

دانشگاه صنعتی امیر کبیر (پلی تکنیک تهران)

گروه مستقل رباتیک

پایاننامه کارشناسی ارشد گرایش هوش مصنوعی و رباتیک

یادگیری تقلیدی مبتنی بر مدلهای احتمالاتی در کاربردهای رباتیک

نگارش علی جوادی

استاد راهنما دکتر شیری قیداری دکتر نیک آبادی

تیر ۱۳۹۷

به نام خدا تعهدنامه اصالت اثر تاریخ:



اینجانب علی جوادی متعهد می شوم که مطالب مندرج در این پایان نامه حاصل کار پژوهشی اینجانب تحت نظارت و راهنمایی اساتید دانشگاه صنعتی امیر کبیر 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 پایان نامه 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 پایان نامه 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 پایان نامه بدون موافقت کتبی دانشگاه صنعتی امیرکبیر ممنوع است. نقل مطالب با ذکر مآخذ بلامانع است.

على جوادى

امضا

چکیده

امروزه استفاده از رباتها در صنایع مختلف و در وظایف متنوع به سرعت رو به گسترش است. استفاده از رباتها باعث کاهش هزینه ی نیروی انسانی و افزایش دقت در عملکرد میگردد. یکی از وظایف معمول که در صنایع مختلف بر عهده ی رباتها قرار داده می شود، وظیفه ی انتقال جسم از یک موقعیت اولیه به موقعیت هدف است. برای آموزش انجام وظایف به ربات الگوریتمهای متنوعی ارائه شده است که دو دسته ی مهم از این الگوریتمها، الگوریتمهای یادگیری تقویتی هستند. در یادگیری تقلیدی یک عامل طوری الگوریتمهای یادگیری تقویتی هستند. در یادگیری تقلیدی یک عامل طوری آموزش داده می شود تا یک کار را از روی نمایشها با یادگیری نگاشتی بین مشاهدات و اعمال انجام دهد. در الگوریتمهای یادگیری تقویتی نحوه انجام وظیفه از طریق تعیین پاداش و جریمه به ازای هر یک از اعمال ربات، به ربات آموزش داده می شود.

در این پایاننامه، برای بهبود عملکرد الگوریتم یادگیری تقویتی Q-Learning، از مسیرهای راهنما استفاده شده است که این ایده برگرفته از روشهای یادگیری تقلیدی میباشد. جهت استخراج مسیرهای راهنما میتوان از ویدئوهایی استفاده کرد که در آن یک جسم توسط انسان به داخل حفره هدایت میشود.

اخیرا، برای حل مسائلی که طراحی مدل دینامیکی ربات با پیچیدگی بسیاری همراه است و یا مدل دینامیکی دارای عدم قطعیت میباشد، الگوریتم PILCO مورد استفاده قرار می گیرد. این الگوریتم که از استنتاج و مدلهای احتمالاتی برای یادگیری کنترل استفاده می کند، پیش از این در مسائل مختلفی مورد استفاده قرار گرفته است، اما از آن در حل مسئله انتقال جسم به درون یک حفره توسط بازوی رباتیک شبیهسازی شده بهره برده نشده است. با توجه به وجود ذات عدم قطعیت در این مسئله، الگوریتم PILCO می تواند به عنوان راهکاری مناسب برای حل این مسئله در محیطهای پیوسته مورد استفاده قرار بگیرد. از اینرو در این پایاننامه روشی بر مبنای الگوریتم به نحوی تغییر داده شده است که محدودیتهای فضای حالت را نیز در روند محاسباتش در نظر می گیرد. جهت ارزیابی روشهای پیشنهادی، شبیهسازیها در محیطهای گسسته و پیوسته انجام شده است و لازم به ذکر است که جهت انتقال جسم از نقطه اولیه به نقطه هدف از یک ربات شبیهسازی شده استفاده شده است. نتایج شبیهسازی در این دو محیط گسسته و پیوسته حاکی از عملکرد بهتر شبیهسازی شده است به الگوریتمهای اصلی یعنی الگوریتمهای و Q-Learning و بیوسته حاکی از عملکرد بهتر الگوریتمهای پیشنهادی نسبت به الگوریتمهای اصلی یعنی الگوریتمهای و Q-Learning و بیوسته حاکی از عملکرد بهتر

واژههای کلیدی:

یادگیری تقلیدی، مسیرهای راهنما، یادگیری تقویتی،مدلهای احتمالاتی، رباتیک

فهرست مطالب

1	1– مقدمه
۴	<u>ن</u> صل دوم
۴	۴– مروری بر ادبیات و کارهای پیشین
۵	۱-۲– مرور کلی روشهای یادگیری تقلیدی و یادگیری تقویتی
	۱–۱–۲ روشهای یادگیری تقلیدی
۶	١-١-١-٢ مقدمه
11	٣-١-١-٢- بازنمايي ويژگي
17	۳–۱–۱-۲ یادگیری مستقیم
١٧	۴-۱-۱-۲ یادگیری غیرمستقیم
	٢-١-٢ روشهای یادگیری تقویتی
٣٢	۱-۲-۱-۲- حل فرآیندهای تصادفی مارکف
٣٩	۲-۲ روشهای یادگیری تقویتی با جستجوی سیاست
۴٠	١–٢–٢ مقدمه
	٢-٢-٢ جستجوى سياست فارغ از مدل
	۳–۲–۲ جستجوی سیاست مبتنی بر مدل
	۱-۳-۲-۲ مدلهای پیشرو احتمالاتی
	۲-۳-۲ پیش بینیهای بلندمدت با یک مدل معلوم
	۳-۲-۲-۳ بهروز رسانیهای سیاست

۶۶	۴-۳-۲-۲ الگوریتمهای جستجوی سیاست مبتنی بر مدل با کاربردهای رباتیکی
	۵-۳-۲-۲ ویژگیهای مهم روشهای مبتنی بر مدل
	٣-٢- جمع بندى
	فصل سوم
۸١	٣– روشهای پیشنهادی
	۳-۱ الگوريتم Q-Learning
	۳-۲ روش پیشنهادی: استفاده از مسیرهای راهنما در Q-Learning
ለ۶	۱-۲-۳ ایجاد مسیرهای راهنما از روی ویدئوهای ضبط شده
۸۸	۲-۲-۳ مقداردهی اولیه ماتریس ${ m Q}$ با استفاده از مسیرهای راهنما
۵ ۸۸	۳-۳ استفاده از مدل PILCO در حل مسئله انتقال جسم به درون حفره در محیطهای پیوست
	۱-۳-۳ یادگیری سیاست با محدودیتهای فضای حالت
۹٠	۱-۱-۳-۳ مدل دینامیکی احتمالاتی
	۲-۱-۳-۳ برنامهریزی بلندمدت از طریق استنتاج تقریبی
۹۵	۳-۱-۳-۳ یادگیری کنترلکننده از طریق جستجوی غیرمستقیم سیاست
۹۵	۴-۱-۳-۳ برنامهریزی با محدودیتهای فضای حالت
٩٨	۴-۳- جمعبندی
۹۹	فصل چهارم
۹۹	۴– تحلیل و ارزیابی نتایج
١٠٠	١-۴- معرفي محيطهاي شبيهسازي
١٠٠	١-١-١- و ثاكرهاي محبط گسسته

1 • 1	۲-۱-۴ ویژگیهای محیط پیوسته
	۲-۴- نحوه آمادهسازی ویدئوهای مسیر راهنما
	۳-۴- تنظیم پارامترها و معیار ارزیابی
	۱-۳-۴ پارامترها و معيار ارزيابي الگوريتم Q-Learning
1.5	۳-۳-۲ پارامترها و معيار ارزيابي الگوريتم PILCO
	۴-۴- نتایج آزمایش
11.	۱-۴-۴- بررسی تاثیر اندازه جسم گرفته شده توسط ربات
117	۲-۴-۴ بررسی تاثیر قابلیت چرخاندن جسم گرفته شده توسط ربات
118	۳-۴-۴ بررسی تاثیر تعداد مسیرهای راهنما
	۴-۴-۴ بررسی تاثیر اندازه فضای حالت
بط پیوسته	۵-۴-۴ بررسی تاثیر در نظر گرفتن محدودیتهای فضای حالت در محب
179	۵-۴- جمع بندی
18	فصل پنجم
14	۵– نتیجهگیری و پیشنهادها
171	۱ –۵– نتیجه گیری
171	۱ –۵– نتیجه گیری
	مراجع

فهرست اشكال

۲	شکل ۱-۱ نمایی از قرار دادن قطعه کامپیوتری درون محفظه توسط ربات در مراحل مونتاژ کامپیوتر
	شکل ۲-۱ فلوچارت یادگیری تقلیدی
۱٧	شکل ۲-۲ مثالی از یادگیری سلسلهمراتبی عملها
	شکل ۲-۳ روشهای یادگیری از منابع مختلف
	شكل ٢-۴ فرم كلى الگوريتمهاى حل فرآيندهاى تصميم ماركف
۴۶	شکل ۲-۵ چرخه کلی در یادگیری تقویتی مبتنی بر مدل
۴٧	شکل ۲-۶ خطاهای مدل در جستجوی سیاست مبتنی بر مدل
	شکل ۲-۲ مدل گرافی برای رگرسیون خطی بیزی
۵٩	شکل ۲-۸ محاسبه توزیع پیشبینی کننده تقریبی با استفاده از خطیسازی
۶۰	شکل ۲-۹ محاسبه توزیع پیشبینیشده تقریبی با استفاده از انتقال بیرنگ و بو
۶۱	شکل ۲-۱۰ محاسبه توزیع پیشبینیشده تقریبی با استفاده از تطبیق ممان
۶۹	شکل ۲-۱۱ روشهای جستجوی سیاست مبتنی بر مدل با استنتاج تصادفی
	شکل ۲-۱۲ ترکیب مقدم پارامتری و فرآیندهای گوسی
٧٣	شکل ۲-۱۳ پیشبینیهای تقریبی با فرآیندهای گوسی در ورودیهای نامطمئن
٧۴	شکل ۲-۱۴ موفقیتهای یادگیری PILCO
	شکل ۲-۱۵ مقایسه عملکرد PILCO با روشهای دیگر
	شکل ۳-۱ چهار نمای مختلف از یکی از ویدئوهای ضبط شده مسیر راهنما
	شکل ۳-۲ پیش,بینی GP در یک ورودی نامطمئن
	شکل ۳-۳ تابع هزینهای که محدودیتها (برای مثال موانع) را با "نامطلوب" کردن آنها در نظر میگیرد.
١٠١.	شکل ۴-۱ نمایی از یک محیط شبیهسازی گسسته
۱۰۲.	شکل ۴-۲ نمایی از یک محیط شبیهسازی پیوسته
١٠٣.	شکل ۴-۳ نمایی از ویدئوهای مسیر راهنما

شکل ۴-۴ کوتاهترین مسیر از نظر فاصله اقلیدسی از موقعیتهای اولیه مختلف به هدف.....

فهرست جداول

۵١	ول ۲-۲ رویکردهای جستجوی سیاست مبتنی بر مدل	جد
٧۶	ول ۲-۲ مرور الگوریتمهای جستجوی سیاست مبتنی بر مدل با کاربردهای رباتیکی	جد
مبتنی بر مدل۷۸	ول ۲-۳ ویژگیهای پیشبینیهای قطعی و تصادفی تراژکتوری در جستجوی سیاست	جد
111	ول ۴-۱ پارامترهای آزمایش بررسی تاثیر اندازه جسم گرفته شده توسط ربات	جد
117	ول ۴-۲ نتایج حاصل از آزمایش بررسی تاثیر اندازه جسم گرفته شده توسط ربات	جد
ت	ول ۴-۳ پارامترهای آزمایش بررسی تاثیر قابلیت چرخاندن جسم گرفته شده توسط رب	جد
۱۱۵	ول ۴-۴ نتایج حاصل از بررسی تاثیر قابلیت چرخاندن جسم گرفته شده توسط ربات	جد
\ \ \ \	ول ۴-۵ پارامترهای آزمایش بررسی تاثیر تعداد مسیرهای راهنما	جد
١١٨	ول ۴-۶ نتایج بررسی تاثیر تعداد مسیرهای راهنما در انتقال جسم به اندازه یک سلول	جد
ندن ۱۱۹	ول ۴-۷ نتایج بررسی تاثیر تعداد مسیرهای راهنما در انتقال جسم بدون قابلیت چرخ	جد
17	ول ۴-۸ نتایج بررسی تاثیر تعداد مسیرهای راهنما در انتقال جسم با قابلیت چرخاندن	جد
171	ول ۴-۹ پارامترهای آزمایش بررسی تاثیر اندازه فضای حالت	جد
177	ول ۴-۲۰ نتایج بررسی تاثیر اندازه فضای حالت	جد
177	ول ۱۱-۴ پارامترهای آزمایش بررسی تاثیر در نظر گرفتن محدودیتهای فضای حالت	جد
179	ول ۴-۱۲ نتایج بررسی تاثیر در نظر گرفتن محدودیتهای فضای حالت	جد

فصل اول ۱- مقدمه

مقدمه

در دهههای اخیر و با گسترش صنایع مختلف، به منظور کاهش هزینه نیروی انسانی و افزایش دقت عملکرد استفاده از ابزارهای خودکار برای انجام وظایف مختلف بسیار مورد توجه قرار گرفته است. رباتها به عنوان یکی از مهمترین ابزارهای انجام خودکار وظایف در صنایع مختلف مانند خودروسازی، ساخت تجهیزات پزشکی، ساخت تجهیزات کامپیوتر و ... مطرح هستند. رباتها قادرند بسیاری از وظایف را با دقت بالاتری نسبت به انسان انجام دهند. نمونهای از کوچکترین وظایف که میتواند توسط یک ربات انجام شود در شکل ۱-۱ آمده است. در این شکل، وظیفه قرار دادن یک قطعه کامپیوتری درون یک محفظه در مراحل مونتاژ یک سیستم کامپیوتری بر عهده ربات قرار داده شده است. پیدا کردن محل مناسب قطعه و انتقال جسم به آن نقطه از اساسی ترین چالشهایی است که ربات با آن روبرو است.



شکل ۱-۱ نمایی از قرار دادن قطعه کامپیوتری درون محفظه توسط ربات در مراحل مونتاژ کامپیوتر.[1]

برای آموزش ربات به منظور انجام وظایف مختلف، الگوریتمهای متنوعی ارائه شده است. دو دسته از مهمترین الگوریتمها برای این منظور، الگوریتمهای یادگیری تقلیدی و یادگیری تقویتی هستند. روش کار در الگوریتمهای یادگیری تقلیدی به این صورت است که نمونههایی از انجام درست وظیفه به ربات نمایش داده می شود تا ربات بتواند نحوه انجام وظیفه را از روی نمونهها بیاموزد. در الگوریتمهای یادگیری تقویتی، نحوه

انجام وظیفه از طریق تعیین پاداش و جریمه به ازای هر یک از اعمال ربات به ربات آموزش داده می شود. در این پایان نامه، الگوریتم یادگیری تقویتی Q-Learning [2]، با استفاده از تزریق مسیرهای راهنما به سمت یک الگوریتم یادگیری تقلیدی سوق داده شده است. انجام این کار باعث بهبود عملکرد الگوریتم و کاهش زمان یادگیری می گردد. در بخش دوم از روش پیشنهادی ارائه شده در این پایان نامه، الگوریتم یادگیری تقویتی یادگیری می PILCO آق به نحوی تغییر داده شده است که امکان در نظر گرفتن محدودیتهای محیطی از جمله دیوارها در تابع هزینه الگوریتم وجود داشته باشد. افزودن این قابلیت به الگوریتم، عملکرد ربات را از جنبه درصد موفقیت رساندن جسم به هدف بهبود داده است. هر دو روش پیشنهادی، در وظیفه انتقال جسم از نقطه اولیه به درون حفره مورد بررسی و ارزیابی قرار گرفته اند، با این تفاوت که الگوریتم Q-Learning در محیطهای پیوسته مورد استفاده قرار گرفته است.

روند ادامه ی بخشهای این پایان نامه به این صورت است که در فصل دوم مروری بر کارهای پیشین صورت گرفته در حوزههای یادگیری تقلیدی و یادگیری تقویتی خواهیم داشت و دسته بندی موجود در این حوزهها را شرح خواهیم داد. سپس در فصل سوم پس از شرح الگوریتمهای پایه ی Q-Learning و PILCO و روش پیشنهادی مربوط به پیشنهادی مربوط به استفاده از مسیرهای راهنما در الگوریتم PILCO و روش پیشنهادی مربوط به ترزیق دانش محدودیتهای محیط در الگوریتم PILCO توضیح داده می شود. در فصل چهارم محیطهای شبیه سازی مورد استفاده معرفی شده و نتایج آزمایشها در دو محیط پیوسته و گسسته با پارامترهای مختلف بیان می گردد. در فصل پنجم نیز نتیجه گیری و جمع بندی نهایی پایان نامه ارائه خواهد شد.

فصل دوم

۲- مروری بر ادبیات و کارهای پیشین

مروری بر ادبیات و کارهای پیشین

در این فصل مروری بر کارهای پیشین مرتبط با روشهای یادگیری خواهیم داشت. ابتدا در بخش ۱-۲ مرور کلی بر روشهای یادگیری تقویتی ارائه میشود. سپس با توجه مدل استفاده شده در این پایاننامه، در بخش ۲-۲ روشهای یادگیری تقویتی با جستجوی سیاست و به خصوص روشهای جستجوی سیاست مبتنی بر مدل با تفصیل بیشتری شرح داده میشوند.

۱-۲- مرور کلی روشهای یادگیری تقلیدی و یادگیری تقویتی

۲-۱-۱- روشهای یادگیری تقلیدی

تکنیکهای یادگیری تقلیدی برای تقلید رفتار انسان در یک کار معین مناسب هستند. یک عامل (ماشین فراگیر) آموزش داده می شود تا یک کار را از روی نمایشها با یادگیری نگاشتی بین مشاهدات و اعمال انجام دهد. ایده آموزش با استفاده از تقلید برای سالیان دراز وجود داشته است، اما این رشته اخیرا به دلیل یمیشرفتهای صورت گرفته در محاسبات و حسگرها و همچنین به دلیل نیاز فزاینده به برنامههای کاربردی هوشمند مورد توجه قرار گرفته است. الگوی یادگیری با استفاده از تقلید در حال محبوب شدن است چون آموزش کارهای پیچیده با در اختیار داشتن حداقل دانشی از دانش خبره در مورد آن کارها را آسان می کند. روشهای یادگیری تقلیدی به طور بالقوه مسئله آموزش یک کار را به تهیه نمایش کاهش می دهند بدون اینکه نیاز به برنامه نویسی صریح یا طراحی توابع پاداش ویژه آن کار باشد. حسگرهای مدرن، به سرعت قادر به جمع آوری و انتقال حجم زیادی از دادهها هستند و پردازندههای با قدرت محاسباتی بالا امکان پردازش سریع را فراهم می کنند تا به موقع دادههای دریافتی را از حسگرها به اعمال نگاشت کنند. این موضوع، بستری مناسب را فراهم می کنند تا به موقع دادههای دریافتی را از حسگرها به اعمال نگاشت کنند. این موضوع، بستری مناسب کامپیوتر و بازیهای کامپیوتری فراهم می کند که به مشاهده و عکسالعمل بلادرنگ نیاز دارند. با این حال، از آنجایی که یادگیری با استفاده از تقلید مجموعه چالشهای خودش را دارد الگوریتهای تخصصی برای یادگیری موثر و قوی مدلها نیاز است. در این بخش، روشهای یادگیری تقلیدی را مورد بررسی قرار می دهیم و گزینههای موثر و قوی مدلها نیاز است. در این بخش، روشهای یادگیری تقلیدی را مورد بررسی قرار می دهیم و گزینههای

را معرفی کرده و همچنین چالشهای ویژه مسئله تقلید را مورد تاکید قرار می دهیم. همچنین، روشهای طراحی و ارزیابی کارهای یادگیری تقلیدی دسته بندی و مرور می شوند. از آنجایی که در منابع علمی روشهای یادگیری در رباتیک و بازی ها محبوب هستند و مجموعه گسترده ای از مسائل و متدلوژی ها را پیش رو قرار می دهند، توجه ویژه ای به این زمینه ها می شود. همچنین در این بخش، به طور گسترده در مورد ترکیب رویکردهای یادگیری تقلیدی با استفاده از روشها و منابع مختلف و نیز در مورد درآمیختن دیگر روشهای یادگیری حرکت برای بهبود تقلید بحث می کنیم.

1-1-1- مقدمه

در سالهای اخیر، اساسا نیاز به عاملهای هوشمند که قابلیت تقلید رفتار انسان را دارند افزایش یافته است. پیشرفتها در رباتیک و فناوری ارتباطات کاربردهای بالقوه زیادی به وجود آورده است که به هوش مصنوعیای نیاز دارند که نه تنها تصمیمات هوشمند می گیرد بلکه همچنین قادر است اعمال حرکتی را در موقعیتهای گوناگون اجرا کند. خیلی از جهت گیریهای آینده در فناوری متکی به این توانایی عاملهای هوش مصنوعی است که هنگام مواجهه با موقعیت مشابه، همانند انسان رفتار کنند. خودروهای خودران، رباتهای امدادرسان و تعامل انسان – کامپیوتر مثالهایی از این رشتهها هستند. خصوصا برای مورد آخر به خاطر گرایش اخیر به واقعیت مجازی مصرف کننده و سیستمهای ضبط حرکت، فرصتها برای کاربردهای جدید در حال افزایش هستند. در این کاربردها و خیلی از کارهای رباتیکی، ما با مسئله اجرای یک عمل با توجه به حالت فعلی عامل و محیط پیرامونش مواجه هستیم. تعداد سناریوهای ممکن در یک کاربرد پیچیده خیلی زیادتر از آن حدی است که با برنامهنویسی صریح پوشش داده شود و از همین رو یک عامل موفق باید قادر باشد سناریوهای دیدهنشده را مدیریت کند. در حالی که چنین کاری می تواند به عنوان یک مسئله بهینه سازی فرموله سازی شود، به طور گسترده پذیرفته شده که داشتن دانش قبلی که توسط فرد خبره فراهم شده است موثرتر و مقرون به صرفهتر از جستجو بدون دانش قبلی برای یک راه حل است[4]. علاوه بر این، بهینهسازی از طریق آزمون و خطا به توابع پاداشی نیاز دارد که مخصوص هر کار طراحی شدهاند. میتوان تصور کرد که حتی برای کارهای ساده، تعداد اعمال ممکن که عامل می تواند انجام دهد به طور نمایی زیاد می شود. تعریف پاداش برای چنین مسئلهای مشکل و در خیلی از موارد نامشخص است. یکی از راههای طبیعی تر و شهودی تر ابلاغ دانش توسط فرد خبره، فراهم کردن نمایشهایی برای رفتار مطلوب است که نیاز است یادگیرنده از آنها تقلید کند. برای (انسان) معلم انتقال دانشش از طریق نمایش آسان تر از این است که آن را طوری که یادگیرنده خواهد فهمید بیان کند [5]. در این بخش روشهایی مرور می شوند که با استفاده از آنها به عاملهای مصنوعی آموزش داده می شود تا سلسله اعمال پیچیده را از طریق تقلید انجام دهند.

یادگیری تقلیدی یک رشته تحقیقاتی میان رشته ای است. به طور ویژه تر، در این بخش روشهای هوش مصنوعی را مرور می کنیم که برای یادگیری سیاستهایی استفاده می شوند که مسائل را مطابق با نمایشهای انسان حل می کنند. این بخش با تمرکز بر روی روشهای یادگیری، یادگیری برای هر عامل هوشمند اعم از اینکه یک ربات فیزیکی باشد یا یک عامل نرم افزاری (از قبیل بازی ها، شبیه سازی ها و غیره) را مورد بحث قرار می دهد. منابع علمی مرور شده کاربردهای متنوع را مورد بررسی قرار می دهند، اما خیلی از روشهایی که استفاده می شوند کلی هستند و می توانند به کارهای عمومی (کلی) یادگیری حرکت اعمال شوند. فرآیند یادگیری برای ایجاد بازنمایی ویژگی، تقلید مستقیم و یادگیری غیرمستقیم دسته بندی می شود. روشها و منابع یادگیری برای هر فرآیند و همچنین معیارهای ارزیابی و کاربردهای مناسب این روشها مرور می شوند.

یادگیری تقلیدی به فراگیری مهارتها و رفتارهای یک عامل از طریق مشاهده یک مربی که کار مورد نظر را به نمایش می گذارد گفته می شود. با الهام از عصب شناسی، یادگیری تقلیدی بخش مهمی از هوش ماشین و تعامل انسان – کامپیوتر است و از ابتدا به آن به عنوان بخش اساسی در آینده رباتیک نگاه شده است [6]. یک الگوی محبوب دیگر، یادگیری از طریق آزمون و خطاست. اما فراهم کردن مثالهای خوب برای یادگیری از آنها، فرآیند پیدا کردن یک مدل عمل مناسب را تسریع کرده و از گیر افتادن عامل در کمینه محلی جلوگیری می کند. علاوه بر این، یادگیرنده می تواند خودش به خوبی به یک راه حل مناسب برسد یعنی راه حلی که به هدف معین قابل سنجشی می رسد، اما با روشی که یک انسان با همین کار مواجه می شود به طور قابل توجهی فرق می کند. گاهی اوقات مهم است که اعمال یادگیرنده قابل باور بوده و طبیعی به نظر برسند. در خیلی از زمینههای رباتیک و همچنین تعامل انسان – کامپیوتر که کارایی یادگیرنده فقط بستگی به فهم و درک انسان مشاهده کننده از آن دارد این امری ضروری است. بنابراین مطلوب است که به یادگیرنده رفتار مطلوب را از مجموعهای از نمونههای دارد این امری ضروری است. بنابراین مطلوب است که به یادگیرنده رفتار مطلوب را از مجموعهای از نمونههای دارد این امری ضروری است. بنابراین مطلوب است که به یادگیرنده به خاطر تغییرات در کار از قبیل محل جمع آوری شده آموزش داد. اما اغلب اوقات این مورد پیش می آید که به خاطر تغییرات در کار از قبیل محل

اشیا یا نمایشهای ناکافی و نامساعد، تقلید مستقیم حرکت فرد خبره کفایت نمیکند. بنابراین، تکنیکهای یادگیری تقلیدی باید قادر باشند که به گونهای یک سیاست را از نمایشهای (در اختیار قرار) داده شده یاد گرفته که بتوانند آن را به سناریوهای دیده نشده تعمیم دهند. بدین ترتیب، عامل به جای اینکه به طور کامل از معلم کپی کند، یاد می گیرد که کار را انجام دهد.

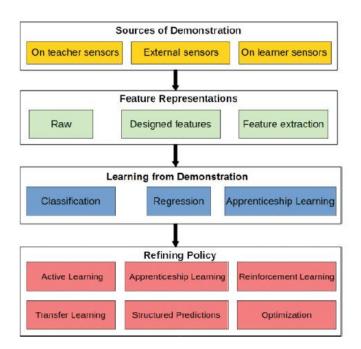
اهمیت رشته یادگیری تقلیدی از ارتباطش با انواع کاربردها از قبیل تعامل انسان-کامپیوتر و رباتهای امدادگر نشأت می گیرد. از این روش برای آموزش آرایهای از کارهای مختلف به رباتهای با اسکلتها و درجه آزادیهای مختلف استفاده شده است. بعضی از مثالها، مسائل ناوبری هستند که معمولا از رباتهای شبیه وسایل نقلیه ربا درجه آزادی نسبتا کمتر) استفاده می کنند که شامل وسایل نقلیه پروازی[7] و وسایل نقلیه زمینی[8] می شوند. کاربردهای دیگر، روی رباتهای با درجات آزادی بالاتر از قبیل رباتهای انسانها[9] و بازوهای رباتیکی [10] تمرکز دارند. رباتهای انسانهای با درجه آزادی بالا می توانند کارهای گسسته از قبیل ایستادن و کارهای چرخهای از قبیل قدمزنی را یاد بگیرند[11]. اگرچه اکثریت کاربردها رباتیک را هدف قرار می دهند، یادگیری تقلیدی برای کارهای شبیهسازی نیز به کار رفته[11] و حتی در بازیهای کامپیوتری استفاده می شود[12].

یادگیری تقلیدی با استخراج اطلاعات درباره رفتار مربی و محیط پیرامون شامل هر شیء دستکاری شده و همچنین، یادگیری یک نگاشت بین وضعیت و رفتار مشاهده شده کار میکند. الگوریتمهای یادگیری ماشین مرسوم، به عاملهای با ابعاد بالا و با درجه آزادی زیاد مقیاس نمیشوند[13]. بنابراین الگوریتمهای به خصوصی نیاز است تا بازنماییها و پیشبینیهای مناسب تولید کنند که قادر به تقلید کارکردهای حرکتی در انسان باشند.

مشابه با یادگیری با ناظر مرسوم که نمونهها جفت ویژگیها و برچسبها را نشان میدهند، در یادگیری تقلیدی نیز نمونهها جفت حالات و اعمال را نشان میدهند که حالت، وضع فعلی عامل شامل موقعیت و سرعت مفاصل مربوطه و در صورت وجود هدف، وضعیت هدف (از قبیل موقعیت، سرعت، اطلاعات هندسی و غیره) را بازنمایی میکند. از همین رو، فرآیندهای تصمیم مارکفی (MDP) به طور طبیعی خودشان را معطوف مسائل یادگیری تقلیدی کرده و به طور رایج برای بازنمایی نمایشهای خبره استفاده میشوند. ویژگی مارکفی تعیین میکند که حالت بعدی فقط به حالت و عمل قبلی وابسته است، موضوعی که نیاز به درنظرگرفتن حالات قبل تر را در

بازنمایی حالت کم می کند [14]. یک چارچوب کاری یادگیری تقلیدی معمول، با به دست آوردن نمایشهای نمونه از یک خبره شروع می شود، نمایشهایی که بعدا به صورت جفت حالت — عمل کدگذاری می شوند. سپس، این نمونه ها برای آموزش یک سیاست استفاده می شوند. اما اغلب اوقات برای رسیدن به رفتار مورد نیاز، یادگیری یک نگاشت مستقیم بین حالت و عمل کافی نیست. این موضوع می تواند بر اثر برخی از مسائل رخ دهد که خطا در به دست آوردن نمایشها، مغایرت در اسکلت مربی و یادگیرنده (مسئله تناظر) یا نمایشهای ناکافی از این قبیل هستند. علاوه بر این، کاری که توسط یادگیرنده اجرا می شود ممکن است به دلیل تغییرات در محیط، موانع یا اهداف کمی از کار نمایش داده شده متفاوت باشد. بنابراین، غالبا یادگیری تقلیدی گام دیگری را دربرمی گیرد که نیاز دارد که یادگیرنده، عمل آموخته شده را اجرا کرده و مطابق با عملکرد خودش در آن کار، سیاست آموخته شده را دوباره بهینه کند. این خودبهبودی می تواند با توجه به یک پاداش قابل سنجش انجام شده یا از نمونهها آموخته شود. خیلی از این رویکردها زیر چتر گسترده یادگیری تقویتی قرار می گیرند.

