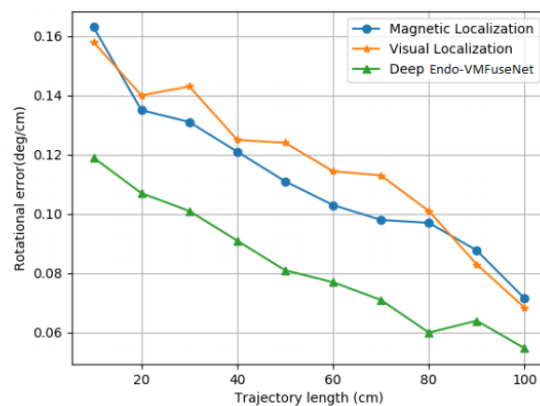


(a) Trajectory length vs translation error



(b) Trajectory length vs rotation error

شکل ۱۱-۲ مقایسه نتایج روش موقعیت‌یابی ترکیبی یادگیری عمیق با روش‌های موقعیت‌یابی براساس حسگر مغناطیسی و تصویر [۱۱۱].

عملکرد رویکرد عمیق Endo-VMFuseNet با استفاده از برآورد میانگین خطای میانگین مربعات ریشه ($RMSE^1$) برای حرکات انتقالی و چرخشی مورد تجزیه و تحلیل قرار گرفت. برای مسیری با پیچیدگی‌های مختلف، مانند مسیرهای بدون مانع با انتقال و چرخش‌های آهسته و اسکن‌های جامع و مسیرهای پیچیده با حرکات چرخشی و انتقالی سریع، آزمایشات عمیق در مقایسه با سیستم موقعیت‌یابی مبتنی بر تصویر و موقعیت‌یابی مغناطیسی انجام شده است. میانگین $RMSE$ ‌های انتقالی و چرخشی برای Endo-VMFuseNet عمیق، موقعیت‌یابی مبتنی بر تصویر به تنهایی و موقعیت‌یابی مغناطیسی به تنهایی در طول مسیرهای مختلف به ترتیب در شکل ۱۱-۲ نشان داده شده است. نتایج نشان می‌دهد که Endo-VMFuseNet عمیق به وضوح بهتر از موقعیت‌یابی مبتنی بر تصویر به تنهایی و موقعیت‌یابی مغناطیسی به تنهایی عمل می‌کند. استفاده موثر از معماری LSTM در معماری Endo-VMFuseNet امکان یادگیری از آرایه سنسور ناهمزمان و غیر کالبره شده را فراهم کرده است. نتایج همچنین نشان می‌دهد که Endo-VMFuseNet با تفسیر اطلاعات موقعیت‌یابی از داده‌های مغناطیسی فعلی و اطلاعات موقعیت‌یابی مبتنی بر تصویر قبلی ذخیره شده توسط حافظه پنهان داخلی واحدهای LSTM، قادر به مدیریت داده‌های ناهمزمان (داده‌های مغناطیسی ۵۰ هرتز و داده‌های تصویری ۳۰ هرتز) است.

۲-۲. محیط شبیه‌سازی

ابزار ML-Agents از سال ۲۰۱۷ آغاز به کار کرده است و امکانات زیادی در زمینه یادگیری تقویتی عمیق در اختیار محققان، توسعه دهندگان و برنامه نویسان بازی قرار داده است. این ابزار از ابتدا با هدف کمک به محققان و توسعه دهندگان در تبدیل بازی‌ها و شبیه‌سازی‌های ایجاد شده با محیط توسعه Unity به محیط‌های یادگیری عمیق ایجاد شد که در آن عامل‌ها می‌توانند با استفاده از پیشرفته‌ترین الگوریتم‌های یادگیری تقویتی عمیق،

¹ Root-mean-square deviation

استراتژی‌های تکاملی و ژنتیکی، و سایر روش‌های یادگیری عمیق (شامل بینایی کامپیوتر، تولید داده‌های مصنوعی) صرفاً از طریق یک API ساده شده پایتون آموزش داده شوند.

۲-۲-۱. پلتفرم یونیتی (Unity)

یونیتی (Unity) یک پلتفرم توسعه سه بعدی بلادرنگ است که از یک موتور رندر^۱ تصاویر و یک موتور فیزیک و همچنین یک رابط کاربری گرافیکی^۲ به نام Unity Editor تشکیل شده است. یونیتی به طور گسترده در صنایع بازی، AEC^۳ (معماری، مهندسی، ساخت و ساز)، خودرو و صنایع فیلم پذیرفته شده است و توسط جامعه بزرگی از توسعه دهندگان بازی برای ساخت انواع شبیه سازی‌های تعاملی، از بازی‌های کوچک موبایل و مبتنی بر مرورگر گرفته تا بازی‌های کنسولی و بازی‌های AR/VR^۴ استفاده می‌شود. استفاده از موتور یونیتی می‌تواند پیشرفت تحقیقاتی را در زمینه‌های عقب مانده و حیاتی تسریع کند.

پیشینه تمرکز یونیتی بر توسعه یک موتور همه منظوره برای پشتیبانی از انواع پلتفرم‌ها، انواع توسعه‌دهندگان و انواع بازی‌ها، موتور یونیتی را به یک کاندیدای ایده‌آل پلتفرم شبیه‌سازی برای تحقیقات هوش مصنوعی تبدیل می‌کند. انعطاف‌پذیری زیرساخت موتور، امکان ایجاد برنامه‌هایی از دنیای شبکه دوبعدی ساده تا بازی‌های استراتژیک پیچیده سه بعدی، پازل‌های مبتنی بر فیزیک و بازی‌های رقابتی چند عامله را ممکن می‌سازد. برخلاف بسیاری از پلتفرم‌های تحقیقاتی، زیرساخت موتور به هیچ سبک خاصی از محیط بازی یا شبیه‌سازی محدود نمی‌شود و یونیتی را به یک پلتفرم همه منظوره تبدیل می‌کند. علاوه بر این، محیط ویرایشگر یونیتی، نمونه سازی سریع، توسعه بازی و محیط‌های شبیه سازی شده را امکان پذیر می‌کند.

یک پروژه یونیتی از مجموعه ای از دارایی‌ها^۵ تشکیل شده است. صحنه‌ها^۶ نوع خاصی از دارایی هستند که محیط یک پروژه را مشخص می‌کنند. صحنه‌ها ترکیبی سلسله مراتبی از اشیاء بازی^۷ هستند که با اشیاء واقعی (چه فیزیکی و چه کاملاً منطقی) در محیط مطابقت دارد. رفتار و عملکرد هر کدام از اشیاء بازی توسط اجزای متصل به آن تعیین می‌شود. این اجزا^۸ میتواند انواع مختلفی داشته باشند همچون دوربین‌ها، مش‌ها، رندررها، ساختارهای غیر منعطف و بسیاری اجزای دیگر که توسط ویرایشگر یونیتی ارائه می‌شوند. همچنین می‌توان اجزای سفارشی^۹ را با استفاده از اسکریپت‌های C# یا افزونه‌های خارجی تعریف کرد [۱۱۶].

^۱ Render

^۲ GUI (Graphical user interface)

^۳ Architecture, Engineering & Construction

^۴ Augmented reality (AR) and Virtual Reality (VR)

^۵ Asset

^۶ Scene

^۷ GameObjects

^۸ Components

^۹ Custom components

موتور یونیتی پیچیدگی لازم را جهت ایجاد محیط‌های یادگیری چالش برانگیز امکان‌پذیر می‌کند.

