



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)  
دانشکده مهندسی کامپیوتر

پایان نامه کارشناسی

# استفاده از یادگیری تقویتی برای گل زدن در موقعیت تک به تک در فوتبال

نگارش

آراد فیروزکوهی

استاد راهنما

دکتر احسان ناظر فرد

فروردین ۱۴۰۳

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)  
دانشکده مهندسی کامپیوتر

پایان نامه کارشناسی

# استفاده از یادگیری تقویتی برای گل زدن در موقعیت تک به تک در فوتبال

نگارش

آراد فیروزکوهی

استاد راهنما

دکتر احسان ناظر فرد

فروردین ۱۴۰۲

# صفحه فرم ارزیابی و تصویب پایان نامه - فرم تأیید اعضاء کمیته دفاع

در این صفحه فرم دفاع یا تأیید و تصویب پایان نامه موسوم به فرم کمیته دفاع - موجود در پرونده آموزشی - را قرار دهید.

## نکات مهم:

- نگارش پایان نامه/رساله باید به **زبان فارسی** و بر اساس آخرین نسخه دستورالعمل و راهنمای تدوین پایان نامه های دانشگاه صنعتی امیرکبیر باشد.(دستورالعمل و راهنمای حاضر)
- رنگ جلد پایان نامه/رساله چاپی کارشناسی، کارشناسی ارشد و دکترا باید به ترتیب مشکی، طوسی و سفید رنگ باشد.
- چاپ و صحافی پایان نامه/رساله بصورت **پشت و رو(دورو)** بلامانع است و انجام آن توصیه می شود.



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)

به نام خدا

## تعهدنامه اصالت اثر

تاریخ: فروردین ۱۴۰۲

اینجانب **آراد فیروزکوهی** متعهد می‌شوم که مطالب مندرج در این پایان‌نامه حاصل کار پژوهشی اینجانب تحت نظارت و راهنمایی اساتید دانشگاه صنعتی امیرکبیر بوده و به دستاوردهای دیگران که در این پژوهش از آنها استفاده شده است مطابق مقررات و روال متعارف ارجاع و در فهرست منابع و مآخذ ذکر گردیده است. این پایان‌نامه قبلاً برای احراز هیچ مدرک هم‌سطح یا بالاتر ارائه نگردیده است. در صورت اثبات تخلف در هر زمان، مدرک تحصیلی صادر شده توسط دانشگاه از درجه اعتبار ساقط بوده و دانشگاه حق پیگیری قانونی خواهد داشت.

کلیه نتایج و حقوق حاصل از این پایان‌نامه متعلق به دانشگاه صنعتی امیرکبیر می‌باشد. هرگونه استفاده از نتایج علمی و عملی، واگذاری اطلاعات به دیگران یا چاپ و تکثیر، نسخه‌برداری، ترجمه و اقتباس از این پایان‌نامه بدون موافقت کتبی دانشگاه صنعتی امیرکبیر ممنوع است. نقل مطالب با ذکر مآخذ بلامانع است.

**آراد فیروزکوهی**

امضا

تقدیم به پدر بزرگوار و مادر مهربانم

آن دو فرشته ای که از خواسته هایشان گذشتند، سختی ها را به جان خریدند و  
خود را سربلای مشکلات و ناملایمات کردند تا من به جایگاهی که اکنون  
در آن ایستاده ام برسم.

# سپاس‌گزاری

به مصداق «من لم يشكر المخلوق لم يشكر الخالق» بسی شایسته است از استاد فرهیخته و فرزانه جناب آقای دکتر رضا صفابخش که با کرامتی چون خورشید، سرزمین دل را روشنی بخشیدند و گلشن‌سرای علم و دانش را با راهنمایی‌های کارساز و سازنده بارور ساختند تقدیر و تشکر نمایم.

آراد فیروزکوهی  
فروردین ۱۴۰۲

## چکیده

چکیده محتویات

### واژه‌های کلیدی:

یادگیری تقویتی، یادگیری تقویتی عمیق، شبیه‌ساز دوبعدی فوتبال، تک به تک با دروازه‌بان



# فهرست مطالب

صفحه

عنوان

۱	.....	۱ مقدمه
۲	.....	۱-۱ مقدمه
۲	.....	۲-۱ تعریف مساله
۲	.....	۳-۱ اهداف
۲	.....	۴-۱ ابزارها و نرم افزارهای مورد استفاده
۲	.....	۵-۱ ساختار پایان نامه
۳	.....	۲ یادگیری تقویتی
۴	.....	۱-۲ مقدمه
۵	.....	۲-۲ اصول یادگیری تقویتی
۵	.....	۱-۲-۲ عامل، محیط، حالت، عمل و پاداش
۵	.....	۲-۲-۲ فرایند تصمیم گیری مارکوف
۶	.....	۳-۲-۲ سیاست، تابع ارزش، و تابع ارزش عمل
۷	.....	۳-۲ الگوریتم های پایه یادگیری تقویتی
۷	.....	۱-۳-۲ برنامه نویسی پویا
۸	.....	۲-۳-۲ یادگیری به کمک نمونه برداری مونته کارلو
۹	.....	۳-۳-۲ یادگیری به کمک تفاوت زمانی
۱۲	.....	۴-۲ استفاده از تقریب گر تابع
۱۳	.....	۵-۲ الگوریتم یادگیری Q عمیق
۱۳	.....	۱-۵-۲ بازیابی تجربه
۱۴	.....	۲-۵-۲ استفاده از شبکه هدف
۱۴	.....	۳-۵-۲ فرآیند یادگیری
۱۵	.....	۶-۲ الگوریتم بهبود گرادیان سیاست معین عمیق
۱۸	.....	۷-۲ جمع بندی
۱۹	.....	۳ محیط فوتبال روبوکاپ

۲۰	۱-۳ معرفی لیگ
۲۰	۱-۱-۳ اهداف لیگ
۲۰	۲-۱-۳ ویژگی‌های محیط و قوانین بازی
۲۱	۲-۳ کد پایه ایجنت
۲۲	۳-۳ معرفی رفتارهای ممکن
۲۲	۱-۳-۳ رفتارهای سطح پایین
۲۳	۲-۳-۳ رفتارهای سطح متوسط
۲۳	۳-۳-۳ رفتارهای سطح بالا
۲۳	۴-۳ حالت دید کامل
۲۴	۵-۳ معرفی حالت پناستی
۲۵	۶-۳ کار با مربی تمرینی برای تولید محیط قابل تکرار
۲۵	۷-۳ جمع‌بندی
۲۷	۴ پیاده‌سازی و آماده‌سازی محیط استاندارد
۲۸	۱-۴ زیرساخت پایتونی تهاجم نصف زمین
۲۸	۲-۴ کد پایه پایرس
۲۹	۳-۴ کد پایه جی‌آرپی‌سی
۳۱	۴-۴ محیط استاندارد جیم
۳۲	۱-۴-۴ توصیف رابط و توابع موجود
۳۴	۲-۴-۴ نحوه اداام با فضای ربوکاپ
۳۵	۳-۴-۴ فضای حالت و فضای عمل عامل
۳۷	۴-۴-۴ طراحی پاداش
۳۸	۵-۴ پیاده‌سازی یادگیری تقویتی
۳۸	۱-۵-۴ پیاده‌سازی شبکه کیو عمیق
۳۸	۲-۵-۴ پیاده‌سازی سایر الگوریتم‌ها به کمک کتابخانه
۳۹	۶-۴ جمع‌بندی
۴۰	۵ مقایسه، آزمایش‌ها، و نتایج
۴۱	۱-۵ ارزیابی الگوریتم‌ها

۴۱	۲-۵ مقایسه سرعت یادگیری با گسسته سازی‌ها متفاوت . . . . .
۴۱	۳-۵ بررسی تعمیم پذیری . . . . .
۴۱	۴-۵ تسریع یادگیری با تلفیق با رفتار سطح بالای شوت . . . . .
۴۲	۶ جمع‌بندی، نتیجه‌گیری و پیشنهادات . . . . .
۴۳	۱-۶ کارهای آینده . . . . .
۴۴	منابع و مراجع . . . . .

## فهرست تصاویر

شکل	صفحه
۱-۲	تعامل عامل و محیط . . . . . ۵
۲-۲	تاثیر پارامتر $\lambda$ بر اهمیت پاداش‌های آینده. محور افقی نماینده تعداد گام‌ها، و محور عمودی نماینده وزن این بازگشت متناظر با آن است. . . . . ۱۰
۳-۲	روش‌های مختلف استفاده از تقریب‌گر تابع . . . . . ۱۲
۴-۲	شبکه‌های عصبی بازیگر و نقاد در الگوریتم DDPG و مقایسه با DQN . . . . . ۱۶
۱-۳	نمونه دید یک بازیکن. در این مثال، بازیکن ۳، بازیکن ۵ و ۸ و توپ را می‌بیند و از موقعیت سایر بازیکنان اطلاعی ندارد. . . . . ۲۱
۲-۳	ارتباط بین محیط و عامل، و روش محاسبه پاداش . . . . . ۲۴
۳-۳	نمونه‌ای از بازی در حالت پنالیتی . . . . . ۲۵
۱-۴	پرچم‌های کنار زمین، که برای موقعیت‌یابی استفاده می‌شوند. . . . . ۲۹
۲-۴	پروتکل‌های استفاده‌شده برای ارتباط بین اجزای مسابقه . . . . . ۳۰
۳-۴	نحوه کارکرد و اتصال کد پایه جی‌آرپی‌سی به سرور مسابقات و نمایشگر بازی . . . . . ۳۱
۴-۴	نمونه‌ای از محیط‌های آماده شده در پلتفرم جیم . . . . . ۳۲
۵-۴	نحوه کلی ارتباط کد پایه، تصمیم‌گیرنده جی‌آرپی‌سی و جیم . . . . . ۳۴
۶-۴	ترتیب فراخوانی توابع برای اجرای جیم با سرور ربوکاپ . . . . . ۳۵
۷-۴	تقسیم‌بندی فضای عمل به ۶ زاویه و ۲ سرعت . . . . . ۳۶
۸-۴	استفاده از کتابخانه Stable Baselines 3 . . . . . ۳۸

صفحه

## فهرست جداول

جدول

# فصل اول

## مقدمه

۱-۱ مقدمه

۲-۱ تعریف مساله

۳-۱ اهداف

۴-۱ ابزارها و نرم افزارهای مورد استفاده

۵-۱ ساختار پایان نامه

## فصل دوم

### یادگیری تقویتی



## ۱-۲ مقدمه

به طور عمومی، یادگیری ماشین به دسته‌ای از الگوریتم‌ها و روش‌های محاسباتی گفته می‌شود که به ماشین‌ها امکان یادگیری از داده‌ها و تجربه‌های خود را می‌دهند. یادگیری ماشین به دو دسته‌ی اصلی تقسیم می‌شود: یادگیری نظارت‌شده و یادگیری بدون نظارت. در یادگیری نظارت‌شده، مدل به کمک داده‌های برچسب‌خورده آموزش داده می‌شود، و سپس برای پیش‌بینی خروجی‌های جدید از این مدل استفاده می‌شود. در یادگیری بدون نظارت، مدل بدون داده‌های برچسب‌خورده آموزش داده می‌شود، و باید خودش مفاهیم و الگوهای موجود در داده‌ها را کشف کند. یادگیری تقویتی در این دو دسته قرار نمی‌گیرد، و به عنوان یک دسته جداگانه در نظر گرفته می‌شود.

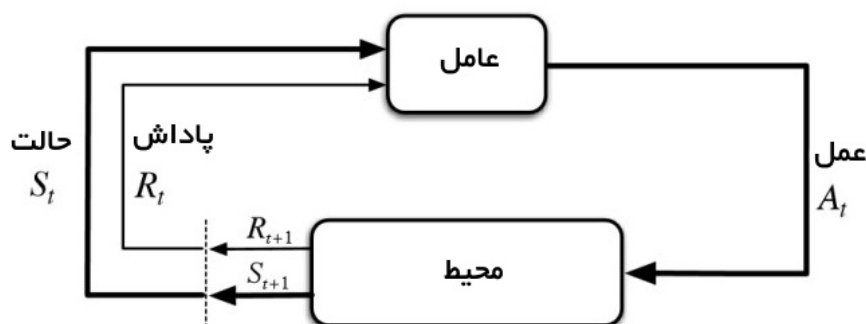
یادگیری تقویتی یک روش یادگیری ماشین است که به عامل اجازه می‌دهد تا رفتار بهینه را از طریق تعامل و آزمون و خطا با محیط یاد بگیرد. در این روش، عامل مشاهدات خود را از محیط (به کمک سنسور) دریافت کرده، و بر اساس آن تصمیم خود را اخذ می‌کند. پس از انجام هر عمل، محیط به عامل پاداشی می‌دهد که نشان‌دهنده‌ی عملکرد عامل در آن حالت است. هدف این است که عامل مجموع پاداش‌های دریافتی خود را بیشینه کند، که در صورت صحیح بودن تعریف سیگنال پاداش، معمولاً به معنای رسیدن به هدف مطلوب است.