شکل ۲-۱، چارچوب کاری یک فرآیند یادگیری تقلیدی را نشان می دهد. این فرآیند، با ضبط اعمالی که قرار است از آنها آموخته شود شروع می شود که می توان از طریق روشهای حسی مختلف به آن دست یافت. سپس، داده حسگرها پردازش شده تا ویژگیهایی استخراج شوند که حالت و محیط پیرامون انجام دهنده کار را توصیف می کنند. از ویژگیها برای یادگیری یک سیاست استفاده شده تا با استفاده از آن، رفتار نمایش داده شده را تقلید کرد. سرانجام، با اجازه دادن به عامل برای عمل کردن به سیاست آموخته شده و بهبود آن براساس عملکردش، این سیاست می تواند به تر شود. این گام می تواند به ورودی مربی نیاز داشته و یا نداشته باشد. در نظر گرفتن بهبود سیاست به عنوان یک گام بعد از یادگیری می تواند قابل درک باشد، اما در بسیاری از موارد همراه با یادگیری از نمایشها اتفاق می افتد.



شكل ۲-۱ فلوچارت يادگيري تقليدي.[15]

کاربردهای یادگیری تقلیدی به خاطر ماهیت میان رشته ایشان با چالشهای متنوعی روبرو هستند که عبارتند از:

-شروع این چالشها با فرآیند به دست آوردن نمایشهاست که چه داده را از حسگرهای روی یادگیرنده یا مربی ضبط کرده و چه از اطلاعات بصری استفاده شود، سیگنالهای ضبط شده در معرض نویز و خطاهای حسگر هستند. مشکلات مشابه در طول اجرا، زمانی که عامل محیط را حس می کند پیش می آیند. حس کردن نویزی یا غیرقابل اعتماد منجر به رفتار اشتباه می شود، اگرچه مدل به اندازه کافی آموزش داده شده باشد. مرجع [16] روشهای مختلف جمع آوری نمایشها و چالشهای هر رویکرد را مورد بررسی قرار داده است.

-مسئله دیگر که مربوط به نمایش است، مسئله یا مشکل تناسب میباشد[17]. تناسب، تطبیق قابلیتها، اسکلت و درجه آزادی یادگیرنده و مربی با هم است. نیاز است هر تفاوتی در اندازه یا ساختار بین مربی و یادگیرنده در طول آموزش جبران شود. غالبا در این مورد، یادگیرنده میتواند شکل حرکت را از نمایشها یاد گرفته، سپس مدل را با آزمون و خطا بهتر کرده تا به هدفش برسد.

-یک چالش مرتبط، مسئله مشاهده پذیری است که در آن، سینماتیک مربی برای یادگیرنده شناخته شده نیست [18]. اگر نمایشها، توسط یک مربی برگزیده تهیه نشده باشند، یادگیرنده ممکن است از قابلیتها و اعمال احتمالی مربی آگاه نباشد. یادگیرنده فقط می تواند آثار اعمال مربی را مشاهده کرده و با استفاده از سینماتیک خودش سعی در تکرار آنها داشته باشد.

از آنجایی که تکنیکهای سنتی یادگیری ماشین در درجات آزادی بالا به خوبی عمل نمیکنند، فرآیند یادگیری یادگیری با مشکلات عملی نیز مواجه است[13]. به خاطر ماهیت بلادرنگ بسیاری از کاربردهای یادگیری تقلیدی، الگوریتمهای یادگیری از لحاظ قدرت محاسباتی و محدودیتهای حافظه محدود هستند، به خصوص، در کاربردهای رباتیکی که به کامپیوترهای جاسازی شده ا برای اجرای پردازش بلادرنگ نیاز است.

-علاوه بر این، غالبا رفتارهای پیچیده میتوانند به صورت یک تراژکتوری از اعمال کوچک وابسته دیده شوند که فرض i.i.d بودن که در بیشتر کارهای یادگیری ماشین اتخاذ شده را نقض میکند. سیاست آموخته شده باید قادر باشد که رفتارش را مبتنی بر اعمال قبلی وفق داده و در صورت لزوم ویرایشهایی انجام دهد.

-سیاست باید قادر باشد تا خودش را به تغییرات در کار مورد نظر و محیط پیرامون وفق دهد. ماهیت پیچیده کاربردهای یادگیری تقلیدی این امر را تحمیل می کند که عاملها باید قادر باشند تا کار مورد نظر را حتی تحت شرایطی که با نمایش آموزشی فرق می کند به طور قابل اعتماد اجرا کنند.

-کارهایی که تعامل انسان-کامپیوتر را دربرمی گیرند مجموعه جدیدی از چالشها را پیش رو قرار می دهند. طبیعتا ایمنی، نگرانی اصلی در چنین کاربردهایی است[19] و اندازه گیریهایی باید انجام شود تا مانع از صدمه زدن به همکاران شده و ایمنی آنها تضمین شود. علاوه بر این، چالشهای دیگر به مکانیک ربات مربوط می شوند که توانایی آن برای عکس العمل نشان دادن به نیروی انسانها و وفق دادن خود به اعمالشان از این موارد هستند.

۲-۱-۱-۲ بازنمایی ویژگی

قبل از یادگیری یک سیاست، مهم است که حالت قابل مشاهده به فرمی بازنمایی شود که برای آموزش، مناسب و کارا باشد. این بازنمایی، بردار ویژگی نامیده میشود. یک ویژگی ممکن است شامل اطلاعاتی درباره یادگیرنده،

-

¹ on-board

محیطش، اشیای قابل دستکاری و عاملهای دیگر در آزمایش باشد. ویژگیهای آموزش باید مناسب باشند، به این معنی که اطلاعات کافی را منتقل کنند تا یک سیاست یکپارچه برای حل مسئله مورد نظر را تشکیل دهند. همچنین مهم است که ویژگیها بتوانند از لحاظ زمانی و محدودیت محاسباتی کاربردهای یادگیری تقلیدی به طور کارا استخراج شده و پردازش شوند.

هنگام ضبط داده، سوال مهم این است: چه چیزی را باید تقلید کرد؟ در بیشتر کاربردهای واقعی، غالبا محیط برای بازنمایی کلیت آن خیلی پیچیده است، چون معمولا مقدار زیادی اطلاعات نامرتبط یا زاید دارد. بنابراین، ضروری است در نظر بگیریم که چه جنبههایی از نمایشها را میخواهیم در اختیار یادگیرنده قرار دهیم.

٣-١-١-٢ يادگيري مستقيم

حال، روشهای مختلف یادگیری یک سیاست از نمایشها را مورد بررسی قرار میدهیم. بعد از در نظر گرفتن این که چه چیزی را باید گرفت، این فرآیند به این سوال که چطور باید یاد گرفت مربوط می شود.

سرراست ترین راه برای یادگیری یک سیاست از نمایشها تقلید مستقیم است، یعنی این که یک مدل باناظر از نمایش یاد می گیرد که در آن عملی که توسط خبره در اختیار قرار گرفته، به صورت برچسب برای نمونه داده شده عمل می کند. از همین رو، مدل قابلیت این را دارد که عمل مناسب را به هنگام مواجهه با یک وضعیت پیشبینی کند. روشهای یادگیری باناظر به دو دسته دستهبندی و رگرسیون تفکیک می شوند.

• دستهبندی

دستهبندی، یک کار محبوب در یادگیری ماشین است که در آن به طور خودکار مشاهدات به مجموعه متناهی x از کلاسها دستهبندی میشوند. یک دستهبند h(x) برای پیشبینی کلاس y ای که یک مشاهده مستقل x = y به آن مربوط است استفاده میشود که $y \in Y$ و $y \in Y$ و $y \in Y$ مجموعه متناهی از کلاسها و $y \in Y$ برداری از $y \in Y$ است. در دستهبندی باناظر، $y \in Y$ با استفاده از یک مجموعه داده با $y \in Y$ برداری از $y \in Y$ و روزش داده میشود که $y \in Y$ برداری از $y \in Y$ هستند. $y \in Y$ میشود که $y \in Y$ هستند.

رویکردهای دستهبندی، زمانی استفاده میشوند که اعمال یادگیرنده میتوانند به کلاسهای گسسته دستهبندی شوند [16]. این رویکرد برای کاربردهایی از قبیل ناوبری [20] و شبیه سازهای پرواز [7] مناسب است که در آن

به عمل می توان به عنوان یک تصمیم نگاه کرد. در مرجع [20]، یک مدل ترکیبی گوسی (GMM) برای پیشبینی تصمیمات ناوبری آموزش داده شده است. در مرجع [21]، متادستهبند ها برای یادگیری انجام بازیهای کامپیوتری استفاده شدهاند. دستهبند پایهای که در این مرجع استفاده شده یک شبکه عصبی است. در بازی اتومبیل,انی کارت'، دستورات دسته بازی ٔ آنالوگ به ۱۵ دسته گسستهسازی شده و مسئله به یک مسئله دستهبندی ۱۵ کلاسه کاهش می یابد. بنابراین شبکه عصبیای که استفاده شده ۱۵ نود خروجی دارد. بازی ماریو بروز ٔ از یک کنترل کننده گسسته استفاده می کند. اعمال با فشردن یک یا تعداد بیش تری از چهار دکمه انتخاب میشوند. بنابراین در شبکه عصبی مربوطه، عمل مربوط به یک فریم با چهار نود خروجی بازنمایی می شود. این امر، دستهبند را قادر می سازد تا برای یک نمونه چندین کلاس را انتخاب کند. اگرچه نتایج امیدوارکننده هستند، استدلال می شود که استفاده از تکنیک کنترل بهینه معکوس[22] به عنوان دستهبند یایه ممکن است مفید باشد. در مرجع [23]، این آزمایشها با استفاده از رگرسیون تکرار شده تا ورودی آنالوگ در سوپر تاکس کارت $^{\circ}$ یاد گرفته شود. برای بازی ماریو بروز، چهار دستهبند SVM جایگزین شبکه عصبی شده تا مقدار هر یک از چهار کلاس دودویی پیشبینی شود. همچنین، دستهبندی می تواند برای اتخاذ تصمیماتی که اعمال سطح پایین تری را شامل می شوند استفاده شود. در مرجع [5]، در یک شبیه سازی چندعامله فوتبال، تصمیمات سطح بالا به وسیله یک دستهبند پیشبینی میشوند. تصمیماتی از قبیل "نزدیک شدن به توپ" و "دریبل کردن به سمت گل" می توانند به طور قطعی و با استفاده از اعمال سطح پایین تر اجرا شوند. یک مطالعه تجربی انجام شده تا ارزیابی کند که کدام دستهبندها برای کار یادگیری تقلیدی مناسبتر هستند. چهار دستهبند مختلف بر حسب دقت و زمان یادگیری مقایسه میشوند. نتایج نشان میدهند که تعدادی از دستهبندها می توانند پیش بینی هایی با دقت قابل مقایسه انجام دهند. اما زمان یادگیری نسبت به تعداد نمایشها می تواند تا حد زیادی تغییر کند[5]. در مرجع [24]، شبکههای عصبی بازگشتی (RNN) به منظور یادگیری تراژکتوریهایی از نمایشها برای دستکاری شی استفاده شدهاند. RNN ها هنگام در نظر

¹ Meta-classifier

² kart

³ Joystick

⁴ Mario Bros

⁵ Super Tux Kart

گرفتن عمل بعدی، حافظه اعمال قبلی را با هم ادغام میکنند. ذخیرهسازی حافظه، شبکه را قادر میکند تا رفتار اصلاحی مانند بازیابی از شکست را یاد بگیرد به این شرط که مربی چنین سناریویی را نمایش داده باشد.

• رگرسیون

روشهای رگرسیون، برای یادگیری اعمال در یک فضای پیوسته استفاده می شوند. برخلاف دستهبندی، روشهای رگرسیون ورودی حالت را به یک خروجی عددی که یک عمل را نشان می دهد نگاشت می کنند. بنابراین، این روشها برای اعمال حرکتی سطح پایین مناسب هستند تا تصمیمات سطح بالاتر، به خصوص زمانی که اعمال به صورت مقادیر پیوسته بازنمایی می شوند مانند ورودی از یک دسته بازی [23]. رگرسور $y \in \mathbb{R}$ یک نمونه مستقل y را به جای مجموعهای از کلاسها به یک مقدار پیوسته y نگاشت می کند که $y \in \mathbb{R}$ مجموعه اعداد حقیقی است. به طور مشابه، رگرسور با استفاده از مجموعهای از نمونههای برچسب خورده $y \in \mathbb{R}$ است. $y \in \mathbb{R}$ آموزش داده می شود که $y \in \mathbb{R}$ و $y \in \mathbb{R}$ است.

تکنیکی که به طور رایج استفاده می شود رگرسیون وزن دهی شده محلی الست. LWR برای یادگیری تراژکتوریها مناسب است چون این حرکات از دنبالهای از مقادیر پیوسته تشکیل شده اند. مثالهایی از این تبیل حرکات کارهای ضربه زدن [10] و گام برداشتن [25] هستند که در مورد اول نیاز است که عامل یک تراژکتوری حرکت را اجرا کرده تا از یک نقطه عبور کرده و یک هدف را بزند و در مورد دوم عامل باید یک تراژکتوری تولید کند که منجر به حرکت پایدار و آرام می شود. یک کاربرد فراگیرتر، تنیس روی میز است. در مرجع [26]، از رگرسیون بیزی خطی استفاده شده تا به یک بازوی ربات آموزش داده شود که بازی پیوسته تنیس روی میز را انجام دهد. نیاز است که این عامل، با دقت در یک فضای سه بعدی پیوسته و در وضعیتهای مختلف از قبیل هنگام زدن توپ، بازیابی موقعیت بعد از زدن توپ و آماده شدن برای حرکت بعدی رقیب حرکت کند. نمونه دیگری که به طور رایج برای رگرسیون استفاده می شود شبکههای عصبی مصنوعی (ANNs) است. شبکههای عصبی از این لحاظ با دیگر تکنیکهای رگرسیون متفاوت هستند که به زمان آموزش زیاد و تعداد نمونههای آموزش زیاد نیاز دارند. رویکردهای شبکه عصبی غالبا از مطالعات زیستشناسی و عصبشناسی الهام گرفته و قصد تقلید فرآیند تقلید و یادگیری انسانها و حیوانات را دارند[4]. استفاده از رگرسیون همراه با یک

¹ Locally weighted regression

² Batting

سیستم دینامیکی از مقدمات حرکتی، کاربردهایی برای یادگیری حرکات گسسته و ریتمیک تولید کرده است[27]، اما بیش تر رویکردها برروی تقلید مستقیم بدون بهینهسازی تمرکز بیش تری کردهاند[28]. در چنین کاربردهایی یک سیستم پویا درجه آزادی را بازنمایی میکند، به طوری که هر درجه آزادی هدف و محدودیتهای متفاوت دارد[29].

سیستمهای پویا می توانند با روشهای یادگیری ماشین احتمالاتی ترکیب شده تا فواید هر دو رویکرد به دست آید. این امر استخراج الگوهایی که برای کار مورد نظر مهم هستند و تعمیم به سناریوهای مختلف را امکان پذیر می کند در حالی که توانایی تطبیق و تصحیح تراژ کتوریهای حرکت در زمان واقعی را حفظ می کند[30]. در مرجع [30]، تخمین پارامترهای سیستمهای دینامیکی به صورت یک مسئله رگرسیون ترکیب گوسی (GMR) بازنمایی شده است. این رویکرد یک مزیت بر رویکردهای مبتنی بر LWR دارد چنانچه امکان یادگیری توابع فعالسازی به همراه اعمال حرکتی را فراهم می کند. روش پیشنهاد شده برای یادگیری حرکت مبتنی بر زمان و حرکت نامتغیر با زمان استفاده شده است. در مرجع [31]، از یک روش مبتنی بر GMM مشابه در یک چارچوبکاری وظیفه-پارامتربندیشده ٔ استفاده شده است که اجازه شکل دهی حرکت ربات به صورت تابعی از پارامترهای کار را می دهد. نمایشهای انسان کدگذاری شده تا پارامترهایی که به کار دم دست 7 مرتبط هستند را بازتاب کرده و موقعیت، سرعت و محدودیتهای نیروی آن کار را مشخص کنند. این کدگذاری به چارچوبکاری امکان داده تا حالتی که ربات باید در آن باشد را استنتاج کرده و مطابق آن حرکت ربات را بهینه کند. این رویکرد در زمینه همکاری ربات و انسان (HRC) استفاده شده و قصد بهینهسازی مداخله انسان و همچنین تلاش ربات را دارد. در مرجع [32]، یادگیری عمیق با سیستمهای دینامیکی ترکیب شده است. مقدمات حرکت پویا (DMPs) در اتوانکدرهایی که بازنماییهای حرکت از داده نمایش را یاد گرفته جاسازی شدهاند. اتوانکدرها در لایه مخفی ویژگیها را به صورت غیرخطی به فضای نهفته ٔ با بعد کمتر نگاشت می کنند. در این رویکرد، واحدهای مخفی به DMP ها محدود شده تا لایه مخفی به بازنمایی پویای سیستم محدود شود.

¹ Dynamic system

² task-parameterized

³ At hand

⁴ Latent

در هر دو روش دستهبندی و رگرسیون با توجه به منابع مدلهای یادگیری، یک تصمیم طراحی می تواند اتخاذ شود. یادگیرندههای تنبل از قبیل kNN و kNN نیازی به آموزش ندارند اما نیاز دارند که به هنگام پیشبینی کردن تمام نمونههای آموزشی حفظ شوند. از طرف دیگر، مدلهای آموزش داده شده از قبیل ANN و SVM به زمان آموزش نیاز دارند اما به محض ایجاد یک مدل، دیگر به نمونههای آموزشی احتیاجی نیست و فقط مدل ذخیره می شود که این امر مانع از هدر رفتن حافظه می شود. هم چنین این مدلها می توانند منجر به زمانهای پیشبینی خیلی کوتاه شوند.

• مدلهای سلسله مراتبی

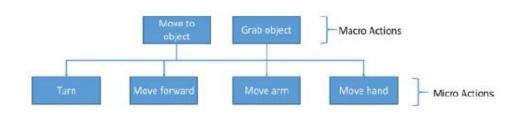
به جای استفاده از یک مدل برای بازتولید رفتار انسان، از یک مدل سلسلهمراتبی می تواند استفاده شود که اعمال آموخته شده را به تعدادی از فازها تجزیه می کند. یک مدل دستهبندی می تواند استفاده شود تا تصمیم گرفته شود که کدام عمل یا زیرعمل از مجموعهای از اعمال ممکن باید در زمان مورد نظر انجام شود. سپس از یک مدل متفاوت برای تعیین جزئیات عمل انتخاب شده استفاده می شود که در آن، هر عمل سطح پایین ممکن یک مدل مشخص و معین دارد. شکل ۲-۲ مثالی از اعمال سلسلهمراتبی را نشان می دهد. در مرجع [33]، از یک مدل مشخص و معین دارد. شکل ۲-۲ مثالی از اعمال سلسلهمراتبی را نشان می دهد. در مرجع [43]، از یک رویکرد سلسلهمراتبی برای یادگیری تقلیدی در دو مسئله مختلف استفاده شده است. مسئله اول، ایرهاکی است که در مقابل یک رقیب بازی می شود و هدف، شوت کردن یک گوی درون دروازه رقیب و در عین حال حفاظت از دروازه خودی است. بازی دوم، ماربل میز ۲ است. هزارتو ۳ می تواند حول محورهای مختلف کج شده تا توپ به سمت انتهای هزارتو حرکت کند. هر کدام از دو کار ذکر شده مجموعهای از اعمال سطح پایین به نام مقدمات حرکتی دارند که امکانهای بازی کردن برای عامل را تشکیل می دهند (مثل شوت مستقیم، دفاع از عملی که باید انجام شود استفاده شده است. با مشاهده حالت بازی، این دستهبند در نمایش های تهیهشده توسط انسان خبره دنبال مشابه ترین نمونهها می گردد و عمل مقدماتی که در آن نقطه توسط انسان انتخاب شده است را بازیابی می کند. مرحله بعدی تعیین هدف عمل انتخاب شده است برای مثال سرعت توپ یا موقعیت گوی هنگامی که عمل مقدماتی کامل می شود. سپس، هدف در یک مدل رگرسیون استفاده شده تا موقعیت گوی هنگامی که عمل مقدماتی کامل می شود. سپس، هدف در یک مدل رگرسیون استفاده شده تا

¹ air hockey

² marble maze

³ Maze

پارامترهای عملی که هدف مورد نظر را بهینه می کند پیدا شود. هدف، از k نزدیک ترین نمایش همسایه که در مرحله قبل پیدا شدند استنتاج می شود. اهداف در آن نمایشها ورودی یک مدل رگرسیون وزن دار برای انجام عمل مقدماتی هستند. با سبکی مشابه در مرجع [34]، از یک دسته بند برای تصمیم گیری در یک کار مرتبسازی استفاده شده که شامل اعمال بزرگ (wait, sort left, sort right and pass) است. هر عمل بزرگ، اعمال حرکتی گذرا آ از قبیل برداشتن یک توپ و حرکت دادن و قرار دادن آن توپ را شامل می شود.



شكل ٢-٢ مثالى از يادگيرى سلسلهمراتبى عملها. [15]

۴-۱-۱-۲- یادگیری غیرمستقیم

در این بخش، در مورد روشهای غیرمستقیم یادگیری سیاستها بحث می شود که می تواند یا یادگیری مستقیم را کامل کرده یا جایگزین آن شود. سیاست می تواند با استفاده از نمایشها، تجربه یا مشاهده بهتر شود تا یا دقیق تر باشد یا در شرایط دیده نشده عمومی تر و قوی تر باشد. غالبا یادگیری مستقیم به تنهایی برای باز تولید رفتار قوی و شبیه انسان در عاملهای هوشمند کافی نیست. این محدودیت می تواند به دو عامل اصلی نسبت داده شود: ۱) خطاهای در نمایش و ۲) تعمیم ضعیف. به خاطر محدودیتهای تکنیکهای به دست آوردن داده که عبار تند از مسئله تناسب، خطای حسگر و تاثیرات فیزیکی در نمایشهای حرکتی[16]، تقلید مستقیم می تواند منجر به کارایی ناپایدار یا نادقیق شود (این مشکل، به ویژه در کارهایی از قبیل زدن گوی یا رسیدن به جایی یا چیزی نمود پیدا می کند که نیاز به حرکت دقیق در فضای پیوسته دارند). برای مثال، در مرجع به جایی یا چیزی نمود پیدا می کند که نیاز به حرکت دقیق در فضای پیوسته دارند). برای مثال، در مرجع به جایی یا چیزی نمود پیدا می کند که نیاز به حرکت دقیق در فضای پیوسته دارند). برای مثال، در مرجع به جایی یا چیزی نمود راه رفتن با تقلید مستقیم نمایشها را دارد می افتد چون نمایشها به طور دقیق

¹ Macro

² Temporal

ویژگیهای فیزیکی از قبیل وزن و مرکز جرم ربات که در این کار درگیر هستند را در نظر نمی گیرند. اما بهبود سیاست از طریق آزمون و خطا این عوامل را در نظر گرفته و یک حرکت پایدار تولید می کند.

در حالی که تعمیمپذیری مقولهای مهم در کارهای یادگیری ماشین است یک مورد خاص تعمیمپذیری در کاربردهای یادگیری تقلیدی مورد تاکید قرار گرفته است. رایج است که نمایشهای انسان به صورت دنبالهای از اعمال در اختیار قرار میگیرند. وابستگی هر عمل به بخش قبلی دنباله، فرض i.i.d بودن نمونههای آموزشی که در یادگیری باناظر برای تعمیمپذیری مهم است را نقض میکند[21]. علاوه بر این، چون انسانهای خبره فقط نمونههای درست را در اختیار قرار میدهند یادگیرنده آماده مدیریت خطاها در تراژکتوری نیست. اگر یادگیرنده در هر نقطهای، از عملکرد بهینه در تراژکتوری فاصله بگیرد (چیزی که در هر کار یادگیری ماشین انتظار میرود)، در معرض یک موقعیت دیدهنشده قرار میگیرد که مدل برای وفق یافتن با آن آموزش داده نشده است. یک مثال واضح در مرجع [35] آورده شده است که در آن از یادگیری باناظر استفاده شده تا یک سیاست برای راندن یک ماشین یاد گرفته شود. با توجه به این که نمایشهای انسان فقط شامل رانندگی خوب بدون تصادف یا حتی نزدیک خطر است، زمانی که خطایی رخ میدهد و ماشین از تراژکتوریهای نمایش داده بدون تصادف یا حتی نزدیک خطر است، زمانی که خطایی رخ میدهد و ماشین از تراژکتوریهای نمایش داده شده منحرف میشود یادگیرنده نمیداند که چگونه جبران کند.

• یادگیری تقویتی

یادگیری تقویتی یک سیاست را یاد گرفته تا از طریق آزمون و خطا یک مسئله را حل کند.

در یادگیری تقویتی، یک عامل به صورت یک MDP مدل شده که حرکت در یک فضای حالت را یاد می گیرد. یک الله MDP متناهی شامل چهارتایی (S,A,T,R) است که S مجموعه متناهی حالتها، A مجموعه اعمال محموعه اعمال TP_{sa} محموعه احتمالات انتقال حالت و R یک تابع پاداش است. TP_{sa} شامل مجموعهای از احتمالات است که TP_{sa} شامل مجموعهای از احتمالات است که TP_{sa} احتمال رسیدن در حالت TP_{sa} به شرط عمل TP_{sa} است. تابع یاداش رسیدن در حالت TP_{sa} به شرط عمل TP_{sa} احتمال رسیدن در حالت TP_{sa} به شرط عمل TP_{sa} به شرط عمل در یک حالت معلوم و قرار گرفتن در یک پاداش آنی برای انجام یک عمل در یک حالت معلوم و قرار گرفتن در یک حالت جدید برمی گرداند TP_{sa} که گام زمانی است. این پاداش با گذشت زمان به اندازه ضریب کاهش می یابد و هدف عامل، بیشینه کردن امید ریاضی پاداش کاهش یافته در هر گام زمانی است.

یادگیری تقویتی با یک سیاست تصادفی شروع شده و پارامترهایش را براساس پاداشهایی که از اجرای این سیاست به دست میآورد تغییر میدهد. یادگیری تقویتی میتواند به تنهایی برای یادگیری یک سیاست در انواع کاربردهای رباتیک استفاده شود. اما اگر یک سیاست از نمایش یاد گرفته شود، یادگیری تقویتی می تواند برای تنظیم پارامترها استفاده شود. تهیه کردن نمونههای مثبت و منفی برای اموزش یک سیاست، با کاهش فضای جستجوی دردسترس به یادگیری تقویتی کمک می کند[4]. بهبود سیاست با استفاده از یادگیری تقویتی گاهی اوقات ضروری است برای نمونه در صورتی که بین مربی و یادگیرنده مغایرتهای فیزیکی وجود داشته باشد یا برای کاهش خطاهای موجود در به دست آوردن نمایشها. یادگیری تقویتی همچنین میتواند برای آموزش سیاست برای وضعیتهای دیده نشده که در نمایشها پوشش داده نشدهاند مفید باشد. اعمال یادگیری تقویتی به سیاست آموخته شده به جای یک سیاست تصادفی میتواند به طور قابل توجهی سرعت فرآیند یادگیری تقویتی را زیاد کرده و از خطر همگرایی سیاست به کمینه محلی جلوگیری کند. علاوه بر این، یادگیری تقویتی می تواند سیاستی را برای اجرای یک کار پیدا کند که به نظر (انسان) شاهد طبیعی نیست. در کاربردهایی که یادگیرنده با یک انسان تعامل می کند، برای کاربر مهم است که اعمال عامل را به طور شهودی تشخیص دهد. این امر در مواردی که رباتها در محیطهای تثبیت شده (از قبیل خانهها و ادارات) قرار داده میشوند تا با کاربران (انسان) آموزش دیده نشده تعامل کنند رایج است[36]. با اعمال یادگیری تقویتی به سیاستی که از نمایشهای انسان آموخته شده است میتوان از مشکل رفتار ناشناخته پرهیز کرد. در روشهای یادگیری تقلیدی، زمانی که شایستگی کار اجرا شده قابل ارزیابی است غالبا برای بهبود سیاست آموخته شده یادگیری تقویتی با یادگیری از نمایشها ترکیب میشود.

در پژوهشهای اولیه، آموزش با استفاده از نمایشهای اعمال موفق برای بهبود و افزایش سرعت یادگیری تقویتی استفاده شده است تا یک سیاست برای انجام یک بازی در محیط پویای دو بعدی آموزش داده شود. روشهای مختلف برای بهبود سیاست یادگیری تقویتی مورد بررسی قرار گرفتهاند. نتایج نشان میدهند که آموزش یادگیرنده با نمایشها امتیازش را بهتر کرده و به جلوگیری از افتادن یادگیرنده در کمینه محلی کمک میکند. همچنین مشاهده شده است که بهبود از آموزش با دشواری کار افزایش می یابد. راهحلهای مسائل ساده به سادگی و بدون نیاز به نمایشهایی از خبره می توانند

استنتاج شوند. اما چنانچه پیچیدگی کار افزایش مییابد فایده یادگیری از نمایشها مهمتر میشود و حتی برای یادگیری موفق در کارهای دشوارتر ضروری میشود[38].