۱-۲-۲-۲. ویژگی‌های محیط

پیچیدگی سنسور - موتور یونیتی رندر گرافیکی با کیفیت بالا را امکان‌پذیر می‌کند. این موتور از نورپردازی پیش ساخته و نورپردازی زمان واقعی پشتیبانی می‌کند. همچنین توانایی تعریف سایه بان‌های سفارشی به صورت برنامه نویسی یا از طریق یک زبان برنامه نویسی بصری را دارد. بنابراین، می‌توان به سرعت تصاویری از محیط رندر کرد که کاملاً شبیه عکس برداری واقعی از محیط باشد و این تصاویر را به عنوان داده‌های آموزشی^۱ در مدل یادگیری ماشین استفاده کرد. بعلاوه با استفاده از سایه‌بان‌های سفارشی میتوان اطلاعات عمق، پوشش اشیاء، مادون قرمز و تصاویر با نویز تزریق شده رندر کرد. موتور یونیتی وسیله‌ای برای تعریف سیگنال‌های صوتی ارائه می‌کند که می‌تواند به عنوان اطلاعات مشاهده شده در عوامل یادگیری به کار گرفته شود. پس قابلیت شبیه سازی سیستم‌های تشخیصی مبتنی بر پرتو^۲ مانند Lidar را نیز دارد.

پیچیدگی فیزیکی - پدیده‌های فیزیکی در محیط‌های یونیتی را می‌توان با موتورهای Nvidia PhysX یا Havok Physics شبیه سازی کرد. که اینکار، امکان تحقیق در محیط‌هایی با بدنه صلب، بدنه نرم، دینامیک ذرات، دینامیک سیالات و همچنین فیزیک ragdoll را فراهم می‌کند. در صورت تمایل، ماهیت توسعه پذیر بودن این پلتفرم امکان استفاده از موتورهای فیزیکی دیگر^۳ را فراهم می‌کند. برای مثال، پلاگین‌هایی برای یونیتی وجود دارد که موتورهای فیزیک Bullet و MuJoCo را به عنوان جایگزین PhysX به کار می‌گیرد.

پیچیدگی منطق کار - موتور یونیتی یک سیستم اسکریپت نویسی قدرتمند و انعطاف پذیر از طریق C# ارائه می‌دهد. این سیستم پیاده سازی هر شکلی از گیم پلی، شبیه سازی و کنترل دینامیکی را امکان‌پذیر می‌سازد. علاوه بر زبان اسکریپت نویسی، اشیاء بازی و سیستم اجزاء^۴ امکان مدیریت همزمان چندین نمونه از عوامل، سیاست ها و محیط ها را فراهم می‌کند و امکان تعریف وظایف پیچیده سلسله مراتبی و همچنین وظایفی که برای حل آنها به یادگیری نیاز است را ممکن می‌سازد.

پیچیدگی اجتماعی - ماهیت زبان اسکریپت نویسی یونیتی و سیستم اجزاء، نمایش سناریوهای چند عامله را بسیار ساده و راحت می‌کند. در واقع، از آنجایی که این پلتفرم برای پشتیبانی از توسعه بازی‌های ویدیویی چند نفره طراحی شده است، چندین اجزای انتزاعی مفید مانند سیستم شبکه چند نفره از قبل طراحی شده است.

¹ Training Data

² ray-cast

³ 3rd party

⁴ component

شبیه سازی سریع و توزیع شده - در موتور یونیتی شبیه سازی فیزیک و رندر کردن فریم ها به صورت ناهمزمان انجام می‌شود. لذا بدون نیاز به افزایش نرخ فریم در فرآیند رندر، می‌توان سرعت شبیه سازی را تا حد زیادی افزایش داد. همچنین اگر شبیه سازی نیاز به رندر کردن نیازی نداشته باشد می‌توان شبیه سازی‌های یونیتی را بدون رندر نیز اجرا کرد که در این حالت یونیتی منابع پردازشی بسیار کمتری استفاده می‌کند. در سناریوهایی که رندر کردن لازم است، مانند یادگیری از پیکسل ها، می‌توان نرخ فریم و سرعت منطق بازی را کنترل کرد. همچنین کنترل زیاد روی کیفیت رندر این امکان را فراهم می‌کند که هر زمان نیاز باشد نرخ فریم تا حد زیادی افزایش یابد. اما قابلیت‌های اضافه شده موتور یونیتی در شبیه سازی‌های مقیاس بزرگ، سربار اضافی به همراه دارد. به طوری که میزان حافظه رم نیاز برای شبیه سازی یونیتی از محیط‌های دیگر پلتفرم ها مانند ALE بیشتر است.

کنترل انعطاف پذیر - کنترل بیشتر جنبه‌های شبیه سازی به صورت برنامه نویسی امکان پذیر است. این ویژگی محققان را قادر می‌سازد تا برنامه‌های آموزشی، سناریوهای متخاصم یا سایر روش‌های پیچیده ای که محیط را در طول فرآیند آموزش تغییر می‌دهد، تعریف کنند. به عنوان مثال، اشیاء بازی را می‌توان به صورت مشروط ایجاد و به صورت بلادرنگ از بین برد. در نهایت موتور یونیتی از طریق پارامترهای بیرونی شبیه سازی و یک رابط برنامه نویسی کاربردی (API¹) پایتون²، کنترل بیشتری در شبیه سازی را امکان پذیر می‌کند.

۲-۲-۳. نرم افزار ویرایشگر Unity

ویرایشگر یونیتی (شکل ۱) یک رابط کاربری گرافیکی است که برای ایجاد محتوای ۲ بعدی و ۳ بعدی استفاده می‌شود. نرم افزار یونیتی در سیستم عامل‌های ویندوز، مک و لینوکس موجود است و مزایای زیادی برای تحقیقات هوش مصنوعی فراهم می‌کند. از جمله آنها می‌توان به مواردی نظیر اشاره کرد:

- ایجاد صحنه‌های سفارشی: این امکان جهت آموزش هوش مصنوعی در محیط‌های مختلف با تجربه‌های گوناگون به کار می‌رود.
- نمایش متخصص: به این وسیله آموزش هوش مصنوعی به عهده یک متخصص گذاشته می‌شود. متخصص می‌تواند شبیه سازی را شروع کند و عامل‌های مختلف در صحنه را کنترل کند. بنابراین داده‌هایی تولید می‌شود که بر طبق آن آموزش انجام می‌شود. در حقیقت این ویژگی ما را قادر می‌سازد تا از الگوریتم‌های یادگیری تقلیدی^۳ استفاده کنیم. این الگوریتم ها سرعت یادگیری ماشین را در محیط‌های پیچیده بسیار بالا می‌برد.

¹ Application Programming Interface

² Python

³ Imitation Learning

- نمایش در مقیاس بزرگ: یکی از قدرتمندترین ویژگی‌های یونیتی توانایی ساخت محیط شبیه سازی در بیش از ۲۰ پلتفرم شامل انواع رایانه، انواع موبایل، کنسولها و ... است. این ویژگی توانایی اندازه گیری کارایی از کاربران متنوع را فراهم می‌کند. همچنین در ویرایشگر یونیتی می‌توان همزمان چندین محیط را به صورت موازی ایجاد کرد که این کار منجر به افزایش نمای سرعت یادگیری می‌شود.

۴-۲-۲. ابزار ML-Agent در Unity

ابزار ML-Agents یک پروژه منبع باز است که محققان و توسعه دهندگان را قادر می‌سازد محیط‌های شبیه سازی شده را با استفاده از ویرایشگر یونیتی ایجاد کرده و از طریق یک API پایتون با آنها تعامل داشته باشند. ML-Agents یک جعبه ابزار (SDK) ارائه می‌دهد که شامل تمام توابع لازم برای تعریف محیط‌ها در ویرایشگر یونیتی است. این توابع شامل اسکریپت‌های اصلی C# برای ساخت یک روند یادگیری است.