یادگیری ماشین با سایر رویکردهای یادگیری ماشین مانند یادگیری نظارت‌شده و یادگیری بدون نظارت تفاوت‌های زیادی دارد، و معمولاً بر مسائلی به کار می‌رود که تمرکز بر تصمیم‌گیری یا پیدا کردن سیاست بهینه بدون نیاز به داده‌های برچسب‌خورده است. در عوض، عامل با اکتشاف و بهره‌برداری ابتدا انتقالات بین حالت‌ها و پاداش‌های مرتبط با آن‌ها را یاد می‌گیرد، و سپس سیاستی را یاد می‌گیرد که مجموع پاداش‌ها را بیشینه کند.

## ۲-۲ اصول یادگیری تقویتی

### ۱-۲-۲ عامل، محیط، حالت، عمل و پاداش

در یادگیری تقویتی، عامل<sup>۱</sup> تصمیم گیرنده برای رسیدن به هدف خود با محیط<sup>۲</sup> تعامل دارد. محیط معمولاً به صورت مجموعه‌ای از حالت‌ها<sup>۳</sup> و عمل‌هایی<sup>۴</sup> که عامل می‌تواند انجام دهد مدل می‌شود. در هر گام، عامل مشاهده‌ای از حالت محیط را دریافت می‌کند و بر اساس آن تصمیمی اتخاذ می‌کند. پس از انجام عمل، محیط به عامل پاداش<sup>۵</sup> می‌دهد که نشان‌دهنده عملکرد عامل در آن حالت است.



شکل ۱-۲ تعامل عامل و محیط

در اکثر مواقع، یادگیری تقویتی در مسائلی استفاده می‌شود، که بتوان آن را به دنباله‌ای از گام‌ها تقسیم کرد که قطعاً به یک حالت پایانی می‌رسد. هر یک از این دنباله‌ها را قسمت<sup>۶</sup> می‌نامند.

### ۲-۲-۲ فرایند تصمیم‌گیری مارکوف

فرایند تصمیم‌گیری مارکوف<sup>۷</sup> مدلی است برای تصمیم‌گیری در محیط‌هایی که به صورت مارکوف هستند. محیط‌های دارای خاصیت مارکوف، محیط‌هایی هستند که حالت بعدی به صورت کامل به حالت فعلی و عمل انجام شده وابسته است.

یک فرایند تصمیم‌گیری مارکوف، چهارتایی  $(S, A, P, R)$  است که در آن:

<sup>۱</sup> Agent

<sup>۲</sup> Environment

<sup>۳</sup> State

<sup>۴</sup> Action

<sup>۵</sup> Reward

<sup>۶</sup> Episode

<sup>۷</sup> Markov Decision Process (MDP)

- $S$  مجموعه‌ی تمام حالت‌های ممکن محیط است.
  - $A$  مجموعه‌ی تمام عمل‌های ممکن است.
  - $P$  تابع انتقال<sup>۸</sup> است که به ازای هر حالت و عمل، توزیع احتمال حالت بعدی را مشخص می‌کند.
  - $R$  تابع پاداش<sup>۹</sup> است که به ازای هر حالت و عمل، پاداش مورد انتظار را مشخص می‌کند.
  - فاکتور تخفیف  $\gamma$  نیز معمولاً به عنوان یک پارامتر دیگر در نظر گرفته می‌شود که نشان‌دهنده‌ی اهمیت پاداش‌های آینده نسبت به پاداش‌های فعلی است.
- همانطور که گفته شد، در هر گام عامل با اخذ تصمیم خود، محیط را به حالت جدیدی می‌برد و پاداشی دریافت می‌کند. به مجموع کل پاداش‌هایی که عامل در یک قسمت دریافت می‌کند، خروجی<sup>۱۰</sup> گفته می‌شود.

$$G_t = R_{t+1} + \gamma \times R_{t+2} + \gamma^2 \times R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k \times R_{t+k+1} \quad (1-2)$$

## ۳-۲-۲ سیاست، تابع ارزش، و تابع ارزش عمل

سیاست<sup>۱۱</sup> یک تابع از حالت‌ها به عمل‌ها است که نشان‌دهنده‌ی رفتار عامل در هر حالت محیط است. در واقع مسئله یادگیری تقویتی را می‌توان به ((یافتن سیاستی که مجموع پاداش‌ها را بیشینه می‌کند)) تعبیر کرد.

$$\pi(a|s) = \mathbb{P}\{A_t = a | S_t = s\} \quad (2-2)$$

تابع ارزش<sup>۱۲</sup>  $V(s)$  یک تابع از حالت‌ها است که نشان‌دهنده‌ی میزان پاداش مورد انتظار از یک

<sup>8</sup>Transition Function

<sup>9</sup>Reward Function

<sup>10</sup>Return

<sup>11</sup>Policy

<sup>12</sup>Value Function

حالت تا به انتهای قسمت است.

$$v_{\pi}(s) = \mathbb{E}_{\pi}\{G_t | S_t = s\} = \sum_{a \in A} \pi(a|s) \times \{R_s^a + \gamma \times \sum_{s' \in S} p(s'|s, a) \times v_{\pi}(s')\} \quad (3-2)$$

تابع ارزش عمل <sup>۱۳</sup>  $Q(s, a)$  نیز مشابه تابع ارزش است، با این تفاوت که به جای حالت، از یک حالت و یک عمل مشخص محاسبه می‌شود. رایج است که به این تابع، تابع کیو گفته شود.

$$q_{\pi}(s, a) = \mathbb{E}_{\pi}\{G_t | S_t = s, A_t = a\} = R_s^a + \gamma \times \sum_{s' \in S} p(s'|s, a) \times v_{\pi}(s') \quad (4-2)$$

به این معادلات، که از کلیدی‌ترین روابط در یادگیری تقویتی هستند، معادله بلمن <sup>۱۴</sup> گفته می‌شود.

## ۳-۲ الگوریتم‌های پایه یادگیری تقویتی

### ۱-۳-۲ برنامه‌نویسی پویا

همان‌طور که در معادله ۳-۲ مشخص است، می‌توان تابع ارزش را به صورت بازگشتی محاسبه کرد. این روش برنامه‌نویسی پویا <sup>۱۵</sup> نام دارد.

### یادگیری به کمک تکرار ارزش

در روش تکرار ارزش <sup>۱۶</sup>، ابتدا تابع ارزش را به صورت تصادفی مقداردهی می‌کنیم و سپس آن را به صورت بازگشتی به روزرسانی می‌کنیم. قابل اثبات است که این روش به تابع ارزش بهینه همگرا می‌شود. به صورت شهودی نیز، می‌توان دید که ارزش از سمت حالت‌های پایانی به سمت حالت‌های ابتدایی به روزرسانی می‌شود.

<sup>13</sup>Action Value Function

<sup>14</sup>Bellman Equation

<sup>15</sup>Dynamic Programming

<sup>16</sup>Value Iteration

## یادگیری به کمک تکرار سیاست

در روش تکرار سیاست<sup>۱۷</sup>، ابتدا سیاست را به صورت تصادفی مقداردهی می‌کنیم و سپس تابع ارزش را برای آن محاسبه می‌کنیم. سپس سیاست را به صورت بازگشتی به روزرسانی می‌کنیم. این فرآیند را تا زمانی که سیاست تغییر نکند ادامه می‌دهیم. قابل اثبات است که این روش به سیاست بهینه همگرا می‌شود.

در عمل، با توجه به اینکه معمولاً به تابع انتقال دسترسی نداریم، نمی‌توانیم از روش‌های برنامه‌نویسی پویا به صورت مستقیم استفاده کنیم. در واقع نیاز به راه‌حل‌های مستقل از مدل داریم که به کمک نمونه‌برداری و کاوش محیط، سیاست بهینه را یاد بگیرند.

## ۲-۳-۲ یادگیری به کمک نمونه‌برداری مونته کارلو

در روش‌های یادگیری به کمک نمونه‌برداری مونته کارلو<sup>۱۸</sup>، به جای استفاده از مدل، از نمونه‌برداری برای تخمین ارزش استفاده می‌شود. کافیست ابتدا یک قسمت را به طور کامل اجرا کنیم، و سپس ارزش هر حالت را، به سمت خروجی قسمت، به روزرسانی کنیم:

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha \times (G_t - Q(s_t, a_t)) \quad (5-2)$$

که در این فرمول،  $\alpha$  نرخ یادگیری<sup>۱۹</sup> است.

لازم به ذکر است که در حین اجرای قسمت، از سیاست اپسیلون-حریصانه<sup>۲۰</sup> استفاده می‌شود. در این سیاست، با احتمال  $\epsilon$  عملی تصادفی انجام می‌شود و با احتمال  $1 - \epsilon$  عملی که ارزش بیشتری دارد انجام می‌شود. دلیل استفاده از این سیاست، نیاز به کاوش محیط و جلوگیری از گیر کردن در حالت‌های محلی است.

از معایب این روش، می‌توان به نیاز به اجرای کامل قسمت‌ها و نیاز به زمان برای یادگیری اشاره کرد. به همین دلیل، این روش برای مسائلی که قسمت‌های طولانی دارند، مناسب نیست. از مشکلات دیگر این روش، عدم استفاده از ویژگی مارکوف محیط است.

<sup>17</sup>Policy Iteration

<sup>18</sup>Monte Carlo Sampling

<sup>19</sup>Learning Rate

<sup>20</sup> $\epsilon$ -greedy

## ۳-۳-۲ یادگیری به کمک تفاوت زمانی

همان‌طور که گفته شد، استفاده از یادگیری مونته کارلو باعث می‌شود که عامل در حین انجام قسمت، از تجربه‌ی قبلی خود استفاده نکند و به روز رسانی ارزش‌ها فقط پس از اتمام قسمت انجام شود. به این روش، یادگیری آفلاین<sup>۲۱</sup> گفته می‌شود. در روش یادگیری به کمک تفاوت زمانی<sup>۲۲</sup>، عامل در حین انجام قسمت، از تجربه‌ی خود استفاده می‌کند و ارزش‌ها را به صورت آنلاین به روزرسانی می‌کند. در واقع به کمک معادله بلمن، مقادیر کیو به سمت مقدار کیو بعدی به روزرسانی می‌شوند.

پیاده‌سازی این دو الگوریتم، به دو روش دید رو به جلو و دید رو به عقب انجام می‌شود. در حالت دید رو به جلو، مشابه با یادگیری مونته کارلو، پس از رسیدن به پایان قسمت، ارزش‌ها به روزرسانی می‌شوند. در حالت دید رو به عقب، ارزش‌ها به صورت آنلاین به روزرسانی می‌شوند. به این صورت که پس از دریافت یک پاداش، عامل مقادیر کیو  $n$  حالت قبلی خود را به روزرسانی می‌کند.

### تی‌دی صفر

در ساده‌ترین حالت، عامل در حین انجام عمل، با دیدن یک گام در آینده یا گذشته، ارزش عمل فعلی را به روزرسانی می‌کند. به این روش تی‌دی صفر<sup>۲۳</sup> گفته می‌شود.

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha \times (R_{t+1} + \gamma \times Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)) \quad (۶-۲)$$

در واقع، به کمک رابطه بلمن (که فرم ارزش عمل آن در رابطه ۴-۲ دیده می‌شود) میزان صحیح بودن ارزش عمل فعلی، با ارزش عمل بعدی و پاداش فعلی مقایسه می‌شود و ارزش عمل فعلی به روزرسانی می‌شود.

<sup>۲۱</sup>Offline Learning

<sup>۲۲</sup>Temporal Difference Learning

<sup>۲۳</sup>TD(0)

## تی‌دی لامبدا رو به جلو

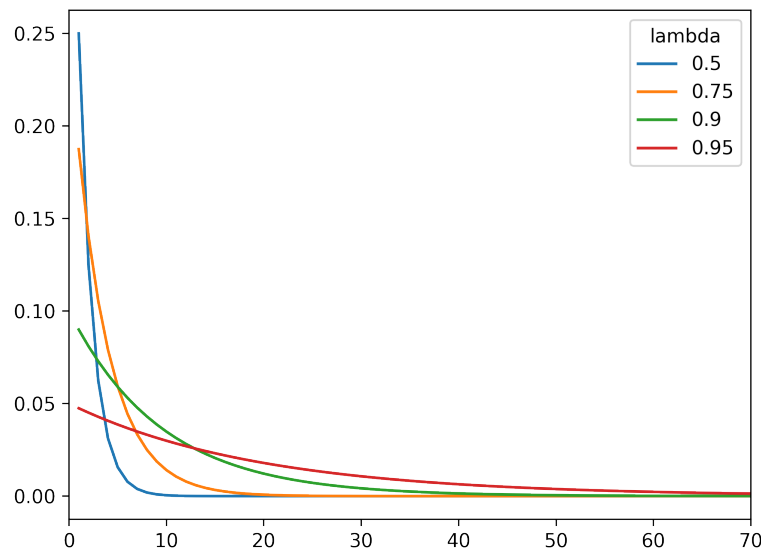
در روش تی‌دی  $\lambda$ ،<sup>۲۴</sup> با دید رو به جلو به جای استفاده از یک گام در آینده، از یک ترکیب خطی از پاداش‌ها و بازگشت چندین گام استفاده می‌شود. به این ترکیب خطی،  $\lambda$ -بازگشت<sup>۲۵</sup> گفته می‌شود.