در مرجع [39]، GMR برای آموزش یک ربات در کار در دست گرفتن شی استفاده شده است. چون سناریوهای دیده نشده از قبیل موانع و موقعیت متغیر شی در این کاربرد مورد انتظار هستند یادگیری تقویتی استفاده شده تا راههای جدید برای انجام این کار مورد بررسی قرار گیرد. سیستم آموزش داده شده یک سیستم پویاست که روی تراژکتوریهای تقلید شده میرایی از اجرا می کند. این امر به ربات اجازه داده تا به آرامی به هدفش رسیده و از اینکه یادگیری تقویتی منجر به نوسانات شود جلوگیری می کند. استفاده از میرایی در سیستمهای پویا به هنگام ترکیب یادگیری تقلیدی و یادگیری تقویتی یک رویکرد رایج است[13].

یک کاربرد شگفتانگیز یادگیری تقلیدی و تقویتی، آموزش یک عامل برای انجام بازی تختهای "GO" است که با انسانهای خبره رقابت می کند[40]. یک شبکه عصبی کانولوشنی عمیق با استفاده از بازیهای قبلی آموزش داده شده است. سپس یادگیری تقویتی برای اصلاح وزنهای شبکه و بهبود سیاست استفاده شده است. یک رویکرد متفاوت برای ترکیب یادگیری از نمایشها با یادگیری تقویتی در مرجع [41] استفاده شده است. به جای استفاده از نمایشها برای آموزش یک سیاست اولیه، از آنها برای به دست آوردن دانش قبلی برای شکل دهی پاداش استفاده شده است[42]. از یک تابع پاداش برای تشویق دستیابی به اهداف کوچکتر در کار (از قبیل نقاط عطفی که در نمایشها به آنها دست یافته شده است) استفاده شده است. این تابع پاداش با تابع پاداش اولیه ترکیب شده تا هزینه اعمال عامل را در اختیارش قرار دهد. این الگوی استفاده از نمایشهای خبره

روشهای جستجوی سیاست زیرمجموعهای از یادگیری تقویتی هستند که به طور طبیعی مناسب کاربردهای روشهای جستجوی سیاست رباتیکی هستند چون به MDP های با ابعاد بالا مقیاس می شوند [14]. بنابراین، روشهای جستجوی سیاست برای ادغام شدن با روشهای یادگیری تقلیدی مناسب هستند. یک روش گرادیان سیاست در مرجع [44] استفاده شده تا یک سیاست موجود که می تواند از طریق یادگیری با ناظر یا برنامهنویسی صریح تولید شود را

برای به دست آوردن یک تابع پاداش مشابه رویکردهای یادگیری تقویتی معکوس است[43].

¹ Damping

² Milestone

³ Explicit programming

بهبود دهد. یک رویکرد مشابه در مرجع [45] در یک سیستم پویا استفاده شده است که قبلا برای یادگیری باناظر از نمایشها استفاده شده بود [46]. این امر منجر به یک سری از کارها شده است که از سیستم پویای معرفی شده در مرجع [46] بهره میبرند تا از نمایشها یاد گرفته و متعاقبا از یادگیری تقویتی برای خود بهسازی استفاده میکنند[13]. از این چارچوب کاری برای آموزش تعدادی از کاربردها به بازوهای رباتیکی از قبیل قرار دادن توپ در فنجان، زدن توپ با پارو [10] و انجام بازی تنیس[26] استفاده شده است. به جای استفاده از یادگیری تقویتی برای اصلاح سیاستی که با استفاده از نمایشها آموزش داده شده است، نمایشها می توانند برای راهنمایی جستجوی سیاست استفاده شوند. در مرجع [47]، از برنامهنویسی پویای تفاضلی استفاده شده تا از نمایشهای انسان نمونههای هدایت کننده به جستجوی سیاست کمک کرده تا مناطق با پاداش بالا از فضای جستجو را مورد بررسی قرار دهند.

در مرجع [48]، شبکههای عصبی بازگشتی در جستجوی هدایتشده سیاست گنجانده شده تا پرداختن به مسائل مشاهدهپذیر جزئی را آسان کند. حافظههای پیشین به فضای حالت افزوده شده و هنگام پیشبینی عمل بعدی در نظر گرفته میشوند. یک رویکرد باناظر از تراژکتوریهای نمایشداده شده استفاده کرده تا در مورد اینکه کدام حافظهها باید ذخیره شوند تصمیم گیری کرده، در حالی که از یادگیری تقویتی استفاده شده تا سیاستی را که شامل مقادیر حالت حافظه میشود بهینه کند.

یک روش متفاوت برای استفاده از یادگیری تقویتی در یادگیری تقلیدی، استفاده از یادگیری تقویتی برای تهیه نمایشهای تقلید مستقیم است. این رویکرد به (انسان) مربی احتیاج ندارد چون سیاست از آغاز و با استفاده از آزمون و خطا آموخته شده و سپس برای تولید نمایشهای آموزش استفاده میشود. یک دلیل برای تولید نمایشها و آموزش یک مدل باناظر به جای استفاده مستقیم از سیاست یادگیری تقویتی این است که روش یادگیری تقویتی در زمان واقعی (به طور بلادرنگ) کار نمی کند[49]. وضعیت دیگر هنگامی است که سیاست یادگیری تقویتی در یک محیط کنترل شده آموخته میشود. در مرجع [50]، یادگیری تقویتی برای یادگیری انواع کارهای رباتیکی در یک محیط کنترل شده استفاده شده است. اطلاعاتی از قبیل موقعیت اشیای هدف در طول این فاز در دسترس است. سپس یک شبکه عصبی کانولوشنی عمیق، با استفاده از نمایشهای برگرفته از

¹ ball paddling

² Differential

سیاست یادگیری تقویتی آموزش داده می شود. شبکه عصبی، نگاشت ورودی بصری به اعمال را یاد گرفته و در نتیجه، اجرای کارها را بدون اطلاعاتی که در فاز یادگیری تقویتی به آنها نیاز است یاد می گیرد. این روش، از نمایشهای انسان تقلید می کند همان طور که انسانها برای تهیه نمایشها از دانش خبرهای که در فرآیند آموزش در گیر نشده است بهره می برند.

بهبنهسازی

رویکردهای بهینهسازی نیز می توانند برای پیدا کردن راه حلی برای یک مسئله معلوم استفاده شوند.

با فرض یک تابع هزینه $A \to \mathbb{R}$ که منعکس کننده عملکرد یک عامل است، به طوری که A مجموعهای از پارامترهای ورودی و \mathbb{R} مجموعه اعداد حقیقی است، روشهای بهینه سازی سعی در پیدا کردن پارامترهای ورودی و $f(x_0) \leq f(x)$ است.

مشابه یادگیری تقویتی، تکنیکهای بهینهسازی می توانند با شروع از یک راه حل تصادفی و بهبود پیاپی آن به منظور بهینهسازی تابع برازش، برای پیدا کردن راه حلهای مسائل استفاده شوند. الگوریتمهای تکاملی (EAs) روشهای محبوب بهینهسازی هستند که به طور گسترده برای یافتن تراژکتوریهای حرکتی در کارهای رباتیکی استفاده شدهاند[51]. الگوریتمهای تکاملی برای تولید تراژکتوریهای حرکتی برای رباتهای با درجه آزادی بالا و پایین استفاده شدهاند[52]. روشهای محبوب هوش گروهی از قبیل بهینهسازی ازدحام ذرات (PSO) بالا و پایین استفاده شدهاند [52]. روشهای محبوب هوش گروهی از قبیل بهینهسازی باری ناوبری وسیله بدون سرنشین استفاده شدهاند. این تکنیکها رفتار موجودات زنده را شبیهسازی کرده تا یک راه حل سراسری بهینه را در فضای جستجو بیابند. مشابه یادگیری تقویتی، الگوریتمهای تکاملی می توانند با یادگیری تقلیدی آمیخته شده تا تراژکتوریهایی که از نمایشها آموخته شده را بهبود داده یا سرعت فرآیند بهینهسازی را افزایش دهند. در مرجع [11]، یک الگوریتم ژنتیک (GA) برای بهینهسازی تراژکتوریهای حرکتی نمایش داده شده استفاده شده است. تراژکتوریها به عنوان یک جمعیت آغازی برای الگوریتم ژنتیک استفاده شدهاند. تراژکتوریهای می کنند ضبط شده به صورت کروموزومهایی کدگذاری شدهاند که از ژنهایی که مقدمات حرکتی را بازنمایی می کنند ضبط شده به صورت کروموزومهایی کدگذاری شده اند که از ژنهایی که مقدمات حرکتی را بازنمایی می کنند

¹ Swarm intelligence

² Particle Swarm Optimization

³ Ant Colony Optimization

تشکیل شدهاند. الگوریتم ژنتیک دنبال کروموزومی می گردد که یک تابع برازش که موفقیت کار را ارزیابی کرده بهینه می کند. تصویر کردن کردن تراژکتوریهای حرکتی به بعد پایین تر، تغییر قابل توجه بین حرکت بهینه شده و حرکتی که مستقیما از دست کاری جنبشی آموخته شده را نشان می دهد[11].

به طور مشابه در مرجع [55]، الگوریتمهای تکاملی بعد از آموزش عاملها در یک شبیه سازی فوتبال استفاده شده اند. یک راه حل ممکن (کروموزوم)، به صورت مجموعه ای از قوانین اگر – آنگاه بازنمایی می شود. به دلیل جایگشتهای متناهی مشاهدات و اعمال، قوانین متناهی هستند. برای ارزیابی مناسب بودن یک راه حل از یک تابع وزن دار از تعداد اهداف و مقیاسهای عملکرد دیگر استفاده می شود. اگرچه الگوریتم تکاملی اندازه جمعیت کمی داشته و تقاطع 7 را به کار نگرفته، اما نتایج امیدوار کننده ای روی قوانین آموخته شده از نمایشها نشان داده است.

در مرجع [35] نیز، از الگوریتمهای تکاملی برای بهینهسازی چندین هدف در یک بازی مسابقه سرعت استفاده شده است. این الگوریتمها یک راهحل بهینهشده (کنترلکننده) را از یک جمعیت اولیه از تراژکتوریهای رانندگی استنتاج کردهاند. ارزیابی کنترلکنندههای استنتاج شده مشخص کرده است که آنها به سبک رانندگی بازیکنانی که از روی آنها مدلسازی شدهاند وفادار میمانند. این امر در مورد مقیاسهای کمی از قبیل سرعت و پیشرفت و در مورد مشاهدات ذهنی^۳ از قبیل راندن در وسط جاده صدق می کند.

در مرجع [56]، وزنهای شبکه عصبی به صورت ژنومی که باید بهینهسازی شود در نظر گرفته شده است. جمعیت اولیه از طریق آموزش شبکه با نمونههای نمایش داده شده فراهم شده تا وزنها را مقداردهی اولیه کنند. از نمایشها نیز برای ایجاد یک مقدار برازش که مربوط به فاصله میانگین مربع خطا از خروجیهای مورد نظر (اعمال انسان) است استفاده شده است.

در مرجع [57]، از PSO استفاده شده تا مسیر بهینه برای یک وسیله هوایی بدون سرنشین (UAV) از طریق یافتن بهترین نقاط کنترلی روی یک منحنی B-spline پیدا شود. نقاط اولیهای که حکم ذرات اولیه PSO را دارند از اسکلتسازی 7 تهیه می شوند. یک گونه اجتماعی PSO در مرجع [58] معرفی شده است که از یادگیری

¹ Kinesthetic

² Crossover

³ Subjective observations

⁴ Skeletonization

حیوانات در طبیعت که از طریق مشاهده همتایانشان صورت می گیرد الهام گرفته است. هر ذره با یک راهحل تصادفی شروع کرده و یک تابع برازش برای ارزیابی هر راهحل استفاده می شود. سپس ذرات تقلیدگر (همه به جز ذرهای با بهترین برازش) رفتارشان را از طریق مشاهده ذرات نمایشدهنده (ذرات با عملکرد بهتر) تغییر می دهند. همانند طبیعت هر تقلیدگر می تواند از چندین نمایش دهنده یاد گرفته و یک نمایش دهنده می تواند براي آموزش بيش تر از يک تقليدگر استفاده شود. الگوريتمهاي تكاملي تعاملي (IEAs) [59] الگوي متفاوتي را به کار می گیرند. به جای استفاده از ورودی انسان برای شروع یک جمعیت اولیه از راهحلها و بهینهسازی آنها، الگوریتم تکاملی تعاملی از ورودی انسان برای برآورد برازش (مناسب بودن) راهحلها استفاده می کند. برای پرهیز از ارزیابی تعداد خیلی زیادی از راهحلهای بالقوه، یک مدل روی نمونههای باناظر آموزش داده شده تا ارزیابی کاربر (انسان) را تخمین بزند. در مرجع [60]، برای یادگیری ناوبری ربات جستجوی مبتنی بر برازش با یادگیری سیاست مبتنی بر ارجحیت (PPL) ترکیب شده است. ارزیابیهای کاربر از PPL جستجو را به مناطقی دور از کمینه محلی هدایت کرده در حالی که جستجوی مبتنی بر برازش دنبال یک راهحل می گردد. با یک ماهیت مشابه، در مرجع [61] یک ربات آموزش شده تا حرکت بازوی انسان را تقلید کند. با استفاده از نمایشها به عنوان جمعیت اولیه، از نمود اختلاف در درجه ازادی بین (انسان) نمایشدهنده و ربات جلوگیری می شود. اما به جای استفاده از ورودی انسان برای ارزیابی ذهنی یک راه حل، شباهت حرکت ربات به نمایشهای انسان به صورت کیفی ارزیابی میشود. از بازنمایی مفصل مستقل از توالی نمایشدهنده و یادگیرنده برای تشکیل یک تابع برازش استفاده شده است. از PSO برای یافتن زوایای مفصل استفاده شده تا این مقیاس شباهت بهینه شود. یک روش متفاوت برای ادغام نمایشها در مرجع [62] پیشنهاد شده است. با الهام از یادگیری تقویتی معکوس، از یک چارچوب کاری تنظیم کننده مربعی خطی معکوس (ILQR) استفاده شده تا تابع هزینهای که توسط (انسان) نمایشدهنده بهینه شده یاد گرفته شود. سپس به جای روشهای گرادیان از PSO به منظور یافتن راه حلی برای تابع آموخته شده استفاده شده است.

• یادگیری انتقالی^۱

یادگیری انتقالی یک الگوی یادگیری ماشین است که از دانش یک کار یا یک حیطه برای بهبود یا تقویت یادگیری کار دیگر استفاده می کند.

با فرض حیطه مبدا D_s و کار مبدا T_s ، یادگیری انتقالی به این صورت تعریف میشود که یادگیری یک کار $T_s \neq T_t$ یا $D_s \neq D_t$ میدف $T_s \neq T_t$ یا $D_s \neq D_t$ و D_s بهبود میدهد به طوری که D_t و یا استفاده از دانش $D_s \neq D_t$ بهبود میدهد به طوری که $D_s \neq D_t$ یا $D_s \neq D_t$ به صورت یک فضای ویژگی $D_s \neq D_t$ و توزیع احتمال حاشیه یا $D_s \neq D_t$ به صورت یک فضای ویژگی $D_s \neq D_t$ با نقرار است که $D_s \neq D_t$ باست. شرط $D_s \neq D_t$ است. شرط $D_s \neq D_t$ در صورتی برقرار است که $D_s \neq D_t$ یا $D_s \neq D_t$ است.

یک یادگیرنده می تواند اشکال متنوع دانش در مورد یک کار از قبیل بازنماییهای ویژگی یا پارامترهای مفید برای یادگیری مدل را از عامل دیگر فراگیرد. یادگیری انتقالی به یادگیری تقلیدی و کاربردهای رباتیکی مرتبط است چون به دست آوردن نمونهها دشوار و پرهزینه است. بهره بردن از دانش کاری که قبل از این در آن سرمایه گذاری کرده تا آن را یاد بگیریم می تواند کارآمد و نتیجه بخش باشد.

سیاستی که در یک کار یاد گرفته شده، می تواند برای نصیحت (آموزش) یک یادگیرنده در کار دیگری که شباهتهایی با کار اول دارد استفاده شود. در مرجع [64]، این رویکرد بر روی دو کار شبیه ساز فوتبال روبو کاپ پیاده سازی شده است. کار اول دور نگه داشتن توپ از تیم دیگر است و کار دوم گل زدن است. واضح است که مهارتهایی که برای اجرای کار اول یاد گرفته شده اند می توانند در کار دوم نیز استفاده شوند. در این مورد، نصیحت به صورت یک قانون که در مورد حالت و یک یا تعداد بیش تری عمل است فرموله شده است. برای ایجاد نصیحت، سیاست مربوط به کار اول با استفاده از یادگیری تقویتی یاد گرفته شده است. سپس سیاست یاد گرفته شده، توسط یک کاربر (برای پرهیز از اختلافات در فضاهای حالت و عمل) به فرم نصیحتی نگاشت شده که برای آغاز سیاست مربوط به کار دوم استفاده می شود. بعد از دریافت نصیحت، یادگیرنده به اصلاح

¹ Transfer learning

² Domain

سیاست از طریق یادگیری تقویتی ادامه داده و در صورتی که از طریق تجربه اثبات کند که نصیحت داده شده نادقیق یا نامربوط است می تواند آن را تغییر داده یا نادیده بگیرد.

غالبا در یادگیری انتقالی، ورودی انسان برای نگاشت دانش از یک حیطه به حیطه دیگر نیاز است. اما در برخی موارد، رویه نگاشت می تواند خودکار شود [65]. برای مثال در مرجع [66]، یک تابع نگاشت برای انجام بازی ارائه شده است. این تابع به صورت خودکار بین حیطههای مختلف نگاشت کرده تا از تجربه قبلی آموخته شود. عامل قادر بوده که بازیهایی که قبلا انجام شده و به کار فعلی مرتبط هستند را شناسایی کند. ممکن است عامل قبلا همان بازی یا مشابهش را انجام داده باشد و قادر است که یک کار مبدا مناسب را انتخاب کرده تا از آن بیاموزد بدون اینکه آن کار به طور واضح معین شده باشد. آزمایشات نشان می دهند که رویکرد یادگیری انتقالی سرعت یادگیری بازی از طریق یادگیری تقویتی را (در مقایسه با یادگیری بدون دانش قبلی) افزایش داده و بعد از تکمیل تکرارهای یادگیری به عملکرد بهتری دست می یابد. همچنین نتایج حاکی از آنند که فایده استفاده از یادگیری انتقالی وابسته به تعداد نمونههای آموزشی است که از کارهای مبدا منتقل می شوند. حتی اگر عامل با انتقال منفی مواجه شود [63] برای مثال بر اثر بیش برازش به کار مبدا، می تواند از طریق یادگیری از تجربه و تصحیح مدلش در کار فعلی جبران (بازیابی) کرده تا در یک زمان مناسب همگرا شود [66].

در مرجع [67]، شکل دهی پاداش و یادگیری انتقالی ترکیب شده تا انواع کارهای محک (معیار) ایاد گرفته شوند. از آنجایی که شکل دهی پاداش متکی به دانش قبلی است تا بر روی تابع پاداش تاثیر بگذارد، یادگیری انتقالی می تواند از یک سیاست آموخته شده برای یک کار سود برده و شکل دهی پاداش را برای یک کار مشابه انجام دهد. در مرجع [67]، یادگیری انتقالی از یک نسخه ساده مسئله به نسخه پیچیده تر آن اعمال شده است (برای مثال از ماشین کوهنوردی دوبعدی به سه بعدی و از بازی ماریو بدون دشمن به یک بازی با دشمن).

یادگیری کار آموزی^۲

در خیلی از کاربردهای هوش مصنوعی از قبیل بازیها یا کارهای پیچیده رباتیک، اندازه گیری (سنجیدن) موفقیت یک عمل دشوار است. در این گونه موارد، نمونههای نمایش داده شده می توانند به عنوان یک قالب برای

¹ Benchmark

² Apprenticeship

عملکرد مورد نظر استفاده شوند. در مرجع [43]، یادگیری کارآموزی (یا یادگیری تقویتی معکوس) پیشنهاد شده تا زمانی که تابع پاداش صریحی در دسترس نباشد یک سیاست آموخته شده را بهبود دهد.

در چنین کاربردهایی، هدف تقلید رفتار مربیها (انسان) به این شرط است که مربی در حال بهینهسازی یک تابع ناشناخته باشد.

یادگیری تقویتی معکوس (IRL) از نمونههای آموزشی استفاده کرده تا تابع پاداشی که توسط خبره بهینه می شود را یاد گرفته و از آن برای بهبود مدل آموزش داده شده استفاده کند.

از همین رو، یادگیری تقویتی معکوس عملکردی مشابه خبره به دست می آورد. بدون هیچ تابع پاداشی، عامل μ_E سلامی ویژگی μ_E می MDP/R (S,A,T) مدل می شود. در عوض، بعد از اینکه امید ریاضی های ویژگی MDP/R (S,A,T) به صورت یک $\{S_0^{(i)}, S_1^{(i)}, \cdots\}_{i=1}^m$ به ریاضی استنتاج شدند سیاست مدل می شود. با فرض μ_E به ریاضی ویژگی سیاست خبره μ_E به صورت معادله (۱-۲) مشخص می شود

$$\hat{\mu}_E = \frac{1}{m} \sum_{i=1}^m \sum_{t=0}^\infty \gamma^t \phi(s_t^{(i)})$$
 (۱-۲)معادله

در مرجع [48]، یک رویکرد بیشینه آنتروپی برای یادگیری تقویتی معکوس به کار گرفته شده تا ابهام را کاهش دهد. ابهام در کارهای یادگیری تقویتی معکوس به وجود می آید چون خیلی از توابع پاداش می توانند با سیاست یکسان بهینه شوند. این امر هنگام یادگیری تابع پاداش، یک مشکل را نمایان می کند مخصوصا زمانی که نمایشها ناقص و معیوب باشند. روش پیشنهادی بر روی کار یادگیری انتخابهای مسیر راننده نشان داده شده است که در آن ممکن است نمایشها زیربهینه و غیرقطعی باشند. این رویکرد در مرجع [68] به یک چارچوب

_

¹ Suboptimal

کاری یادگیری عمیق بسط داده شده است. توابع هدف بیشینه آنتروپی، یادگیری مستقیم وزنهای شبکه و در نتیجه استفاده از شبکههای عمیق که با گرادیان تصادفی نزولی آموزش داده شدهاند را امکانپذیر می کنند[68]. معماری عمیق بیش تر بسط داده شده تا به جای استفاده از ویژگیهای از پیش استخراج شده، ویژگیها از طریق لایههای کانولوشنی یاد گرفته شوند. این یک گام مهم در مسیر خودکار کردن فرآیند یادگیری است. یکی از چالشهای اصلی در یادگیری تقویتی از راه آزمون و خطا نیاز به دانش انسان در طراحی بازنماییهای ویژگی و با ویژگی و توابع هدف است[14]. با استفاده از یادگیری عمیق برای یادگیری خودکار بازنماییهای ویژگی و با استفاده از یادگیری تقویتی معکوس برای استنتاج توابع پاداش از نمایشها نیاز به ورودی و طراحی انسان کمینه میشود. الگوی یادگیری تقویتی معکوس مزیتی که بر اشکال دیگر یادگیری از نمایشها دارد این است که تابع هزینه کار از محیط جدا شده است. چون هدف نمایشها یاد گرفته شدهاند نه خود نمایشها، نیاز نیست که نمایش دهنده و یادگیرنده، اسکلت و محیط پیرامون یکسان داشته باشند، بنابراین چالشهایی از قبیل مسئله تناسب را کاهش میدهد. بنابراین، تهیه نمایشهایی که کلی بوده و برای یک ربات یا محیط خاص ساخته نشدهاند آسان تر است.

علاوه بر این، یادگیری تقویتی معکوس می تواند به جای یادگیری تقویتی سنتی به کار گرفته شود اگرچه تابع پاداشی وجود داشته باشد (به این شرط که نمایشها در درسترس باشند). برای مثال در مرجع [69]، از یادگیری کارآموزی استفاده شده تا یک تابع پاداش از نمایشهای خبره در بازی ماریو استنتاج شود. در حالی که اهداف در یک بازی مثل ماریو می توانند از قبل تعریف شوند (از قبیل امتیاز کشتن دشمنان و جمع آوری سکهها یا زمان کامل کردن سطح)، معلوم نیست که چگونه یک کاربر خبره این اهداف را اولویت بندی می کند. بنابراین در تلاش برای تقلید رفتار انسان، تابع پاداشی که از نمایشها استخراج شده به تابع پاداشی که به طور دستی طراحی شده ترجیح داده می شود.

• یادگیری فعال

یادگیری فعال الگویی است که در آن، مدل قادر است از یک خبره برای جواب بهینه به یک حالت یا وضعیت معلوم سوال کرده و از این نمونههای فعال برای بهبود سیاستش استفاده کند.

یک دستهبند $D_K(x^{(i)},y^{(i)})$ داده داده و برای پیشبینی برچسبهای $D_K(x^{(i)},y^{(i)})$ مجموعه داده بدون برچسب $D_C(x^{(i)}) \subset D_U$ استفاده شده است. یک زیرمجموعه و داده بدون برچسب $D_U(x^{(i)})$ استفاده شده است. یک زیرمجموعه و ناول کند. از نمونههای فعال یادگیرنده انتخاب شده تا در مورد برچسبهای صحیح آنها $y^{*(i)}$ از خبره سوال کند. از نمونههای در $D_C(x^{(i)},y^{*(i)})$ برای آموزش $D_C(x^{(i)},y^{*(i)})$ با هدف کمینه کردن $D_C(x^{(i)},y^{*(i)})$ تعداد نمونههای در $D_C(x^{(i)},y^{*(i)})$ استفاده شده است:

یادگیری فعال روشی مفید برای تطبیق مدل به وضعیتهایی است که در نمونههای آموزشی اصلی پوشش داده نشدهاند. از آنجایی که یادگیری تقلیدی دربرگیرنده تقلید تراژکتوری کامل یک حرکت است ممکن است خطا در هر مرحلهای از اجرا رخ دهد. ساخت مجموعههای آموزشی منفعل که میتوانند از این مشکل پرهیز کنند خیلی دشوار است.

یک رویکرد برای تصمیم گیری در مورد زمان سوال کردن از خبره، استفاده از تخمینهای اطمینان است تا بخشهایی از مدل آموخته شده که نیاز به بهبود دارند شناسایی شوند. هنگام اجرای اعمال آموخته شده اطمینان در این پیشبینی تخمین زده شده و یادگیرنده می تواند تصمیم بگیرد که درخواست نمایشهای جدید برای بهبود این قسمت از کاربرد را داده یا در صورتی که اطمینان کافی باشد از سیاست فعلی استفاده کند. با تناوب کردن بین اجرای سیاست و بهروز رسانی آن با نمونههای جدید، کم کم یادگیرنده اطمینان پیدا کرده و یک سیاست کلی به دست آورده که بعد از مدتی، دیگر نیازی به درخواست برای بهروز رسانیهای بیش تر ندارد. در مرجع [20]، از بهبود سیاست مبتنی بر اطمینان برای یادگیری ناوبری استفاده شده و در مرجع [34]، از همین روش برای یک کار مرتبسازی کلان ۲ استفاده شده است.

در مرجع [70]، یادگیری فعال معرفی شده تا عامل را قادر ساخته که در هر گامی در تراژکتوری، با در نظر گرفتن همه گامهای قبلی، از خبرهها سوال کند. این مسئله به یادگیری فعال i.i.d کاهش یافته (تبدیل شده) و استدلال شده که به طور قابل توجهی تعداد نمایشهای مورد نیاز را کاهش میدهد.

در مرجع [71]، یادگیری فعال در کارهای مشارکتی انسان- ربات مطرح شده است. انسان و ربات به صورت فیزیکی تعامل کرده تا در یک کار نامتقارن (به عبارت دیگر، انسان و ربات نقشهای متفاوتی دارند) به یک

¹ Passive

² Macro sorting

هدف مشترک برسند. یادگیری فعال بین دوره های تعامل رخ داده و انسان از طریق یک واسط کاربری گرافیکی (GUI) به ربات بازخورد می دهد. بازخورد ضبط شده و به پایگاه داده نمونههای آموزشی اضافه شده که برای آموزش مدل ترکیبی گوسی که اعمال ربات را کنترل می کند استفاده می شود. تعامل فیزیکی بین انسان و ربات منجر به رفتار وابسته متقابل می شود. بنابراین با هر تکرار تعامل، اعمال بههم پیوسته دو طرف به تراژکتوری حرکت هموارتری همگرا می شود. تجزیه و تحلیل کیفی آزمایشات نشان می دهد که در صورتی که انسان به اعمال ربات وفق یابد تعامل بین آنها می تواند بهتر شود و در صورتی که ربات هم به نوبه خود با هر نوبت تعامل به عمل انسان وفق یابد تعامل به طور قابل توجهی بهتر می شود.

در مرجع [72]، به جای اینکه یادگیرنده یک سوال بفرستد مربی اصلاحات را آغاز می کند. در حین اجرای کار توسط ربات، مربی رفتار یادگیرنده را مشاهده کرده و به طور جنبشی موقعیت مفاصل ربات را تصحیح می کند. یادگیرنده حرکتش که به کمک مربی انجام داده را از طریق سنسورهایش دنبال کرده و از این تراژکتوریها برای اصلاح مدل استفاده می کند، مدلی که به تدریج آموخته شده تا امکان اضافه شدن نمایشها در هر نقطهای را فراهم کند.

• پیشبینیهای ساختارمند

به سبکی مشابه، DAGGER [23] اجتماع نمونه را به کار گرفته تا برای وضعیتهای دیده نشده تعمیم داده شود. اما اساسا این رویکرد متفاوت است. DAGGER مسئله یادگیری تقلیدی را به صورت یک مسئله پیشبینی ساختارمند که از مرجع [73] الهام گرفته شده فرموله میکند. یک عمل به صورت دنبالهای از پیشبینیهای وابسته تلقی میشود. چون یک عمل وابسته به حالت قبلی است یک خطا منجر به یک حالت دیده نشده میشود که یادگیرنده نمی تواند از آن بازیابی کند که این امر منجر به خطاهای مرکب میشود. دیده نشده می دهند لازم و کافی DAGGER نشان می دهد که تجمیع نمونههایی که خطاهای یادگیری اولیه را پوشش می دهند لازم و کافی است. بنابراین، یک رویکرد تکراری پیشنهاد شده که از یک سیاست بهینه استفاده می کند تا هر مرحله از اعمالی را که با استفاده از سیاست فعلی پیشبینی شده تصحیح کند، بنابراین نمونههایی اصلاح شده ایجاد می شوند

¹ round

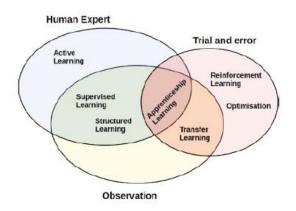
² Kinesthetically

که برای بهروز رسانی سیاست به کار میروند. همانطور که الگوریتم تکرار میشود استفاده از سیاست بهینه کم میشود تا زمانی که سیاست یاد گرفته شده به عنوان مدل نهایی استفاده میشود.