ویژگی‌های این جعبه ابزار شامل مجموعه‌ای از محیط‌های نمونه، الگوریتم‌های پیشرفته یادگیری تقویتی (SAC^1) و (PPO^2)، الگوریتم‌های یادگیری تقلیدی ($GAIL^3$) و شبیه‌سازی رفتاری (BC) است. پشتیبانی از Self-Play در بازی‌های متقارن و نامتقارن، و همچنین امکان گسترش الگوریتم‌ها و سیاست‌ها با روش‌های ICM و LSTM از دیگر ویژگی‌های این جعبه ابزار است.

در ابزار ML-Agents یک مفهوم به نام عامل‌ها وجود دارد. مؤلفه عامل نشان دهنده اشیای بازی موجود در یک صحنه هستند که می‌توانند مشاهدات را جمع‌آوری کنند، اقداماتی انجام دهد و پاداش دریافت کند. عامل می‌تواند مشاهدات را با استفاده از انواع حسگرها مبتنی بر اشکال مختلف داده مانند تصاویر رندر شده، پخش پرتو و ... در ابعاد مختلف برداری، جمع‌آوری کند. هر جزء عامل حاوی یک سیاست است که با یک نام رفتاری برچسب گذاری شده است.

هر تعداد عامل می‌تواند سیاستی با نام رفتاری مشابه داشته باشد که این عوامل همگی یک سیاست را اجرا می‌کنند و داده‌های تجربی را در طول آموزش به اشتراک می‌گذارند. به علاوه، می‌تواند هر تعداد نام رفتاری برای سیاست‌ها در یک صحنه وجود داشته باشد که ساخت سناریوهای چند عاملی را با گروه‌ها یا عامل‌های فردی، ممکن می‌سازد. یک سیاست می‌تواند برآمده از مکانیسم‌های تصمیم‌گیری مختلف از جمله ورودی بازیکن، اسکریپت‌های از پیش کدگذاری شده، مدل‌های شبکه عصبی تعبیه شده داخلی باشد یا اینکه از طریق تعامل با API پایتون مشخص شود. این امکان برای عامل‌ها وجود دارد تصمیم‌ها را از سیاست‌های ثابت یا پویا (همانطور که توسط توسعه دهنده محیط تعریف شده است)، درخواست کنند.

¹ Actor-Critic

² Proximal Policy Optimization

³ Generative Adversarial Imitation Learning

⁴ Agents

تابع پاداش، که برای ارائه سیگنال یادگیری به عامل استفاده می‌شود، می‌تواند در هر زمانی در طول شبیه سازی با استفاده از سیستم اسکرپت یونیتی تعریف یا تغییر یابد. به همین ترتیب، شبیه سازی را می‌توان در یک حالت انجام شده در سطح یک عامل فردی یا محیط قرار داد. این حالت یا از طریق فراخوانی اسکرپت یونیتی یا با رسیدن به حداکثر تعداد گام از پیش تعریف شده اتفاق می‌افتد.

در هر شبیه سازی یک آکادمی^۱ وجود دارد که برای پیگیری مراحل شبیه سازی و مدیریت عامل ها استفاده می‌شود. آکادمی همچنین دارای قابلیت تعریف پارامترهای محیطی است که می‌توان از آن برای تغییر پیکربندی محیط در زمان اجرا استفاده کرد. به طور خاص، جنبه‌های فیزیک محیطی و بافت‌ها، اندازه‌ها و وجود اشیاء از طریق پارامترهای قابل دسترسی کنترل می‌شوند. پس می‌توانند مجدداً نمونه‌برداری شده و در طول آموزش تغییر کنند. به عنوان مثال، گرانش در محیط می‌تواند در هر بازه ثابتی نوسان کند یا موانع اضافی می‌تواند زمانی که یک عامل به مهارت خاصی برسد، ایجاد شود. این قابلیت، ارزیابی یک عامل را در محیط‌های مختلف و در شرایط یادگیری و تست فراهم می‌کند و ایجاد سناریوهای یادگیری برنامه آموزشی را تسهیل می‌کند.

۲-۲-۵. معیارهای عملکرد

برای یک شبیه سازی ضروری است که محیط قادر به ارائه سرعت شبیه سازی بیشتری نسبت به زمان واقعی باشد. در ابزار شبیه سازی ML-Agents امکان افزایش سرعت تا صد برابر زمان واقعی وجود دارد. با این حال در عمل افزایش سرعت، ممکن مطابق با منابع پردازشی موجود و همچنین پیچیدگی‌های محیط متفاوت باشد. در موتور یونیتی، منطق بازی می‌تواند مستقل از رندر کردن فریم ها اجرا شود. به این ترتیب، محیط‌هایی که به مشاهدات بصری متکی نیستند، مانند محیط‌هایی که از پرتوها استفاده نمی‌کنند، می‌توانند از شبیه سازی با سرعت بیشتری نسبت به محیط‌هایی که استفاده می‌کنند بهره ببرند.

۲-۲-۶. اجزای اصلی ابزار ML-Agents

ابزار ML-Agents از ۴ جزء اصلی سطح بالا تشکیل شده است:

- **محیط یادگیری** - شامل صحنه یونیتی و تمامی شخصیت‌های بازی می‌باشد. صحنه یونیتی محیطی را فراهم می‌کند که در آن عامل ها محیط را مشاهده می‌کنند، در آن عمل‌های متفاوت انجام می‌دهند و همچنین از آن محیط یاد می‌گیرند. نحوه تنظیم صحنه یونیتی برای استفاده به عنوان یک محیط یادگیری به هدف آموزش بستگی دارد. صحنه ها هم جهت آموزش و هم جهت آزمایش استفاده می‌شوند. پس برای یک مسئله یادگیری تقویتی خاص با دامنه محدود می‌توان از یک صحنه همزمان برای آموزش و آزمایش عوامل آموزش دیده استفاده کرد. همچنین به منظور آموزش عامل ها در یک بازی پیچیده یا شبیه سازی، می‌توان از صحنه‌های مجزا بهره برد که ممکن است کارآمدتر و کاربردی‌تر

^۱ Academy

باشد. جعبه ابزار ML-Agents به وسیله بسته پایتون این امکان را می‌دهد که با تعریف عامل‌ها و رفتارهای آنها، هر صحنه یونیتی را به یک محیط یادگیری تبدیل کنید.

- **رابط برنامه نویسی کاربردی پایتون** - شامل یک رابط پایتون سطح پایین برای تعامل و دستکاری یک محیط یادگیری است. برخلاف محیط یادگیری، API پایتون بخشی از یونیتی نیست، بلکه یک ماژول خارجی است و از طریق ارتباط دهنده^۱ با یونیتی ارتباط برقرار می‌کند. این API در یک بسته اختصاصی پایتون موجود است و توسط فرآیند آموزش پایتون برای برقراری ارتباط و کنترل آکادمی در طول آموزش استفاده می‌شود. با این حال، می‌توان آن را برای اهداف دیگری نیز مانند استفاده از یونیتی به عنوان موتور شبیه سازی الگوریتم‌های یادگیری ماشین، بکار برد.
- **ارتباط دهنده خارجی** - محیط یادگیری را به API پایتون متصل می‌کند.
- **موتور آموزش دهنده پایتون** - که شامل تمام الگوریتم‌های یادگیری ماشینی است که عامل‌ها را آموزش می‌دهد. این الگوریتم‌ها در پایتون پیاده‌سازی شده‌اند و فقط با API سطح پایین پایتون رابط کاربری دارد.