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha \times (G_t^\lambda - Q(s_t, a_t)) \quad (۷-۲)$$

که در این رابطه،  $G_t^\lambda$  به صورت زیر محاسبه می‌شود:

$$G_t^\lambda = (1 - \lambda) \times \sum_{n=1}^{\infty} \lambda^{n-1} \times G_t^{(n)} \quad (۸-۲)$$

که در آن  $G_t^{(n)}$  مقدار بازگشتی بعد از  $n$  گام است، و  $\lambda$  یک پارامتر بین صفر و یک است که نشان‌دهنده اهمیت پاداش‌های آینده نسبت به پاداش‌های فعلی است. می‌توان میزان اهمیت پاداش‌های را به ازای مقادیر مختلف  $\lambda$  مشاهده کرد.



شکل ۲-۲ تاثیر پارامتر  $\lambda$  بر اهمیت پاداش‌های آینده. محور افقی نماینده تعداد گام‌ها، و محور عمودی نماینده وزن این بازگشت متناظر با آن است.

<sup>۲۴</sup>TD( $\lambda$ )

<sup>۲۵</sup> $\lambda$ -return

## تی‌دی لامبدا رو به عقب

در روش تی‌دی  $\lambda$  با دید رو به عقب، از مفهومی به نام آثار شایستگی<sup>۲۶</sup> استفاده می‌شود. این مفهوم، نشان‌دهنده‌ی این است که هر عملی که در گذشته انجام شده و موجب تغییر ارزش عملی شده است، به چه میزان روی این تغییر سهمیم بوده‌است.

ابتدا به ازای هر حالت و عمل، یک متغیر آثار شایستگی  $E(s, a)$  صفر تعیین می‌شود و سپس در هر گام، این متغیر به صورت زیر به روزرسانی می‌شود:

$$E(s, a) = \gamma \times \lambda \times E(s, a) + 1(s = s_t, a = a_t) \quad (۹-۲)$$

در واقع کل آثار شایستگی در لامبدا ضرب شده، و یکی به آثار شایستگی عمل و حالت فعلی اضافه می‌شود. در این حالت، آثار شایستگی معادل با میزان نزدیکی زمانی هر حالت و عمل، به حالت و عمل فعلی است.

سپس ارزش عمل به صورت زیر به روزرسانی می‌شود:

$$Q(s, a) = Q(s, a) + \alpha \times \delta \times E(s, a) \quad (۱۰-۲)$$

که در این رابطه،  $\delta$  خطای تخمین است که به صورت زیر محاسبه می‌شود:

$$\delta = R_{t+1} + \gamma \times (Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)) \quad (۱۱-۲)$$

می‌توان نشان داد که پس از انجام یک قسمت، به‌روزرسانی‌های انجام شده توسط این روش، معادل با به‌روزرسانی‌های انجام شده توسط تی‌دی  $\lambda$  با دید رو به جلو است.

با ترکیب این روش و سیاست اپسیلون-حریصانه، می‌توان به یک روش یادگیری تقویتی کارا دست یافت که سارسا-لامبدا<sup>۲۷</sup> نام دارد.

<sup>۲۶</sup>Eligibility Traces

<sup>۲۷</sup>SARSA( $\lambda$ )

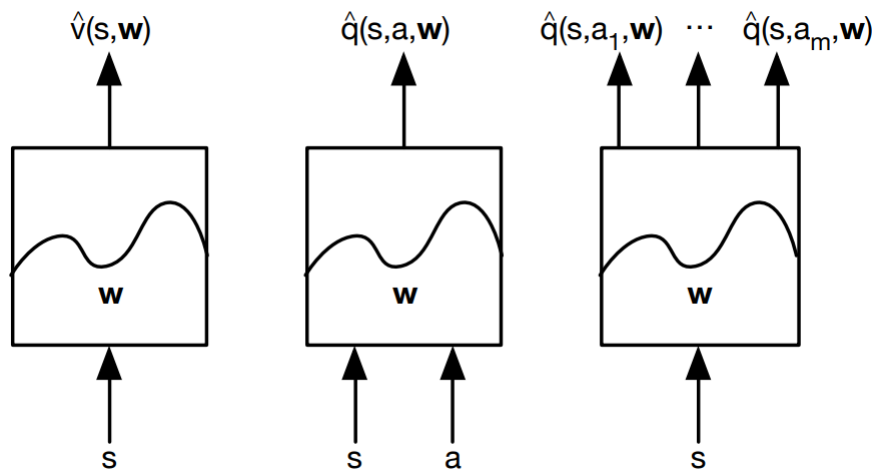


## ۴-۲ استفاده از تقریب‌گر تابع

با توجه به زیاد بودن تعداد حالت‌های ممکن در اکثر مسائل، یادگیری با روش‌های توضیح‌داده‌شده بسیار کند و حتی غیرمعقول می‌باشد. از این رو، به جای نگهداری مقادیر  $Q$  در یک جدول به ازای هر زوج حالت-عمل از تقریب‌گرهای تابع<sup>۲۸</sup> برای تخمین این مقادیر استفاده می‌کنیم.

از تقریب‌گرهای تابع مختلفی می‌توان برای این کار استفاده کرد. از جمله‌ی این تقریب‌گرها می‌توان به شبکه‌های عصبی، درخت‌های تصمیم، و توابع پایه‌ای مثل چندجمله‌ای‌ها اشاره کرد.

در واقع ابتدا از حالت، برخی ویژگی‌های مهم را استخراج کرده، سپس از این ویژگی‌ها برای تخمین ارزش‌ها به کمک پارامترهایی که قرار است یاد گرفته شوند استفاده می‌کنیم. رایج است که در روابط، به این پارامترها با  $\theta$  یا  $w$  اشاره شود.



شکل ۳-۲ روش‌های مختلف استفاده از تقریب‌گر تابع

حال کافیست که معادلات ۲-۶، ۲-۷، و ۲-۱۰ را به صورت تقریبی برای تقریب‌گر تابع بنویسیم، و به جای به روزرسانی مقادیر  $Q$ ، پارامترهای تقریب‌گر تابع را به روز کنیم. با توجه به نیاز به به روزرسانی، لازم است از تقریب‌گرهای مشتق‌پذیر استفاده کنیم. همانطور که در عکس ۳-۲ می‌توان دید، در حالتی که از تقریب‌گرهای تابع برای پیش‌بینی مقادیر  $Q$  استفاده می‌شود، رایج است از عبارت  $Q(s, a; \theta)$  یا  $Q(s, a; w)$  برای نشان دادن این تقریب‌گرها استفاده کرد.

<sup>28</sup>Function Approximators

## ۵-۲ الگوریتم یادگیری Q عمیق

این الگوریتم، یکی از مستقیم‌ترین راه‌های استفاده از تقریب‌گرهای توابع، برای ترکیب یادگیری عمیق و تقویتی است. در این الگوریتم، هدف این است که مشابه با حالت سوم شکل ۲-۳، پارامترهای یک شبکه عصبی را به‌گونه‌ای تنظیم کنیم که مقادیر ارزش-عمل را بتواند تخمین بزند. از دستاوردهای این الگوریتم، عملکرد در حد انسان در چندین بازی آتاری<sup>۲۹</sup> است که توسط تیم دیپ‌مایند<sup>۳۰</sup> در سال ۲۰۱۵ به‌وقوع پیوست.

شبکه‌های عصبی الهام گرفته از ساختار و عملکرد مغز انسان هستند که برای یادگیری از داده‌ها و تصمیم‌گیری‌های پیچیده استفاده می‌شوند. این شبکه‌ها از واحدهای پردازشی به نام پرسپترون‌ها تشکیل شده‌اند که در لایه‌های مختلف قرار گرفته‌اند و از طریق وزن‌هایی به هم متصل می‌شوند. یادگیری در شبکه‌های عصبی اغلب از طریق فرایندی به نام نزول گرادیان انجام می‌گیرد که در آن وزن‌های شبکه به صورت تکراری تنظیم می‌شوند تا خطا بین پیش‌بینی‌های شبکه و داده‌های واقعی به حداقل برسد. این فرایند شامل محاسبه گرادیان یا شیب تابع خطا نسبت به وزن‌ها و به‌روزرسانی وزن‌ها در جهت مخالف گرادیان برای کاهش خطا است.

$$\theta_{t+1} = \theta_t - \alpha \nabla J(\theta_t) \quad (۱۲-۲)$$

در این فرمول،  $\theta$  پارامترهای شبکه،  $\alpha$  نرخ یادگیری،  $J$  تابع خطا، و  $\nabla J$  گرادیان تابع خطا نسبت به پارامترها هستند.

برای اینکه یادگیری شبکه عصبی با ثبات بالاتری رخ دهد، این الگوریتم از دو تکنیک بازیابی تجربه و استفاده از شبکه هدف بهره می‌برد.

### ۱-۵-۲ بازیابی تجربه

برای رفع مشکلات داده‌های همبسته و توزیع‌های غیر ایستا در یادگیری آنلاین، DQN از یک مکانیزم بازیابی تجربه استفاده می‌کند. این شامل ذخیره تجربیات عامل در هر گام زمانی در یک بافر بازیابی و سپس نمونه‌برداری تصادفی ریزدسته‌ها<sup>۳۱</sup> از این بافر برای آموزش شبکه است. این رویکرد به شکستن

<sup>۲۹</sup>Atari

<sup>۳۰</sup>DeepMind

<sup>۳۱</sup>Mini-batches

همبستگی بین نمونه‌های پیاپی کمک می‌کند و فرآیند یادگیری را پایدار می‌سازد.

## ۲-۵-۲ استفاده از شبکه هدف

DQN یک شبکه دوم به نام شبکه هدف را برای استقرار بیشتر آموزش معرفی می‌کند. شبکه هدف یک کپی از شبکه Q است، اما وزن‌های آن کمتر به‌روز می‌شوند. این جداسازی نوسان ارزش‌های هدف در به‌روزرسانی یادگیری Q را کاهش می‌دهد و خطر حلقه‌های بازخورد خودتقویتی را کاهش می‌دهد.

## ۳-۵-۲ فرآیند یادگیری

در ابتدا وزن‌های شبکه ارزش-عمل را به صورت تصادفی تنظیم می‌کنیم و آن را  $\theta$  می‌نامیم. شبکه هدف را با وزن‌های یکسان مقداردهی می‌کنیم و آن را  $\theta^-$  می‌نامیم. و بافر بازیابی را به طول  $N$  که یک ابرپارامتر است، مقداردهی می‌کنیم.

سپس در هر حالت، پاداش‌های خروجی شبکه اصلی را می‌گیریم و به صورت اپسیلون-حریصانه عمل را انتخاب می‌کنیم؛ یعنی به احتمال  $\epsilon$  رفتار تصادفی داریم، و با احتمال  $1 - \epsilon$  رفتار با بالاترین ارزش-عمل پیش‌بینی شده را برمی‌داریم، یعنی:  $a_t = \arg\max_a Q(s_t, a; \theta)$ . ترکیب چهارتایی حالت قبل عمل، عمل انتخاب شده، پاداش دریافتی، و حالت بعد عمل  $(s_t, a_t, r_t, s_{t+1})$  را به بافر اضافه می‌کنیم. یک گروه با اندازه مشخص را با نمونه‌برداری تصادفی از بافر انتخاب می‌کنیم، و به صورت زیر یادگیری را انجام می‌دهیم:

$$y_j = \begin{cases} r_j & \text{اگر } s_{j+1} \text{ حالت پایانی قسمت باشد} \\ r_j + \gamma \times \max_{a'} Q(s_{j+1}, a'; \theta^-) & \text{در غیر این صورت} \end{cases} \quad (۱۳-۲)$$

در این معادله  $y_j$  تخمینی از میزان پاداش دریافتی واقعی است، که مشابه با روش‌های TD، از ترکیب پاداش دریافتی و ارزش-عمل بعدی به دست می‌آید. بنابراین می‌توان از این مقدار، مشابه با برجسب در یادگیری نظارت‌شده، برای آموزش شبکه استفاده کرد. میزان خطای شبکه به صورت زیر تعریف می‌شود:

$$L(\theta) = \mathbb{E}[(y_j - Q(s_j, a_j; \theta))^2] \quad (۱۴-۲)$$

که با استفاده از این خطا، می‌توان گرادیان تابع خطا نسبت به وزن‌های شبکه را به دست آورد و با استفاده از این گرادیان، می‌توان وزن‌های شبکه را به‌روز کرد:

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} L(\theta) \quad (15-2)$$

در نهایت کافی‌ست به صورت تناوبی و هر چند گام یک بار وزن‌های شبکه هدف را برابر با کپی وزن‌های شبکه اصلی قرار دهیم. با تکرار این مراحل، به شبکه عصبی دست می‌یابیم که قدرت پیش‌بینی ارزش‌های عمل را دارد.

