement learning with factored states and actions. *Journal of Machine Learning Research*, 5:10631088, 2004.
- [16] Lucas Jen An application of SARSA temporal difference learning to Super Mario URL: <http://x3ro.de/downloads/MarioSarsa.pdf>
- [17] van Seijen, H., van Hasselt, H., Whiteson, S. and Wiering, M. (2009). A theoretical and empirical analysis of Expected Sarsa, 2009 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning pp. 177184. URL: <http://goo.gl/Oo1lu>

## **Abstract**

In this research, we explore the hypothesis that Reinforcement Learning applications are not amenable to conduct close analysis. With the combination of Deep Reinforcement Learning and neural network, the snake game agent is trained. The dependences of this game are on python, Keras, sci-kit image, and tensor flow. For the training of the agent in the game, SARSA is used to control the moving capabilities of the snake as per the user's desire. Deep Q-learning is used in the learning module and the learning rate of the SARSA algorithm shows the effective results in the different agent iterations.



College of Science  
School of Mathematics, Statistics, and Computer Science

# Exploration of Reinforcement Learning to Play Snake Game

**Ramtin Rahvarian**

Supervisor: Dr. Hediye Sajedi

A thesis submitted in partial fulfillment of the requirements for  
the degree of B.Sc. in Computer Science

February 2022