در مرجع [74]، یک الگوریتم به نام SIMILE پیشنهاد شده که محدودیتهای روشهای معرفی شده در مرجع [73] و [73] را از طریق تولید یک سیاست پایا که به تجمیع داده نیاز ندارد کم کرده است. SIMILE، نیاز به خبره برای تامین عمل در هر مرحله از تراژکتوری را از طریق فراهم کردن "بازخورد خبره مجازی" کم کرده که همواری تراژکتوری تصحیح شده را کنترل کرده و به اعمال خبره همگرا می شود.

در نظر گرفتن اعمال قبلی در فرآیند یادگیری یک نکته مهم در یادگیری تقلیدی است چون تعداد زیادی از کاربردها به اجرای تراژکتوریهایی از مقدمات حرکتی وابسته متکی هستند. یک روش کلی برای گنجاندن حافظه در یادگیری استفاده از RNN هاست[75]. RNN ها میان لایههای مخفی یک حلقه بازخورد ایجاد کرده تا خروجیهای قبلی شبکه را در نظر بگیرند و از همین رو برای کارهای با تراژکتوریهای ساختارمند مناسب هستند[76].

در نهایت، شکل ۲-۳ یک نمودار ون را نشان می دهد که منابع دادهای که توسط روشهای یادگیری مختلف به کار گرفته شدهاند را به طور خلاصه شرح داده است. یک عامل می تواند از نمایشهای اختصاصی مربی و با مشاهده اعمال عاملهای دیگر یا از طریق آزمون وخطا یاد بگیرد. یادگیری فعال به یک اوراکل اختصاصی نیاز دارد که بتواند از آن درخواست نمایش کند، در حالی که روشهای دیگر که از نمایشها بهره می برند می توانند آنها را از یک خبره اختصاصی یا با مشاهده رفتار مورد نیاز از عاملهای دیگر به دست بیاورند. روشهای یادگیری تقویتی و بهینهسازی با آزمون و خطا یاد گرفته و از نمایشها استفاده نمی کنند. یادگیری انتقالی از تجربه کارهای قدیمی یا از دانش عاملهای دیگر استفاده کرده تا یک سیاست جدید یاد بگیرد. یادگیری کارآموزی از نمایشهای یک خبره یا مشاهده استفاده کرده تا یک تابع پاداش یاد بگیرد. سپس یک سیاست که تابع پاداش را بهینه می کند می تواند از طریق تجربه آموخته شود.



شكل ٢-٣ روشهاي يادگيري از منابع مختلف.[15]

۲-۱-۲ روشهای یادگیری تقویتی

عمل کردن در شرایطی که عدم قطعیت وجود دارد نقش مهمی در مطالعه مسائل تصمیم گیری دنبالهای دارد و در رشتههای مختلفی از قبیل برنامه ریزی نظری تصمیم گیری، تحقیق در عملیات، یادگیری تقویتی و اقتصاد مورد بررسی قرار گرفته است. این هدف که با روشی معقول و به صورت بهینه عمل کرد، مقوله اصلی در بسیاری از رشتههای کمابیش مرتبط به هوش مصنوعی است. بخش زیادی از مسائلی که تصمیم گیری دنبالهای بهینه را دربرمی گیرند می توانند به صورت فرآیندهای تصمیم مارکف (MDP)[77] بیان شوند. در این مدل، یک محیط به صورت مجموعهای از حالتها و عملهایی که برای کنترل حالت سیستم انجام می شوند مدل می شود. هدف، کنترل سیستم به نحوی است که بهترین عملکرد حاصل شود. خیلی از مسائل از قبیل مسائل برنامه ریزی (تصادفی)، یادگیری کنترل ربات و مسائل انجام بازی در قالب یک MDP مدل شوند.

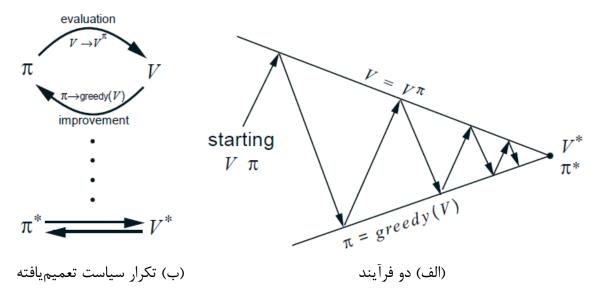
۱-۲-۱-۲- حل فرآیندهای تصادفی مارکف

حل یک MDP معلوم به معنای محاسبه سیاست بهینه π^* است که به هر حالتی در فضای حالت بهترین عمل را انتساب می کند. برای دستیابی به این هدف، تعداد زیادی روشهای برنامهنویسی پویا و یادگیری تقویتی وجود دارد. فرق اصلی بین این روشها این است که آیا دانش (کامل) در مورد تابع انتقال T و تابع

-

¹ Sequential decision making

پاداش R را فرض می کنند یا خیر. در دو بخش بعدی به طور مختصر برخی از الگوریتمهای پایه را شرح می دهیم. فرم کلی همه این الگوریتمها شامل دو فرآیند است (به شکل ۲-۴ رجوع شود)، یکی تخمین مقادیر برای حالتها و عملهاست و دیگری بهبود سیاستهاست.



شکل ۲-۴ فرم کلی الگوریتمهای حل فرآیندهای تصمیم مارکف. (الف) همگرایی تدریجی تابع مقدار و سیاست به نسخههای بهینه. (ب) الگوریتمهای معرفی شده در بخش 1-1-1-1- می توانند به صورت مقداردهی اولیه حلقه تکرار سیاست تعمیمیافته (\mathbf{GPI}) دیده شوند. گام ارزیابی سیاست، عملکرد سیاست \mathbf{V}^{π} را تخمین میزند. گام بهبود سیاست، سیاست \mathbf{T} را براساس تخمینهای در \mathbf{V}^{π} بهبود میدهد. [78]

• برنامهنویسی پویا: تکنیکهای حل مبتنی بر مدل

تعداد زیادی تکنیک حل مبتنی بر مدل که به برنامهنویسی پویا (DP) معروف هستند وجود دارد. ایده اصلی، محاسبه توابع مقدار و استنتاج سیاستهای بهینه (که معمولا با یک روش تکراری صورت میپذیرد) از آنهاست. دو روش DP کلاسیک، تکرار سیاست و تکرار مقدار هستند. آنها را به ترتیب به طور مختصر شرح خواهیم داد. روش تکرار سیاست (PI) یک سیاست اولیه دلخواه π_0 آغاز میشود. سپس دنبالهای از تکرارها در آن می آیند که در آنها سیاست فعلی بعد از اینکه بهبود داده شد ارزیابی میشود. گام اول، گام ارزیابی سیاست، V^{π_k} را با استفاده مکرر از معادله(۲-۲) محاسبه می کند.

$$V_{n+1}(s) := R(s) + \gamma \max_{a \in A} [\sum_{s' \in S} T(s, a, s') V_n(s')]$$
 (۲-۲) معادله

گام دوم، گام بهبود سیاست، با استفاده از π_k را از π_{k+1} را از π_{k+1} محاسبه می کند. برای هر حالت، با استفاده از معادله(۲-۲)، عمل بهینه معین می شود.

$$\pi^*(s) := argmax_a[R(s) + \gamma \sum_{s' \in S} T(s, a, s') V^*(s')]$$
 معادله (۳-۲) معادله

در صورتی که برای همه حالتها $\pi_k(s)=\pi_k(s)$ باشد، سیاست پایدار است و الگوریتم تکرار سیاست می تواند متوقف شود. شکل کامل الگوریتم در الگوریتم ۱-۲ آورده شده است. هنگامی که MDP متناهی است یا به عبارت دیگر مجموعههای حالت و عمل متناهی هستند، تکرار سیاست بعد از تعداد متناهی تکرار همگرا می شود. هر سیاست $\pi_k=\pi^*$ اکیدا از سیاست π_k بهتر است مگر اینکه در موردی که $\pi_k=\pi^*$ است که در این صورت الگوریتم متوقف می شود. و چون برای یک MDP متناهی تعداد سیاستهای مختلف متناهی است، تکرار سیاست در زمان متناهی همگرا می شود. در عمل، معمولا بعد از تعداد کمی تکرار همگرا می شود. اگر چه تکرار سیاست، سیاست بهینه برای یک MDP معلوم را در زمان متناهی محاسبه می کند اما نسبتا ناکاراست. توابع مقدار برای همه سیاستهای میانی به ویژه، گام ارزیابی سیاست از لحاظ محاسباتی پرهزینه است. توابع مقدار برای همه سیاستهای میانی به ویژه، گام ارزیابی میاسیه می شوند که این به ازای هر تکرار شامل چندین جابجایی در فضای حالت کامل است.

```
Require: V(s) \in \mathbb{R} and \pi(s) \in A(s) arbitrarily for all s \in S
  {POLICY EVALUATION}
  repeat
     \Delta := 0
     for each s \in S do
        v := V(s)
        V(s) := R(s) + \gamma \cdot \sum_{s'} T(s, \pi(s), s') V(s')
        \Delta := \max(\Delta, |v - V(s)|)
  until \Delta < \sigma
   {POLICY IMPROVEMENT}
   policy-stable := true
   for each s \in S do
      b := \pi(s)
      \pi(s) := R(s) + \gamma \arg\max[\sum_{s'} T(s, a, s') V(s')]
      if b \neq \pi(s) then policy-stable := false
   if policy-stable then stop; else go to POLICY EVALUATION
```

الگوريتم ٢-١ تكرار سياست[79].

در مقابل، تکرار مقدار (VI) تا همگرایی، هر تابع مقدار را محاسبه نمی کند، بلکه بهروز رسانی های مورد نیاز را در حین کار انجام می دهد. در اصل، نسخه مختصر شده گام ارزیابی سیاست را با گام بهبود سیاست ترکیب می کند که اساسا معادله (Y-Y) است که به یک قانون بهروز رسانی تبدیل شده است:

$$V^{\pi}(s) := R(s) + \gamma \sum_{s' \in S} T(s, \pi(s), s') V^{\pi}(s')$$
 (۴-۲) معادله $V_{n+1}(s) := R(s) + \max_{a} [\sum_{s'} T(s, a, s') V_{t}(s')]$ (۵-۲) معادله (۶-۲) معادله معادله (۶-۲)

براساس معادلات معادله(\mathcal{A} - \mathcal{V}) و معادله(\mathcal{F} - \mathcal{V}) ، الگوریتم VI (به الگوریتم V-V رجوع شود) می تواند به صورتی که در ادامه گفته می شود بیان شود: با شروع از یک تابع مقدار V_0 برای همه حالتها، به طور مکرر مقدار هر حالت مطابق با معادله(V-V) به روز رسانی شده تا توابع مقدار بعدی (V-V) به روز رسانی شده تا توابع مقدار بعدی (V-V) به دست آیند.

تضمین می شود که VI در حد به V^* همگرا می شود، به عبارت دیگر، معادله بهینگی Bellman برای هر حالت صدق می کند. برای همه حالتها $S \in S$ ، یک سیاست قطعی π می تواند محاسبه شود.

$$\pi^*(s) := \max_{a} [R(s) + \gamma \sum_{s' \in S} T(s, a, s') V^*(s')]$$
 (۷-۲)معادله

```
Require: initialize V arbitrarily (e.g. V(s) := 0, \forall s \in S)

repeat
\Delta := 0
for each s \in S do
v := V(s)
for each a \in A(s) do
Q(s, a) := R(s) + \gamma. \sum_{s'} T(s, a, s') V(s')
V(s) := \max_{a} Q(s, a)
\Delta := \max(\Delta, |v - V(s)|)
until \Delta < \sigma
```

الگوريتم ٢-٢ تكرار مقدار [80].

تکرار سیاست تغییریافته (MPI) [81] یک روش بینابینی بین تکرار سیاست و تکرار مقدار را اتخاذ کرده است. تکرار مقدار اساسا هر دو گام تخمین سیاست و بهبود سیاست را در یک بهروز رسانی ترکیب می کند. از طرفی دیگر، تکرار سیاست دو گام را جدا کرده و هر دو آنها را در حد محاسبه می کند. MPI دو گام جدا را نگه می دارد، اما هر دو گام لزوما در حد محاسبه نمی شوند. در اینجا دیدگاه کلیدی این است که برای بهبود سیاست، به یک سیاست به طور دقیق ارزیابی شده نیاز نیست تا آن را بهبود داد. برای مثال، گام تخمین سیاست می تواند تقریبی باشد که بعد از آن یک گام بهبود سیاست می تواند بیاید. به طور کلی، هر دو گام می توانند کاملا مستقل و با روشهای مختلف انجام شوند. برای مثال، به جای اعمال مکرر قانون بهروز رسانی Bellman از معادله (۲۶) گام تخمین سیاست می تواند با استفاده از یک رویه نمونه برداری از قبیل تخمین افته هستند.

تعداد زیادی از نسخههای دیگر تکرار سیاست و مقدار پیشنهاد شدهاند. یک دسته بزرگ تحت عنوان روشهای ناهمگام از لحاظ ترتیبی - و همچنین بخشهایی خاص در آن ترتیب - که بهروز رسانیها انجام

میشوند با VI و VI استاندارد تفاوت دارند. روشهای ناهمگام مبتنی بر تخمین مقدار، مقادیر حالتها یا جفت حالت-عمل را نه در یک جابجایی در سراسر فضای حالت بلکه در عوض، مقادیر را در مناطق مخصوصا هدف قرار داده شده بهروز رسانی می کنند برای مثال، در امتداد تراژکتوریهای شبیهسازی شده که در رویکردهای یادگیری تقویتی رایج است. روشهای تکرار مقدار ناهمگام بهروز رسانیها را با یک ترتیب غیرثابتی انجام می دهند. تکنیک برنامهنویسی پویای بلادرنگ که توسط بارتو و همکاران [82] ارائه شده، جستجوی پیشرو را با برنامهنویسی پویا ترکیب می کند. برای بسیاری از این روشهای ناهمگام، همگرایی می تواند تحت شرایط کلی از قبیل اینکه اغلب همه حالتها باید به طور نامتناهی بهروز رسانی شوند تضمین شود.

• یادگیری تقویتی: تکنیکهای حل فارغ از مدل

در مقایسه با الگوریتمهایی که در بخش قبل در موردشان بحث شد، روشهای فارغ از مدل به در دسترس بودن مدلهای انتقال و پاداش از پیش معلوم، به طور مختصر مدل MDP، متکی نیستند. نبود یک مدل این نیاز را ایجاد می کند که باید از MDP نمونهبرداری کرد تا دانش آماری در مورد مدل نامعلوم جمعآوری شود. خیلی از تکنیکهای یادگیری تقویتی فارغ از مدل وجود دارند که محیط را با انجام عملها مورد بررسی قرار می دهند و از این راه، مشابه تکنیکهای مبتنی بر مدل همان نوع توابع مقدار حالت و حالت-عمل را تخمین میزنند. یکی از پایهای ترین و محبوب ترین روشها برای تخمین توابع Q – مقدار به شیوه فارغ از مدل، الگوریتم میزنند. یکی از پایهای ترین و محبوب ترین روشها برای تخمین توابع Q – مقدار به شیوه فارغ از مدل، الگوریتم Q – میدنان در سال Q – ارائه شده است (الگوریتم Q – طرح کلی آن را نشان می دهد).

```
Require: discount factor \gamma, learning parameter \alpha initialize Q arbitrarily (e.g. Q(s,a) := 0, \forall s \in S, \forall a \in A) for each episode do s is initialized as the starting state repeat choose an action a \in A(s) based on an exploration strategy perform action a observe the new state s' and received reward r Q(s,a) := Q(s,a) + \alpha \left[ R(s) + \gamma \cdot \max_{a' \in A(s')} Q(s',a') - Q(s,a) \right]s := s'
```

الگوريتم T-T [2] Q-Learning

یک جنبه مهم الگوریتمهای فارغ از مدل این است که نیاز به اکتشاف وجود دارد. چون مدل ناشناخته (نامعلوم) است، یادگیرنده باید عملهای مختلفی را امتحان کرده تا نتایجشان را ببیند. یک الگوریتم یادگیری باید بین اکتشاف و بهرهبرداری توازن برقرار کند، به عبارت دیگر برای اینکه پاداش زیادی به دست آورد یادگیرنده باید دانش فعلیش در مورد عملهای خوب را به کار گیرد، اگرچه گاهی اوقات باید عملهای مختلف را امتحان کرده

 ϵ - تا محیط را برای عملهای احتمالا بهتر مورد بررسی قرار دهد. پایهای ترین استراتژی اکتشاف سیاست و محیط را برای عملهای احتمالا بهتر مورد بهترین عمل فعلیش را با احتمال ϵ عمل دیگر (که به طور تصادفی انتخاب شده است) را با احتمال ϵ برمی گزیند.

الگوریتمهایی از قبیل Q-learning وگونههایی از قبیل SARSA [83] نیز، در صورتی که همه مقادیر مجزا در یک جدول پشتیبان ذخیره شده باشند همگراییشان به تابع مقدار (و سیاست) بهینه تضمین شده است. به دلیل تجرید^۱، موقعیتهایی پیش می آیند که این نتایج همگرایی برقرار نیستند.

الگوریتمهایی از قبیل Q-learning از خود راهاندازی آبرای محاسبه مقادیر استفاده می کنند: تخمین یک مقدار با مقدار تخمینزده شده حالت بعدی بهروز رسانی می شود. الگوریتمهای دیگری که از تخمینهای بدون تبعیض تر استفاده می کنند تکنیکهای Monte Carlo هستند. این روشها، تعداد تکرار انتقالها و پاداشها را نگه داشته و مقادیرشان را بر مبنای این تخمینها قرار می دهند.

روشهای غیرمستقیم توازنی بین یادگیری مبتنی بر مدل و یادگیری فارغ از مدل برقرار می کنند. این روشها اساسا فارغ از مدل هستند، اما به طور موازی با یادگیری تقویتی فارغ از مدل، یک مدل انتقال و پاداش یاد می گیرند و از این مدل برای یادگیری کاراتر تابع مقدار استفاده می کنند. نمونهای از این روش مدل Dyna می گیرند و از این مدل برای یادگیری ۱۹۹۱ ارائه شده است. روش دیگری که در آن سودمند بودن یادگیری مدل اثبات شده است پاکسازی اولویت بندی شده آ[85] است.

۲-۲- روشهای یادگیری تقویتی با جستجوی سیاست

جستجوی سیاست زیررشتهای در یادگیری تقویتی است که روی پیدا کردن پارامترهای خوب برای پارامتربندی یک سیاست معلوم تمرکز می کند. این روش مناسب رباتیک است چون می تواند از عهده فضاهای عمل و حالت با ابعاد بالا که یکی از چالشهای اصلی در یادگیری ربات است برآید.

² bootstrapping

¹ Abstraction

³ prioritized sweeping

جستجوی سیاست فارغ از مدل، یک رویکرد کلی برای یادگیری سیاستهای مبتنی بر تراژکتوریهای نمونهبرداری شده است. روشهای فارغ از مدل را براساس استراتژیهای ارزیابی سیاست، بهروز رسانی سیاست واکتشاف آنها دستهبندی کرده و یک دید یکپارچه به الگوریتمهای موجود ارائه میدهیم. غالبا یادگیری یک سیاست آسان تر از یادگیری یک مدل پیشرو دقیق است و از همین رو، روشهای فارغ از مدل به کرات در عمل استفاده میشوند. اما برای هر تراژکتوری نمونهبرداری شده تعامل با ربات ضروری است، امری که در عمل می تواند زمان بر و چالش برانگیز باشد.

جستجوی سیاست مبتنی بر مدل، این مشکل را ابتدا با یادگیری یک شبیهساز دینامیک ربات از دادهها حل می کند. سپس، شبیهساز تراژکتوریهایی تولید می کند که برای یادگیری سیاست استفاده می شوند.

۲-۲-۱ مقدمه

از رباتهای نظافت کار خانه گرفته تا صندلی چرخدارهای رباتیکی و رباتهای عمومی حمل و نقل، تعداد و انواع رباتهای که در زندگی همهروزه ما استفاده میشوند به سرعت در حال افزایش هستند. اکثرا کنترل کنندههای این رباتها توسط یک مهندس طراحی و تنظیم میشوند. برنامهنویسی رباتها کاری خسته کننده است که به سالها تجربه و درجهای بالا از تخصص نیاز دارد. کنترل کنندههای برنامهنویسی شده حاصله مبتنی بر فرض مدلهای دقیق از رفتار ربات و محیط پیرامونش هستند. در نتیجه، زمانی که یک ربات باید با وضعیتهای جدید تطبیق پیدا کرده یا زمانی که ربات/محیط نمیتوانند به اندازه کافی دقیق مدل شوند هارد کد کردن کنترل کنندهها محدودیتهای خودش را دارد. از همین رو، بین رباتهایی که در حال حاضر از آنها استفاده میشود و رباتهای کاملا خودمختار فاصله وجود دارد. در یادگیری ربات، برای حل یک کار رباتیکی از روشهای یادگیری ماشین به منظور استخراج خودکار اطلاعات مرتبط از دادهها استفاده میشود. با استفاده از قدرت و انعطاف پذیری تکنیکهای مدرن یادگیری ماشین، حوزه کنترل ربات بیشتر میتواند خودکار شود و فاصله تا انعطاف پذیری تکنیکهای مدرن یادگیری ماشین، حوزه کنترل ربات بیشتر میتواند خودکار شود و فاصله تا رباتهای خودمختار (برای مثال برای همکاری کلی در کارهای خانه، مراقبت از سالمندان و خدمات عمومی)

¹ Hardcode

روشهای جستجوی سیاست فارغ از مدل، از تعاملها با ربات واقعی به منظور تولید تراژکتوریهای نمونه استفاده میکنند. در حالی که نمونهبرداری تراژکتوریها در شبیهسازی کامپیوتر نسبتا سرراست و مستقیم است، معمولا هنگام کار با رباتها تولید هر نمونه به سطحی از نظارت انسانی نیاز دارد. در نتیجه، تولید تراژکتوری با سیستم واقعی به طور قابل ملاحظهای وقتگیرتر از کار با سیستمهای شبیهسازی شده است. علاوه بر این، تعاملات ربات واقعی موجب فرسودگی رباتهای غیرصنعتی میشود. اما با وجود تعداد نسبتا زیاد تعاملها با ربات که برای جستجوی سیاست فارغ از مدل مورد نیاز است، غالبا یادگیری یک سیاست آسان تر از یادگیری مدلهای پیشرو دقیق است و از همین رو، جستجوی سیاست فارغ از مدل در مقایسه با روشهای میشود.

روشهای جستجوی سیاست مبتنی بر مدل سعی دارند که مشکل کارا نبودن از لحاظ نمونه را با استفاده از تراژکتوریهای مشاهده شده به منظور یادگیری یک مدل پیشرو از دینامیک ربات و محیط پیرامونش حل کنند. سپس، این مدل پیشرو برای شبیهسازیهای درونی محیط و دینامیک ربات استفاده میشود که مبتنی بر آن سیاست آموخته میشود. روشهای جستجوی سیاست مبتنی بر مدل این پتانسیل را دارند که به تعاملهای کمتری با ربات نیاز داشته باشند و به طور کارا به وضعیتهای دیده نشده تعمیم یابند[88]. در حالی که ایده استفاده از مدلها در زمینه یادگیری ربات از دهه ۱۹۸۰ شناخته شده است[87]، به خاطر وابستگی زیادش به کیفیت مدلهای آموخته شده محدود شده است. در عمل مدل آموخته شده دقیق نیست، بلکه فقط یک تقریب کمابیش دقیق از دینامیک واقعی است. از آنجایی که سیاست آموخته شده ذاتا مبتنی بر شبیهسازیهای درونی با مدل آموخته شده است بنابراین مدلهای نادقیق می توانند منجر به استراتژیهای کنترلی شوند که در برابر خطاهای مدل مقاوم نیستند. در بعضی از موارد، مدلهای آموخته شده ممکن است کنترلی شوند که در برابر خطاهای مدل مقاوم نیستند. در بعضی یا جرمهای منفی باشند. اغلب اوقات، این آثار غیرمحتمل توسط الگوریتم جستجوی سیاست به کار گرفته شده که منجر به کیفیت پایین سیاست آموخته شده میشود. این اثر می تواند با استفاده از مدلهایی که خطاهای مدل را به طور صریح در نظر می گیرند کاهش یابد[3].

۲-۲-۲ جستجوی سیاست فارغ از مدل

روشهای جستجوی سیاست فارغ از مدل (PS)، سیاست را مستقیما و مبتنی بر تراژکتوریهای نمونهبرداری (وشهای جستجوی سیاست فارغ از مدل (PS)، سیاست را مستقیما و مبتنی بر تراژکتوریهای آنی به دست آمده برای شده $au^{[i]}$ به روز رسانی می کنند که تراژکتوریها و آرگیری و آرگیری به روز رسانی کنند که هنگام دنبال کردن سیاست جدید، تراژکتوریهای با پاداشهای بیش تر محتمل تر شوند و از همین رو، میانگین حاصل افزایش پیدا می کند.

$$J_{\theta} = \mathbb{E}[R(\tau)|\theta] = \int R(\tau)p_{\theta}(\tau)d au$$
 هعادله(۸-۲)

یادگیری سیاست غالبا آسان تر از یادگیری مدل ربات و محیط پیرامونش است و از همین رو، غالبا روشهای جستجوی سیاست مبتنی بر مدل استفاده میشوند. رویکردهای جستجوی سیاست فارغ از مدل را مبتنی بر استراتژیهای ارزیابی سیاستشان، بهروز رسانیشان[88] و اکتشافشان[89] دستهبندی می کنیم.

repeat

Explore: Generate trajectories $\tau^{[i]}$ using policy π_k

Evaluate: Assess quality of trajectories or actions

Update: Compute π_{k+1} given trajectories $\tau^{[i]}$ and evaluations

until Policy converges $\pi_{k+1} \approx \pi_k$

الگوريتم ٢-٢ جستجوى سياست فارغ از مدل[90].

استراتژی اکتشاف معین میکند که چگونه تراژکتوریهای جدید برای گام ارزیابی سیاست بعدی ایجاد شوند. استراتژی اکتشاف برای جستجوی سیاست فارغ از مدل کارا ضروری است چون به تنوع در تراژکتوریهای تولید شده نیاز داریم تا بهروز رسانی سیاست را معین کنیم اما اکتشاف بیش از حد نیز احتمال دارد به ربات صدمه بزند. بنابراین، بیشتر روشهای فارغ از مدل از یک سیاست تصادفی برای اکتشاف استفاده میکنند که فقط به صورت محلی اکتشاف میکند. استراتژیهای اکتشاف میتوانند به استراتژیهای اکتشاف مبتنی بر گام از عمل گام و استراتژیهای اکتشاف مبتنی بر گام از عمل

-

¹ Average return

اکتشافی در هر گام زمانی استفاده می کند، اکتشاف مبتنی بر اپیزود مستقیما بردار پارامتر θ سیاست را فقط در آغاز اپیزود تغییر می دهد.

استراتژی ارزیابی سیاست در مورد چگونگی ارزیابی کیفیت تراژکتوریهای اجرا شده تصمیم گیری می کند. در اینجا نیز می توانیم بین ارزیابیهای مبتنی بر اپیزود و ارزیابیهای مبتنی بر گام تفاوت قائل شویم. استراتژیهای ارزیابی مبتنی بر گام تراژکتوری au را به تک تک گامهایش au را به تک تک گامهایش au را به تک تک گامهایش کنند. در مقابل، ارزیابی مبتنی بر اپیزود مستقیما از حاصلهای کل تراژکتوریها استفاده کرده تا کیفیت پارامترهای سیاست استفاده شده au را ارزیابی کنند.

در آخر، استراتژی بهروز رسانی سیاست از کیفیت برآورد استراتژی ارزیابی استفاده کرده تا بهروز رسانی سیاست را تعیین کند. استراتژیهای بهروز رسانی می توانند براساس روش بهینه سازی ای که توسط الگوریتم جستجوی سیاست به کار گرفته شده است دسته بندی شوند. در حالی که بیش تر استراتژیهای رایج بهروز رسانی مبتنی بر استنتاج بر گرادیان صعودی هستند که منجر به روشهای گرادیان سیاست می شوند [91]، رویکردهای مبتنی بر استنتاج از بیشینه سازی امید ریاضی استفاده کرده [89] و رویکردهای نظریه اطلاعات [92] از دیدگاههای نظریه اطلاعات برای بهروز رسانی سیاست استفاده می کنند. روشهای مهم دیگری از قبیل رویکردهای انتگرال مسیر و بهینه سازی تصادفی نیز وجود دارند.

جستجوی سیاست فارغ از مدل می تواند به سیاستهایی با تعداد متوسط پارامتر یعنی تا چند صد پارامتر اعمال شود. بیش تر کاربردها از بازنماییهای سیاست خطی از قبیل کنترلکنندههای خطی یا مقدمات حرکت دینامیکی استفاده می کنند.

۲-۲-۳ جستجوی سیاست مبتنی بر مدل

روشهای جستجوی سیاست فارغ از مدل ذاتا مبتنی بر نمونهبرداری تراژکتوریها با استفاده از ربات هستند تا سیاستهای خوب پیدا را کنند. نمونهبرداری کردن تراژکتوریها در شبیهسازی کامپیوتر نسبتا سرراست و مستقیم است. اما هنگام کار با سیستمهای مکانیکی از قبیل رباتها، هر نمونه مربوط به تعامل مستقیم با ربات میشود. میشود، چیزی که نیاز به زمان آزمایش قابل توجه داشته و موجب فرسودگی در رباتهای غیرصنعتی میشود. بسته به کار می تواند یادگیری مدل آسان تر بوده یا یادگیری مستقیم سیاست آسان تر باشد. روشهای جستجوی

سیاست مبتنی بر مدل سعی دارند مشکل کارا نبودن از لحاظ داده را با استفاده از داده مشاهده شده به منظور یادگیری مدل پیشرو دینامیک ربات حل کنند. سپس، این مدل پیشرو برای شبیهسازیهای درونی دینامیک ربات استفاده می شود که مبتنی بر آن سیاست یاد گرفته می شود.