حال می‌توان اهمیت دو تکنیک گفته شده را بهتر فهمید: در صورت عدم استفاده از شبکه هدف، باید همزمان از شبکه اصلی برای پیش‌بینی آینده (در فرمول ۲-۱۳) و انتخاب عمل استفاده کنیم، که می‌تواند به مشکلاتی مانند نوسانات در آموزش و حلقه‌های بازخورد خودتقویتی منجر شود. همچنین در صورت عدم استفاده از بافر تجربه، باید صرفاً از تجارب اخیر خود استفاده کنیم، که در این صورت هبستگی بین نمونه‌ها و توزیع‌های غیر ایستا می‌تواند مشکل‌ساز شود.

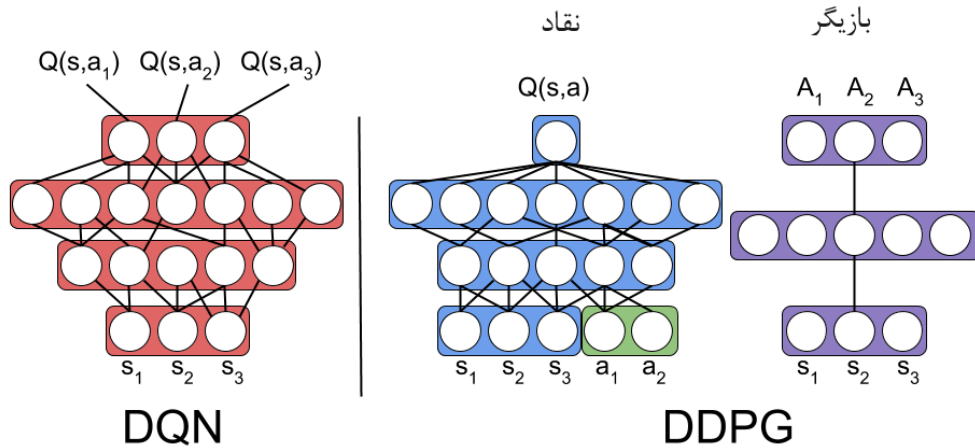
## ۶-۲ الگوریتم بهبود گرادیان سیاست معین عمیق

الگوریتم بهبود گرادیان سیاست معین عمیق<sup>۳۲</sup> یا به اختصار DDPG، یک الگوریتم یادگیری تقویتی است که برای حل مسائل پیوسته و فضای عمل پیوسته طراحی شده است. همانطور که در بخش‌های قبل دیدیم، برای انتخاب عمل معمولاً از سیاست افسیلون-حریصانه استفاده می‌شود، اما این روش برای فضای عمل پیوسته مناسب نیست، چرا که به‌دست‌آوردن بیشینه در فضای پیوسته به مراتب دشوارتر از حالت گسسته است. این الگوریتم از نوع یادگیری خارج از سیاست است؛ به این معنا که عامل با دنبال کردن سیاستی که سیاست انتخابی خود نیست به بهبود عملکرد خود می‌پردازد.

DDPG از دو شبکه عصبی استفاده می‌کند: یک شبکه برای تخمین ارزش عمل که به آن نقاد می‌گوییم، و یک شبکه برای تخمین سیاست و انتخاب عمل که به آن بازیگر می‌گوییم. شبکه نقاد مشابه با شبکه Q در DQN است، با این تفاوت که حالت و عمل را ورودی گرفته و ارزش عمل را خروجی می‌دهد. و شبکه سیاست یک شبکه عصبی است که خروجی آن عمل است. این الگوریتم مشابه با DQN از بافر تجربه برای ذخیره تجربیات عامل، و از شبکه‌های هدف برای هر دو شبکه نقاد و بازیگر استفاده

<sup>32</sup>Deep Deterministic Policy Gradient

می‌کند. البته در این الگوریتم، به جای کپی وزن‌های شبکه به شبکه هدف، از میانگین وزن‌دار استفاده می‌شود تا به روزرسانی آرام‌تر انجام شود و ثبات بیشتری داشته باشد.



شکل ۲-۴ شبکه‌های عصبی بازیگر و نقاد در الگوریتم DDPG و مقایسه با DQN

رایج است که شبکه ارزش عمل را با  $Q(s, a; \phi)$  و شبکه سیاست را با  $\mu(s; \theta)$  نشان دهیم. فرآیند یادگیری این الگوریتم به صورت زیر است:

۱. ابتدا وزن‌های شبکه‌های نقاد و بازیگر را به صورت تصادفی مقداردهی می‌کنیم و آن‌ها را به ترتیب  $\phi$  و  $\theta$  می‌نامیم.

۲. بافر تجربه را به طول  $N$  مقداردهی می‌کنیم.

۳. در هر گام زمانی، حالت را به شبکه بازیگر می‌دهیم و عملی که این شبکه پیشنهاد می‌دهد را انجام می‌دهیم.

۴. پاداش دریافتی و حالت بعد عمل را به بافر تجربه اضافه می‌کنیم.

۵. یک ریزدسته از بافر تجربه را انتخاب کرده و با استفاده از فرمول زیر، مقدار واقعی ارزش حالت را (طبق الگوریتم TD) محاسبه می‌کنیم:

$$y_j = \begin{cases} r_j & \text{اگر } s_{j+1} \text{ حالت پایانی قسمت باشد} \\ r_j + \gamma \times Q(s_{j+1}, \mu(s_{j+1}; \theta^-); \phi^-) & \text{در غیر این صورت} \end{cases} \quad (2-16)$$

۶. خطای شبکه نقاد را به صورت زیر محاسبه می‌کنیم:

$$L(\phi) = \mathbb{E}[(y_j - Q(s_j, a_j; \phi))^2] \quad (۱۷-۲)$$

و وزن‌های شبکه نقاد را به‌روز می‌کنیم:

$$\phi_{t+1} = \phi_t - \alpha \nabla_{\phi} L(\phi) \quad (۱۸-۲)$$

۷. خطای شبکه بازیگر را به صورت زیر محاسبه می‌کنیم:

$$L(\theta) = -\mathbb{E}[Q(s, \mu(s; \theta); \phi)] \quad (۱۹-۲)$$

و وزن‌های شبکه بازیگر را به‌روز می‌کنیم:

$$\theta_{t+1} = \theta_t + \alpha \nabla_{\theta} L(\theta) \quad (۲۰-۲)$$

۸. در نهایت هر چند گام یک بار، وزن‌های شبکه‌های هدف را به‌روز می‌کنیم:

$$\theta_{t+1}^- = \tau \theta_t + (1 - \tau) \theta_t^- \quad (۲۱-۲)$$

$$\phi_{t+1}^- = \tau \phi_t + (1 - \tau) \phi_t^- \quad (۲۲-۲)$$

همانطور که مشاهده‌شد، این الگوریتم تعداد زیادی پارامتر دارد که باید به‌درستی تنظیم شوند تا به عملکرد بهینه برسد. از جمله این پارامترها می‌توان به نرخ یادگیری ( $\alpha$ ) که می‌تواند برای شبکه‌های نقاد و بازیگر متفاوت باشد، نرخ تخفیف ( $\gamma$ )، فرکانس به‌روزرسانی شبکه‌ها (هر چند گام یک بار ریزدسته انتخاب می‌کنیم و مراحل ۵ تا ۸ را طی می‌کنیم)، فرکانس به‌روزرسانی وزن‌های هدف (هر چند گام یک بار وزن‌های شبکه‌های هدف را به‌روزرسانی می‌کنیم)، اندازه ریزدسته، طول بافر تجربه، و ضریب

به روزرسانی شبکه‌های هدف  $\tau$  اشاره کرد.

## ۷-۲ جمع‌بندی

# فصل سوم

## محیط فوتبال ريوکاپ



در این بخش به معرفی مفاهیم پایه محیط فوتبال روبوکاپ و توصیف چگونگی عملکرد آن می‌پردازیم. این مفاهیم شامل قوانین بازی، رفتارهای ممکن، کد پایه ایجنت، حالت پناستی، و روش اجرای بازی می‌باشند.

## ۳-۱ معرفی لیگ

### ۳-۱-۱ اهداف لیگ

لیگ روبوکاپ مجموعه مسابقاتی سالانه است که قصد دارد با کمک فوتبال، به پیشرفت زمینه‌های رباتیک و هوش مصنوعی کمک کند. علت انتخاب فوتبال به عنوان محیط مسابقه، این است که فوتبال یکی از محیط‌هایی است که می‌تواند مسائل مختلفی از جمله تصمیم‌گیری، هماهنگی، بینایی ربات، و ارتباط را در بر داشته باشد. یکی از اهداف بلندمدت لیگ، ساختن ربات‌هایی است که بتوانند تا سال ۲۰۵۰، تیمی از انسان‌ها را به چالش بکشند.

یکی از لیگ‌های این مسابقات، لیگ شبیه‌ساز فوتبال دو بعدی<sup>۱</sup> است. همانطور که از اسم لیگ پیداست، مسابقه حالت دو بعدی دارد، به این منظور که فضا حالت دید از بالا دارد، و همه حرکات بازیکنان و توپ روی سطح زمین انجام می‌شود. تمرکز اصلی این لیگ، تصمیم‌گیری و استراتژی، و ساختن الگوریتم‌های مناسب با محیط‌های چندعامله با دید ناقص است.

### ۳-۱-۲ ویژگی‌های محیط و قوانین بازی

مشابه با فوتبال واقعی، هر بازی متشکل از دو تیم ۱۱ نفره است که هر کدام از این نفرات یک عامل مستقل می‌باشد. بازی در یک مستطیل به ابعاد ۱۱۵ در ۶۸ انجام می‌شود، و هر تیم یک دروازه‌بان دارد که در هر طرف زمین قرار دارد. بازی در ۶۰۰۰ گام<sup>۲</sup> انجام می‌شود که به دو نیمه تقسیم شده، و هر گام ۱۰۰ میلی‌ثانیه طول می‌کشد. سرور مسابقات در هر گام، اطلاعات مربوط به وضعیت بازی را به تمامی عامل‌ها می‌فرستد، و عامل‌ها باید تصمیم‌گیری خود را بر اساس این اطلاعات انجام دهند.

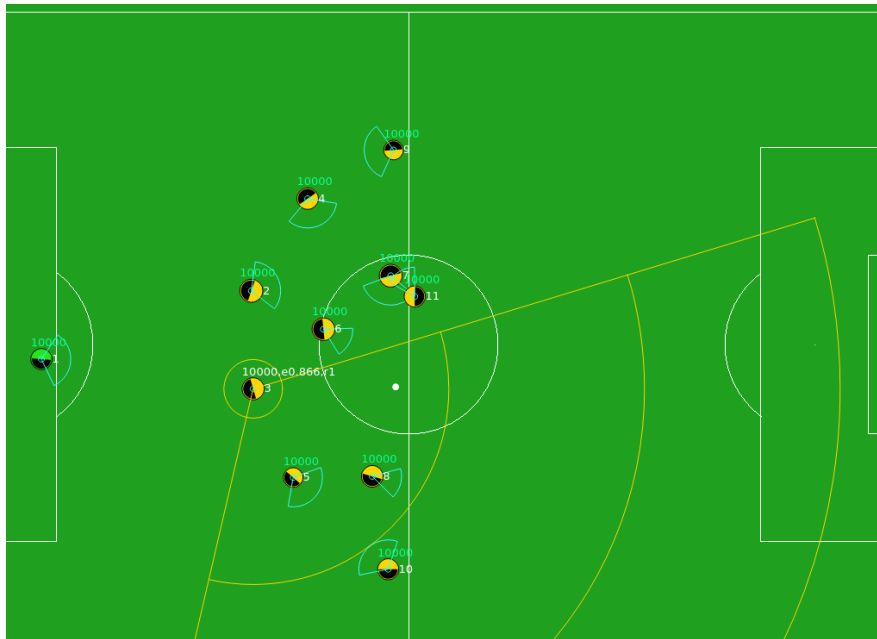
بازیکنان دید محدودی دارند. با توجه به زاویه گردنی که تنظیم کرده‌اند، یک قطاع از زمین را مشاهده می‌کنند و فقط اطلاعات مربوط به بازیکنان و اشیاء داخل این قطاع را دریافت می‌کنند در شکل ۳-۱ نمونه دید یک بازیکن نشان داده شده است. اطلاعات بازیکنانی که داخل دید نیستند، باید از حافظه بازیکن و یا از ارتباط با سایر بازیکنان به دست آید. شایان به ذکر است که مشاهدات هر بازیکن، دارای

<sup>1</sup>2D Soccer Simulation

<sup>2</sup>Cycle

کمی خطای اندازه‌گیری است و به طور کامل دقیق نیست. به طور مشابه، ضربات بازیکن به توپ نیز کمی خطا دارند و کامل دقیق نیستند.