۲-۶-۱. محیط یادگیری

محیط یادگیری شامل دو مؤلفه است که به سازماندهی صحنه یونیتی کمک می‌کند:

- **عامل‌ها** - که به یک شی از یونیتی (هر شی در صحنه) متصل است و مشاهدات آن را ایجاد می‌کند، اقداماتی را که دریافت می‌کند انجام می‌دهد و در صورت لزوم پاداش (مثبت / منفی) اختصاص می‌دهد. هر عامل به یک رفتار پیوند خورده است.
- **رفتار**^۲ - ویژگی‌های خاص عامل را تعریف می‌کند، مانند تعداد اقداماتی که عامل می‌تواند انجام دهد. یک رفتار را می‌توان به عنوان عملکردی در نظر گرفت که مشاهدات و پاداش‌ها را از عامل دریافت می‌کند و اعمال را بر می‌گرداند. یک رفتار می‌تواند یکی از سه نوع باشد: یادگیری، اکتشافی یا استنتاج.

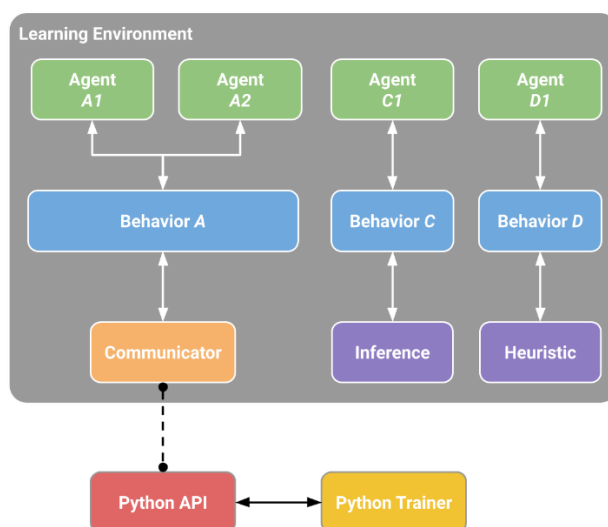
هر محیط آموزشی همیشه یک عامل برای هر شخصیت در صحنه خواهد داشت. در حالی که هر عامل باید به یک رفتار پیوند خورده باشد، این امکان وجود دارد که عامل‌هایی که مشاهدات و اقدامات مشابهی دارند، رفتار مشابهی داشته باشند و این بدان معنا نیست که در هر کدام، ارزش‌های مشاهده و عمل یکسانی خواهند داشت. در یک محیط، می‌تواند همزمان چندین عامل و چندین رفتار وجود داشته باشد. محیط یادگیری از طریق آکادمی

^۱ Communicator

^۲ Behavior

(که در نمودار نشان داده نشده است) تضمین می کند که همه عامل ها علاوه بر کنترل تنظیمات محیطی، همگام هستند. به طور کلی آکادمی موارد زیر را تضمین می کند:

- جمع آوری مشاهدات
- انتخاب عمل با استفاده از سیاست
- انجام عمل
- ریست کردن یادگیری اگر یادگیری انجام شد یا به حداکثر مراحل خود رسید.



شکل ۱۲-۲ محیط یادگیری.

۲-۶-۲-۲ معماری رفتار -آکادمی

داشتن عامل هایی که به وسیله ی شبکه عصبی کنترل می شوند، مفهوم اساسی در ساخت محیط بازی یا شبیه سازی به وسیله ML-Agents به شمار می آید. هر کدام از این عالم ها به یک رفتار پیوند خورده اند که می توان به طور زیر طبقه بندی کرد:

- **رفتار استنتاجی** - این نوع معماری شامل مدل های از پیش آموزش دیده است که با استفاده از الگوریتم های مختلف یادگیری عمیق آموزش داده شده اند. رفتار استنتاجی برای آموزش، نیازی به ارتباط با API پایتون در زمان اجرا ندارد، زیرا قبلاً روی یک مجموعه از پیش موجود فضای مشاهده - عمل آموزش داده شده است.
- **رفتار ابتکاری** - این نوع معماری معمولاً دارای هیچ مدل شبکه عصبی نیست. بلکه شامل یک راه حل ابتکاری است که در این رفتار گنجانده شده است. به جای استفاده از یادگیری عمیق، رفتار ابتکاری بیشتر بر پایه برنامه نویسی پویا و یادگیری مبتنی بر حالت، توسعه می یابند.

- **رفتار یادگیری** - این رفتار نوع دیگری است که از الگوریتم‌های یادگیری عمیق استفاده می‌کند و در هنگام آموزش عامل در یونیتی استفاده می‌شود. رفتار یادگیری برای آموزش عامل‌ها در زمان اجرای آموزش از طریق ارتباط دهنده و API پایتون به یونیتی متصل می‌شود. ارتباط دهنده، نقش مهمی در برقراری ارتباط روی پورت ۵۰۴۰ ایفا می‌کند. به وسیله ارتباط دهنده می‌توان برای آموزش از مدل‌های Tensorflow بهره برد. همچنین مدل‌های آموزش دیده از رفتار یادگیری، به عنوان مدل‌های تولید شده و از پیش آموزش دیده شده برای رفتار استنتاجی نیز استفاده می‌شوند.

۳-۶-۲-۲. موتور استنتاج یونیتی: Barracuda

جعبه ابزار ML Agents برای اجرای مدل‌های شبکه عصبی از پیش آموزش دیده با استفاده از رفتار یادگیری، از موتور استنتاج یونیتی یا Barracuda بهره می‌برد. این یک کتابخانه سبک وزن به تبدیل مدل‌های شبکه عصبی از کتابخانه Tensorflow به فایل‌های سریالی با پسوند nn. کمک می‌کند. موتور استنتاج یونیتی از دو فرمت Barracuda و ONNX پشتیبانی می‌کند. تنظیمات موتور استنتاج یونیتی را می‌توان در داخل اسکریپت پارامترهای رفتاری مشاهده کرد. یک ویژگی مهم موتور استنتاج در بخش مدل وجود دارد که می‌توان به وسیله آن دستگاه پردازشی CPU یا GPU^۱ را انتخاب کرد. این قابلیت موتور استنتاج را قادر می‌سازد تا متناسب با پروژه، بهترین منابع پردازشی را جهت آموزش استفاده کند.

۷-۲-۲. سناریوهای آموزشی انعطاف پذیر

با استفاده از ابزار ML-Agent تعریف کردن چندین سناریو جهت آموزش عامل‌ها امکان پذیر است:

- **تک عامل:** یک عامل واحد، با سیگنال پاداش خاص خود. این سناریو روش سنتی آموزش یک عامل مانند بازی‌های تک نفره است.
- **تک عامل همزمان:** چندین عامل مستقل با سیگنال‌های پاداش مستقل با پارامترهای رفتاری یکسان. یک نسخه موازی از سناریوی آموزشی سنتی است که می‌تواند روند تمرین را تسریع کند. این سناریو زمانی مفید است که چندین نسخه از یک شخصیت در محیطی دارید که باید رفتارهای مشابهی را بیاموزد. به عنوان مثال آموزش چند ربات بازو برای باز کردن یک در به صورت همزمان باشد.
- **تقابل دو عامل:** دو عامل در تعامل باهم اما با سیگنال‌های پاداش معکوس. در بازی‌های دو نفره، تقابل دو عامل می‌تواند به یک عامل اجازه دهد تا به طور فزاینده‌ای ماهرتر شود، در حالی که همیشه حریف کاملاً همسانی داشته باشد.