الگوریتمهای جستجوی سیاست مبتنی بر مدل معمولا تنظیم ٔ زیر را در نظر می گیرند: حالت X مطابق دینامیک مارکفی زیر تکامل می یابد

$$x_{t+1} = f(x_t, u_t) + \omega, \ x_0 \sim p(x_0)$$
 (۹-۲) معادله

که در این رابطه، f یک تابع غیرخطی، u یک سیگنال کنترلی (عمل) و ω نویز جمعپذیر (افزودنی) است که در این رابطه، f در نظر گرفته می شود. علاوه بر این، یک تنظیم اپیزودیک (چند بخشی) در نظر گرفته شده غالبا گوسی i.i.d در نظر گرفته می سیاست به حالت اولیه χ_0 بازگردانده می شود. توزیع حالت اولیه است که در آن ربات بعد از اجرای یک سیاست به حالت اولیه χ_0 بازگردانده می شود. توزیع گوسی χ_0 در نظر گرفته می شود. علاوه بر این، مسائل با افق متناهی را در نظر غالبا یک توزیع گوسی χ_0 در نظر گرفته می شود. علاوه بر این، مسائل با افق متناهی را در نظر می گیریم، به این معنی که هدف از جستجوی سیاست پیدا کردن یک سیاست پارامتربندی شده χ_0 است که امید ریاضی پاداش بلندمدت را بیشینه می کند

$$\pi_{\theta}^* \in \operatorname*{argmax} J_{\theta} \operatorname*{argmax} \sum_{t=1}^T \gamma^t \mathbb{E}[r(x_t, u_t) | \pi_{\theta}], \quad \gamma \in [0,1]$$
 معادله (۱۰-۲)

که در این رابطه r سیگنال پاداش آنی، γ ضریب کاهش و سیاست π_{θ} با پارامترهای θ پارامتربندی شده θ^* است. بنابراین، پیدا کردن π_{θ}^* در معادله(۲-۱۰) معادل با پیدا کردن پارامترهای سیاست بهینه مربوطه π_{θ}^* است.

برای برخی مسائل، روشهای یادگیری تقویتی مبتنی بر مدل نوید این را میدهند که با یادگیری مدلی از نگاشت انتقال در معادله(۲-۹)، در مقایسه با یادگیری تقویتی فارغ از مدل به تعاملهای کمتری با ربات نیاز دارند، در عین حال با استفاده از مدلی که از دادههای مشاهده شده آموخته شده به صورت کارا به وضعیتهای دیده نشده تعمیم مییابند[86].

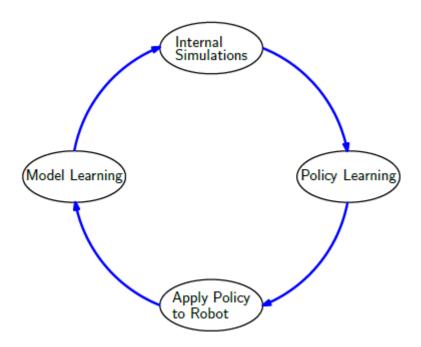
¹ Set-up

² Additive

³ Episodic

ایده کلی یادگیری تقویتی مبتنی بر مدل در شکل ۲-۵ نشان داده شده است. مدل آموخته شده برای شبیه سازیهای درونی استفاده می شود، به عبارت دیگر، پیشبینیها در مورد چگونگی رفتار ربات واقعی و محیط پیرامونش است در صورتی که ربات سیاست فعلی را دنبال کند. براساس این شبیه سازیهای درونی، کیفیت سیاست با استفاده از معادله(۲-۱۰) ارزیابی شده و مطابق آن بهبود می یابد. سپس، سیاست به روز رسانی شده دوباره با استفاده از معادله(۲-۱۰) ارزیابی شده و بهبود می یابد. این چرخه ارزیابی ابهبود سیاست هنگامی که سیاست یاد گرفته می شود پایان می یابد، به این معنی که دیگر تغییر نمی کند و به بهینه (محلی) رسیده است. به محض یادگیری یک سیاست، آن سیاست به ربات اعمال شده و داده جدید ضبط می شود. این مجموعه داده بعد از ترکیب با داده از قبل جمعآوری شده، برای بهروز رسانی و اصلاح مدل دینامیکی آموخته شده استفاده می شود. از لحاظ نظری، این چرخه برای همیشه ادامه می یابد. توجه داشته باشید که فقط اعمال سیاست به تعامل با ربات نیاز دارد. شبیه سازی های درونی و یادگیری سیاست فقط از مدل کامپیوتری یادگرفته شده دینامیک ربات استفاده می کنند.

در حالی که ایده استفاده از مدلها در زمینه یادگیری ربات، از کار Aboaf [87] در دهه ۱۹۸۰ شناخته شده است، به خاطر وابستگی شدیدش به کیفیت مدل آموخته شده محدود شده است، موضوعی که در شکل ۲-۵ مشهود است: سیاست آموخته شده ذاتا مبتنی بر شبیهسازیهای درونی است که با استفاده از مدل آموخته شده انجام شده است. زمانی که مدل دقیقا به دینامیک درست ربات مربوط میشود، نمونهبرداری از مدل آموخته شده معادل با نمونهبرداری از ربات واقعی است.



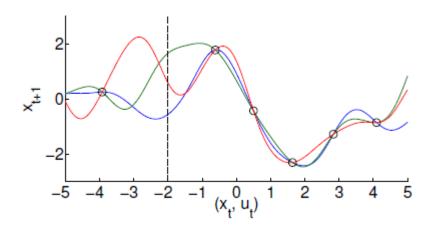
شکل ۲-۵ چرخه کلی در یادگیری تقویتی مبتنی بر مدل. مدل یاد گرفته شده برای شبیهسازیهای درونی استفاده میشود (تمرین ذهنی) که براساس آن سیاست بهبود می یابد. سپس سیاست آموخته شده به ربات اعمال میشود. داده حاصل از این تعامل برای اصلاح مدل استفاده میشود و این چرخه ادامه می یابد تا زمانی که مدل و سیاست آموخته شده همگرا شوند.[90]

اما، در عمل مدل آموخته شده دقیق نیست بلکه فقط یک تقریب کمابیش دقیق از دینامیک واقعی است. برای مثال، در مناطقی که داده آموزشی پراکنده است، همانطور که در شکل ۲-۶ نشان داده شده است کیفیت مدل یاد گرفته شده می تواند نامناسب باشد. چندین تابع محتمل وجود دارند که می توانند مقادیر تابع مشاهده شده را تولید کرده باشند (دایرههای سیاه). در مناطق با داده آموزشی پراکنده، مدلها و پیشبینیهایشان به طور قابل توجهی فرق می کنند. هر مدل منفرد منجر به پیشبینیهای خود رای میشود که این پیشبینیها به نوبه خود می توانند منجر به استراتژیهای کنترلی شوند که در برابر خطاهای مدل مقاوم نیستند. این رفتار می تواند اثری شدید در رباتیک داشته باشد، برای مثال می تواند منجر به تخمین جرمهای منفی یا ضرایب اصطکاک منفی شود. این آثار غیرمحتمل غالبا توسط سیستم یادگیری به کار گرفته می شوند چون انرژی را در سیستم قرار می دهند که باعث می شود سیستم به "حرکت دائمی" اعتقاد پیدا کند. بنابراین، به جای انتخاب سیستم قرار می دهند که باعث می شود سیستم به "حرکت دائمی" اعتقاد پیدا کند. بنابراین، به جای انتخاب

¹ Overconfident

² Perpetuum mobile

یک مدل منفرد (برای مثال مدل بیشینه درستنمایی)، باید عدم اطمینانمان در مورد تابع پنهان f را با یک توزیع احتمالاتی p(f) نشان دهیم تا در مقابل اینچنین خطاهای مدل مقاوم باشدp(f). با در نظر گرفتن عدم قطعیت مدل، احتمال اثر حرکت دائمی به طور چشمگیر کمتر است.



شکل ۲-۶ خطاهای مدل در جستجوی سیاست مبتنی بر مدل. در این مثال، شش مقدار تابع مشاهده شدهاند (دایرههای سیاه). سه تابع نشان داده شدهاند که می توانند این مشاهدات را تولید کرده باشند. هر تابع منفرد (برای مثال تابع بیشینه درستنمایی) پیشبینیهای کمابیش دلخواهی در مناطق با داده آموزشی پراکنده تولید می کند که یک مکان نمونه از آن با خطچین نشان داده شده است. به جای انتخاب یک تقریبزننده تابع منفرد، عدم اطمینانمان در مورد تابع زیرین را با یک توزیع احتمالاتی نشان مىدهيم تا در برابر اين چنين خطاهاى مدل مقاوم باشد.[90]

از آنجایی که سیاست آموخته شده ذاتا به کیفیت مدل پیشرو آموخته شده (که اساسا به عنوان شبیهساز ربات به کار می رود) متکی است، خطاهای مدل فقط باعث سیاستهای ضعیفتر نمی شوند، بلکه غالبا فرآیند یادگیری را به شدت به یک طرف متمایل می کنند ٔ از همین رو، منابع علمی جستجوی سیاست مبتنی بر مدل اکثرا روی ساخت مدل تمرکز می کند یا به عبارت دیگر توضیح می دهد که چه نوع مدلی برای دینامیک پیشرو استفاده شده و چگونه آموزش داده شده است.

¹ Bias

رویکردهای سر و کله زدن با مدلهای غیرقطعی. سر و کله زدن با مدلهای دینامیکی نادقیق یکی از بزرگترین چالشها در یادگیری تقویتی مبتنی بر مدل است چون خطاهای کوچک در مدل میتواند منجر به خطاهای بزرگ در سیاست شود[86]. بیدقتیها ممکن است از یک دسته مدل بیشازحد محدودکننده یا از نبود مجموعه دادگان به اندازه کافی قوی برای آموزش مدلها نشأت بگیرند، چیزی که میتواند منجر به مدلسازی بد دینامیک پیشرو درست شود.

علاوه بر این، معمولا نویز سیستم منشأ دیگری از عدم قطعیت را اضافه می کند.

به طور معمول، یک فرض معادل قطعیت انجام می شود و مدل بیشینه درستنمایی برای برنامه ریزی انتخاب می شود [93]. اما این فرض معادل قطعیت در بیش تر موارد جالب توجه نقض شده و می تواند منجر به خطاهای بزرگ در سیاست شود. علاوه بر این، همان طور که قبلا در مرجع [87] اشاره شده است، خیلی از رویکردها مشتقات بازگشتی مورد نظر (امید ریاضی بازگشت) را با پس – انتشار مشتقات از طریق مدل های پیشرو آموخته شده از سیستم به دست می آورند. به ویژه، این گام در معرض خطاهای مدل است چون بهینه سازهای مبتنی بر گرادیان، پارامترهای سیاست را در امتداد گرادیان هایشان بهبود می بخشند. سیاست آموخته شده باید قوی باشد تا خطاهای مدل را جبران کند به طوری که زمانی که به سیستم واقعی اعمال می شود منجر به عملکرد خوب شود. یادگیری مدل های دینامیکی درست برای ساختن سیاستهای قوی حیاتی است و یکی از بزرگ ترین چالش ها در جستجوی سیاست مبتنی بر مدل باقی مانده است.

در مرجع [94]، نویسندگان خطای ناشناخته (نامعلوم) دینامیکی را با استفاده از رگرسیون وزندار میدان دریافتی [94] مدلسازی کردهاند. مدلسازی صریح اختلالات ناشناخته منجر به نیرومندی فزاینده کنترل کنندههای یاد گرفته شده میشود. ایده طراحی کنترل کنندهها در مواجهه با مدلهای پیشرو نادقیق (تلقی شده به عنوان ایده آل)، به طور تنگاتنگ مرتبط با کنترل مقاوم در رباتیک کلاسیک است.

¹ Deal

^۳ فرض می شود که سیاست بهینه برای مدل آموخته شده به سیاست بهینه برای دینامیک درست مربوط می شود. عدم قطعیت در مورد مدل آموخته شده نادیده گرفته می شود.

³ Receptive-field weighted regression

کنترل مقاوم سعی دارد با وجود خطاهای مدلسازی (به طور معمول محدود) به عملکرد تضمین شده یا پایداری برسد. برای مثال، شکل دهی چرخه \mathcal{H}_{∞} [96] تضمین می کند که سیستم نزدیک به رفتار مورد انتظارش می ماند اگرچه اختلالات (محدود) به سیستم وارد شوند.

در کنترل تطبیقی، عدم قطعیتهای پارامتر معمولا به وسیله توزیعهای احتمالاتی نامحدود (بی کران) توصیف می شود [97]. عدم قطعیت پارامتر مدل به طور معمول در طراحی الگوریتمهای کنترل تطبیقی استفاده نمی شود. در عوض، تخمینهای پارامترها به عنوان پارامترهای درست تلقی می شوند [98]. یک رویکرد به طراحی کنترل کنندههای تطبیقی که عدم قطعیت در مورد پارامترهای مدل را در نظر می گیرد کنترل تطبیقی تصادفی تصادفی است [97]. هنگامی که عدم قطعیت پارامتر از طریق کاوش کاهش می یابد، کنترل تطبیقی تصادفی به اصل کنترل دو گانه منجر می شود [99]. کنترل تطبیقی دو گانه عمدتا برای سیستمهای خطی مورد بررسی قرار گرفته است [98]. در مرجع [100]، بسطی از کنترل تطبیقی دو گانه به مورد سیستمهای غیر خطی با کنترل های همگر آپیشنهاد شده است. یک قانون کنترلی واریانس کمینه به دست می آید و عدم قطعیت در مورد پارامترهای مدل جریمه شده تا تخمینشان را بهبود دهد، امری که نیاز به شناسایی قبلی سیستم را حذف می کند.

رویکردهای یادگیری تقویتی که به طور صریح مشکل دقیق نبودن مدلها در رباتیک را حل می کنند اخیرا معرفی شدهاند[93]. برای مثال، در مرجع [101]، ایده کلیدی استفاده از آزمایش واقعی برای ارزیابی یک سیاست است، اما سپس از یک مدل ابتدایی سیستم برای تخمین مشتق ارزیابی نسبت به پارامترهای سیاست استفاده می کند (و بهبودهای محلی را پیشنهاد می کند). به ویژه، الگوریتم پیشنهاد شده به صورت مکرر مدل استفاده می کند (و بهبودهای محلی را پیشنهاد می کند). به ویژه، الگوریتم پیشنهاد شده به صورت مکرر مدل $f^{[i+1]}(x_t,u_t)=f^{[i]}(x_t,u_t)+x_{t+1}^{[i]}-f^{[i]}(x_t^{[i]},u_t^{[i]})$ به روز رسانی می کند که در این رابطه $f^{[i+1]}(x_t,u_t)=f^{[i]}(x_t,u_t)$ در اینجاه $f^{[i+1]}(x_t,u_t)=f^{[i]}(x_t,u_t)$ در اینجاه تراژکتوری مشاهده شده را پیشبینی می کند به عبارت دیگر می شود که مدل به روز رسانی شده $f^{[i+1]}(x_t,u_t)$ دقیقا تراژکتوری مشاهده شده را پیشبینی می کند به عبارت دیگر $f^{[i+1]}(x_t,u_t)=x_{t+1}^{[i]}$. این الگوریتم، گرادیانهای سیاست را در امتداد تراژکتوری ای از حالتها و

¹ Stochastic

² Affine

³ Real-life

⁴ Crude

کنترلها در سیستم واقعی ارزیابی می کند. در مقابل، یک رویکرد مبتنی بر مدل معمول ۱، مشتقات را در امتداد تراژکتوریای که به وسیله مدل پیشبینی شده است ارزیابی می کند، چیزی که به هنگام دقیق نبودن مدل به تراژکتوری سیستم واقعی مربوط نمی شود. توجه داشته باشید که رویکرد مرجع [101] مستقیما به دینامیک انتقال تصادفی یا سیستمهای با حالتهای پنهان تعمیم نمی یابد. علاوه بر این، یک مدل پارامتری تقریبی از دینامیک اصلی باید از پیش معلوم باشد.

چالشهای بزرگ در جستجوی سیاست مبتنی بر مدل. سه چالش کلی وجود دارد که در روشهای جستجوی سیاست مبتنی بر مدل باید مورد بررسی قرار گیرد: چه مدلی را باید یاد بگیرد، چگونه از این مدل برای پیشبینیهای بلندمدت چگونه این سیاست را بهروزرسانی کند. این سه چالش به سه عنصر در شکل ۲-۵ مربوط میشوند که نیازی به تعامل فیزیکی با ربات ندارند: یادگیری مدل، شبیهسازیهای درونی و یادگیری سیاست.

در زیربخش اول، مرور کلی مختصری درباره دو مدلی که غالبا در جستجوی سیاست مبتنی بر مدل استفاده میشوند[3] ارائه میدهیم که عبارتند از رگرسیون (بیزی) وزندار محلی (LWBR) و فرآیندهای گوسی (GPs) (GPs) در زیربخش دوم، در مورد دو روش کلی طرز استفاده از این مدلها برای پیشبینیهای بلندمدت بحث میکنیم: استنتاج تصادفی یا به عبارت دیگر نمونهبرداری و استنتاج تقریبی قطعی. در زیربخش سوم، مختصرا در مورد چند گزینه بهروز رسانی سیاست بحث میکنیم.

بعد از معرفی این مفاهیم کلی، در زیربخش چهارم، در مورد الگوریتمهای جستجوی سیاست مبتنی بر مدلی بحث می کنیم که مدلهای آموخته شده و الگوریتمهای استنتاجی که در جدول ۱-۲ نشان داده شدهاند را ترکیب می کنند. ما روی مورد اپیزودیک تمرکز می کنیم که در آن توزیع حالت آغازین $p(x_0)$ معلوم است. به ویژه، جزئیات چهار روش جستجوی سیاست مبتنی بر مدل را بیان می کنیم. ابتدا با PEGASUS شروع می کنیم که یک مفهوم کلی برای نمونه برداری کارای تراژکتوری در MDP های تصادفی برای مدلی معلوم است است [103]. سپس، دو روش ترکیب PEGASUS با LWBR با منظور یادگیری کنترل کنندههای مقاوم برای پرواز هلی کوپترها [104] ارائه می دهیم. سپس، رویکردی برای استفاده از الگوریتم PEGASUS برای برای پرواز هلی کوپترها [104] ارائه می دهیم. سپس، رویکردی برای استفاده از الگوریتم PEGASUS برای

¹ Typical

نمونهبرداری تراژکتوریها از مدلهای پیشرو GP [105] در زمینه یادگیری یک کنترلکننده بالون را ارائه میدهیم. سرانجام به عنوان رویکرد چهارم، ایدههای چارچوب کاری جستجوی سیاست GP [3] را به طور کلی شرح میدهیم که استنتاج تقریبی قطعی کارا برای پیشبینیهای بلندمدت را با مدلهای دینامیکی GP برای یادگیری کنترل سیستمهای مکانیکی و بازوهای ربات ترکیب میکند.

جدول ۲-۱ رویکردهای جستجوی سیاست مبتنی بر مدل بر اساس مدل آموخته شده و روش تولید تراژکتوریها گروهبندی میشوند. به خاطر سادگی نمونهبرداری تراژکتوریها، بیش تر روشهای جستجوی سیاست این رویکرد را دنبال میکنند. پیشبینیهای قطعی تراژکتوری فقط می توانند در موارد خاصی انجام شوند که استنتاج تقریبی به فرم بسته ممکن است. گرچه آنها از لحاظ ریاضیاتی بیش تر از نمونهبرداری تراژکتوری درگیر هستند، اما از واریانسهای بزرگ نمونهها رنج نمیبرند. علاوه بر این، می توانند امکان محاسبات تحلیلی گرادیان را فراهم کنند که در هنگام وجود صدها پارامتر سیاست حیاتی است. [90]

پیشبینی تراژکتوری		مدل پیشرو آموخته شده
قطعی	تصادفی	
_	[104][106][107]	مدلهای خطی (محلی)
[3][108]	[105]	فرآیندهای گوسی

۱-۳-۲-۲ مدلهای پیشرو احتمالاتی

برای کاهش اثر خطاهای مدل، مدلهای احتمالاتی که عدم قطعیت در مورد دینامیک انتقال اصلی را بیان میکنند به مدلهای قطعی (برای مثال مدلهای بیشینه درستنمایی از دینامیک انتقال) که یک فرض معادل قطعیت را ایجاب میکنند ترجیح داده میشوند.

دو مدل احتمالاتی غیرپارامتری را در زیر به طور مختصر معرفی می کنیم که غالبا برای یادگیری دینامیک پیشرو در زمینه یادگیری تقویتی و رباتیک استفاده می شوند: رگرسیون بیزی وزندار محلی (LWBR) [104] و فرآیندهای گوسی[3].

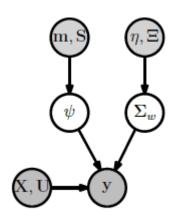
رگرسیون بیزی وزندار محلی

ابتدا با مدل رگرسیون خطی شروع می کنیم که در آن دینامیک انتقال به صورت زیر بیان می شود:

$$x_{t+1} = [x_t, u_t]^T \psi + \omega, \qquad \omega \sim \mathcal{N}(0, \Sigma_\omega)$$
 (۱۱-۲)معادله

در اینجا، ψ پارامترهای مدل رگرسیون خطی بیزی هستند و $\omega \sim \mathcal{N}(0,\Sigma_\omega)$ نویز گوسی u نویز گوسی نامعلوم u که ورودی u که ورودی u که ورودی این مدل نسبت به پارامترهای نامعلوم u که ورودی u که ورودی از را وزن دار می کند خطی است.

در رگرسیون خطی بیزی، توزیعهای مقدم روی پارامترها ψ و روی واریانس نویز Σ_{ω} در نظر می گیریم. به طور $1/\sigma_i^2$ معمول، توزیع مقدم روی ψ ، گوسی با میانگین m و کوواریانس S و توزیع مقدم روی عناصر قطری Σ_{ω}^2 از Σ_{ω}^{-1} یک توزیع گاما با پارامترهای مقیاس و شکل S و آست، به طوری که مدل کلی مزدوج است، به شکل ۲-۷ نگاه کنید که در آن ورودیهای آموزشی با S و اهداف آموزشی با S نشان داده شدهاند. در معادله (۱۱-۲ نگاه کنید که در آن ورودیهای آموزشی با کردن تخمینهای بیشینه درستنمایی یا توزیعهای پسین پارامترها S نسبتا سرراست و مستقیم است. اما خود مدل خیلی گویا نیست چون یک رابطه خطی کلی بین ورودیها S و حالت بعدی S بین فرض می کند.



شکل ۷-۲ مدل گرافی برای رگرسیون خطی بیزی. یک مقدم گوسی با پارامترهای \mathbf{m} و \mathbf{S} که روی پارامترهای \mathbf{v} در نظر گرفته شده است و مقدمهای گاما با پارامترهای \mathbf{v} و \mathbf{v} که روی عناصر قطری ماتریس دقت \mathbf{v} در نظر گرفته شدهاند. ورودیها و اهداف آموزشی به ترتیب با \mathbf{v} و \mathbf{v} نشان داده شدهاند. [90]

ایده رگرسیون خطی وزندار محلی (LWR)، استفاده از ویژگیهای خوب مدل رگرسیون خطی و اما در نظر گرفتن دسته کلی تر از توابع است: توابع خطی محلی. LWR یک تقریب خطی محلی از تابع اصلی پیدا می کند [109]. به این منظور، هر ورودی آزمون (x_t, u_t) با یک ضریب وزندهی b_i تجهیز می شوند که معین می کند نقطه آموزشی (x_i, u_i) چقدر به (x_t, u_t) نزدیک است. یک مثال برای چنین وزنی، وزندهی با شکل می کند نقطه آموزشی (x_i, u_i) چقدر به (x_i, u_i) نزدیک است. یک مثال برای چنین وزنی، وزندهی با شکل گوسی است (x_i, u_i) و (x_i, u_i) و (x_i, u_i) و (x_i, u_i) خیلی بزرگ تر از (x_i, u_i) باید برای هر نقطه مورد سوال محاسبه بزرگ تر از (x_i, u_i) باید برای هر نقطه مورد سوال محاسبه شوند، فقط ذخیره پارامترها (x_i, u_i) کافی نیست بلکه کل مجموعه داده آموزشی (x_i, u_i) نیاز است، چیزی که منجر به رویکرد غیر پارامتری می شود.

مشابه رگرسیون خطی بیزی، می توانیم مقدمهایی روی پارامترها و کوواریانس نویز در نظر بگیریم. برای سادگی، یک ماتریس کوواریانس نویز ψ فرض می کنیم یک ماتریس کوواریانس نویز Σ_ω و یک توزیع مقدم گوسی با میانگین صفر روی پارامترها ψ فرض می کنیم . $\mathcal{N}(\psi|0,S)$

برای هر نقطه مورد سوال (x_t, u_t) ، مطابق قضیه بیز به صورت زیر یک توزیع پسین روی پارامترها ψ محاسبه می شود:

$$p(\psi|X,U,y) = \frac{p(\psi)p(y|X,U,\psi)}{p(y|X,U)} \propto p(\psi)p(y|X,U,\psi)$$
 (۱۲-۲)معادله

 (x_t,u_t) برای راحتی در نمادگذاری $\tilde{X}\coloneqq [X,U]$ را تعریف می کنیم. روابط میانگین و کوواریانس پسین $\tilde{X}\coloneqq [X,U]$ به ترتیب به صورت زیر است:

$$\mathbb{E}[\psi | \tilde{X}, y] = S\tilde{X}B(B\tilde{X}^TS\tilde{X}B + \Sigma_\omega)^{-1}y = S\tilde{X}B\Omega^{-1}y$$
 معادله(۱۳-۲) معادله

$$\Omega^{-1} = (B\tilde{X}^T S\tilde{X}B + \Sigma_{\omega})^{-1}$$

$$cov[\psi|\tilde{X},y] = S - S^T \tilde{X} B \Omega^{-1} B \tilde{X}^T S$$
 (۱۴-۲) معادله

$$b_i = exp(-\|(x_i, u_i) - (x_t, u_t)\|^2/\kappa^2)$$
 (10-7)

در روابط بالا، $B = diag(b_1, \cdots, b_n)$ است و y اهداف آموزشی هستند.

توزیع پیشبینی کننده $p(x_{t+1})$ برای یک جفت حالت-کنترل معلوم (x_t, u_t) به ترتیب به صورت زیر هستند:

$$\mu_{t+1}^{x} = [x_t, u_t]^T \mathbb{E}[\psi | \tilde{X}, y] = [x_t, u_t]^T S \tilde{X} B \Omega^{-1} y$$
 (۱۶-۲) معادله

$$\Sigma_{t+1}^{x} = [x_t, u_t]^T cov[\psi | \tilde{X}, y][x_t, u_t] + \Sigma_{\omega}$$
 (۱۷-۲) معادله

در عمل، میانگین و کوواریانس پسین روی پارامترها ψ میتواند از طریق اعمال لمهای معکوسسازی ماتریس[109] و استفاده از پراکندگی خیلی کاراتر محاسبه شود.

حال نگاهی به کوواریانس پیشبینی کننده می کنیم در زمانی که $[x_t,u_t]$ از مجموعه آموزشی [X,U] خیلی دور است: ماتریس وزن تقریبا صفر است، چیزی که منجر به واریانس پسین روی پارامترهای مدل ψ می شود که معادل با عدم قطعیت مقدم S است (به معادله(۲-۱۴) رجوع شود). از همین رو، برخلاف مورد رگرسیون وزن دار محلی غیربیزی، واریانس پیشبینی کننده در $[x_t,u_t]$ غیرصفر است.

- رگرسیون فرآیند گوسی

یک GP، تماما با یک تابع میانگین m(.) و یک تابع/کرنل کوواریانس نیمهمعین مثبت k(.,.) مشخص GP، تماما با یک تابع میانگین m(.) و یک تابع مقدم میانگین m(.) و تابع کوواریانس به صورت می شود. فرضیات استاندارد در مدلهای m(.) است. m(.) هستند که m(.) است.

$$k(\tilde{x}_p, \tilde{x}_q) = \sigma_f^2 \exp\left(-\frac{1}{2}(\tilde{x}_p - \tilde{x}_q)^T \Lambda^{-1}(\tilde{x}_p - \tilde{x}_q)\right) + \delta_{pq}\sigma_\omega^2$$
 (۱۸-۲) معادله

در معادله(۱۸-۲)، $\Lambda:=diag([\ell_1^2,\cdots,\ell_D^2])$ در معادله $\Lambda:=diag([\ell_1^2,\cdots,\ell_D^2])$ در معادله المخصه $X=[\tilde{x}_1,\cdots,\tilde{x}_n]$ وابسته است و σ_f^2 واریانس مقدم تابع پنهان T است. با فرض T ورودی آموزشی

مربوطه ℓ_i واریانس سیگنال σ_f^2 و واریانس $y=[y_1,\cdots,y_n]^T$ و واریانس $y=[y_1,\cdots,y_n]^T$ و واریانس نویز σ_ω^2) با استفاده از بیشینه سازی شواهد یاد گرفته می شوند σ_ω^2 0.

 x_{t+1} مدل پیشبینی شده و حالت بعدی پیشبینی شده GP بسین یک مدل پیشبینی شده است و حالت بعدی پیشبینی شده دارای توزیع گوسی است

$$p(x_{t+1}|x_t, u_t) = \mathcal{N}(x_{t+1}|\mu_{t+1}^x, \Sigma_{t+1}^x)$$
 (۱۹-۲)معادله

$$\mu_{t+1}^x = \mathbb{E}_f[f(x_t,u_t)]$$
 , $\Sigma_{t+1}^x = var_f[f(x_t,u_t)]$ (۲۰-۲) معادله

که میانگین و واریانس پیشبینی GP به ترتیب به صورت زیر هستند

$$\mu_{t+1}^{x} = k_{*}^{T} K^{-1} y = k_{*}^{T} \beta$$
 معادله(۲۱-۲)

$$\Sigma_{t+1}^{x} = k_{**} - k_{*}^{T} K^{-1} k_{*}$$
 (۲۲-۲)

 K_{ij} عناصر کرنل با عناصر $\beta:=K^{-1}y$ هستند و $k_*:=k(\widetilde{x}_t,\widetilde{x}_t)$ که $k_*:=k(\widetilde{X},\widetilde{x}_t)$ هستند و $k(\widetilde{x}_t,\widetilde{x}_t)$ است.