هر عامل باید به صورت مستقل تصمیم‌گیری کند و ارتباط بسیار ناچیزی با سایر عاملین دارد. در صورتی که دو عامل بخواهند ارتباط برقرار کنند، گیرنده باید از قبل توجه خود را به فرستنده تنظیم کرده‌باشد، و فرستنده حداکثر ۸ بایت می‌تواند ارسال کند.



شکل ۱-۳ نمونه دید یک بازیکن. در این مثال، بازیکن ۳، بازیکن ۵ و ۸ و ۱۰ و توپ را می‌بیند و از موقعیت سایر بازیکنان اطلاعی ندارد.

هر عامل یک رزرو انرژی و یک منبع انرژی دارد. در صورت پر نبودن انرژی بازیکن، انرژی به نرخ ثابتی از رزرو به منبع منتقل می‌شود. در صورت اتمام انرژی منبع، رفتارهای بازیکن از جمله حرکت کردن کندتر می‌شوند. رفتارهای بیشتری (مانند تکل زدن، خطا و ضربه آزاد، کارت زرد و قرمز و ...) نیز در محیط به کمک مدل‌های ریاضیاتی پیاده‌سازی شده‌است که خارج از دامنه این پروژه می‌باشد.

## ۲-۳ کد پایه ایجنت

در سال‌های اول مسابقات، هر تیمی کد عامل خود را از صفر می‌نوشت، که میزان در دسترس بودن لیگ را بسیار پایین آورده‌بود. با توجه به یکسان بودن بخش عمده‌ای از نیازمندی‌های تیم‌ها، همچون نیاز به اتصال به سرور، نیاز به تفکیک وظایف بازیکنان به دروازه‌بان و مدافع و مهاجم، نیاز به موقعیت‌یابی اشیا

در زمین، توابع هندسی، و ...، هیدهیسا آکیام از تیم هلیوس در سال ۲۰۱۲ تصمیم به ساختن یک کد پایه به صورت متن باز محض استفاده سایر تیم‌ها گرفت. این کد پایه، که کد پایه ایجنت نام دارد، زیربنای ۱۳ تیم از ۱۵ شرکت کننده سال اخیر ربوکاپ بوده، و نقطه شروع اکثر کسانی است که قصد فعالیت در این فضا را دارند. این کد پایه همچنان در حال به‌روزرسانی و تقویت شدن است<sup>۳</sup>، و خود منشا سایر کد پایه‌های به اشتراک گذاشته شده همچون کد پایه گلایدر و کد پایه سائرس است.

با استفاده از این کد پایه، توسعه‌دهندگان می‌توانند سطح کدزدن خود را از سطوح پایین، به سطح استراتژی و تاکتیک و تصمیم‌گیری منتقل کنند. به طور مثال بدون استفاده از یک کد پایه، برای حرکت عامل به مرکز زمین باید کد اتصال به سرور و موقعیت‌یابی را پیاده‌سازی کنیم. سپس درگیر مسائلی همچون محاسبه شتاب بازیکن، نیرویی که باید اعمال شود، بهینه‌ترین مسیر حرکت (چرخیدن و دویدن یا دویدن مورب) و ... شوند. در حالی که این عمل، به صورت یک دستور سطح بالا در کد پایه ایجنت قابل اجراست. در بخش بعدی به تمام رفتارهای ممکن در این کد پایه و سطح‌بندی رفتارها می‌پردازیم.

### ۳-۳ معرفی رفتارهای ممکن

هر عامل در هر لحظه، باید تصمیم‌گیری خود را به سرور بفرستد. این تصمیم‌گیری می‌تواند شامل انجام یکی از رفتارهای ممکن باشد. در هر لحظه، عامل می‌تواند گردن خود را بچرخاند تا اطراف خود را ببیند، و همزمان یکی از پنج رفتار ضربه به توپ، حرکت بدن، چرخش بدن، تکل، یا گرفتن توپ را انجام دهد. رایج است که رفتارهای ممکن و پیاده‌سازی شده را به سه طبقه تقسیم‌بندی کنیم که در پایین‌ترین سطح، رفتارهای سطح سرور، و در بالاترین سطح، رفتارهای سطح استراتژیک و فکری قرار می‌گیرند.

### ۱-۳-۳ رفتارهای سطح پایین

در پایین‌ترین سطح، رفتارها مستقیماً معادل با رفتارهای مورد پذیرش سرور می‌باشند. این رفتارها شامل اعمال نیرو روی توپ در یک راستا (نسبت به بدن بازیکن) و نیروی خاص، چرخاندن بازیکن به یک راستا، اعمال نیروی حرکت بازیکن در راستا و نیرو، تکل زدن در یک راستا، و گرفتن توپ (در صورتی که بازیکن دروازه‌بان باشد) می‌باشند.

<sup>۳</sup><https://github.com/helios-base/helios-base>

### ۲-۳-۳ رفتارهای سطح متوسط

در این سطح، رفتارها ساده‌سازی شده‌اند تا استفاده از آن‌ها برای تصمیم‌گیری راحت‌تر باشد. به طور مثال رفتار حرکت به سمت یک نقطه خاص از زمین، رفتار ضربه زدن به توپ با سرعت دلخواه در راستای غیرنمایی، رفتار ضربه زدن توپ به صورت چندضرب به یک نقطه خاص، و... را می‌توان از رفتارهای این سطح معرفی کرد.

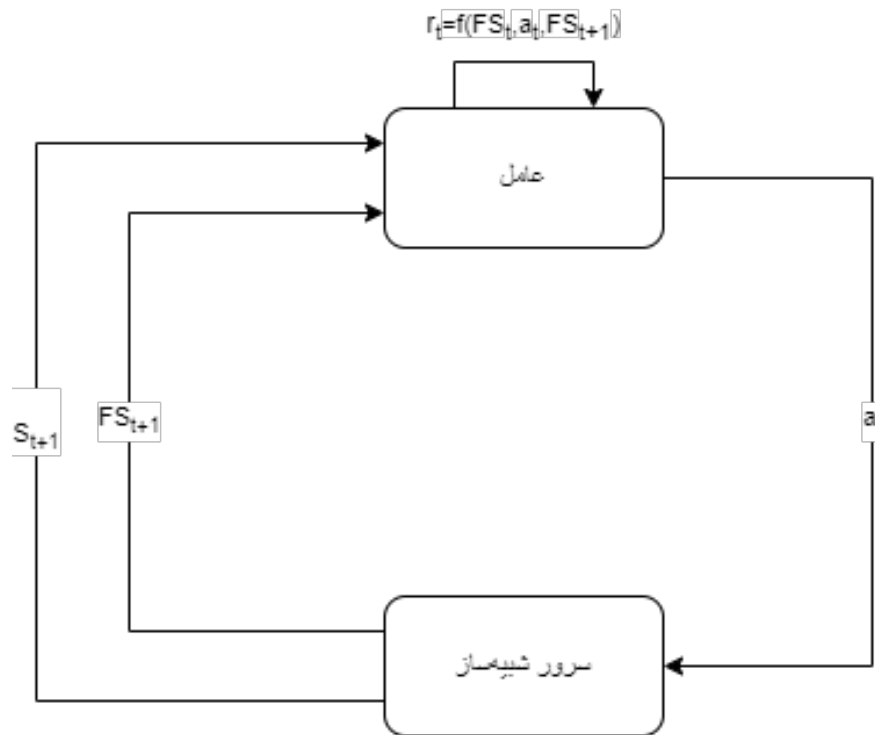
### ۳-۳-۳ رفتارهای سطح بالا

در این سطح، رفتارها حالت استراتژی و فکر کردن دارند، و به صورت انتظاری و مشابه با فوتبال واقعی می‌باشند. به طور مثال، حرکت به سمت نقطه‌ی فرمیشن، حرکت برای قطع توپ، زدن شوت (در صورت ممکن بودن شوت موفق)، و... از رفتارهای سطح بالا می‌باشند.

### ۴-۳ حالت دید کامل

همانطور که گفته‌شد، محیط اطلاعات کامل را در اختیار عامل قرار نمی‌دهد. با توجه به اینکه قصد تغییر کد سرور مسابقات و افزودن پاداش به سرور را نداریم، این محاسبات درون خود عامل باید صورت بگیرند. خوش‌بختانه سرور مسابقات قابلیت ارسال محیط در حالت دید کامل<sup>۴</sup> را دارد که در این حالت، همه اطلاعات بدون نویز به عامل ارسال می‌شوند. داخل کد ایجنت می‌توان انتخاب کرد که در صورت دریافت حالت دید کامل آن‌را جایگزین دید ناقص کند، یا اینکه هر دو را کنار هم نگاه‌دارد. ما در فرآیند یادگیری تقویتی از حالت دوم استفاده می‌کنیم تا با کمک دید کامل محاسبات پاداش را انجام دهیم، و با حالت دید ناقص تصمیم‌گیری کنیم. در عکس ۲-۳ می‌توان نحوه مشاهده کرد که پاداش داخل عامل محاسبه‌شده، و تابعی از حالت دید کامل (Fullstate یا به مختصر FS) است، و محیط هر دو حالت دید را برای عامل ارسال می‌کند. این نمودار در واقع حالت اصلاح‌شده ۱-۲ است که خاص‌منظوره این پروژه می‌باشد.

<sup>۴</sup>Fullstate WorldModel



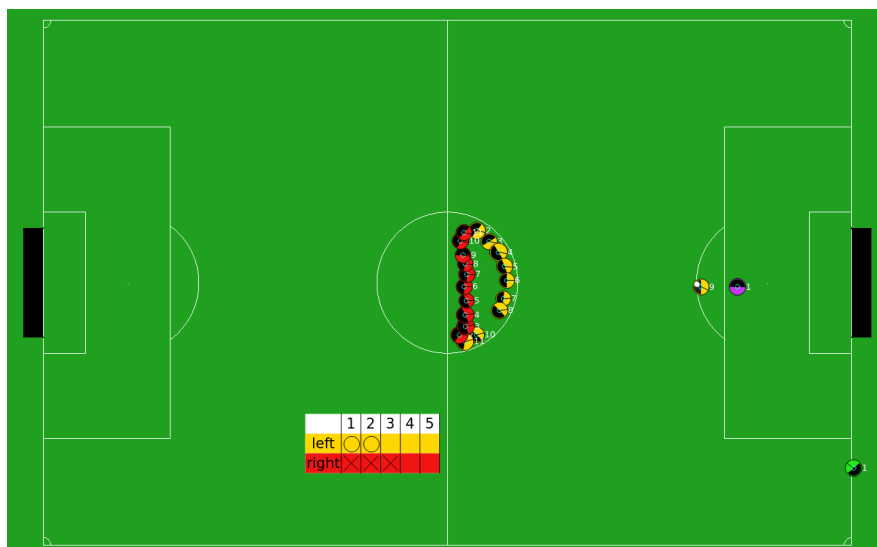
شکل ۳-۲ ارتباط بین محیط و عامل، و روش محاسبه پاداش

### ۳-۵ معرفی حالت پنالتی

در صورت تساوی بازی در ده دقیقه زمان عادی، و تداوم تساوی پس از وقت اضافی، بازی به حالت پنالتی می‌رود. با توجه به دو بعدی بودن محیط، در صورتی که مشابه با فوتبال واقعی، پنالتی به صورت تک ضرب باشد، می‌توان برای آن استراتژی قطعی ارائه داد. به همین منظور، پنالتی در فضای لیگ دوبعدی، به صورت تک به تک با دروازه‌بان است.

بازیکن مهاجم، ۱۵ ثانیه فرصت دارد تا توپ را به گل تبدیل کند. دروازه‌بان باید در این مدت تلاش کند جوری موقعیت‌گیری کند که حریف نتواند شوت منجر به گل داشته‌باشد. در صورتی که زمان مهاجم تمام شود، توپ به بیرون محیط بازی برود، یا توپ توسط دروازه‌بان گرفته‌شود، دروازه‌بان برنده می‌شود. در صورتی که توپ وارد گل شود، مهاجم برنده می‌شود.

با توجه به تک‌عامله بودن، و محدود بودن زمان محیط، می‌توان از روش‌های یادگیری تقویتی برای یادگیری استراتژی‌های بهینه برای پنالتی استفاده کرد. در شکل ۳-۳ یک نمونه از حالت پنالتی را مشاهده می‌کنید.