^۱ Graphics processing unit

- چند عاملی تعاونی: چندین عامل تعاملی با یک سیگنال پاداش مشترک با پارامترهای رفتاری یکسان یا متفاوت. در این سناریو، همه عامل ها باید با هم همکاری کنند تا کاری را انجام دهند که به تنهایی قابل انجام نیست. به عنوان مثال می توان به محیط هایی اشاره کرد که در آن هر عامل تنها به قسمتی از اطلاعات دسترسی دارد، که برای انجام کار یا حل مشترک یک مسئله باید به اشتراک گذاشته شود.
- چند عاملی رقابتی: چندین عامل تعاملی با سیگنال های پاداش معکوس با پارامترهای رفتاری یکسان یا متفاوت. در این سناریو، عامل ها باید با یکدیگر رقابت کنند تا یا در یک مسابقه برنده شوند یا مجموعه محدودی از منابع را به دست آورند. به عنوان مثال تمامی ورزش های تیمی در این سناریو قرار می گیرند.
- زیست بوم: چندین عامل تعاملی با سیگنال های پاداش مستقل با پارامترهای رفتاری یکسان یا متفاوت. این سناریو را می توان به عنوان ایجاد دنیای کوچکی در نظر گرفت که در آن عامل ها با اهداف متفاوت همه با هم تعامل دارند، مانند طبیعت که ممکن است در آن حیوانات مختلف وجود داشته باشد، یا یک شبیه سازی رانندگی مستقل در یک محیط شهری.

۲-۷-۱. سیگنال های پاداش

در یادگیری تقویتی، هدف نهایی عامل، کشف رفتاری (سیاست) است که پاداش را به حداکثر می رساند. شما باید یک یا چند سیگنال پاداش را در اختیار عامل قرار دهید تا در طول آموزش از آنها استفاده کند. به طور معمول، یک پاداش توسط محیط برای رسیدن به هدف تعریف می شود. این نوع از پاداش، بیرونی است زیرا آنها خارج از الگوریتم یادگیری تعریف می شوند. با این حال، پاداش ها را می توان خارج از محیط نیز تعریف کرد، تا عامل را تشویق کند تا به روش های خاصی رفتار کند، یا به یادگیری پاداش بیرونی واقعی کمک کند. این پاداش ها به عنوان سیگنال های پاداش ذاتی (مانند GAIL، Curiosity و RND) بشمار می روند. کل پاداشی که عامل یاد می گیرد تا به حداکثر برسد، می تواند ترکیبی از سیگنال های پاداش بیرونی و درونی باشد. جعبه ابزار ML-Agents اجازه می دهد تا سیگنال های پاداش به صورت ماژولار تعریف شود.

۲-۳. یادگیری عمیق

۲-۳-۱. مقدمه ای بر یادگیری تقویتی

کنترل سیستم های مکانیکی به سطح بالایی از درک سیستم متکی است [۱۱۷]، که عمدتاً از طریق مدل های دینامیک بیان می شود [۱۱۸]، [۱۱۹]. با این حال، این روش برای سیستم های دارای درجه پیچیدگی یا عدم اطمینان بالا محدودیت های زیادی دارد [۱۲۰]. به طور خاص، با در نظر گرفتن الگوریتم های کنترل مبتنی بر مدل، مانند

تکنیک‌های تنظیم کننده خطی-درجه دوم (LQR^1) و پیشگام (BS^2)، استفاده از معادلات دیفرانسیل-جبری حرکت چالش برانگیز است. اینها به شدت به دقت مدل متکی هستند، که مانع کاربرد آنها در سیستم‌های واقعی می‌شود [۱۲۱]. بنابراین، توسعه کنترل کننده‌ها برای سیستم‌های غیرخطی بسیار چالش برانگیز، وقت گیر و در بسیاری از موارد کاری غیرقابل اجرا است. روش مهندسی سنتی شامل استخراج تحلیلی سیستم حاکم بر معادلات و تنظیم دستی پارامترهای سیستم کنترل متناسب با برخی پارامترهای فیزیکی اندازه گیری شده است. بنابراین تخصص پیشرفته در ریاضیات کاربردی، نظریه سیستم‌های دینامیکی، روش‌های محاسباتی، مدل سازی ریاضی، چارچوب‌های بهینه سازی و تنظیم الگوریتمی پارامترهای کنترل به کمک اپراتور مورد نیاز است.

۲-۳-۲. توضیح روش یادگیری تقویتی

یک روش جایگزین استفاده از روشهای داده-محور یادگیری ماشین است. با توجه به حجم عظیم داده، می‌توان با بهره گیری از تکنیک‌های مناسب ML مدل‌هایی را برای سیستم‌های پیچیده ایجاد کرد [۱۲۲]. علاوه بر این، شبیه سازی ها و داده‌های مصنوعی برای آموزش امکان استفاده از شبکه‌های عصبی در برنامه‌های متنوع با درجات مختلف پیچیدگی را فراهم می‌کند [۱۲۳]. به عنوان مثال، نویسندگان زیادی سعی در یادگیری تقلیدی برای به دست آوردن مهارت‌های حرکتی از تظاهرات انسانی دارند [۱۲۴]. با این حال، نتایج این روش‌ها در رباتیک نزدیک به بهینه است. راه حل استفاده از روش یادگیری تقویتی (RL) است. این یک روش ML است که بر روی بهینه سازی فرایندهای تصمیم گیری مارکوف تمرکز دارد. هدف RL یادگیری نحوه کنترل سیستم برای به حداکثر رساندن پاداش‌های تجمعی است. برخلاف یادگیری تحت نظارت یا بدون نظارت، این روش نمونه‌هایی را تولید می‌کند که بدون برچسب هستند، بنابراین فقط سیگنال‌های پاداش تنک را در نظر می‌گیرند [۱۲۵]. به طور کلی، بهینه سازی سیاست RL به سیگنال پاداش بستگی دارد. پاداش هنگامی می‌تواند موثر واقع شود که نشان دهد کدام سیاست و عملکرد منفرد بهتر از دیگران است [۱۲۶]. RL برای حل مشکلات کنترل بهینه، با اطلاعات محدود در مورد حالت دینامیکی به دانش قبلی کمی نیاز دارد [۱۲۷]. به طور کلی، روش‌های RL از نظر محاسباتی ساده بوده و نمایانگر روشی است که مستقیماً برای کنترل سیستم‌های غیرخطی [۱۲۸]، با ترکیب ویژگی‌های کنترل کننده‌های بهینه و تطبیقی، قابل استفاده است [۱۲۹]. در دهه گذشته، RL موفقیت زیادی در کنترل فضای کوچک گسسته ی عمل-حالت کسب کرده است. اخیراً، از طریق تقریب تابع، برای حل مشکلات مربوط به فضای بزرگ و پیوسته استفاده شده است [۱۳۰]، به این ترتیب پتانسیل قابل استفاده در رشته‌های متنوعی مانند مدیریت انرژی [۱۳۱]، وسایل نقلیه خودران [۱۳۲]، و تحلیل روند مالی [۱۳۳] را دارا است. بنابراین، یادگیری عمیق به روش اصلی برای دستیابی به RL موثر، تبدیل شده است [۱۳۴]. ترکیب این دو روش اصطلاحاً (DRL^3) نامیده می‌شود که در حال حاضر یک رویکرد امیدوار کننده برای هوش مصنوعی در حالت عمومی است [۱۳۵].

¹ Linear Quadratic Regulator

² Backstepping

³ Deep RL

اخیراً، DRL در رباتیک مورد توجه بسیاری قرار گرفته است، جایی که عامل به کنترل کننده ربات تبدیل می شود و خود ربات و محیط اطراف آن محیط هستند [۱۳۶]. DRL در رباتیک یک مسئله ی کنترلی است که در آن یک ربات به عنوان عامل، با انتخاب متوالی اقدامات از طریق دنباله ای از مراحل زمانی، در یک محیط تصادفی عمل می کند [۱۳۷]. این چارچوب به یک عامل رباتیک این امکان را می دهد تا مهارت هایی را برای حل وظایف از طریق آزمون و خطا به طور مستقل بیاموزد. با این حال، ابعاد مسئله کنترلی رباتیک در فضاهای حالت پیوسته اغلب دارای ابعاد بالا هستند. به همین دلیل، متأسفانه راه حل گسسته سازی فضای عمل محدودیت های زیادی ایجاد می کند.