توجه داشته باشید که در مناطق دور از داده آموزشی، عدم قطعیت پیشبینی کننده در معادله (۲۲-۲) به عدم $k_{**}> \Sigma_{\omega}=0$ است $\Sigma_{\omega}=0$ است مقدم در مورد تابع مربوط می شود، به عبارت دیگر برای سیستمهای قطعی که Ω است Ω به دست می آوریم. بنابراین، Ω یک مدل غیرپارامتری غیرمنحط است.

۲-۳-۲-۲ پیشبینیهای بلندمدت با یک مدل معلوم

در بخشهای زیر فرض می کنیم که یک مدل برای دینامیک انتقال معلوم است. مشروط به این مدل، بین دو رویکرد تولید پیشبینیهای بلندمدت تمایز قائل می شویم: رویکردهای مبتنی بر نمونهبرداری -Monte رویکرد تولید پیشبینی قطعی Carlo

- پیشبینی مبتنی بر نمونهبرداری تراژکتوری: PEGASUS

PEGASUS (جستجو و ارزیابی خوب بودن سیاست با استفاده از سناریوها^۱) یک چارچوب کاری مفهومی برای MDP نمونهبرداری تراژکتوری در MDP های تصادفی است[103]. ایده کلیدی، انتقال MDP تصادفی به یک شونهبرداری تراژکتوری در PEGASUS های تصادفی داخلی قطعی تقویت شده است. به این منظور، PEGASUS دسترسی به یک شبیهساز بدون مولد عدد تصادفی داخلی را فرض می کند. زمانی که از این مدل نمونهبرداری می کنیم، PEGASUS پیشاپیش از بیرون اعداد تصادفی را فراهم می کند. از این طریق، PEGASUS واریانس نمونهبرداری را به شدت کاهش می دهد. از همین رو، نمونهبرداری با پیروی از رویکرد PEGASUS، در جستجوی سیاست فارغ از مدل نیز به صورت رایج استفاده می شود.

■ نمونهبرداری تراژکتوری و ارزیابی سیاست

فرض کنید که یک مدل پیشرو از سیستم پیشرو داده شده است که میتواند برای نمونهبرداری تراژکتوریها استفاده شود. اگر انتقالات حالت تصادفی باشند، محاسبه پاداش بلندمدت مورد انتظار که رابطه آن به صورت معادله(۲-۲۳) است

$$J_{\theta} = \sum_{t=0}^{T} \gamma^{t} \mathbb{E}[r(x_{t})|\pi_{\theta}], \quad x_{0} \sim p(x_{0}), \quad \gamma \in [0,1]$$
 معادله (۲۳-۲)

به تعداد زیادی نمونه تراژکتوری برای محاسبه تقریب $ilde{J}$ به $ilde{J}$ نیاز خواهد داشت که معادله آن به صورت زیر است و و در آن نمونههای $x_0^{[i]}$ از $x_0^{[i]}$ هستند.

$$ilde{J}_{ heta} = rac{1}{m} \sum_{i=1}^m J_{ heta}ig(x_0^{[i]}ig)$$
 معادله(۲۴-۲) معادله

حتی محاسبه گرادیانهای سیاست قابل اعتماد به نمونههای بیش تری برای مشتقات قوی نیاز خواهد داشت. $\lim_{m\to\infty} \tilde{J}_{\theta} = J_{\theta} \quad \text{ lim} \quad J_{\theta} = J_{\theta}.$ اما، به عنوان حد تعداد نامتناهی نمونه، $J_{\theta} = J_{\theta}$ را به دست می آوریم.

_

¹ Policy Evaluation-of-Goodness And Search Using Scenarios

از آنجایی که دنباله اعداد تصادفی ثابت است، تکرار آزمایش یکسان منجر به نمونه تراژکتوری یکسان میشود. PEGASUS میتواند به صورت تولید m تراژکتوری Monte Carlo و به دست آوردن میانگین پاداششان توصیف شود، اما تولید عدد تصادفی از پیش تعیین شده است. میتوان نشان داد که حل MDP قطعی تقویت شده معادل با حل MDP تصادفی اصلی است[103]. الگوریتم PEGASUS در الگوریتم ۵-۲ خلاصه شده است.

```
Init: g(x,u,w), reward r, random numbers w_0, w_1, ..., initial state distribution p(x_0), policy \pi_\theta for i=1,...,m do x_0^{[i]} \sim p(x_0) \qquad \rhd \text{Sample "scenario" from initial state} distribution \text{for } t=0,...,T-1 \text{ do} x_{t+1}^{[i]}=g(x_t^{[i]},\pi_\theta(x_t),w_t) \qquad \rhd \text{Succ. state in augmented} MDP end for end for
```

الگوريتم ۲-۵ الگوريتم PEGASUS براى نمونهبردارى تراژکتورىها [90].

- پیش بینیهای بلندمدت قطعی

به جای انجام نمونهبرداری تصادفی، یک توزیع احتمالاتی $p(\tau)$ روی تراژ کتوریها $au=(x_0,\cdots,x_T)$ می تواند با استفاده از تقریبهای قطعی از قبیل خطیسازی[110]، روشهای نقطه سیگما (برای مثال انتقال بیرنگ و

¹ Sigm-point

بو[111]) و یا تطبیق ممان 7 نیز محاسبه شود. این روشهای استنتاج رایج، توزیعهای پیشبینی کننده سنگین 7 را با گوسیها تقریب میزنند.

با فرض یک توزیع احتمالاتی گوسی توام با فرض یک توزیع احتمالاتی گوسی توام با فرض یک توزیع حالت، $p(x_t,u_t)=\mathcal{N}([x_t,u_t]|\mu_t^{xu},\Sigma_t^{xu})$ ورخ عصل توانع حالت بعدی $x_{t+1}=f(x_t,u_t)+w$ است.

$$p(x_{t+1}) = \iiint p(x_{t+1}|x_t, u_t) p(x_t, u_t) dx_t du_t dw$$
 (۲۵-۲) معادله

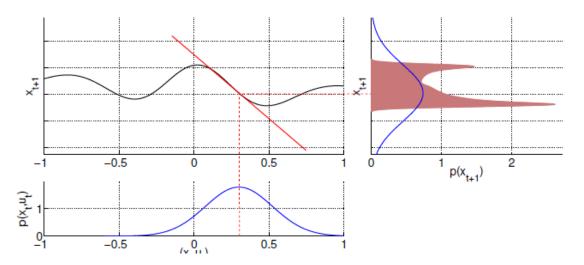
اگر تابع انتقال f غیرخطی باشد، $p(x_{t+1})$ غیرگوسی است و باید به تکنیکهای استنتاج تقریبی متوسل شویم. یک تقریب آسان از توزیع پیشبینی کننده سنگین $p(x_{t+1})$ ، گوسی $p(x_{t+1})$ ، گوسی $p(x_{t+1})$ است. میانگین شویم. یک تقریب آسان از توزیع پیشبینی کننده می توانند به روشهای مختلفی محاسبه شوند. در ادامه، μ_{t+1}^x و کوواریانس Σ_{t+1}^x این توزیع پیشبینی کننده می توانند به روشهای مختلفی محاسبه شوند. در ادامه، سه رویکردی که به طور رایج استفاده می شوند را به صورت خلاصه شرح می دهیم: خطی سازی، انتقال بی رنگ و بو و تطبیق ممان.

خطی سازی. یک روش محاسبه μ_{t+1}^x و μ_{t+1}^x این است که تابع انتقال $f \approx F$ را به صورت محلی حول خطی سازی. یک روش محاسبه کواریانس پیشبینی کننده را با نگاشت توزیع ورودی گوسی به واسطه سیستم خطی شده تخمین زد. با خطی سازی، کوواریانس و میانگین پیشبینی کننده را به ترتیب با رابطه های $\mu_{t+1}^x = F\Sigma_t^{xu}F^T + \Sigma_w$ و $f(\mu_t^{xu})$ و دست می آوریم. شکل ۲-۸ ایده خطی سازی را نشان می دهد.

¹ Unscented transformation

² Moment matching

³ Unwieldy



شکل ۲-۸ محاسبه توزیع پیشبینی کننده تقریبی با استفاده از خطیسازی. یک توزیع گوسی $p(x_t,u_t)$ (پنل سمت چپ پایین) باید به واسطه یک تابع غیرخطی (نمودار مشکی در پنل سمت چپ بالا) نگاشت شود. توزیع پیشبینی کننده درست با ناحیه سایه خورده در پنل راست نشان داده شده است. برای به دست آوردن یک تقریب گوسی از توزیع سایه خورده سنگین، تابع غیرخطی در میانگین توزیع ورودی خطیسازی شده است (پنل سمت چپ بالا، خط قرمز). سپس، این گوسی به واسطه این تقریب خطی نگاشت شده و توزیع پیشبینی کننده تقریبی گوسی $p(x_{t+1})$ حاصل شده که در پنل سمت راست به رنگ آبی نشان داده شده است. [90]

خطی سازی از لحاظ مفهومی سرراست و از لحاظ محاسباتی کاراست. توجه داشته باشید که این رویکرد، توزیع ورودی گوسی $p(x_t,u_t)$ را دست خورده می گذارد اما تابع انتقال f را تقریب می زند. یک اشکال بالقوه، این است که برای اجرای خطی سازی، تابع انتقال f باید مشتق پذیر باشد. علاوه بر این، خطی سازی به راحتی می تواند واریانس های پیش بینی کننده را کم تخمین بزند f می تواند باعث شود که سیاست ها بیش از حد مهاجم f باشند که باعث صدمه به سیستم های ربات واقعی می شود.

انتقال بیرنگ و بو. ایده کلیدی پشت انتقال بیرنگ و بو $p(x_t,u_t)$ بازنمایی توزیع $p(x_t,u_t)$ با مجموعهای از نقاط سیگما $(X_t^{[i]},\mathcal{U}_t^{[i]})$ است که به طور قطعی انتخاب شدهاند. برای این نقاط سیگما، مقادیر دقیق تابع

[ٔ] فرضیات مشتق پذیری می تواند در رباتیک مشکل ساز باشد. برای مثال، اتصالات در جابجایی و دستکاری می توانند این فرض را نادرست جلوه دهند.

² Underestimate

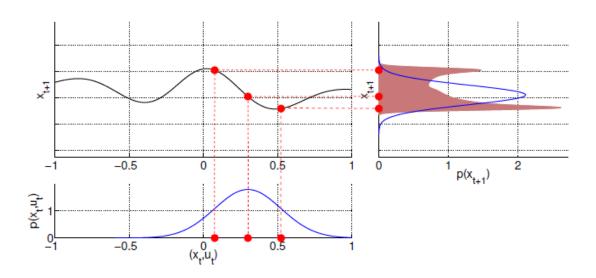
³ Aggressive

مربوطه محاسبه میشوند. میانگین μ_{t+1}^x و کوواریانس Σ_{t+1}^x توزیع پیشبینی کننده $p(x_{t+1})$ از نقاط سیگمای نگاشتشده وزندار محاسبه میشوند. میانگین و واریانس به ترتیب از رابطههای زیر به دست می آیند

$$\mu_{t+1}^{x} = \sum_{i=0}^{2d} w_m^{[i]} f(\mathcal{X}_t^{[i]}, \mathcal{U}_t^{[i]})$$
 معادله(۲۶-۲) معادله

$$\Sigma_{t+1}^{x} = \sum_{i=0}^{2d} w_c^{[i]} (f(X_t^{[i]}, \mathcal{U}_t^{[i]}) - \mu_{t+1}^{x}) (f(X_t^{[i]}, \mathcal{U}_t^{[i]}) - \mu_{t+1}^{x})^T$$

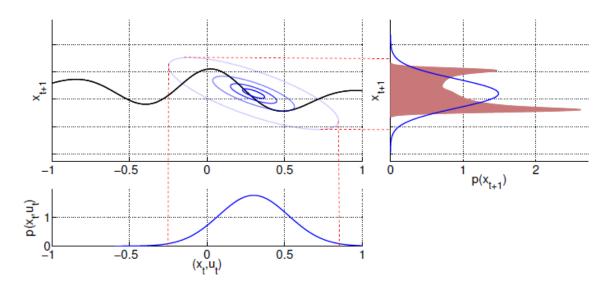
که در این روابط، p بعد (x,u) است، (x,u) است، (x,u) نقاط سیگما هستند به عبارت دیگر، نمونههایی از توزیع (x,u) بقار بیش و ب



شکل ۹-۲ محاسبه توزیع پیشبینی شده تقریبی با استفاده از انتقال بی رنگ و بو. یک توزیع گوسی $p(x_t,u_t)$ (پنل سمت چپ پایین) باید به واسطه یک تابع غیرخطی نگاشت شود (نمودار مشکلی در پنل سمت چپ بالا). توزیع پیشبینی کننده درست در ناحیه سایه خورده در پنل سمت راست نشان داده شده است. برای به دست آوردن یک تقریب گوسی از توزیع سایه خورده سنگین، توزیع ورودی با سه نقطه سیگما نشان داده شده است (نقاط قرمز در پنل سمت چپ پایین). سپس، نقاط سیگما به واسطه تابع غیرخطی نگاشت می شوند (پنل سمت چپ بالا) و میانگین و کوواریانس نمونه اشان، توزیع پیشبینی کننده تقریبی $p(x_{t+1})$ را نتیجه می دهد که در پنل سمت راست نشان داده شده است. [90]

باید توجه داشت که انتقال بیرنگ و بو توزیع گوسی $p(x_t, u_t)$ را با نقاط سیگما تقریب میزند، نقاطی که سپس به واسطه تابع انتقال اصلی f نگاشت میشوند. انتقال بیرنگ و بو به مشتق پذیری نیاز ندارد و انتظار میرود که نسبت به خطیسازی صریح، تقریبهای دقیق تری از توزیع پیش بینی کننده $p(x_{t+1})$ نتیجه دهد [112].

تطبیق ممان. ایده تطبیق ممان، محاسبه دقیق میانگین و کوواریانس پیشبینی کننده $p(x_{t+1})$ و تقریب تطبیق ممان. ایده تطبیق ممان، این میانگین و کوواریانس دقیق به آن تعلق دارد. در اینجا، نه توزیع توام $p(x_{t+1})$ با یک گوسی است که آن میانگین و کوواریانس دقیق به آن تعلق دارد. در اینجا، نه توزیع توام $p(x_{t+1})$ و نه توزیع انتقال f تقریب زده میشوند. تقریب تطبیق ممان، از این لحاظ بهترین تقریب تکمدی از توزیع پیشبینی کننده درست و تقریب تکمدی را کمینه می کند[113]. شکل ۲-۱۰ ایده تطبیق ممان را نشان می دهد.



شکل 1-t محاسبه توزیع پیشبینی شده تقریبی با استفاده از تطبیق ممان. یک توزیع گوسی $p(x_t,u_t)$ (پنل سمت چپ پایین) باید به واسطه یک تابع غیرخطی نگاشت شود (نمودار مشکی در پنل سمت چپ بالا). توزیع پیشبینی کننده درست با ناحیه سایه خورده در پنل راست نشان داده شده است. برای به دست آوردن یک تقریب گوسی از توزیع سایه خورده سنگین، میانگین و کوواریانس توزیع سایه خورده به صورت تحلیلی محاسبه می شوند. این ممانهای مر تبه اول و دوم، به طور کامل توزیع پیشبینی کننده تقریبی گوسی $p(x_{t+1})$ را معین می کنند که نمودار آن به رنگ آبی در ینل راست نشان پیشبینی کننده تقریبی گوسی $p(x_{t+1})$ را معین می کنند که نمودار آن به رنگ آبی در ینل راست نشان

داده شده است. خطوط ترازی که در پنل سمت چپ بالا نشان داده شدهاند توزیع توام بین ورودیها و پیشبینی را نشان میدهند.[90]

■ ملاحظات عملی

ممانهای دقیق فقط در موارد خاص می توانند محاسبه شوند چون انتگرالهای مورد نیاز برای محاسبه میانگین و کوواریانس پیشبینی کننده ممکن است لاینحل باشند. علاوه بر این، به طور معمول تقریب دقیق تطبیق ممان از لحاظ محاسباتی پرهزینه تر از تقریبها به وسیله خطی سازی یا نقاط سیگماست.

برخلاف رویکردهای مبتنی بر نمونهبرداری از قبیل PEGASUS، روشهای استنتاج تقریبی قطعی برای برنامه ریزی بلندمدت می توانند به منظور یادگیری چندین هزار پارامتر سیاست استفاده شوند[3]. دلیلی که پیش بینیهای بلندمدت قطعی می توانند سیاستهایی با تعداد زیادی پارامتر را یاد بگیرند این است که گرادیانها می توانند به صورت تحلیلی محاسبه شوند. بنابراین، این تخمینهای گرادیان از واریانسهای بزرگ رنج نمی برند که یک مشکل معمول در تخمین مبتنی بر نمونهبرداری است. با این وجود، غالبا استنتاج قطعی نسبت به رویکردهای نمونهبرداری به تلاش بیش تری برای پیاده سازی نیاز دارد.

۳-۳-۲-۲ بهروز رسانیهای سیاست

با معرفی دو دسته مدل بزرگ و دو روش کلی انجام پیشبینیهای بلندمدت با این مدلها، در ادامه در مورد روشهای به روز رسانی سیاست فارغ از گرادیان و مبتنی بر گرادیان را برمی شماریم.

- بهروز رسانیهای سیاست مبتنی بر مدل بدون اطلاعات گرادیان

روشهای فارغ از گرادیان احتمالا آسان ترین راه بهروز رسانی سیاست هستند چون به محاسبه و تخمین گرادیانهای سیاست نیازی ندارند. اصولا این روشها محدودیتهای مشتق پذیری روی سیاست یا مدل انتقال ندارند. روش بهینه سازی فارغ از گرادیان استاندارد عبارتند از روش Nelder-Mead [114] که یک روش هیوریستیک ساده است و یا تپهنوردی که یک روش جستجوی محلی است که به طور تنگاتنگ به تابکاری شبیه سازی شده آ[115]

-

¹ Hill-climbing

² Simulated annealing

مرتبط است. به خاطر سادگیشان و تلاش محاسباتی کمی که نیاز دارند، به طور رایج در زمینه جستجوی سیاست مبتنی بر مدل[107] و بهخصوص در ترکیب با تولید تراژکتوری مبتنی بر نمونهبرداری استفاده میشوند.

اشکال واضح بهینهسازی فارغ از گرادیان، نرخ همگرایی نسبتا کند آن است. برای همگرایی سریعتر، میتوانیم از بهروز رسانیهای سیاست مبتنی بر گرادیان استفاده کنیم که در بخشهای بعدی معرفی میشوند.

- بهروز رسانیهای سیاست مبتنی بر مدل با اطلاعات گرادیان

انتظار میرود که بهروز رسانیهای سیاست مبتنی بر گرادیان همگرایی سریعتری از بهروز رسانیهای فارغ از گرادیان داشته باشند. بین دو مورد تمایز قائل میشویم: تخمین مبتنی بر نمونه گرادیانهای سیاست و محاسبه تحلیلی گرادیانهای سیاست $dJ_{\theta}(\theta)/d\theta$.

■ گرادیانهای سیاست مبتنی بر نمونهبرداری

زمانی که از نمونه تراژکتوریهای $au^{[i]}$ از مدل آموخته شده استفاده می کنیم تا پاداش بلندمدت مورد انتظار $J_ heta$ در معادله(۲-۱۰) را تخمین بزنیم، می توانیم گرادیان $dJ_ heta/d heta$ را به صورت عددی تقریب بزنیم.

آسان ترین روش برای تخمین گرادیانها، استفاده از روشهای تفاضل متناهی (محدود) است. اما، روشهای تفاضل متناهی به تعداد O(F) ارزیابی پاداش بلندمدت مورد انتظار نیاز دارند که در آن F تعداد پارامترهای سیاست O(F) است. از آنجایی که هر یک از این ارزیابیها مبتنی بر میانگین F نمونه roll-out است، تعداد نمونه تراژکتوریهای مورد نیاز بیش از حد می شود. در تنظیم مبتنی بر مدل، این فقط یک مشکل محاسباتی است اما مشکلی برای فرسوده کردن ربات نیست چون نمونه ها از مدل تولید می شوند نه از خود ربات.

چندین روش برای کاراتر کردن تخمین گرادیان مبتنی بر مدل وجود دارد که عبارتند از: اولا، برای تخمین قوی تر از J_{θ} یا به عبارت دیگر تخمینی با واریانس کوچکتر، رویکرد PEGASUS [103] میتواند استفاده شود. ثانیا، برای تخمین کاراتر گرادیان، هر کدام از روشهای فارغ از مدل برای تخمین گرادیان میتوانند در زمینه مبتنی بر مدل استفاده شوند. تنها تفاوت این است که به جای ربات، از مدل آموخته شده برای تولید تراژ کتوریها استفاده می شود.

■ گرادیانهای تحلیلی سیاست

محاسبه تحلیلی گرادیانهای $dJ_{\theta}/d\theta$ نیاز دارد که سیاست، (امید ریاضی) تابع پاداش و مدل انتقال آموخته شده مشتق پذیر باشند. با وجود این محدودیت، محاسبات تحلیلی گرادیان به دو دلیل، جایگزینی پایا برای گرادیانهای محاسبه مبتنی بر نمونه برداری هستند: اولا، از واریانس نمونه برداری رنج نمی برند، چیزی که به خصوص هنگام محاسبه گرادیانها نمود پیدا می کند. ثانیا، به طور مطلوب تلاش محاسباتی با تعداد پارامترهای سیاست مقیاس می شود که این موضوع امکان یادگیری سیاستها با هزاران پارامتر را می دهد. اما به خاطر اعمال مکرر قاعده زنجیره ای نالیا محاسبه خود گرادیان از لحاظ ریاضیاتی پیچیده تر از یک تخمین مبتنی بر نمونه برداری است.

یک مثال را در نظر بگیرید که پاداش آنی r فقط به حالت بستگی دارد (تعمیم به کنترلهای وابسته به پاداش $x_{t+1}=f(x_t,u_t)=f(x_t,\pi_{\theta}(x_t,\theta))$ سرراست است) و دینامیک سیستم قطعی است، به طوری که $\pi_{\theta}(x_t,u_t)=f(x_t,u_t)=f(x_t,u_t)=f(x_t,u_t)$ در آن f یک تابع انتقال (غیرخطی)، π_{θ} سیاست (قطعی) و π_{θ} پارامترهای سیاست هستند. گرادیان پاداش بلندمدت $\pi_{\theta}(x_t,u_t)=f(x_t,u_t)$ نسبت به پارامترهای سیاست با اعمال مکرر قاعده زنجیرهای به دست می آید:

$$\begin{split} \frac{dJ_{\theta}}{d\theta} &= \sum_{t} \gamma^{t} \frac{dr(x_{t})}{d\theta} = \sum_{t} \gamma^{t} \frac{dr(x_{t})}{dx_{t}} \frac{dx_{t}}{d\theta} \\ &= \sum_{t} \gamma^{t} \frac{dr(x_{t})}{dx_{t}} \left(\frac{\partial x_{t}}{\partial x_{t-1}} \frac{\partial x_{t-1}}{d\theta} + \frac{\partial x_{t}}{\partial u_{t-1}} \frac{\partial u_{t-1}}{d\theta} \right) \end{split} \tag{7.4-7}$$
 معادله(۲۹-۲)

با توجه به این معادلات مشاهده می کنیم که مشتق کلی $dx_t/d\theta$ به مشتق کلی $dx_{t-1}/d\theta$ در گام زمانی قبلی با توجه به این معادلات مشتق $dJ_{\theta}/d\theta$ می تواند با تکرار محاسبه شود.

تعمیم به مدلهای احتمالاتی و MDP های تصادفی. برای بسط مشتقها به MDP های تصادفی و/یا مدلهای تعمیم به مدلهای احتمالاتی و MDP های تصادفی و/یا مدلهای که حالت احتمالاتی، باید تغییرات کمی در گرادیانهای معادلات معادلات معادله (۲۸-۲) معادله (۲۹-۲) ایجاد کنیم: زمانی که حالت $\mathbb{E}[r(x_t)] = \mathbb{E}[r(x_t)] = \mathbb{E}[r(x_t)] = \mathbb{E}[r(x_t)] = \mathbb{E}[r(x_t)]$ با یک توزیع احتمالاتی $p(x_t)$ نشان داده میشود باید امید ریاضی پاداشرهای توزیع حالت محاسبه کنیم، با این فرض که $p(x_t)$ بازنمایی پارامتری دارد.

 Σ_t^x برای مثال، اگر μ_t^x , Σ_t^x و کوواریانس، مشتقات $\mathbb{E}[r(x_t)]$ را نسبت به میانگین $p(x_t) = \mathcal{N}(x_t | \mu_t^x, \Sigma_t^x)$ و کوواریانس $\mathbb{E}[r(x_t)]$ دامه می دهیم: $\mathbb{E}[r(x_t)]$ دامه می دهیم: $\mathbb{E}[r(x_t)]$ دامه می دهیم: با تعریف $\mathcal{E}_t := \mathbb{E}_{x_t}[r(x_t)]$ گرادیان را به صورت زیر به دست می آوریم

$$\begin{split} \frac{dJ_{\theta}}{d\theta} &= \sum_{t} \gamma^{t} \frac{d\varepsilon_{t}}{d\theta} \ , \\ \frac{d\varepsilon_{t}}{d\theta} &= \frac{d\varepsilon_{t}}{dp(x_{t})} \frac{dp(x_{t})}{d\theta} := \frac{d\varepsilon_{t}}{d\mu_{t}^{x}} \frac{d\mu_{t}^{x}}{d\theta} + \frac{d\varepsilon_{t}}{d\Sigma_{t}^{x}} \frac{d\Sigma_{t}^{x}}{d\theta} \end{split} \tag{T--T}$$
 معادله (

 \mathcal{E}_t (کلی) $d\mathcal{E}_t/dp(x_t) = \{d\mathcal{E}_t/d\mu_t^x \ , d\mathcal{E}_t/d\Sigma_t^x \}$ برای گرفتن مشتق (کلی) که در آن، از نمادگذاری مختصر $p(x_t) = \mathcal{N}(x_t | \mu_t^x, \Sigma_t^x)$ برای گرفتن مشتق (کلی) $p(x_t) = \mathcal{N}(x_t | \mu_t^x, \Sigma_t^x)$ استفاده کردهایم. میانگین \mathcal{L}_t^x و کوواریانس) استفاده کردهایم کنترل کننده میانگین \mathcal{L}_t^x و کوواریانس \mathcal{L}_t^x از لحاظ کاربردی به ممانهای \mathcal{L}_t^x و \mathcal{L}_t^x و پارامترهای \mathcal{L}_t^x کنترل کننده وابسته هستند. با اعمال قاعده زنجیرهای به معادله (۲۰-۲) به دست می آوریم:

$$\frac{d\mu_t^x}{d\theta} = \frac{d\mu_t^x}{d\mu_{t-1}^x} \frac{d\mu_{t-1}^x}{d\theta} + \frac{d\mu_t^x}{d\Sigma_{t-1}^x} \frac{d\Sigma_{t-1}^x}{d\theta} + \frac{d\mu_t^x}{d\theta}$$

$$\frac{d\Sigma_t^x}{d\theta} = \frac{d\Sigma_t^x}{d\mu_{t-1}^x} \frac{d\mu_{t-1}^x}{d\theta} + \frac{d\Sigma_t^x}{d\Sigma_{t-1}^x} \frac{d\Sigma_{t-1}^x}{d\theta} + \frac{d\Sigma_t^x}{d\theta}$$
 (٣٢-٢)عادله

توجه داشته باشید که مشتقات کلی $d\mu_{t-1}^x/d heta$ و $d\Sigma_{t-1}^x/d heta$ از گام زمانی t-1 معلوم هستند.

اگر همه این محاسبات بتوانند به فرم بسته انجام شوند، گرادیانهای سیاست $d\tilde{J}_{\theta}/d\theta$ می توانند با اعمال مکرر قاعده زنجیرهای به صورت تحلیلی و بدون نیاز به نمونهبرداری محاسبه شوند. بنابراین، تکنیکهای استاندارد بهینهسازی (برای مثال BFGS یا CG) می توانند برای یادگیری سیاستها با هزاران پارامتر استفاده شوند[3].

- مباحثه

استفاده از گرادیانهای امید ریاضی پاداش بلندمدت J_{θ} نسبت به پارامترهای سیاست θ غالبا منجر به یادگیری سریع تر از بهروز رسانیهای سیاست فارغ از گرادیان می شود. علاوه بر این، روشهای فارغ از گرادیان به طور معمول محدود به دهها پارامتر سیاست هستند. محاسبه گرادیانها می تواند سنگین باشد و به منابع محاسباتی اضافی نیاز دارد. زمانی که گرادیانها را محاسبه می کنیم، گرادیانهای تحلیلی دقیق به گرادیانهای مبتنی بر نمونهبرداری ترجیح داده می شوند چون مورد دوم غالبا از واریانس بزرگ رنج می برد. این واریانسها می توانند حتی منجر به

همگرایی کندتر از بهروز رسانیهای سیاست فارغ از مدل[107] شوند. برای گرادیانهای تحلیلی، فرضیاتی برای مشتق پذیری تابع پاداش f و تابع انتقال f در نظر می گیریم. علاوه بر این، برای گرادیانهای تحلیلی، فقط به روشهای استنتاج تقریبی قطعی از قبیل تطبیق ممان یا خطیسازی متکی هستیم به طوری که فقط تقریب J_{θ} به J_{θ} می تواند محاسبه شود؛ اما فقط با گرادیانهای دقیق $d\tilde{J}_{\theta}/d\theta$.

برای بهروز رسانی سیاست، استفاده از اطلاعات گرادیان را توصیه می کنیم تا ویژگیهای بهتر همگرایی استفاده شوند. به طور ایده آل، گرادیانها به صورت تحلیلی و بدون هیچ تقریبی معین می شوند. از آنجایی که فقط برای سیستمهای خطی می توان به این هدف رسید باید به تقریبها (با استفاده از رویکردهای مبتنی بر نمونهبرداری یا گرادیانهای تقریبی تحلیلی) متوسل شویم. رویکردهای مبتنی بر نمونهبرداری عملا محدود به پارامترهای سیاست نسبتا با ابعاد پایین (k > 50)، استفاده از گرادیانهای پایین (k > 50)، استفاده از گرادیانهای سیاست تحلیلی را (در صورتی که در دسترس باشند) توصیه می کنیم.