شکل ۳-۳ نمونه‌ای از بازی در حالت پنالتی

### ۶-۳ کار با مربی تمرینی برای تولید محیط قابل تکرار

با توجه به اینکه برای یادگیری تقویتی، نیاز به بازگرداندن محیط به حالت اولیه را داریم، می‌توانیم از مربی تمرینی<sup>۵</sup> استفاده کنیم. در صورتی که در تنظیمات سرور، این حالت فعال شده باشد، می‌توان کدی نوشت که در شرایط دلخواه، توپ و بازیکنان را جابه‌جا کنیم، حالت بازی را تنظیم کنیم، انرژی بازیکنان را پر کنیم و ... . در فصل‌های آینده از این امکان، برای ساختن یک رابط استاندارد یادگیری تقویتی برای محیط فوتبال استفاده خواهیم کرد.

### ۷-۳ جمع‌بندی

در این فصل، محیط فوتبال ربوکاپ دوبعدی به عنوان یک ابزار مطالعاتی برای تقویت پیشرفت‌ها در زمینه هوش مصنوعی و رباتیک معرفی شد. این محیط، با ارائه یک پلتفرم مسابقه‌ای مبتنی بر قوانین فوتبال واقعی، فرصت‌هایی برای توسعه و آزمایش الگوریتم‌های تصمیم‌گیری، هماهنگی، و استراتژی در محیط‌های چندعامله فراهم می‌کند.

معرفی این محیط شامل توضیح قوانین بازی، نحوه ارتباط و تصمیم‌گیری عامل‌ها، و همچنین ساختار کد پایه ایجنت بود که توسعه‌دهندگان را قادر می‌سازد تا تمرکز خود را بر روی تاکتیک و تصمیم‌گیری‌های

<sup>5</sup>Trainer

سطح بالای مشابه با فوتبال معطوف دارند.

یکی از نوآوری‌های کلیدی در این فصل، استفاده از حالت دید کامل در کنار دید ناقص است که به عامل‌ها امکان می‌دهد تا پاداش‌ها را بر اساس اطلاعات کامل محیط محاسبه کنند، در حالی که تصمیم‌گیری بر اساس دید ناقص انجام می‌شود. این رویکرد یک قدم مهم در راستای افزایش قابلیت‌های یادگیری تقویتی در محیط‌هایی با دید ناقص است.

همچنین، با معرفی حالت پنالیتی و استفاده از مربی تمرینی برای تولید محیط‌های قابل تکرار، زمینه‌های لازم برای پیاده‌سازی الگوریتم‌های یادگیری تقویتی فراهم شده است. این امکانات به پژوهشگران اجازه می‌دهد تا در یک محیط مسابقه‌ای استاندارد، استراتژی‌های بهینه‌سازی شده را آزمایش و توسعه دهند.

## فصل چهارم

### پیاده‌سازی و آماده‌سازی محیط استاندارد



## ۴-۱ زیرساخت پایتونی تهاجم نصف زمین

در سال ۲۰۱۵، پروژه تهاجم نصف زمین<sup>۱</sup> یا همان اچ‌اف‌او تلاش کرد محیط پایتونی برای یادگیری تقویتی در فوتبال دو بعدی ایجاد کند. کد سرور به گونه‌ای تغییر یافته بود، که بازیکن بتواند به سرور دستور اعمال رفتارها و رفتن به گام بعدی را بدهد.

این روش متکی بر قابلیت‌های لیب‌سی<sup>۲</sup> بود، که از آن برای ارتباط با کد cpp استفاده می‌کرد. لیب‌سی اجازه می‌دهد که اگر کلاس معادل پایتونی و cpp وجود داشته باشد، از سمت پایتون می‌توان به آن دسترسی داشت. در این پروژه، عامل پایتونی به عنوان ایجنت دسترسی به محیط دارد، که با آن از کد cpp محیط را درخواست می‌کند. محیط در دو سطح بالا و پایین قابل درخواست است؛ حالت سطح پایین ۶۰ ویژگی محیط، و حالت سطح بالا ۹ ویژگی را به صورت یک آرایه صفر تا یک به عامل پس می‌دهد. عامل سپس تصمیم خود را اخذ کرده، و به سرور دستور اجرای آن را می‌دهد. معایب استفاده از این محیط عبارت بود از:

- دشواری در تغییر محیط عامل به صورت دلخواه
- تفاوت زیاد بین محیط آماده‌شده در ایجنت، که مورد استفاده اکثر تیم‌هاست، و محیط اچ‌اف‌او
- عقب‌ماندگی محیط به علت تغییر کدهای سرور و به روز نبودن با تغییرات شبیه‌ساز
- دشواری نصب، به علت پیش‌نیازهای قدیمی همچون Qt4
- عدم امکان استفاده از تیم‌های جدید برای حریف یادگیری به علت نسخه قدیمی سرور

## ۴-۲ کد پایه پایرس

با توجه به دشواری استفاده از یادگیری ماشین در CPP، ما در تیم سایرس تصمیم به پیاده‌سازی یک کد پایه معادل با کد پایه Agent2D در پایتون گرفتیم، که نام آن پایرس<sup>۳</sup> است. [؟] این تلاش‌ها از سال ۲۰۱۹ آغاز شد و در نهایت بعد از سه سال و حدود ۲۵ هزار خط کد، پروژه به حالت پایدار و قابل استفاده رسیده است.<sup>۴</sup>

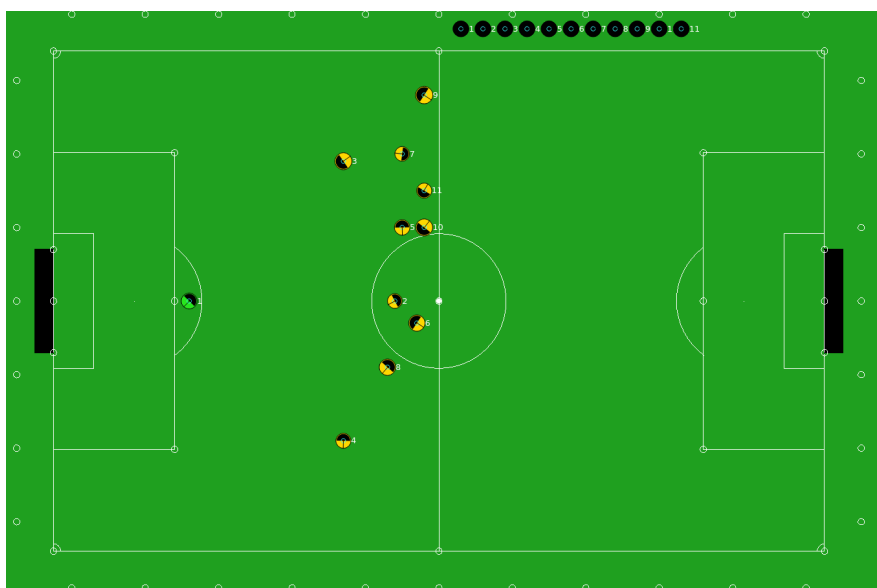
<sup>۱</sup>Half Field Offense (HFO)

<sup>۲</sup>LibC

<sup>۳</sup>Pyrus

<sup>۴</sup><https://github.com/Cyrus2D/Pyrus2D>

امید بر آن بود که با توجه به اتکا بر الگوریتم‌های یادگیری ماشین به جای الگوریتم‌های درخت و گراف، بتوان کندی زبان پایتون را جایگزین کرد. در کد Agent2D، قبل از رسیدن به گام تصمیم‌گیری، محاسبات زیادی رخ می‌دهد تا بازیکن موقعیت خود و سایر بازیکنان را بیابد، سریع‌ترین بازیکنان به توپ، موقعیت آفساید و... را به دست آورد. با توجه به کند بودن پایتون در مقایسه با cpp، همین محاسبات مقدماتی نیز بسیار زمان‌بر بودند و بخش بزرگی از ۱۰۰ میلی‌ثانیه موجود برای هر گام بازی را اشغال می‌کردند. به همین منظور، پایرس جایگزین مناسبی برای پایتونی کردن کد نبود.



شکل ۴-۱ پرچم‌های کنار زمین، که برای موقعیت‌یابی استفاده می‌شوند.

## ۳-۴ کد پایه جی‌آرپی‌سی

با توجه به سریع بودن cpp برای پیش‌پردازش، و کاربردی و راحت بودن پایتون برای یادگیری تقویتی، تصمیم گرفتیم به کمک پروتکل فراخوانی تابع از راه دور<sup>۵</sup>، این دو زبان را به هم متصل کنیم. در این پروژه از چارچوب فراخوانی تابع از راه دور گوگل، یا همان جی‌آرپی‌سی<sup>۶</sup> استفاده می‌کنیم. در این چارچوب، ابتدا یک فایل پروتو<sup>۷</sup> باید تعریف کرد، که حاوی مشخصات اشیاء و امضای توابع است. سپس با استفاده از کامپایلر پروتو می‌توان کدهای سرور و کلاینت را به زبان‌های دلخواه تبدیل کرد. سپس با اضافه کردن کد تولیدشده توسط کامپایلر پروتو می‌توان درون کلاینت از توابع به‌گونه‌ای

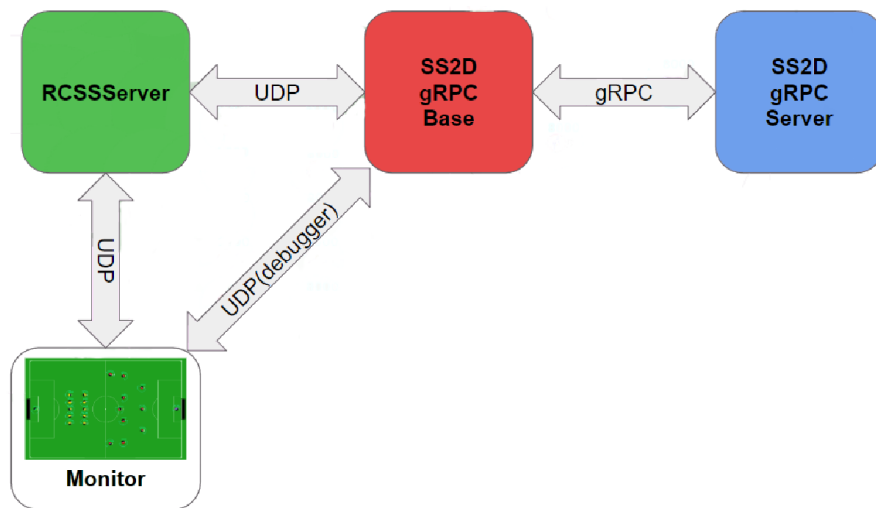
<sup>۵</sup>Remote Procedure Call (RPC)

<sup>۶</sup>gRPC (google Remote Procedure Call)

<sup>۷</sup>Proto

استفاده کرد که گویا داخل خود کد قرار دارند. از مهم‌ترین مزایای این پروتکل ارتباطی، مستقل از زبان بودن آن، و نوشتن بسته‌های پیام به صورت صفر و یکی (در مقابل متنی) است که سربار زمانی نوشتن و خواندن پیام را ناچیز می‌کند. همچنین به علت استفاده از پروتکل اچ‌تی‌تی‌پی<sup>۸</sup> برای ارسال بسته‌ها روی شبکه، می‌توان از ضمانت‌های دریافت پیام نیز استفاده کرد.