قوام در برابر عدم اطمینان در مدل سازی از بیست سال گذشته یکی از مهمترین موضوعات تحقیقاتی در سیستم های کنترل بوده است. این امر عمدتاً ناشی از دشواری هایی است که هنگام انتقال یک مدل به دنیای واقعی پیدا می شود. به طور دقیق تر دلیل این امر ساده سازی بیش از حد در معادلات دینامیکی، پیش فرض های ساده کننده مسئله، دینامیک های مدل نشده و مقادیر تقریبی عددی خصوصیات فیزیکی است. از این منظر، مزیت اصلی RL توانایی یادگیری آن از تعامل با محیط براساس یافتن یک سیاست کنترل بهینه است، نه اینکه توسط یک طراح انسانی برنامه ریزی شود. با این حال، مشکلی که در اینجا ایجاد می شود همگرایی ضعیف RL بر اساس شبکه های عصبی است، که تاکنون یک چالش بزرگ برای ربات های خودمختار است. علاوه بر این، فعل و انفعالات فیزیکی متعددی معمولاً برای این سیستمها غیرقابل انجام است که شاید منجر به فرسودگی و آسیب دیدن سخت افزار شود. بنابراین، هرچه سیستم آسیب پذیرتر باشد، میزان کارایی داده در یادگیری نیز بیشتر است. تکنیک های ML و RL عمیق پتانسیل زیادی برای کاربردهای مهندسی کنترل دارند. اساساً دو نکته اساسی وجود دارد. اول، روش های RL امکان ایجاد کنترل کننده های غیرخطی موثر سیستم های مکانیکی را فراهم می کنند که تا حدی رفتار فیزیکی سیستم دینامیکی را که باید کنترل شود نادیده می گیرند. در حقیقت، با سیستم مکانیکی که باید کنترل شود، به عنوان یک جعبه سیاه که از آن داده های ورودی و خروجی جمع می شود، رفتار می شود. پس این روش امکان استفاده از استراتژی کنترلی طراحی شده را در طیف گسترده ای از سیستم های حالت فراهم می کند [۱۳۸].

۳-۳-۲. الگوریتم های یادگیری تقویتی

برخی از اصطلاحاتی که عناصر یک مساله یادگیری تقویتی را تشریح می کنند در زیر بیان شده است:

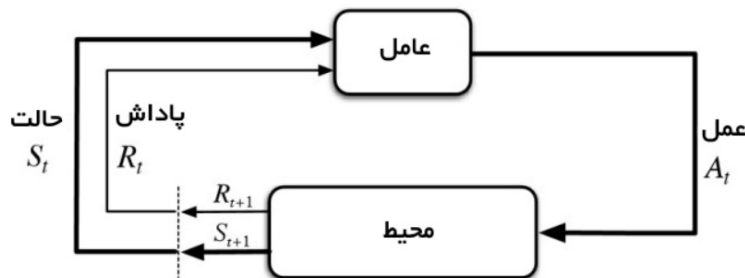
- محیط^۱: جهان فیزیکی که عامل در آن عمل می کند.
- حالت^۲: موقعیت کنونی عامل
- پاداش^۳: بازخورد از محیط

^۱ Environment

^۲ State

^۳ Reward

- سیاست^۱: روشی برای نگاشت حالت عامل به عمل
- ارزش^۲: پاداش آینده که یک عامل با اقدام به یک عمل در یک حالت خاص به آن دست می‌یابد.



شکل ۱۳-۲ الگوریتم‌های یادگیری تقویتی.

به منظور ساخت یک سیاست بهینه، عامل با مساله جستجوی حالت‌های جدید در حالیکه پاداش را نیز به طور هم‌زمان بیشینه می‌کند مواجه است. به این کار موازنه جستجو و استخراج^۳ گفته می‌شود. فرآیندهای تصمیم‌گیری مارکوف^۴ چارچوب‌های ریاضی هستند که برای تشریح محیط در یادگیری تقویتی استفاده می‌شوند و تقریباً همه مسائل این حوزه قابل رسمی شدن با MDPها هستند. یک MDP شامل مجموعه‌ای متناهی از حالت‌های محیط S ، مجموعه‌ای از اعمال ممکن $A(s)$ در هر حالت، یک تابع پاداش دارای مقدار حقیقی $R(s)$ و مدل انتقال $P(s', s|a)$ است. اگرچه، برای محیط‌های جهان واقعی فقدان هرگونه دانش اولیه‌ای از محیط وجود دارد. روش‌های مستقل از مدل یادگیری تقویتی در چنین شرایطی مفید واقع می‌شوند. Q-learning معمولاً به عنوان یک رویکرد مستقل از مدل مورد استفاده قرار می‌گیرد و برای ساخت عامل خودبازی‌کن قابل استفاده خواهد بود. این روش حول محور به‌روز رسانی ارزش‌های Q است که ارزش انجام عمل a را در حالت s نشان می‌دهد. قاعده به‌روز رسانی ارزش، هسته الگوریتم Q-learning محسوب می‌شود.

$$Q(s_t, a_t) \leftarrow (1 - \alpha) \cdot \underbrace{Q(s_t, a_t)}_{\text{مقدار قدیمی}} + \underbrace{\alpha}_{\text{نرخ یادگیری}} \cdot \left(\underbrace{r_t}_{\text{پاداش}} + \underbrace{\gamma}_{\text{فاکتور تنزیل}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{تخمین ارزش آینده مطلوب}} \right)$$

ارزش یاد گرفته شده

در این بخش مروری بر الگوریتم‌های اصلی یادگیری تقویتی استفاده از چارچوب‌های یادگیری عمیق گسترش صورت می‌گیرد.

- روش‌های مبتنی بر ارزش: این روش‌ها احتمال قرار گرفتن در یک حالت معین را تخمین می‌زنند، با استفاده از آن سیاست کنترل تعیین می‌شود. تخمین حالت متوالی با استفاده از معادلات بلمن (معادله امید بلمن و معادله بهینه سازی بلمن) انجام می‌شود. الگوریتم‌های RL مبتنی بر ارزش شامل

¹ Policy

² Value

³ Exploration vs Exploitation trade-off

⁴ Markov Decision Processes (MDP)

(SARSA¹) و Q-Learning هستند، که در اهدافشان متفاوت هستند، این همان مقدار هدف است که Q به صورت بازگشتی و به اندازه یک پله در هر مرحله آپدیت می شوند. SARSA یک روش بر مبنای سیاست است که در آن تخمین ارزش نسبت به یک سیاست به روز می شود در حالی که Q-Learning، یک روش بدون سیاست است که تخمین ارزش را به سمت یک سیاست بهینه هدف به روز می کند. این الگوریتم یک الگوریتم پیچیده است که برای بیان مشکلات مختلف ظاهر چندگانه استفاده می شود اما محدودیت های محاسباتی به مانع بهره بردن از آن می شود.