۴-۲-۲-۲ الگوریتمهای جستجوی سیاست مبتنی بر مدل با کاربردهای رباتیکی

در این بخش، به صورت خلاصه روشهای جستجوی سیاستی را شرح می دهیم که با موفقیت به یادگیری سیاستها برای رباتها اعمال شدهاند. تمایز بین رویکردهایی که امید ریاضی پاداش بلندمدت J_{θ} را با استفاده از روشهای نمونه برداری ارزیابی می کنند و رویکردهایی که آن را با استفاده از روشهای استنتاج تقریبی قطعی ارزیابی می کنند شرح می دهیم.

- پیشبینی مبتنی بر نمونهبرداری تراژکتوری

نمونهبرداری مستقیم از شبیهساز آموختهشده، توسط یک سری از محققان برای حرکت ٔ هلی کوپترها[107] و برای کنترل بالونها[107] مورد بررسی قرار گرفته است. همه رویکردها از الگوریتم PEGASUS [103] استفاده کردهاند تا تراژکتوریها را از مدلهای تصادفی آموختهشده تولید کنند.

-

¹ Maneuvering

انجی و همکاران[106]، مدلهایی برای حرکت در جای هلیکوپتر یاد گرفتهاند که مبتنی بر رگرسیون خطی وزندار محلی است. برای در نظر گرفتن نویز و بیدقتیهای مدل، با اضافه کردن نویز (سیستمی) گوسی i.i.d به دینامیک انتقال، این مدل اصالتا قطعی، تصادفی شده است.

برخلاف مقاله [106]، باگنل و اشنایدر [107] عدم قطعیت در مورد مدل آموخته شده را به وسیله یک توزیع پسین روی مجموعهای متناهی از مدلهای affine محلی به طور صریح شرح داده اند. برای یادگیری سیاست، تراژکتوریها از این ترکیب مدلها نمونه برداری شده اند.

کو و همکاران[105]، مدلهای پارامتری ایدهآلشده را با فرآیندهای گوسی غیرپارامتری ترکیب کردهاند تا دینامیک بالون خودمختار را مدلسازی کنند. GP ها برای مدلسازی اختلاف بین مدل پارامتری غیرخطی بالون و دادهها استفاده شدهاند. هنگام یادگیری سیاست، تراژکتوریها از این مدل آمیخته نمونهبرداری میشوند. در ادامه، در مورد این الگوریتمهای جستجوی سیاست بحث میکنیم.

■ مدلهای پیشرو رگرسیون وزندار محلی و پیشبینی مبتنی بر نمونهبرداری تراژکتوری

در مرجع [107]، رگرسیون بیزی وزندار محلی به منظور یادگیری مدلهای پیشرو برای حرکت در جای یک هلی کوپتر خودمختار استفاده شده است. برای در نظر گرفتن عدم قطعیت مدل، به جای تخمین نقطهای از پارامترهای مدل، یک توزیع پسین احتمالاتی روی پارامترهای ψ مدل و از این رو، روی خود مدل در نظر گرفته شده است. برای یادگیری سیاست ، تراژکتوریها از این ترکیب مدلها نمونهبرداری شدهاند.

تراژکتوریها $\tau^{[i]}$ با استفاده از رویکرد PEGASUS [103] PEGASUS تولید شده اند. در هر گام زمانی، یک مجموعه پارامتر ψ_i مدل از توزیع پسین $p(\psi|X,U,y,x_t,u_t)$ نمونه برداری شده است. بعد از هر انتقال، مدل دینامیکی با انتقال (شبیه سازی شده) مشاهده شده به روز رسانی می شود. بعد از هر تراژکتوری تولید شده $\tau^{[i]}$ ، با حذف تراژکتوری شبیه سازی شده $\tau^{[i]}$ از مدل بازنشانی (راهاندازی مجدد) می شود [107]. برای تخمین گرهای مدل بهینه بیز $\tau^{[i]}$ این رویه معادل با نمونه برداری مدل و نمونه برداری یک تراژکتوری کامل از آن است. الگوریتم $\tau^{[i]}$ به طور مختصر این رویه معادل با نمونه برداری مدل و نمونه برداری یک تراژکتوری کامل از آن است. الگوریتم $\tau^{[i]}$ به طور مختصر

¹ Hybrid

² Bayes-optimal

بیان کرده است که باید چگونه تراژکتوریها را از مدل آموخته شده نمونه برداری کرد در حالی که عدم قطعیت پسین در مورد خود مدل تلفیق شده است. با میانگین گرفتن روی امید ریاضی پاداشهای بلندمدت برای همه تراژکتوریهای تولید شده $au^{[i]}$ عدم قطعیت مدل به صورت ضمنی ادغام می شود.

```
Input: transition model f, posterior distribution over model parameters p(\psi|X,U,y), policy parameters \theta for i=1,\ldots,m do  \text{for } t=0,\ldots,T-1 \text{ do} \qquad \text{Sample trajectory } \tau^{[i]}  Sample local model parameters \psi_i \sim p(\psi|X,U,y,x_t,u_t) Compute control u_t=\pi_\theta(x_t) Generate a sample state transition x_{t+1} \sim p(x_{t+1}|x_t,u_t,\psi_i) Update X, U with simulated transition (x_t,u_t,x_{t+1}) end for Compute J_{\theta,i} Reset the learned model to the original model f end for \tilde{J}_\theta = \frac{1}{m} \sum_{i=1}^m J_{\theta,i}
```

الگوریتم ۲-۶ الگوریتم ارزیابی سیاست و پیشبینیهای T-گامی[107].

در مرجع [107]، روش Nelder-Mead که یک بهینهسازی فارغ از گرادیان است، برای بهروز رسانی پارامترهای θ سیاست استفاده شده است که از بهینهسازی مبتنی بر گرادیان ساده عملکرد بهتری داشته است. رویکرد حاصله از یک کنترل کننده شبکه عصبی با ده پارامتر استفاده کرده تا یک هلی کوپتر را حول یک نقطه ثابت حرکت دهد[107]، شکل ۲-۱۱(الف). با جریمههای زیاد روی حالتهای مربوطه، از برون یابی خارج از محدوده داده جمع آوری شده جلوگیری شده است. این کنترل کننده آموخته شده که مبتنی بر توزیع احتمالاتی روی مدل هاست، به طور چشم گیر کم نوسان تر از کنترل کننده ای استفاده از مدل بیشینه درست نمایی و یا به عبارت دیگر تخمین نقطه ای از پارامترهای مدل، آموخته شده است.





(ب) حرکت در جای هلی کویتر وارونه[107]

(الف) حركت درجاي هلي كوپتر [106]

شکل ۲-۱۱ روشهای جستجوی سیاست مبتنی بر مدل با استنتاج تصادفی که برای یادگیری حرکت در استفاده شدهاند.[90]

انجی و همکاران[106]، مدلهایی مبتنی بر رگرسیون خطی وزندار محلی برای حرکت درجای هلی کوپتر یاد گرفته اند، شکل ۲-۱۱(ب). برخلاف مرجع [107]، تخمین نقطهای از پارامترها للا (برای نمونه با بیشینه درستنمایی یا تخمین بیشینه پسین) در معادله(۲-۱۱) تعیین شده است. برای در نظر گرفتن نویز و بی دقتیهای مدل، با اضافه کردن نویز (سیستمی) گوسی i.i.d به دینامیک انتقال، این مدل اصالتا قطعی، تصادفی شده است. سرعتهای زاویهای که در مختصات هلی کوپتر بیان شدهاند انتگرال گرفته شده و سپس به زوایایی در مختصات جهانی منتقل شدهاند، کاری که مدل را غیرخطی می کند. با این رویکرد، مدلهای مربوط به حرکت درجای هلی کوپتر در وضعیت استاندارد[104] یا وارونه[106]، با استفاده از داده جمعآوری شده از تراژکتوریهای (انسان) خلبانها تعیین شدهاند.

برای یادگیری یک کنترلکننده با این دینامیک انتقال غیرخطی تصادفی، به منظور نمونهبرداری تراژکتوریها از مدل، از روش نمونهبرداری PEGASUS [103] استفاده شده است. با این تراژکتوریهای نمونهبرداری شده، یک تخمین Monte Carlo از امید ریاضی پاداش بلندمدت محاسبه شده است. یک روش تپهنوردی حریصانه برای یادگیری پارامترهای θ سیاست π_{θ} استفاده شده است که با یک شبکه عصبی ساده شده نشان داده شده است π_{θ} استفاده شده است که با یک شبکه عصبی ساده شده نشان داده شده است آلفت [106].

در مورد حرکت در جای هلی کوپتر وارونه، یک (انسان) خلبان هلی کوپتر را وارونه کرده است. سپس، کنترل کننده آموخته شده کنترل کننده آموخته کنترل را در دست گرفته و هلی کوپتر را در موقعیت وارونه پایدار کرده است[106] که مثالی از آن در شکل ۲-۱۱ (ب) نشان داده شده است.

■ مدلهای پیشرو فرآیند گوسی و پیشبینی مبتنی بر نمونهبرداری تراژکتوری

در مرجع [105]، مدلهای پیشرو GP برای مدلسازی دینامیک انحراف بالون خودمختار یاد گرفته شدهاند، هم GP با مدل پارامتری ایدهآلشده دینامیک بالون ترکیب شدهاند، به عبارت دیگر GP شکل ۲-۱۲. مدلهای GP با مدل پارامتری ایدهآلشده دینامیک بالون ترکیب شدهاند، به عبارت دیگر اختلاف بین مدل روی داده تراژکتوری اختلاف بین مدل غیرخطی پارامتری و داده مشاهدهشده را مدلسازی میکند. این مدل روی داده تراژکتوری بالون آموزش داده شده است که به وسیله انسانی تولید شده که بالون را با استفاده از کنترل از راه دور به پرواز در آورده است.

رویکرد PEGASUS ایرای نمونهبرداری تراژکتوریهای بلندمدت استفاده شده است. هر نمونه جدید از طریق بهروز رسانی ماتریس کرنل در مدل آمیخته شده است. کنترلکننده با استفاده از روش بهینهسازی فارغ طریق بهروز رسانی ماتریس کرنل در مدل آمیخته شده است. چهار پارامتر θ کنترلکننده عبارتند از ضریب کشش بهرههای موتور راست/چپ و شیب تابع هموارکننده سیاست. کنترلکننده آموخته شده یک کنترلکننده حلقه باز است، به عبارت دیگر کنترلها برونخط و از قبل محاسبه شده و سپس به بالون اعمال شدهاند. کنترلکننده مبتنی بر دینامیکی پارامتری هایده آلشده بالون است[105] عملکرد بهتری دارد.

¹ Yaw-dynamics

² Drag coefficient

³ Offline



شکل ۲-۱۲ ترکیب مقدم پارامتری و فرآیندهای گوسی برای مدلسازی و یادگیری کنترل یک بالون خودمختار.[105]

- پیشبینیهای قطعی تراژکتوری

در ادامه، به طور مختصر روشهای جستجوی سیاستی را شرح میدهیم که پیشبینیهای قطعی تراژکتوری را برای ارزیابی سیاست انجام میدهند.

■ مدلهای پیشرو فرآیند گوسی و پیشبینی قطعی تراژکتوری

چارچوب کاری جستجوی سیاست PILCO (استنتاج احتمالاتی برای یادگیری کنترل)[3]، از یک مدل پیشرو GP مربوط به دینامیک ربات استفاده کرده تا به طور مداوم خطاهای مدل را در نظر بگیرد. در ترکیب با استنتاج قطعی به وسیله تطبیق ممان برای پیشبینی تراژکتوریهای حالت $J_{\theta}(t)=(p(x_1),\cdots,p(x_T))$ و محالیه تقریب تحلیلی J_{θ} به امید ریاضی پاداش بلندمدت J_{θ} در معادله (۲۰-۲) را محاسبه می کند. علاوه بر این، گرادیانهای $J_{\theta}/d\theta$ امید ریاضی پاداش بلندمدت نسبت به پارامترهای سیاست به صورت تحلیلی محاسبه می شوند. برای یادگیری سیاست، بهینه سازهای استاندارد می توانند استفاده شوند. الگوریتم OPILCO در الگوریتم به طور معمول بدون اطلاعات مفید مقداردهی اولیه می شود، J_{θ} به طور مختصر بیان شده است. این الگوریتم به طور معمول بدون اطلاعات مفید مقداردهی اولیه می شود،

به عبارت دیگر پارامترهای سیاست به طور تصادفی نمونهبرداری میشوند و داده از یک تراژکتوری حالت کوتاه که با اعمال عملهای تصادفی تولید شده است ضبط میشود. در ادامه، به طور مختصر جزئیات محاسبه پیشبینیهای بلندمدت، گرادیانهای سیاست و اعمال چارچوب کاری PILCO به سیستمهای کنترلی و رباتیکی را بیان میکنیم.

```
Init: Sample controller parameters \theta \sim \mathcal{N}(0, I). Apply random control signals and record data. repeat

Learn probabilistic (GP) dynamics model using all data. repeat

Compute p(x_1), \dots, p(x_T) using moment matching and \tilde{f}_{\theta}

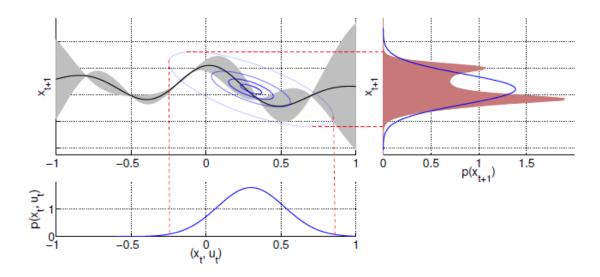
Analytically compute policy gradients \mathrm{d}\tilde{f}_{\theta}/\mathrm{d}\theta

Update parameters \theta (line-search in BFGS or CG). until convergence; return \theta^*

Apply \pi_{\theta}^* to system (single trial/episode) and record data. until task learned
```

الگوريتم ۲-۷ چارچوب كارى جستجوى سياست PILCO [3].

پیشبینیهای بلندمدت. برای پیشبینی یک توزیع $p(\tau|\pi_{\theta})$ روی تراژکتوریها برای یک سیاست معلوم، پیشبینیهای بلندمدت. برای پیشبینیها، عدم قطعیت $p(x_1), \cdots, p(x_T), \cdots, p(x_T)$ مکررا توزیعهای حالت $p(x_1), \cdots, p(x_T)$ را محاسبه می کند. برای این پیشبینیها، عدم قطعیت پسین در مورد مدل پیشرو $p(x_t, u_t)$ آموختهشده ادغام میشود. شکل ۲-۱۳ این سناریو را نشان میدهد: فرض می کنیم که یک توزیع توام گوسی $p(x_t, u_t)$ داده شده است. برای پیشبینی توزیع $p(x_t, u_t)$ حالت بعدی، توزیع توام $p(x_t, u_t)$ در پنل سمت چپ-پایین باید به واسطه توزیع $p(x_t, u_t)$ پسین روی تابع انتقال پنهان نگاشت شود که در پنل سمت راست-بالا نشان داده شده است. به خاطر تابع کوواریانس غیرخطی، استنتاج دقیق لاینحل است. نتیجه نمونهبرداری Monte Carlo گسترده، یک تقریب نزدیک به توزیع پیشبینی کننده است که با توزیع دو مدی سایهخورده در پنل سمت راست نشان داده شده است. PILCO، میانگین و واریانس این توزیع سایه خورده را دقیقا محاسبه کرده و همان طور که با توزیع آبی در پنل سمت راست-بالا نشان داده شده است، توزیع سایه خورده را با یک گوسی با میانگین و واریانس صحیح تقریب میزند [3].

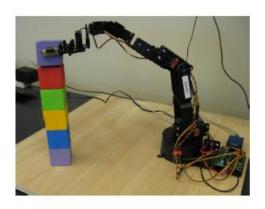


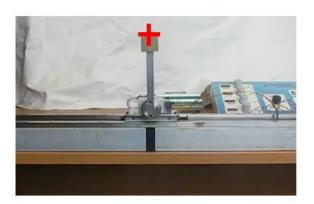
شکل ۱۳-۲ پیشبینیهای تقریبی با فرآیندهای گوسی در ورودیهای نامطمئن. برای تعیین توزیع پیشبینیکننده $p(x_t,u_t)$ و نیاز است که توزیع ورودی $p(x_t,u_t)$ (پنل سمت چپ پایین) را به واسطه توزیع GP پسین (پنل سمت چپ بالا) نگاشت کرد در حالی که صریحا عدم قطعیت مدل را میانگین گرفت (ناحیه سایه خورده). نمونهبرداری جامع Monte-Carlo توزیع دقیق را نتیجه می دهد که با توزیع به رنگ قرمز سایه خورده نشان داده شده است (پنل سمت راست). تقریب تطبیق ممان قطعی، میانگین و و اریانس توزیع پیشبینی کننده دقیق را محاسبه کرده و یک گوسی (منحنی آبی در پنل سمت راست) با دو ممان اول دقیق به آن برازش می کند. منحنی های تراز در پنل سمت چپ بالا توزیع توام بین ورودی ها و پیشبینی را نشان می دهند.

گرادیانهای سیاست تحلیلی. حالتهای پیشبینی شده تخمینهای نقطهای هستند اما با توزیعهای $dJ_{\theta}/d\theta$ سیاست تحلیلی و بیشبینی شده میشوند. هنگام محاسبه گرادیانهای سیاست $p(x_t)$, $t=1,\cdots,T$ محاسبی گوسی PILCO صریحا از طریق محاسبه تحلیلی گرادیانهای سیاست برای مدلهای احتمالاتی، فرمول بندی احتمالاتی را در نظر می گیرد.

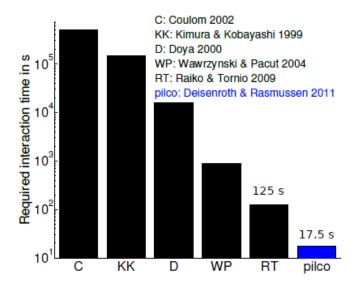
به خاطر گنجاندن صریح عدم قطعیت مدل در پیشبینیهای بلندمدت و محاسبه گرادیان، به طور معمول PILCO خیلی زیاد از خطاهای مدل رنج نمیبرد.

کاربردهای ربات. همانطور که در شکل ۲-۱۵ نشان داده شده است، الگوریتم PILCO به سرعتی بیسابقه در یادگیری روی یک کار محک استاندارد (کار under-actuated نوسان و متعادلسازی cart-pole که در شکل ۱۴-۲ نشان داده شدهاند) رسیده است.





شکل ۱۴-۲ موفقیتهای یادگیری PILCO. سمت چپ: یادگیری خودمختار برای روی هم قرار دادن دستهای از بلوکها با استفاده از بازوی ارزان موجود در بازار [108]. سمت راست: یادگیری خودمختار نوسان و متعادل کردن یک آونگ آزادانه معلق که به یک ارابه متصل است[3].



شکل ۱۵-۲ مقایسه عملکرد PILCO با روشهای دیگر. PILCO به سرعتی بیسابقه در یادگیری کار نوسان cart-pole میرسد. محور افقی به رویکردهای یادگیری تقویتی مراجعه می کند که کار مشابه را حل کردهاند، محور عمودی زمان تعامل مورد نیاز در واحد ثانیه با مقیاس لگاریتمی را نشان می دهد. [90]

به طور ویژه، نوسان cart-pole یاد گرفته شده است که از هر روش یادگیری تقویتی دیگری که از هیچ (یا به عبارت دیگر بدون مقداردهی اولیه با اطلاعات مفید مثلا با نمایشها) یاد می گیرد به زمان که از هیچ (یا به عبارت دیگر بدون مقداردهی اولیه با اطلاعات مفید مثلا با نمایشها) یاد می گیرد به زمان تعامل خیلی کم تری با ربات نیاز دارد. برای مسئله نوسان cart-pole، سیاست غیرخطی آموخته شده یک شبکهی تابع پایه شعاعی با ۵۰ تابع پایه گوسی axis-aligned است. پارامترهای θ سیاست شامل وزنها، موقعیتها و عرضهای توابع پایه هستند که در مجموع π پارامتر سیاست می شوند.

همچنین، PILCO در یک کار روی هم قرار دادن بلوک توسط یک بازوی ربات ارزان با ۵ درجه آزادی، برای یادگیری کارا و موثر کنترل کننده ها از هیچ با موفقیت اعمال شده است[108]، شکل ۲-۱۴. حالت x سیستم یادگیری کارا و موثر کنترل کننده ها از هیچ با موفقیت اعمال شده است. برای دنبال کردن این به صورت مختصات سه بعدی بلوک در محرک (مجری) پایانی بازو تعریف شده است. برای دنبال کردن این مختصات، از یک دوربین RGB-D استفاده شده است. سیاست آموخته شده π_{θ} یک تابع affine از حالت مختصات، از یک دوربین u = Ax + b با پیروی دقیق از الگوریتم ۲-۷، PILCO وی هم قرار دادن برجی از شش بلوک را در کمتر از ۲۰ آزمایش یاد گرفته است. محدودیتهای فضای حالت برای دوری از مانع نیز مستقیما در فرآیند یادگیری گنجانده شده است[108].

- مرور الگوریتمهای جستجوی سیاست مبتنی بر مدل

جدول ۲-۲، رویکردهای جستجوی سیاست مبتنی بر مدلی را خلاصه کرده است که در این بخش ارائه شدهاند. هر الگوریتم مطابق با روش پیشبینی(نمونهبرداری/ قطعی)، مدل پیشرو آموخته شده(LWR/LWBR/GP)، به روز رسانیهای سیاست(فارغ از گرادیان/ مبتنی بر گرادیان) و کاربردهای رباتیکی مربوطهاش فهرست شده است.

باید توجه داشت که در مقاله [106]، خطاهای مدل و کمینه محلی با وارد کردن نویز اضافی به سیستم مدیریت شدهاند تا خطر بیشبرازش کاهش یابد. به طور کلی، نویز در سیستم (چه با تزریق مصنوعی آن[106] و چه با استفاده از مدلهای احتمالاتی[3])، تابع هدف و از همینرو کمینه محلی را هموار می کند.

-

¹ end-effector

جدول ۲-۲ مرور الگوریتمهای جستجوی سیاست مبتنی بر مدل با کاربردهای رباتیکی[90].

كاربرد	بەروز	مدل	پیشبینیها	الگوريتم
	رسانی سیاست	پیشرو		
حرکت در جای	فارغ از گرادیان	LWBR	نمونهبرداری	[107]
هلی کوپتر			(PEGASUS)	
حرکت در جای	فارغ از گرادیان	+ LWR	نمونهبرداري	[106][104]
هلی کوپتر		نويز	(PEGASUS)	
كنترل بالون	فارغ از گرادیان	GP	نمونهبرداري	[103]
			(PEGASUS)	
بازوی ربات، -cart	مبتنی بر	GP	تطبيق ممان	[3][108]
pole	گرادیان			

۵-۳-۲-۲ ویژگیهای مهم روشهای مبتنی بر مدل

در ادامه، در مورد سه موضوع مهم که به جستجوی سیاست مبتنی بر مدل مرتبط هستند بحث می کنیم. به ویژه، ابتدا در مورد فواید و معایب استنتاج تصادفی در برابر استنتاج قطعی بحث می کنیم. سپس، در مورد این بحث می کنیم که چگونه در منابع علمی به عدم قطعیت در مورد مدل آموخته شده پرداخته می شود. سرانجام، نیازمندی ها را مشخص می کنیم در هنگامی که یک سیاست از هیچ یاد گرفته می شود یا به عبارت دیگر یادگیری باید بدون مقداردهی اولیه مناسب رخ دهد. نکته آخر درصورتی مهم است که نه دانش مفید در مورد دینامیک در دسترس باشد و نه مجموعه داده مناسب از نمایش ها. در عوض، ربات باید یاد بگیرد که از داده بالقوه پراکنده و دانش پیشین (مقدم) بی فایده شروع کند.

- پیشبینیهای بلندمدت تصادفی و قطعی

در مورد دو رویکرد کلی مبتنی بر مدل برای محاسبه توزیعها روی تراژکتوریها و پاداش بلندمدت مربوطه بحث کردیم: نمونهبرداری Monte-Carlo از ترفند PEGASUS استفاده کرده و پیشبینیهای قطعی از

خطی سازی، انتقال بی رنگ و بو یا تطبیق ممان استفاده می کنند. فایده نمونه برداری تصادفی این است که نمونه بردار، تخمینی درست از امید ریاضی پاداش بلندمدت J_{θ} در حد تعداد نامتناهی از تراژ کتوری های نمونه برداری شده باز خواهد گرداند. نمونه برداری جامع می تواند از لحاظ محاسباتی کارا نباشد، اما به صورت سر راست می تواند موازی شود. یک معضل مهم تر در رابطه با نمونه برداری، حتی هنگام استفاده از رویکرد راست می تواند موازی است که این روش فقط برای چند ده پارامتر سیاست عملی و کاربردی است.

به عنوان یک جایگزین برای نمونهبرداری تصادفی، پیشبینیهای قطعی فقط یک توزیع تراژکتوری دقیق برای سیستمهای خطی- گوسی را محاسبه می کنند. بنابراین، در سیستمهای غیرخطی، فقط تقریبی از امید ریاضی پاداش بلندمدت بازگردانده میشود. محاسباتی که برای محاسبه توزیعهای پیشبینی کننده نیاز است غیربدیهی هستند و میتوانند از لحاظ محاسباتی پرهزینه باشند. برخلاف نمونهبرداری تصادفی، پیشبینیهای قطعی به صورت سرراست قابل موازی کردن نیستند. از طرفی دیگر، پیشبینیهای قطعی فواید متعددی دارند که میتواند بر معایبش بچربد: اولا، با وجود این حقیقت که پیشبینیهای قطعی از لحاظ محاسباتی پرهزینه تر از تولید یک انتقال ساده هستند نیاز به تعداد زیادی نمونه از لحاظ محاسباتی حتی پرهزینه تر میشود. فایده جالب توجه پیشبینیهای قطعی این است که گرادیانهای نسبت به پارامترهای سیاست میتوانند به صورت تحلیلی بیشبینیهای قطعی میتواند سیاستهای با هزاران پارامتر را یاد بگیرد[3].

جدول ۲-۳ ویژگیهای پیشبینیهای قطعی و تصادفی تراژکتوری را خلاصه کرده است. این جدول فهرست کرده که آیا امید ریاضی پاداش بلندمدت J_{θ} و گرادیانهای مربوطه $dJ_{\theta}/d\theta$ میتوانند دقیقا ارزیابی شده یا فقط میتوانند تقریبا ارزیابی شوند. برای پیشبینیهای تصادفی تراژکتوری یا به عبارت دیگر نمونهبرداری، محاسبات مورد نیاز نسبتا ساده هستند در حالی که محاسبات مربوط به پیشبینیهای قطعی از لحاظ ریاضیاتی پیچیده تر هستند. در آخر، حدود ممکن روی تعداد پارامترهای سیاست که با استفاده از هر یک از روشهای پیشبینی میتوانند یاد گرفته شوند را ارائه می دهیم. برای تولید تراژکتوری تصادفی، J_{θ} می تواند دقیقا در حد

¹ In the limit

تعداد نامتناهی نمونه تراژکتوری ارزیابی شود. گرادیانهای سیاست مربوطه حتی کندتر همگرا میشوند. در عمل که تعداد متناهی نمونه در دسترس است هر دو $J_{ heta}$ و $J_{ heta}/d heta$ نمی توانند دقیقا ارزیابی شوند.

جدول ۲-۳ ویژگیهای پیشبینیهای قطعی و تصادفی تراژکتوری در جستجوی سیاست مبتنی بر مدل[90].

قطعی	تصادفي	
تقریبی	دقیق در حد	$J_{ heta}$
دقیق	دقیق در حد	$dJ_{\theta}/d\theta$
پیچیده	ساده	محاسبات
$1 \le \theta \le ?$	$1 \le \theta \le 50$	تعداد پارامترهای سیاست

- پرداختن به عدم قطعیت مدل

بیان عدم قطعیت در مورد مدل آموختهشده برای جستجوی سیاست مبتنی بر مدل از این نظر مهم است که در مقابل خطاهای مدل مقاوم باشد. هنگام پیشبینی یا تولید تراژکتوریها، دو روش کلی برای پرداختن به عدم قطعیت مدل وجود دارد.

در مقاله [3]، به عدم قطعیت مدل به عنوان نویز ناهمبسته از لحاظ زمانی پرداخته شده است، به عبارت دیگر خطاهای مدل در هر گام زمانی مستقل در نظر گرفته شدهاند. این رویکرد از لحاظ محاسباتی نسبتا ارزان بوده و امکان در نظر گرفتن تعداد نامتناهی مدل در طول میانگین گیری مدل را می دهد[3]. در عوض، نمونه برداری از پارامترهای مدل در ابتدا و تثبیت کردن آنها برای تولید تراژکتوری این فایده را دارد که زمانی که تراژکتوریهای به طور نسبی نمونه برداری شده به عنوان داده آموزشی در نظر گرفته شده تا زمانی که پارامترهای مدل دوباره نمونه برداری شوند به طور خود کار همبستگی زمانی در نظر گرفته می شود [107]. در اینجا منظور از همبستگی زمانی در طول یک تراژکتوری، با حالت در گام زمانی قبلی همبسته است. از طرفی دیگر، فقط تعداد متناهی (محدود) مدل می توانند نمونه برداری شوند.