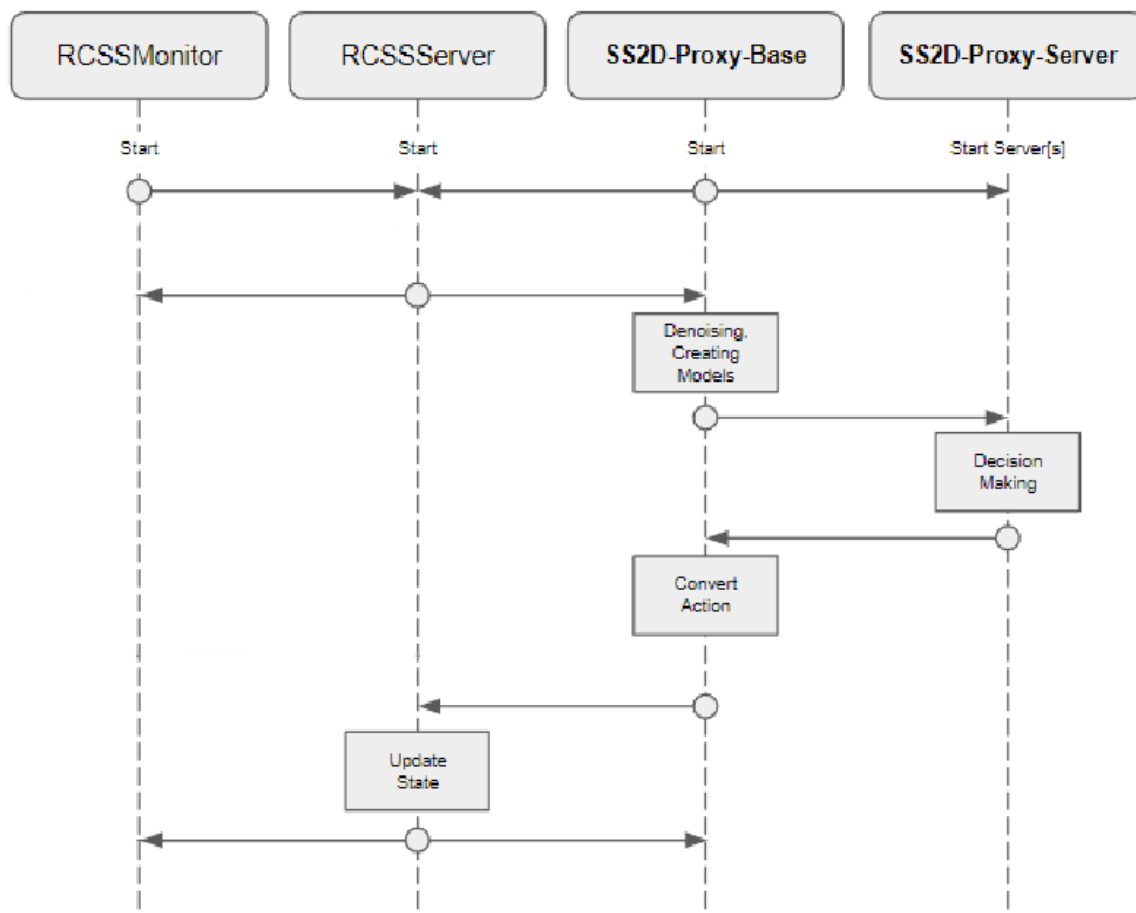
برای تشکیل اتصال بین این دو زبان کد پایه Agent2D را به گونه‌ای تغییر دادیم که پس از اتمام مراحل پیش‌پردازش و آماده شدن مدل دنیای بازیکن<sup>۹</sup> در cpp، اطلاعات آن را با یک درخواست جی‌آرپی‌سی به یک سرور تصمیم‌گیری پایتونی بفرستد. در شکل ۲-۴ و ۳-۴ نحوه‌ی اتصال اجزای مسابقه، و روند عمومی تصمیم‌گیری نمایش داده شده‌است.



شکل ۲-۴ پروتکل‌های استفاده‌شده برای ارتباط بین اجزای مسابقه

<sup>۸</sup>HTTP

<sup>۹</sup>WorldModel



شکل ۴-۳ نحوه کارکرد و اتصال کد پایه جی آر پی سی به سرور مسابقات و نمایشگر بازی

## ۴-۴ محیط استاندارد جیم

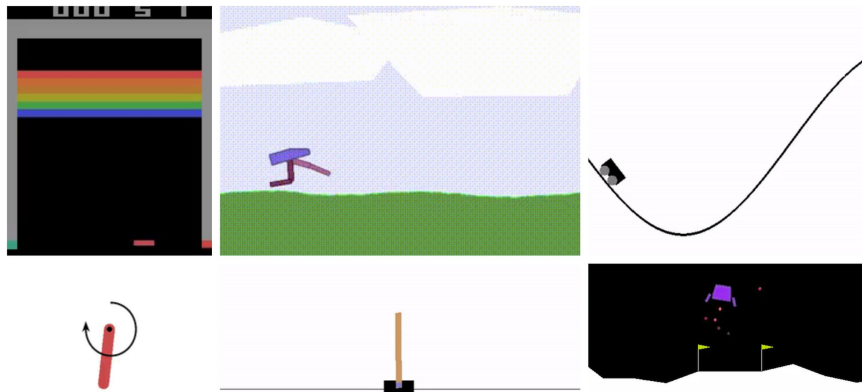
پلتفرم جیم<sup>۱۰</sup> که توسط گروه اپن‌ای‌ای<sup>۱۱</sup> توسعه داده شده‌است، یک محیط استاندارد برای یادگیری تقویتی است. این پلتفرم شامل محیط‌های گوناگون از پیش آماده شده است که می‌توانند برای یادگیری تقویتی استفاده شود. این محیط همچنین به ما قابلیت تعریف محیط‌های جدید را می‌دهد. در سال ۲۰۲۳، مدیریت و نگهداری این پلتفرم به شرکت فاراما<sup>۱۲</sup> واگذار شد و از آن زمان، این پلتفرم به نام پلتفرم جیمنازیوم<sup>۱۳</sup> شناخته می‌شود. تبدیل محیط یادگیری به این نوع محیط، نه تنها به ما قابلیت استفاده از کتابخانه‌های پیشین و کمک گرفتن از ابزارهای از پیش تعبیه شده را می‌دهد، بلکه امکان

<sup>10</sup>Gym

<sup>11</sup>OpenAI

<sup>12</sup>Farama

<sup>13</sup>Gymnasium



شکل ۴-۴ نمونه‌ای از محیط‌های آماده شده در پلتفرم جیم

به اشتراک‌گذاری راحت این محیط به سایر محققان را نیز فراهم می‌کند.

#### ۴-۴-۱ توصیف رابط و توابع موجود

برای تعریف یک محیط جدید در جیم، باید توابع استاندارد آن را پیاده‌سازی کنیم. به این منظور، لازم است با روش کارکرد هر جزء محیط آشنا شویم. لازم به ذکر است که تفاوت‌های جزئی بین نسخه اوپن‌ای و فاراما وجود دارد، که در اینجا به توصیف نسخه فاراما می‌پردازیم.

##### تابع گام برداشتن

این تابع، یک عمل را به عنوان ورودی می‌گیرد، و وضعیت جدید، پاداش این عمل، و اطلاعاتی همچون اتمام بازی، یا تمام شدن وقت عامل را به عنوان خروجی برمی‌گرداند.

##### تابع شروع مجدد

این تابع، محیط را به حالت اولیه باز می‌گرداند، و حالت جدید را به عامل برمی‌گرداند. در صورتی که شروع بازی شامل المان‌های تصادفی باشد، ورودی این تابع باید کاوش تصادفی<sup>۱۴</sup> را نیز به عنوان ورودی بگیرد.

##### انواع فضاهای حالت و عمل

در حین تعریف محیط، باید فضاهای حالت و عمل را نیز تعریف کنیم. هر فضا می‌تواند از انواع زیر باشد:

<sup>14</sup>Random Seed

- گسسته<sup>۱۵</sup>: به ازای ورودی  $n$ ، در این فضا می‌توانیم از اعداد صحیح از  $0$  تا  $n-1$  استفاده کنیم. این حالت بیشتر برای فضای خروجی استفاده می‌شود.
- گسسته چندتایی<sup>۱۶</sup>: این فضا مشابه فضای گسسته است، با این تفاوت که می‌توانیم چند عدد از این فضا را به عنوان خروجی بگیریم. مشابه با گسسته است، ولی جای ورودی یک عدد، یک آرایه از اعداد است؛ هر عدد نشان‌دهنده تعداد حالت‌های ممکن در هر بعد است.
- فضای بسته<sup>۱۷</sup>: نشان‌دهنده یک فضای چندبعدی پیوسته است، که باید حد پایین و بالای هر بعد را مشخص کنیم.
- دوتایی چندبعدی<sup>۱۸</sup>: حالت چندبعدی است، که هر بعد مقدار صفر یا یک می‌تواند به خود بگیرد. به کمک فضاهای فوق، می‌توان فضاهایی از ترکیب حالت‌های فوق تعریف کرد، که به انواع زیر است:
  - دیکشنری<sup>۱۹</sup>: یک دیکشنری پایتونی از فضاهای مختلف است.
  - چندتایی (توپل)<sup>۲۰</sup>: یک چندتایی مرتب از فضاهای مختلف است.
  - دنباله<sup>۲۱</sup>: یک دنباله از فضاهای مختلف است.

## تابع پایان

در اتمام کار با محیط، این تابع فراخوانی می‌شود تا محیط بسته شود، و منابع آن آزاد شوند.

## تابع ترسیم

با توجه به اینکه اکثر محیط‌های جیم، به کمک کتابخانه‌هایی مانند PyGame ترسیم می‌شوند، این تابع برای ترسیم محیط در هر گام استفاده می‌شود. همچنین به عنوان ورودی این تابع می‌توان معین کرد که ترسیم به چه گونه‌ای انجام شود، چرا که ممکن است در حین آموزش برای تسریع عملیات، نیازی به

<sup>15</sup>Discrete

<sup>16</sup>MultiDiscrete

<sup>17</sup>Box

<sup>18</sup>MultiBinary

<sup>19</sup>Dict

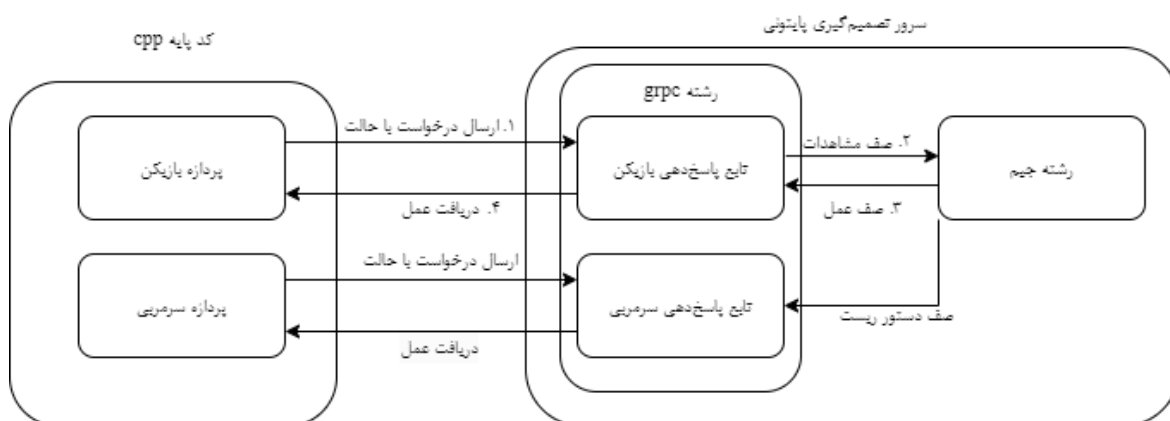
<sup>20</sup>Tuple

<sup>21</sup>Sequence

ترسیم نباشد. از آن جا که در محیط ما، وظیفه ترسیم با نمایش‌گر خارجی انجام می‌شود، این تابع برای ما معنی ندارد.

#### ۴-۴-۲ نحوه اقدام با فضای ربوکاپ

برای پیاده‌سازی محیط جیم در فضای ربوکاپ، از کد پایه gRPC استفاده می‌کنیم. فرآیند سرو کردن درخواست‌های این پروتکل را با کمک کتابخانه‌های استاندارد پایتون به یک رشته<sup>۲۲</sup> جدا منتقل می‌کنیم، که با کمک دو صف با رشته محیط جیم ارتباط دارد. در هنگام دریافت حالت از بخش پایه cpp، رشته سرویس‌دهنده حالت را داخل صف مشاهدات می‌گذارد، و منتظر می‌شود که رشته جیم، تصمیم عامل را بر اساس مشاهده ثبت‌شده در صف عمل‌ها قرار دهد. صف سومی نیز برای سرمربی وجود دارد، که در صورت نیاز به شروع مجدد محیط، از سمت جیم پر می‌شود. کد سرمربی تمرینی در هر لحظه از اجراش چک می‌کند و در صورت خالی نبودن صف، دستورات لازم برای از نو کردن محیط را می‌فرستد. در صورت خالی بودن صف نیز عمل خالی به کد پایه پس می‌فرستد.

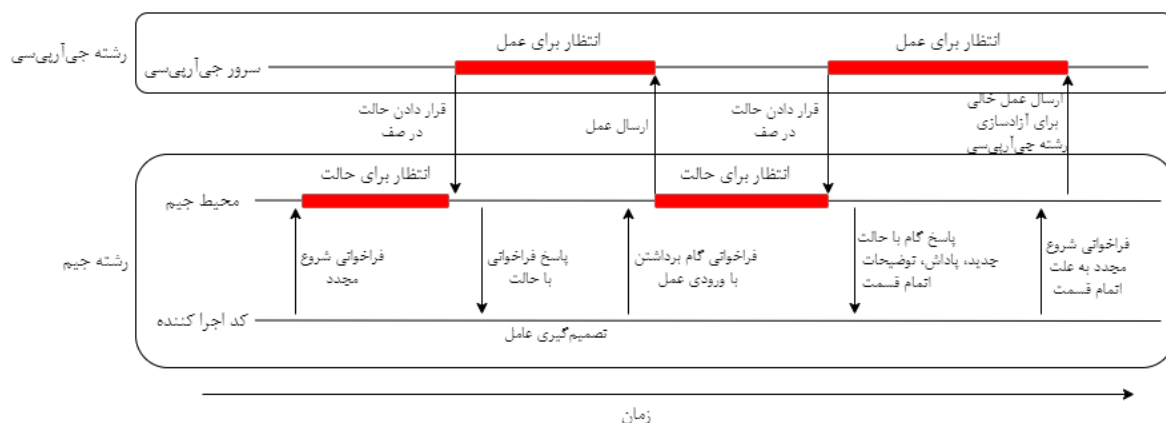


شکل ۴-۵ نحوه کلی ارتباط کد پایه، تصمیم‌گیرنده جی‌آرپی‌سی و جیم

در شکل ۴-۶ به طور دقیق می‌توان ترتیب فراخوانی توابع، و نحوه استفاده از محیط جیم را مشاهده کرد. کد اجرا کننده<sup>۲۳</sup> دقیقاً فرم مشابه سایر محیط‌های جیم را دارد و از این رابط به طور کامل طبیعت می‌کند. بخش‌های قرمز رنگ به معنی قفل بودن رشته‌ی اجرایی و در انتظار ماندن برای محتویات صف است.