- روش های مبتنی بر سیاست: برخلاف روش های مبتنی بر ارزش، روش های مبتنی بر سیاست مستقیماً سیاست را به روز می کنند بدون اینکه به تخمین ارزش نگاه کنند و از نظر همگرایی، حل مسائل با داده های پیوسته در ابعاد بالا و حل سیاست های غیر تصادفی به طور اثربخش، کمی بهتر از روش های مبتنی بر ارزش هستند. آنها به دو روش مبتنی بر گرادیان و بدون گرادیان [۳۱، ۳۲] پارامترها را تخمین می زنند. اینجا بر روی روش های مبتنی بر گرادیان تمرکز می شود که الگوریتم بهینه سازی گرادیان نزولی انتخاب شده است. در اینجا، ما تابع هدف را بهینه می کنیم:

$$J(\pi_{\theta}) = E_{\pi_{\theta}} [f_{\pi_{\theta}}(.)] \quad (۱-۲)$$

که در آن تابع امتیاز [۳۳] برای سیاست π_{θ} توسط $f_{\pi_{\theta}}(.)$ داده می شود. با استفاده از معادله بالا، می توانیم در مورد عملکرد مدل با توجه به وظیفه انجام شده نظر دهیم. الگوریتم RL به سادگی بازده نمونه را برابر با تابع امتیاز قرار می دهد:

$$f_{\pi_{\theta}}(.) = G_t \quad (۲-۲)$$

اصطلاح $b(s)$ از بازده نمونه کم می شود تا واریانس تخمین که معادله زیر را به روز می کند، کاهش دهد:

$$f_{\pi_{\theta}}(.) = G_t - b_t(s_t) \quad (۳-۲)$$

در حالی که از تابع ارزش Q استفاده می کنید، تابع امتیاز می تواند از گرادیان سیاست تصادفی یا گرادیان سیاست قطعی استفاده شده کند که در زیر داده شده است:

$$\nabla_{\theta}(\pi_{\theta}) = E_{s,a} [\nabla_{\theta} \log \pi_{\theta}(a|s) \cdot Q^{\pi}(s,a)] \quad (۴-۲)$$

و

$$\nabla_{\theta}(\mu_{\theta}) = E_s [\nabla_{\theta} \mu_{\theta}(s) \cdot Q^{\mu}(s, \mu_{\theta}(s))] \quad (۵-۲)$$

مشاهده شده است که این روش مطمئناً از نظر محدودیت های زمانی و مکانی محاسباتی بر روش قبلی غلبه دارد. ولی هنوز هم نمی توان آن را به کارهایی شامل تعامل با محیط هایی که به طور پیوسته در حال تکامل هستند و

¹ State-Action-Reward-State-Action

نیاز به عامل تطبیق دارند، تعمیم داد. توجه کنید که به دلیل مسائل ایمنی و محدودیت‌های سخت افزاری، استفاده از سیاست گرادیان در عمل مناسب نیست. بنابراین، بهینه سازی با استفاده از گرادیان سیاست در سیاست‌های تصادفی انجام می‌شود.

روش بازیگر منتقد: روش بازیگر منتقد: این الگوریتم‌ها نمایشی واضح از تخمین سیاست‌ها و حالت‌ها را نگه می‌دارند. تابع امتیاز برای این کار با جایگزینی بازده G_t از معادله روشهای مبتنی بر سیاست با $Q^{\pi_\theta}(s_t, a_t)$ و $b(s)$ پایه با $V^{\pi_\theta}(s_t)$ بدست می‌آید که منجر به معادله زیر می‌شود:

$$f_{\pi_\theta}(\cdot) = Q_{\pi_\theta}(s_t, a_t) - V_{\pi_\theta}(s_t) \quad (6-2)$$

تابع مزیت $A(s, a)$ توسط رابطه زیر بدست می‌آید:

$$A(s, a) = Q(s, a) - V(s) \quad (7-2)$$

روش‌های بازیگر منتقد می‌تواند به عنوان تلاقی روش‌های مبتنی بر سیاست و روش‌های مبتنی بر ارزش توصیف شود، که روش‌های یادگیری تکراری هر دو روش را با هم ترکیب می‌کند.

یکپارچه سازی برنامه ریزی و یادگیری: روشهایی وجود دارد که در آن آنها عامل از تجربیات خود می‌آموزد و می‌تواند روال‌های خیالی را جمع آوری کند. چنین روشهایی با استفاده از روشهای DRL ارتقا یافته‌اند. این روش‌ها در گسترش استفاده از تکنیک‌های RL در سیستم‌های رباتیک نرم به کار رفته‌اند، زیرا وجود درجات آزادی زیاد در فضای حالت، منجر به تعامل هزینه بر با محیط و در نتیجه مصالحه در مورد داده‌های آموزش موجود می‌شود.

۲-۳-۴. یادگیری تقویتی عمیق و کاربرد آن در کپسول آندوسکوپی

این ایده که انسانها با تعامل با محیط یاد می‌گیرند، احتمالاً اولین ایده ای است که وقتی به ماهیت یادگیری می‌اندیشیم، به ذهن ما خطور می‌کند. با استفاده از این ارتباط، اطلاعات زیادی در مورد علت و معلول، در مورد عواقب اعمال و اینکه در جهت دستیابی به اهداف باید چه کاری انجام شود، تولید می‌شود [۱۳۹]. مانند انسانهایی که در حل طیف وسیعی از مشکلات چالش برانگیز از کنترل حرکتی سطح پایین تا کارهای شناختی سطح بالا با تعامل با محیط، سرآمد هستند، اخیراً عوامل مصنوعی می‌توانند در الگوی یادگیری تقویتی، دانش خود را مستقیماً از ورودی‌های خام با دریافت پاداش یا مجازات بسازند [۱۴۰]، [۱۴۱].

با انگیزه موفقیت تکنیک‌های یادگیری تقویتی عمیق در کنترل رباتیک، به ویژه روش‌های Q-Learning عمیق بازیگر منتقد [۱۴۱]، یک راه حل برای کنترل تجهیزات پزشکی پیچیده، به ویژه WCE ها استفاده از یادگیری تقویتی عمیق است. سیستم کنترل مبتنی بر یادگیری و مبتنی بر داده، به مدل سازی پیچیده مبتنی بر فیزیک سیستم نیاز ندارد. این کار، نیاز به دانش تخصصی و کار مهندسی جهت مدل سازی پیچیده حرکت در محیط‌های چالش

برانگیز مانند دستگاه GI با حرکات پرستالتیک را از بین می برد. همچنین سیستم کنترل پیشنهادی قادر به سازگاری با بیماران مختلف و اندامهای موجود در دستگاه گوارش است. شواهد تجربی از چنین سازگاری از طریق آزمایش های گسترده در چندین مورد معده گوشت خوک انجام شده است. این دو مورد مهم از بین مزایای کنترل کننده یادگیری تقویتی عمیق است.

ایده اصلی یادگیری تقویتی این است که یک عامل مصنوعی یاد می گیرد در تعامل با محیط و با بررسی کیفیت اقدام انجام شده از طریق اخذ نمرات پاداش دریافتی، چگونه رفتار خود را بهینه کند. این رویکرد اصولاً در مورد هر نوع مسئله تصمیم گیری متوالی با تکیه بر تجربه گذشته اعمال می شود. محیط ممکن است تصادفی باشد، نماینده فقط ممکن است اطلاعات جزئی در مورد وضعیت فعلی را مشاهده کند، مشاهدات ممکن است در ابعاد بالا باشد (به عنوان مثال، فریم ها و سری های زمانی)، عامل ممکن است آزادانه تجربیات خود را در محیط جمع کند یا برعکس، داده ها ممکن است محدود باشد (به عنوان مثال، عدم دسترسی به شبیه ساز دقیق یا داده های محدود به دلیل هزینه های زیاد یا حریم خصوصی داده ها و غیره). یادگیری تقویت عمیق (DRL) که توانایی استدلال یادگیری تقویت (RL) را با قدرت بازنمایی یادگیری عمیق ترکیب می کند، در سال های اخیر به عوامل بسیار موفق منجر شده است که می توانند از پس مشکلات تصمیم گیری متوالی غیرخطی بآید. DRL به دلیل توانایی یادگیری سطوح مختلف انتزاع از داده ها، در مسائل فضای حالت با ابعاد بالا و کارهای پیچیده با دانش قبلی پایین، کاربرد دارد.