<sup>22</sup>Thread

<sup>23</sup>Driver Code



شکل ۴-۶ ترتیب فراخوانی توابع برای اجرای جیم با سرور روبوکاپ

داخل شکل، یک قسمت که فقط یک گام تصمیم‌گیری دارد را می‌توان دید. در واقعیت مراحل ((گام برداشتن)) تا قبل از ((فراخوانی شروع مجدد)) به تعداد گام‌های قسمت، تکرار می‌شوند. برای محیط پنالتی، در صورتی که توپ قابل ضربه‌زدن نباشد، حرکت صحیح همیشه قطع توپ با حداکثر سرعت برای ضربه مجدد است. بنابراین فقط حالت را در شرایطی به جیم ارسال می‌کنیم که توپ قابل ضربه‌زدن باشد.

### ۴-۴-۳ فضای حالت و فضای عمل عامل

#### فضای حالت

فضای حالت انتخاب‌شده یک بسته ۹ بعدی است، که شامل ویژگی‌های زیر می‌باشد: موقعیت قطبی بازیکن نسبت به مرکز دروازه (اندازه و زاویه)، موقعیت دکارتی توپ (ایکس و ایگرگ)، موقعیت دکارتی دروازه‌بان (ایکس و ایگرگ)، زاویه نسبی حریف نسبت به توپ، موقعیت قطبی دروازه‌بان نسبت به بازیکن (اندازه و زاویه).

علت اهمیت عاملی مثل زاویه بدن دروازه‌بان، قابلیت حرکت سریع‌تر بازیکنان در راستای مستقیم است. می‌توان عواملی هم‌چون سرعت بازیکنان را نیز اثر داد، اما برای همگرایی بهتر و سریع‌تر از این عوامل صرف نظر شده.



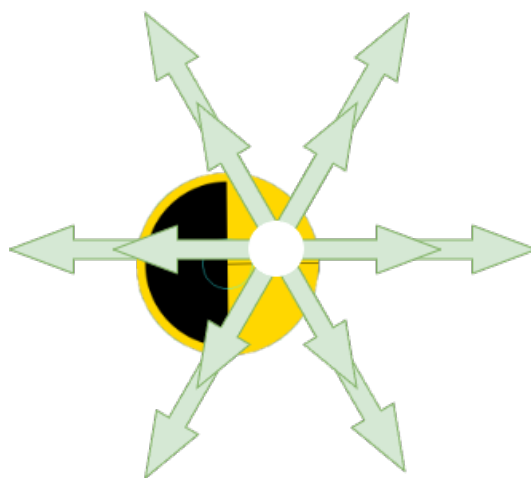
## فضای عمل

همان‌طور که گفته‌شد، عامل به صورت خودکار به دنبال توپ می‌رود و ما با کمک یادگیری تقویتی قرار است این عامل را به یادگیری ضربه‌زدن به دروازه برسانیم. برای این منظور، از رفتار سطح متوسط ضربه با سرعت دلخواه در راستای دلخواه (KickOneStep) استفاده می‌کنیم تا محاسبات مستقل از سرعت توپ قبلی و شتاب توپ شوند.

به این منظور، فضای عمل ما نیاز است سرعت و زاویه ضربه را تعیین کند. با توجه به اینکه برخی الگوریتم‌ها مانند DQN بهترین عمل را از بین اعمال ممکن انتخاب می‌کنند و نیاز به فضای عمل گسسته دارند، و سایر الگوریتم‌ها مانند DDPG نیاز به فضای عمل پیوسته دارند، محیط را در دو حالت گسسته و پیوسته پیاده‌سازی کردیم.

در فضای پیوسته، دو مقدار سرعت توپ و زاویه ضربه به عامل داده می‌شود. سرعت توپ بین صفر تا یک می‌باشد که نشان‌دهنده سرعت توپ نسبت به حداکثر سرعت ممکن است. زاویه ضربه بین منفی یک تا یک می‌باشد که نرمال‌شده زاویه مطلق ضربه بین  $-180$  تا  $180$  درجه است.

در حالت گسسته نیز، خروجی ما سرعت توپ و زاویه‌است، اما به جای مقادیر پیوسته، زوایای ضربه ممکن را به ۱۲ قسمت مساوی تقسیم کرده‌ایم و سرعت توپ را به ۱۰ قسمت مساوی تقسیم کرده‌ایم. می‌توان این تقسیم‌بندی را به راحتی تغییر داد و تاثیرات آن روی سرعت یادگیری و نتیجه نهایی آن را بررسی کرد. در شکل ۴-۷ می‌توان یک نمونه از تقسیم‌بندی فضای عمل را مشاهده کرد. از آنجا که سرعت به ۲ درجه تقسیم شده‌است، حالت‌های ممکن آن معادل با نصف حداکثر سرعت توپ و حداکثر سرعت توپ می‌باشد.



شکل ۴-۷ تقسیم‌بندی فضای عمل به ۶ زاویه و ۲ سرعت

#### ۴-۴-۴ طراحی پاداش

همان‌طور که گفته شد، محیط جیم در حالت‌هایی که توپ قابل ضربه‌زدن است، یا توپ به بیرون رفته و یا توسط دروازه‌بان گرفته‌شده اجرا می‌شود. در این حالت‌ها باید بر اساس شرایط قبلی، عمل و شرایط جدید پاداشی محاسبه کنیم که به عامل کمک می‌کند تا راحت‌تر به هدف خود نزدیک شود.

##### پاداش‌های پایانی

همان‌طور که گفته‌شد، چهار حالت پایانی داریم:

۱. حالت گل زدن: با توجه به رسیدن به هدف، پاداش بزرگی به عامل می‌دهیم که این مقدار، ۱۵۰۰ فرض شده‌است.
۲. حالت بیرون رفتن توپ: از آن‌جا که تنها با ضربات بد این سناریو رخ می‌دهد، پاداش بسیار بزرگ منفی ۵۰۰- دارد.
۳. گرفتن توپ توسط دروازه‌بان: این حالت از آنجا که بهتر از بیرون رفتن توپ است، پاداش منفی کوچک‌تری دارد و مقدار ۲۰۰- دارد.
۴. حالت اتمام زمان: به این حالت، پاداش ۱۵۰- تخصیص داده‌شده، تا عامل در حالتی که هنوز به گل‌زدن نرسیده‌است، مستقیماً توپ را به دروازه‌بان ندهد. در واقع با کوچک‌تر در نظر گرفتن این پاداش، عامل را به کاوش بیشتر وامی‌داریم.

##### پاداش در حالت عادی

در حالت بین دو ضربه، پاداش را بر اساس سه فاکتور حساب می‌کنیم:

- می‌خواهیم عامل را به حرکت به سوی دروازه واداریم، پس از تفاضل فاصله توپ با دروازه در حالت جدید و حالت قبلی به عنوان معیار اصلی استفاده می‌کنیم. در واقع از عکس این فاکتور استفاده می‌کنیم، چرا که کاهش فاصله ویژگی مثبتی است.
- نزدیک شدن زیادی به دروازه‌بان ویژگی خوبی نیست، پس از ضربی از تفاضل فاصله توپ با دروازه‌بان در بین دو حالت استفاده می‌کنیم.

- پاداش ثابت ۱۰- که برای تشویق عامل برای سریع‌تر رسیدن به حالت گل می‌باشد. از این رو که عامل این پاداش را فقط در لحظات ضربه زدن می‌بیند، این فاکتور عامل را به زدن ضربات با قدرت بیشتر تشویق می‌کند.

## ۵-۴ پیاده‌سازی یادگیری تقویتی

### ۱-۵-۴ پیاده‌سازی شبکه کیو عمیق

با کمک کتابخانه PyTorch، یک شبکه عصبی با دو لایه پنهان و تابع فعال‌سازی ReLU پیاده‌سازی شد. این شبکه، ورودی به ابعاد فضای حالت محیط دارد و خروجی به اندازه فضای عمل است.

### ۲-۵-۴ پیاده‌سازی سایر الگوریتم‌ها به کمک کتابخانه

یکی از مهم‌ترین مزایای پیاده‌سازی رابط استاندارد جیم، ایجاد امکان استفاده از پیاده‌سازی‌ها و کتابخانه‌های از پیش موجودی اند که وجود دارند. با توجه به اینکه اکثر الگوریتم‌های یادگیری تقویتی در کتابخانه‌هایی مانند Stable Baselines و OpenAI Baselines پیاده‌سازی شده‌اند، می‌توان از این کتابخانه‌ها برای آموزش عامل‌ها استفاده کرد. در این پروژه، از کتابخانه Stable Baselines برای آموزش عامل‌ها استفاده شده‌است. کافی‌ست ابتدا یک محیط جدید از جیم تعریف کنیم، مدلی از الگوریتم دلخواه بسازیم، و سپس با استفاده از تابع learn مدل را به تعداد گام‌های دلخواه آموزش دهیم.

```
model = DDPG('MlpPolicy', gym_env)
model = model.learn(1_000_000, progress_bar=True)
model.save("DDPG_model_1M")
```

شکل ۴-۸ استفاده از کتابخانه Stable Baselines 3

طبق مستندات این کتابخانه، می‌توان به مدل توابعی را به عنوان ورودی داد تا در حین آموزش به صورت مکرر و پس از تعداد گام دلخواه به صورت تناوبی صدا زده شوند. از این توابع می‌توان برای ارزیابی مدل به صورت تناوبی، و ذخیره‌سازی در حین آموزش استفاده کرد که برای مقایسات فصل بعدی بسیار کاربردی است.

## ۴-۶ جمع‌بندی

## فصل پنجم

مقایسه، آزمایش‌ها، و نتایج

در این فصل به بررسی نتایج حاصل از پیاده‌سازی الگوریتم‌ها و مقایسه آن‌ها می‌پردازیم.

## ۱-۵ ارزیابی الگوریتم‌ها

در این بخش، میانگین پاداش، درصد رخداد هر یک از حالت‌های پایانی بازی (گل، بیرون رفتن از زمین، اتمام زمان، گرفتن دروازه‌بان) و خطای میانگین مدل‌ها را در طی زمان یادگیری ارزیابی کرده، سپس بهترین مدل هر الگوریتم را با هم مقایسه می‌کنیم.

## ۲-۵ مقایسه سرعت یادگیری با گسسته سازی‌ها متفاوت

مقایسه تقسیم‌بندی‌های متفاوت قدرت و تعداد زوایای ممکن، و حالت پیوسته.

## ۳-۵ بررسی تعمیم پذیری

در این بخش، با تست علیه دروازه‌بان تیم‌هایی که مقابل آن‌ها یادگیری رخ نداده، می‌بینیم که آیا این یادگیری تعمیم‌پذیر بوده یا خیر.

## ۴-۵ تسریع یادگیری با تلفیق با رفتار سطح بالای شوت

اگر از رفتار سطح بالای شوت استفاده کنیم، در صورت موجود بودن شوت منجر به گل، این ضربه به صورت خودکار رخ می‌دهد. طبیعتاً این موضوع می‌تواند به تسریع یادگیری منجر شود. پاداش منفی برای شوت بی‌جا در نظر گرفته شده است.

## فصل ششم

جمع‌بندی، نتیجه‌گیری و پیشنهادات

در این فصل مروری بر سیر مطالب عنوان شده خواهیم داشت. سپس مروری بر نتایج و دستاوردهای این پروژه خواهیم داشت و در نهایت پیشنهاداتی برای ادامه‌ی کار در این موضوع ارائه خواهیم کرد.

## ۱-۶ کارهای آینده

- بررسی یادگیری تقویتی در حالت‌های چند عامله.
- بررسی اثر فضای حالت‌های مختلف، اثر قطبی یا دکارتی گرفتن هر ویژگی
- بررسی اثر اضافه کردن سایر ویژگی‌ها همچون سرعت بازیکنان
- آزمایش تاثیر تغییر پاداش بر یادگیری، همچون پاداش برای ضربات محکم‌تر
- استفاده از روش‌های یادگیری متقابل برای آموزش دروازه‌بان و مهاجم به صورت همزمان



## منابع و مراجع