۴-۲. کنترل ربات به وسیله یادگیری تقویتی عمیق

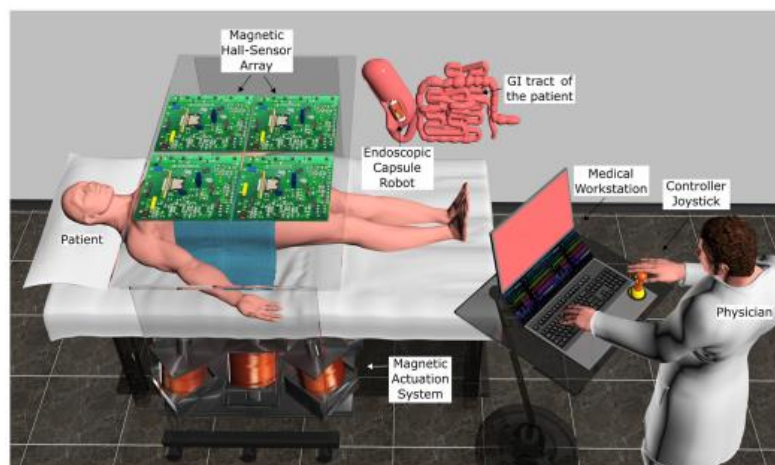
یک موقعیت صحیح و دقیق و کنترل جهت گیری WCE، پزشک را قادر می سازد ربات را به منطقه مربوطه با دقت میلی متر ببرد، داروها به مکان های خاص تحویل داده شود و یا داده های فیزیولوژیکی با وضوح مکانی بالا را اندازه گیری کند. علاوه بر این، بهره گیری از کنترل بهبود یافته، تصویربرداری تشخیصی را نیز ممکن می سازد. به عنوان مثال، WCE های نسل دوم تقریباً در ۱۴٪ از بیماران قادر به تشخیص پولیپ^۱ نیستند [۱۴۲]. به دست آوردن کنترل کامل بر وضعیت ربات کپسول ممکن است میزان خطا را بهبود بخشد، و مهمتر از همه، ربات های کپسول فعال می توانند چندین زاویه دید را تسهیل کنند، در صورت کشف مکان های مورد آسیب، میزان اطلاعات پزشک را بهبود می بخشند [۱۴۳].

۴-۲-۱. تحریک مغناطیسی

تحریک مغناطیسی از راه دور که نوعی از عملگر کپسول بی سیم خارجی است، یک مزیت منحصر به فرد دارد که نیازی به محرک های روی ربات ندارد. این امر پیچیدگی الکترومکانیکی دستگاه را بسیار کاهش می دهد

^۱ polyp

و نیاز به انرژی و فضای اضافی را برطرف می‌کند. سناریوی کاربردی احتمالی کپسول آندوسکوپی قابل کنترل در بیمارستان‌ها را نشان می‌دهد، جایی که پزشک ربات را به مناطق مورد نظر جهت نظارت، تحویل دارو و یا عملیات شبیه بیوپسی می‌برد. آهنرباهای دائمی یا الکترومغناطیسی ممکن است در خارج از بدن بیمار قرار داشته باشد تا نیروی و گشتاور مطلوبی را به کپسولی که در داخل بدن قرار دارد تحمیل کند [۳۰].



شکل ۱۴-۲ سیستم تحریک مغناطیسی.

بسته به محیط، حرکت کپسول ممکن است فرم‌های مختلفی داشته باشد، از جمله کشیدن روی سطوح، حرکت طولی در مایعات یا چرخش به دور محور roll [۸۱]، [۳۲]. تکنیک‌های کنترلی که برای بستن حلقه به پزشک متکی هستند در بعضی موارد به دلیل جاذبه پایدار بین آهنربا خارجی و کپسول امکان‌پذیر است [۳۲]، اما برای دستیابی به ناوبری با دقت بالا یک کنترل‌کننده کاملاً اتوماتیک لازم است. بستن حلقه نیاز به اندازه‌گیری دقیق وضعیت کپسول‌ها دارد و چندین روش غیر دیدی برای این منظور وجود دارد [۱۴۴]، [۷۸]. با این حال، دستکاری مغناطیسی حلقه بسته WCE‌ها همچنان چالش برانگیز است زیرا به طور معمول به مدل‌های دقیق میدان‌های مغناطیسی تولید شده توسط محرک‌ها نیاز دارد تا وابستگی موقعیت ربات‌ها و جهت‌گیری آنها در این زمینه به طور مناسب پیش‌بینی شود [۳۲]. بعلاوه، این تکنیک‌ها اغلب به کالیبراسیون مدل‌های میدانی [۱۴۵] و مدل‌های برهم‌کنش ربات و بافت نیاز دارند.

۲-۴-۲. انواع حرکت کپسول (پیچشی)

از نظر عملکرد مغناطیسی برای ربات‌های کپسولی، در سال‌های اخیر روش‌های مختلفی ارائه شده است. محققان در مقاله [۳۶] از کویل‌های مغناطیسی برای تحریک WCE‌ها در داخل معده پر آب برای اهداف غربالگری استفاده کردند که ۱۰ نوع حرکت اساسی را فراهم می‌کند. کلر و همکارانش کنترل حلقه بسته را ارائه نمی‌کند و رانش‌های بزرگی در قسمت‌های بالای معده دارد. علاوه بر این، از این روش کنترل نمی‌توان برای بیوپسی و تحویل دارو استفاده کرد زیرا به نظر نمی‌رسد بر اساس نتایج ارائه شده دقت کافی را ارائه دهد. جهت پیش‌ران کپسول،

دو قطبی مغناطیسی چرخشی که می‌تواند حرکتی مانند پیچ برای WCE ایجاد کند، توسط پژوهش [۸۲] پیشنهاد شده است. اگرچه نشان داده شده است که چنین دستکاری مغناطیسی در سطوح روده مانند عملی است، اما عملیاتی بودن این روش در یک فضای بیولوژیکی متورم، مانند معده، نشان داده نشده است. روش مقاله [۱۴۳] علاوه بر کنترل موقعیت و جهت گیری در رزولوشن بالا از جهت محدودیت نداشتن ارجح است زیرا که روش پیشراانه ی پیچی به حرکت ۵ درجه آزادی محدود می‌شود.

۲-۴-۳. دلیل استفاده از یادگیری تقویتی

در صورت کنترل نیروهای مغناطیسی به صورت حلقه بسته، حرکت از طریق کنترل جهت می‌تواند به حرکت در سطوح جامد منجر شود. تکنیک‌های کنترل حلقه بسته که ممکن است مبتنی بر کنترل کننده‌های استاندارد مانند تناسبی-انتگرالی و تناسبی-انتگرالی-مشتقی باشد [۳۲]، [۱۴۶] بیش از حد به دقت مدل مورد استفاده برای توصیف تعامل بین ربات و میدان مغناطیسی خارجی بستگی دارد. ماهیت پیچیده و غیرخطی میدان‌های مغناطیسی خارجی تحمیل شده، بسته شدن حلقه را با چنین کنترل کننده‌های وابسته به مدل سنتی دشوار می‌کند. علاوه بر عدم قطعیت در میدان مغناطیسی، عدم قطعیت‌هایی در محیط کار وجود دارد که وظیفه کنترل را بیشتر پیچیده می‌کند. خصوصیات مکانیکی و هندسه متفاوت از یک مکان به مکان دیگر، اختلالات ناشناخته ای نظیر پرستالسیس، فعل و انفعالات ربات و بافت، ناهنجاری‌های روی بافت (مانند تومورهای موجود) یا مواد مغناطیسی اطراف ممکن است تأثیر قابل توجهی بر دینامیک کلی سیستم داشته باشند.