- ویژگیهای برون یابی مدلها

در جستجوی سیاست مبتنی بر مدل فرض می شود که مدلها معلوم هستند یا در گام پیش پردازش آموزش داده شده اند [104]. در اینجا، از انسانها خواسته شده که ربات (برای مثال یک هلی کوپتر یا یک بالون) را حرکت داده تا برای ساخت مدل داده جمع آوری شود. یک جنبه بسیار مهم داده جمع آوری شده این است که مناطقی از فضای حالت را پوشش می دهد که برای یادگیری موفق کار پیشرو مناسب هستند. با این وجود ممکن است که بررسی مناطق خارج از داده آموزشی مدل فعلی بتواند (براساس تابع پاداش) بهینه باشد. اما در این مورد، مدل آموخته شده باید قادر باشد اطمینانش در مورد مناطق دور از داده آموزشی را صادقانه پیش بینی کند. مدل های قطعی (برای مثال LWR یا شبکههای عصبی) نمی توانند اطمینانشان در مورد مناطق دور از داده آموزشی را صادقانه نشان دهند که همین موضوع دلیل این است که غالبا با استفاده از عبارات جریمه بزرگ در تابع پاداش، از بررسی (اکتشاف) الجلوگیری می شود [107]. دو مدل که میزان خطای قابل باوری در خارج از مجموعه آموزشی دارند عبارتند از رگرسیون بیزی وزن دار محلی و فرآیندهای گوسی. بنابراین، در صورتی که مجموعه آموزشی دارند عبارتند از رگرسیون بیزی وزن دار محلی و فرآیندهای گوسی. بنابراین، در صورتی که ربات در برابر اکتشاف دلخواه آغازی مقاوم باشد این مدل ها حتی می توانند به منظور یادگیری از هیچ در زمینه رباتیک استفاده شوند، به عبارت دیگر بدون نیاز به اینکه از یک انسان خبره بخواهند تا داده مناسب برای یادگیری مدل یا سیاست آغازی فطری معقول از را تولید کند [3].

مجموعه دادههای بزرگ

در رباتیک، غیرمتداول نیست که مجموعه دادههای بزرگ با میلیونها داده 7 در دسترس باشد. برای نمونه، ضبط 7 دا ثانیه داده با فرکانس یک کیلوهرتز منجر به مجموعه دادهای با یک میلیون داده می شود. برای مدلهای سراسری 7 از قبیل 7 استاندارد، این اندازههای مجموعه دادهها منجر به زمان محاسبه غیرعملی می شود. برای نمونه، 7 نیاز خواهد داشت تا در طول آموزش به صورت مکرر ماتریس کرنل 7 8 را ذخیره و معکوس کند. یک روش رایج برای کاهش اندازه مجموعه داده زیرنمونه گیری (نمونه گیری از نمونه) 6 است، برای مثال

¹ Exploration

² Reasonably innate starting policy

³ Data point

⁴ Gloabal

⁵ Subsampling

برداشتن فقط هر دهمین یا صدمین داده. این امر غالبا شدنی است چون دینامیک به اندازه کافی هموار است و حولت ربات در 1/1000 ثانیه خیلی تغییر نمی کند. علاوه بر این، تقریبهای پراکنده به GP سراسری وجود دارند[116] که به طور مطلوب تر مقیاس می شوند. اما، حتی برای این روشها میلیونها داده غیرعملی هستند. بنابراین، در صورتی که مجموعههای داده بزرگ باشند مدلهای محلی از قبیل LWR یا مدلهای و ترکیب باید به کار گرفته شوند. ایده مدلهای محلی، آموزش تعداد زیادی مدل با اطلاعات محلی و ترکیب پیشبینیهای این مدلهاست.

٣-٢- جمع بندي

در این بخش ابتدا به شرح کلی روشهای یادگیری تقلیدی و یادگیری تقویتی پرداختیم. در توضیح روشهای یادگیری تقلیدی «بازنمایی ویژگی»، «یادگیری مستقیم» و «یادگیری غیرمستقیم» بیان شدند. در قسمت روشهای یادگیری تقویتی نیز «روشهای یادگیری تقویتی مبتنی بر مدل» و «روشهای یادگیری تقویتی فارغ از مدل» توضیح داده شدند. با توجه به اینکه مدل استفاده شده در این پایان نامه در دستهی «روشهای یادگیری تقویتی جستجوی سیاست مبتنی بر مدل» قرار میگیرد در بخش بعدی این فصل این مدلها به تفصیل شرح داده شدند.

فصل سوم ۳- روشهای پیشنهادی

روشهای پیشنهادی

به طور کلی، در این پژوهش هدف انتقال جسم به درون یک حفره توسط بازوی رباتیک شبیهسازی شده است که از روشهای یادگیری تقلیدی و یادگیری تقویتی برای حل مسئله استفاده شده است. برای شبیهسازی در محیطهای گسسته از الگوریتم Q-Learning استفاده می شود که در این پایان نامه با افزودن مسیرهای راهنما که از درون ویدئوهای ضبط شده استخراج شده اند عملکرد الگوریتم بهبود پیدا کرده است. در این فصل در بخش ۱-۳ به شرح الگوریتم ولاد نظر به عنوان مدل پایه محیطهای گسسته برای حل مسئله مورد نظر می بردازیم. در بخش ۲-۳ نیز افزودن مسیرهای راهنما به الگوریتم Q-Learning به عنوان بخشی از روش می پیشنهادی ارائه شده در این پایان نامه شرح داده خواهد شد.

الگوریتم PILCO این بیچیدگی بسیاری همراه است و یا مدل دینامیکی دارای عدم قطعیت میباشد مورد استفاده قرار می گیرد. این الگوریتم، پیش از این در مسائل مختلفی مورد استفاده قرار گرفته است، اما از آن در حل مسئله انتقال جسم به درون یک حفره توسط بازوی رباتیک شبیهسازی شده بهره برده نشده است. با توجه به وجود ذات عدم قطعیت در این مسئله، الگوریتم PILCO میتواند به عنوان راهکاری مناسب برای حل این مسئله در محیطهای پیوسته مورد استفاده قرار بگیرد. در بخش ۳-۳ استفاده از الگوریتم PILCO در حل مسئله انتقال جسم به درون حفره در محیطهای پیوسته به عنوان بخشی دیگر از روش پیشنهادی ارائه شده در این پایاننامه تشریح می گردد. همچنین، نحوه ایجاد تغییر در تابع هزینه به منظور افزودن قابلیت در نظر گرفتن محدودیت فضای حالت (مانند دیوار یا موانع) به مدل PILCO نیز بیان می گردد.

۱-۳- الگوريتم Q-Learning

Q-learning یک روش یادگیری ساده برای عاملهاست که با تجربه عواقب و نتایج اعمال خود و بدون نیاز به ساخت نقشههای حوزهها، در حوزههای مارکفی کنترلشده به صورت بهینه عمل کنند. این روش معادل روش

افزایشی برای برنامهنویسی پویا است که نیازهای محاسباتی محدودی را اعمال میکند. این روش با بهرهگیری از بهبود یی در پی ارزیابی هایش از کیفیت عمل های خاص در حالات خاص عمل می کند.

یادگیری الگوریتم Q-Learning مشابه روش تفاوت زمانی ساتن عمل می کند. به این صورت که عامل یک عمل را در یک حالت خاص امتحان کرده و عواقب و نتایج آن عمل از لحاظ پاداش یا جریمه آنی که دریافت می کند و تخمین عامل از مقدار حالتی که به آن می رود را ارزیابی می کند. با امتحان کردن همه اعمال در همه حالات به طور مکرر، عامل یاد می گیرد که به طور کلی کدام موارد از لحاظ پاداش بلندمدت کاهش یافته بهتر هستند. Q-learning یک شکل ابتدایی یادگیری است اما به تنهایی می تواند به عنوان پایه و اساس شیوههای خیلی پیچیده تر عمل کند.

یک عامل محاسباتی که در یک جهان گسسته متناهی حرکت میکند و در هر گام زمانی یک عمل را از مجموعه متناهی از اعمال انتخاب می کند را در نظر بگیرید. جهان از یک فرآیند مارکف کنترل شده به همراه یک عامل به عنوان کنترل کننده تشکیل شده است. در گام n عامل می تواند حالت $\chi_n(\epsilon X)$ را ثبت کند و مطابق با آن عمل $a_n(\epsilon A)$ را انتخاب کند. عامل پاداش احتمالاتی r_n را دریافت می کند که مقدار میانگین y_n به (۱-۳) به مطابق معادله (۱-۳) به $R_{\chi_n}(a_n)$ آن $R_{\chi_n}(a_n)$ به به حالت و عمل بستگی دارد و حالت جهان به طور احتمالاتی مطابق معادله تغيير مي كند:

$$Prob[y_n = y | x_n, a_n] = P_{x_n y}[a_n].$$
 (۱-۳) معادله

کاری که عامل با آن مواجه است تعیین یک سیاست بهینه است که یاداش کاهش یافته مورد انتظار کلی را حداکثر می کند. منظور از پاداش کاهش یافته، پاداشهایی است که در s گام دریافت می شود بنابراین از پاداشهایی که در حال حاضر دریافت میشوند با ضریب γ^s بارزش کمتری دارند. تحت سیاست π ، ارزش حالت x به صورت معادله (τ - τ) است:

$$V^{\pi}(x) = R_{x}(\pi(x)) + \gamma \sum_{y} P_{xy}[\pi(x)] V^{\pi}(y)$$
 (۲-۳) معادله

با توجه به این معادله، عامل انتظار دارد برای اجرای عمل π که پیشنهاد کرده پاداش $R_xig(\pi(x)ig)$ را فورا دریافت کند و سپس با احتمال $P_{xy}[\pi(x)]$ به یک حالت که برایش ارزش $V^{\pi}(y)$ دارد حرکت می کند.

¹ Temporal difference

اصول نظری برنامهنویسی پویا به ما اطمینان می دهد که حداقل یک سیاست مانا و بهینه π^* وجود دارد که به صورت معادله(۳-۳) است.

$$V^*(x) = V^{\pi^*}(x) = \max_a \{R_x(\pi(a)) + \gamma \sum_y P_{xy}[\pi(a)]V^{\pi^*}(y)\}$$
 معادله (۳-۳) معادله

برای یک سیاست π ، مقادیر Q به صورت معادله(۳-۴) تعریف میشود.

$$Q^{\pi}(x,a) = R_{x}(a) + \gamma \sum_{y} P_{xy}[\pi(x)] V^{\pi}(y)$$
 (۴-۳)معادله

به عبارت دیگر، مقدار Q امید ریاضی پاداش کاهشیافته برای اجرای عمل a در حالت a و دنبال کردن سیاست a بعد از آن است. در Q-learning، هدف تخمین مقادیر a برای یک سیاست بهینه است. برای راحتی، به ازای هر a و a و a و a تعریف میشود. به راحتی میتوان نشان داد که a^* (a و a و a و a عملی باشد که در آن به بیشینه دست یافته میشود، یک a^* و همچنین اگر a^* عملی باشد که در آن به بیشینه دست یافته میشود، یک سیاست بهینه میتواند به صورت a^* a شکل بگیرد. سودمندی مقادیر a در اینجاست-اگر یک عامل بتواند آنها را یاد بگیرد، به راحتی میتواند تصمیم بگیرد که انجام چه کاری بهینه است. اگر چه ممکن است بیش تر از یک سیاست بهینه یا a^* وجود داشته باشد مقادیر a منحصر به فرد هستند.

در Q-learning، تجربه عامل شامل دنبالهای از مراحل یا قسمتهای مجزاست. در n امین قسمت برای عامل فرآیند زیر اتفاق میافتد:

-

¹ Watkins

- حالت فعلى خودش x_n را مشاهده مى كند،
 - می کند، و اجرا می کند، a_n عمل a_n عمل
 - را مشاهده می کند، y_n حالت بعدی \bullet
 - بازده آنی r_n را دریافت می کند، ullet
- مقادیر Q_{n-1} را با استفاده از ضریب یادگیری $lpha_n$ براساس معادله q_{n-1} تنظیم می کند:

$$Q_n(x,a)=$$

$$\begin{cases} (1-\alpha_n)Q_{n-1}(x,a)+\alpha_n[r_n+\gamma V_{n-1}(y_n)] & if \ x=x_n \ and \ a=a_n, \\ Q_{n-1}(x,a) & otherwise, \end{cases}$$

بهترین کاری که عامل فکر می کند می تواند از حالت y انجام دهد به صورت معادله(y-z) قابل بیان است.

$$V_{n-1}(y_n) \equiv \max_b \{Q_{n-1}(y,b)\}$$
 (۶-۳)معادله

البته، در مراحل اولیه یادگیری، مقادیر Q ممکن است سیاستی را که به طور ضمنی تعریف میکنند به طور دقیق منعکس نکنند. فرض شده است که مقادیر اولیه Q یا $Q_0(x,a)$ ، برای همه حالات و اعمال داده شدهاند.

۳-۲ روش پیشنهادی: استفاده از مسیرهای راهنما در Q-Learning

یکی از اصلی ترین اجزای الگوریتم Q-Learning ماتریس Q است. این ماتریس، میزان ارزش هر عمل را در هر یک از حالتهای ممکن مشخص می کند. در حالت معمول الگوریتم Q-Learning مقداردهی اولیه می شود. این در حالی است که تعیین مقدار اولیه مناسب برای ماتریس Q می تواند تاثیر بسیار قابل توجهی در عملکرد ربات داشته باشد. برای مثال، در حالتی که ماتریس Q با صفر مقداردهی شده باشد، در فضای حالت با ابعاد بزرگ تر از 20×20 ، ربات در زمان معقول نمی تواند جسم را با شروع از موقعیت در فضای حالت با ابعاد میزون از حفره قرار دارد به موقعیت هدف درون حفره برساند. روش پیشنهادی ارائه شده در این پایان نامه، تعیین چارچوبی است که با استفاده از آن بتوان با مقداردهی اولیه مناسب ماتریس Q، الگوریتم Q- لین مسئله، استفاده از مسیرهای راهنما می باشد. این مسیرهای راهنما در قالب مقدار اولیه ماتریس برای حل این مسئله، استفاده از مسیرهای راهنما می باشد. این مسیرهای راهنما در قالب مقدار اولیه ماتریس

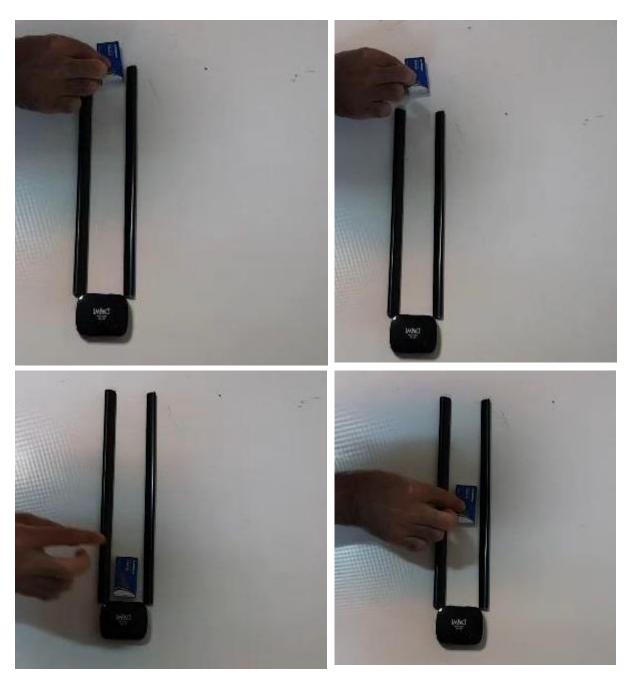
Q در اختیار الگوریتم Q-Learning قرار می گیرند. مسیرهای راهنما می توانند از دو طریق در اختیار الگوریتم Q-Learning قرار بگیرند. در حالت اول، مسیرهای رساندن جسم به هدف به صورت دستی توسط عامل انسانی طراحی شده و در اختیار الگوریتم قرار گرفته اند. استفاده از این روش، در فضاهای حالت با ابعاد بزرگ دارای پیچیدگی زیادی است. در حالت دوم، این مسیرهای راهنما از درون ویدئوهای ضبط شده و با استفاده از الگوریتمهای ردیابی شی استخراج می شوند. در ادامه این فصل، در ابتدا نحوه استفاده از ویدئوهای ضبط شده برای ساخت مسیر راهنما در بخش Y-Y-Y نحوه نفوذ دادن مسیرهای راهنما به عنوان مقدار اولیه ماتریس Q الگوریتم Q-Learning به تفصیل توضیح داده خواهد شد.

۱-۲-۳ ایجاد مسیرهای راهنما از روی ویدئوهای ضبط شده

مسیرهای راهنما می توانند از دو طریق مختلف ایجاد شوند. در روش اول، مسیر انتقال جسم از نقطه اولیه به هدف توسط عامل انسانی مشخص می شود. این روش در محیطهای با ابعاد بزرگ چندان کارایی ندارد و پیچیدگیهای زیادی را برای عامل انسانی به همراه می آورد. با توجه به این مسئله روش دوم ارائه می گردد که در آن مسیر انتقال جسم از نقطه اولیه به هدف به صورت خودکار و توسط ویدئوهای ضبط شده (دمو) تعیین می گردد. در این بخش، هدف تشریح نحوه استفاده از ویدئوهای ضبط شده برای تعیین مسیر راهنما است.

برای اینکه بتوان ویدئوهایی به منظور استخراج مسیر راهنما ایجاد کرد، لازم است که ویدئوها در فضایی مشابه فضای حالت ربات ضبط شوند. با توجه به این موضوع، لازم است محیطی طراحی کنیم که شامل جسم، حفره و فضای حرکت جسم باشد. جسم باید در یک مسیر معقول به درون حفره منتقل شود و ویدئویی از این فرآیند ضبط گردد. در شکل ۳-۱ نمونهای از چهار نمای ویدئوی ضبط شده مشاهده می شود که جسم و حفره در آن تعبیه شده اند.

پس از آماده سازی ویدئوها لازم است تا مسیر حرکت جسم از موقعیت اولیه به هدف به طور خودکار از ویدئوها استخراج شود. برای این منظور، از الگوریتمهای ردیابی شی در محیط دو بعدی استفاده شده است. در این گونه از الگوریتمها محل جسم درون محیط در هر فریم تشخیص داده شده و به صورت مختصات در فضای دو بعدی در اختیار قرار می گیرد.



شکل ۲-۳ چهار نمای مختلف از یکی از ویدئوهای ضبط شده مسیر راهنما. شکل سمت راست بالا زمان شروع انتقال جسم، شکل بالا سمت چپ هنگام عبور دادن جسم از دریچه حفره، شکل سمت راست پایین انتقال جسم به سمت هدف در میانه حفره، شکل سمت چپ پایین قرار گرفتن جسم در نقطه هدف را نشان میدهند.

۲-۲-۲ مقداردهی اولیه ماتریس ${f Q}$ با استفاده از مسیرهای راهنما

پس از استخراج مختصات جسم از فریمهای مختلف ویدئو لازم است که ماتریس Q با توجه به مختصات جسم در ویدئوهای ضبط شده مقداردهی اولیه شود. برای این منظور، در ابتدا مقادیر تمام درایههای ماتریس Q صفر در نظر گرفته می شود و سپس درایههایی از ماتریس Q که جسم در مختصات مربوط به آنها در فریمهای ویدئو مشاهده شده اند ارزش بالاتری متناسب با محل خود دریافت می کنند.

در گام بعدی، با در نظر گرفتن یک میدان پتانسیل پیرامون هر یک از نقاط دارای ارزش (نقاطی که جسم در مختصات مربوط به آنها در ویدئوها مشاهده شده است)، ارزش نقاط دیگر با توجه به میزان قرارگیری آنها در میدان پتانسیل تعیین می گردد. به این صورت که نقاط مجاور مراکز پتانسیل ارزش بیش تری دریافت می کنند و سایر نقاط متناسب با فاصله خود با مراکز پتانسیل ارزش کم تری را به خود اختصاص می دهند.

۳-۳- استفاده از مدل PILCO در حل مسئله انتقال جسم به درون حفره در

محيطهاي ييوسته

در این بخش، توضیح میدهیم که یک بازوی رباتیک شبیهسازی شده چگونه می تواند یاد بگیرد که یک بلوک را به طور کاملا خودمختار داخل یک حفره قرار دهد. از فرضیاتی که در ادامه گفته خواهد شد استفاده می کنیم: ابتدا، از آنجایی که در دست گرفتن شی، تمرکز این کار نیست، فرض می کنیم که بلوک در گیره ربات قرار داده می شود. ثانیا، زوایا و سرعتهای مفاصل بازو اندازه گیری نمی شوند و فقط موقعیت و سرعت مرکز بلوک در گیره ربات اندازه گیری می شود. ثالثا، هیچ مسیر یا مسیر مطلوبی از پیش معلوم نیست. رابعا، فرض می کنیم که موقعیت اولیه و موقعیت هدف بلوک در گیره ربات ثابت است.

به طور ساده دنبال کردن یک مسیر مستقیم بین موقعیت اولیه و هدف ممکن است به دلیل موانع موفقیت آمیز نباشد. برای برنامه ریزی بلندمدت و یادگیری کنترل کننده، این عدم قطعیت باید در نظر گرفته شود. $u \in \mathbb{R}^2$ در هر گام زمانی، ربات از مرکز بلوک در گیرهاش برای محاسبه سیگنال کنترلی با مقدار گسسته $x \in \mathbb{R}^2$ استفاده می کند. چرخش مچ و باز/بسته شدن گیره یاد گرفته نمی شوند. الگوریتم یادگیری تقویتی، از مرکز ۲ بعدی شی و سرعتش به عنوان حالت $x \in \mathbb{R}^4$ استفاده می کند.

۱-۳-۳- یادگیری سیاست با محدودیتهای فضای حالت

در ادامه، چارچوب کاری PILCO را برای یادگیری یک سیاست حلقه-بسته مناسب (کنترل کننده بازخورد در ادامه، چارچوب کاری PILCO را برای یادگیری یک سیاست حلقه-بسته مناسب (کنترل کننده بازخورد حالت یا به طور مختصر شرح می در اینجا، x حالت نامیده می شود که به صورت مختصات و سرعت مرکز بلوک (x_c, y_c, x_c', y_c') در گیره تعریف شده است. ما قصد داریم که این سیاست را از هیچ یا به عبارت دیگر فقط با دانش قبلی خیلی کلی در مورد کار و خود راهحل یاد بگیریم. علاوه بر این، می خواهیم π را در تعداد کمی آزمایش پیدا کنیم یا به عبارت دیگر، به یک روش یادگیری داده-کارا نیاز داریم.

به عنوان معیاری برای قضاوت در مورد عملکرد کنترل کننده π ، می توانیم از امید ریاضی حاصل بلندمدت یک تراژ کتوری (x_0, \cdots, x_T) به هنگام اعمال π استفاده کنیم که به صورت معادله (x_0, \cdots, x_T) است:

$$J^{\pi} = \sum_{t=0}^{T} \mathbb{E}_{x_t}[c(x_t)]$$
 معادله(۲-۳)معادله

در معادله(۷-۳)، T افق پیشبینی و $c(\mathbf{x}_t)$ هزینه لحظهای بودن در حالت \mathbf{x} در زمان t است. در صورتی که به \mathbf{x} در معادله(\mathbf{x}_t) می بیشبینی و \mathbf{x} این پایان نامه از یک تابع هزینه اشباع شده \mathbf{x} در کل این پایان نامه از یک تابع هزینه اشباع شده \mathbf{x} در معرک در معرک (مجری) پایانی از موقعیت هدف \mathbf{x} را جریمه استفاده می کنیم که فواصل اقلیدسی \mathbf{x} با \mathbf{x} پارامتربندی می شود. PILCO سیاست پارامتربندی شده مناسب را با استفاده از الگوریتم \mathbf{x} یاد می گیرد.

-

¹ End effector

۱-۱-۳-۳ مدل دینامیکی احتمالاتی

برای پرهیز از فرضیات همارزی اطمینان (قطعیت) روی مدل آموخته شده، PILCO عدم قطعیتهای مدل را در طول برنامهریزی در نظر می گیرد. از همین رو، یک توزیع (پسین) روی مدلهای دینامیکی ممکن نیاز است. ما از GP ها[102] برای استنتاج این توزیع پسین از تجربه در حال حاضر در دسترس استفاده می کنیم.

با پیروی از [102]، به طور مختصر نمادگذاری و مدلهای پیشبینی استاندارد برای GP ها را معرفی می کنیم که به منظور استنتاج توزیع روی یک تابع پنهان f از مشاهدات نویزی $y_i = f(\mathbf{x}_i) + \varepsilon$ استفاده می شوند، در این پایان نامه، $\varepsilon \sim \mathcal{N}(0, \sigma_{\varepsilon}^2)$ نویز i.i.d سیستم در نظر گرفته شده است. یک GP به طور کامل با یک تابع میانگین و $m(\cdot)$ و یک تابع کوواریانس نیمه معین مثبت $m(\cdot)$ که کرنل نامیده می شود مشخص می شود. در سرتاسر این پایان نامه، یک تابع میانگین مقدم $m(\cdot)$ و کرنل نمایی مجذور (SE) با تعیین ارتباط خود کار در نظر می گیریم که به صورت زیر تعریف شده است:

$$k(x,x') = \alpha^2 \exp\left(-\frac{1}{2}(x-x')^T \Lambda^{-1}(x-x')\right)$$
 هعادله(۸-۳) معادله

در اینجا، α^2 را به عنوان واریانس تابع پنهان f و f را به عنوان واریانس تابع پنهان f و f را تعریف کردهایم که به f را تعریف کردهایم که به f مقیاسهای طول مشخصه f بستگی دارد. با در اختیار داشتن f و اهداف f و اهداف f بستگی دارد. با در اختیار داشتن f و اوریانس سیگنال f و اوریانس سیگنال g و اوریانس سیگنال g و اوریانس سیگنال g و اوریانس های نویز g با بیشینه سازی شواهد g یاد گرفته می شوند g یا بیشینه سازی شواهد g یا بیشینه بیشینه سازی شواهد g یا بیشینه سازی شواهد g یا بیشینه بیشین بیشینه بیشین بیشینه بیشینه بیشینه بیشینه بیشینه بیشینه بیشین بیشینه بیشینه بیشینه بیشینه بیشینه بیشینه بیشین بیشینه بیشینه بیشینه بیشین بیشین بیشینه بیشین بیشی

توزیع پیشبینی کننده پسین $p(f_*|\mathbf{x}_*)$ تابع مقدار $f_*=f(\mathbf{x}_*)$ برای یک ورودی آزمایشی دلخواه اما معلوم $p(f_*|\mathbf{x}_*)$ عستند معادله(۳-۳) و معادله(۳-۳) هستند \mathbf{x}_*

$$m_f(x_*) = \mathbb{E}_f[f_*] = k_*^T (K + \sigma_{\varepsilon}^2)^{-1} y = k_*^T \beta$$
 (۹-۳)معادله

$$\sigma_f^2(x_*) = var_f[f_*] = k_{**} - k_*^T (K + \sigma_{\varepsilon}^2 I)^{-1} k_* + \sigma_{\varepsilon}^2$$
 (۱۰-۳)معادله

¹ Certainty equivalence

² Evidence maximization

 K_{ij} که در آن $\beta:=(K+\sigma_{arepsilon}^2)^{-1}$ ی $k_{**}:=k(\mathrm{X}_*,\mathrm{X}_*)$ $k_*:=k(\mathrm{X},\mathrm{X}_*)$ که در آن $k_*:=k(\mathrm{X},\mathrm{X}_*)$ است.

در سیستم رباتیکی ما، $f: \mathbb{R}^6 \to \mathbb{R}^4$, $(\mathbf{x}_{t-1}, \mathbf{u}_{t-1}) \mapsto \Delta_t := \mathbf{x}_t - \mathbf{x}_{t-1} + \varepsilon_t$ تابع \mathbf{GP} تابع \mathbf{GP} تابع $\mathbf{i}.i.d$ سیستم است. ورودیها و اهداف آموزشی به مدل \mathbf{GP} به ترتیب، چندتاییهای $\varepsilon_t \in \mathbb{R}^4$ هستند. $(\mathbf{x}_{t-1}, \mathbf{u}_{t-1})$

۲-۱-۳-۳ برنامهریزی بلندمدت از طریق استنتاج تقریبی

کمینه کردن و ارزیابی J^π در معادله(۳-۷)، به پیشبینیهای بلندمدت از تغییر و تحول حالت نیاز دارد. برای به دست آوردن توزیعهای حالت $p(\mathbf{x}_1), \cdots, p(\mathbf{x}_T), \cdots, p(\mathbf{x}_T)$ پیشبینیهای یک گامه را پشت سر هم می چینیم. انجام صحیح این کار به نگاشت ورودیهای آزمایشی نامطمئن به واسطه مدل دینامیکی GP نیاز دارد. در ادامه، فرض می کنیم که این ورودیهای آزمایشی دارای توزیع گوسی هستند و نتایج مرجع [3] را به برنامهریزی بلندمدت در سیستمهای تصادفی با ورودیهای کنترلی تعمیم می دهیم.

 $u_{t-1}=u_{t-1}=p(x_{t-1},u_{t-1})$ به توزیع توام $p(x_{t-1},u_{t-1})$ نیاز داریم. برای محاسبه این توزیع، از $p(x_{t-1},u_{t-1})$ به توزیع بیش بینی کننده کنترل $p(x_{t-1},u_{t-1})$ استفاده می کنیم که این یعنی کنترل تابعی از حالت است: ابتدا، توزیع پیش بینی کننده کنترل $p(x_{t-1},u_{t-1})$ و سپس، کوواریانس متقابل $p(x_{t-1},u_{t-1})$ را محاسبه می کنیم. سرانجام، $p(x_{t-1},u_{t-1})$ و سپس، کوواریانس متقابل $p(x_{t-1},u_{t-1})$ را محاسبه می کنیم. این محاسبات به پارامتربندی $p(x_{t-1},u_{t-1})$ با یک توزیع گوسی با میانگین و کواریانس درست تقریب می زنیم. این محاسبات به پارامتربندی $p(x_{t-1},u_{t-1})$ با $u_{t-1}=\pi(x_{t-1})=Ax_{t-1}+b$ با $u_{t-1}=\pi(x_{t-1})=Ax_{t-1}+b$ با اعمال نتایج استاندارد از مدل های گوسی خطی و رابطه زیر حاصل می شوند:

$$p(\mathbf{u}_{t-1}) = \mathcal{N}(\mathbf{u}_{t-1}|\mu_u, \Sigma_u)$$

$$\mu_u = A\mu_{t-1} + \mathbf{b}, \quad \Sigma_u = A\Sigma_{t-1}A^T$$

,

¹ Linear-Gaussian

در این مثال، π یک تابع خطی از \mathbf{x}_{t-1} است و بنابراین توزیع توام مطلوب $\mathbf{x}_{t-1},\mathbf{u}_{t-1}$ دقیقا گوسی است و به صورت معادله(۱۰-۳) مشخص می شود

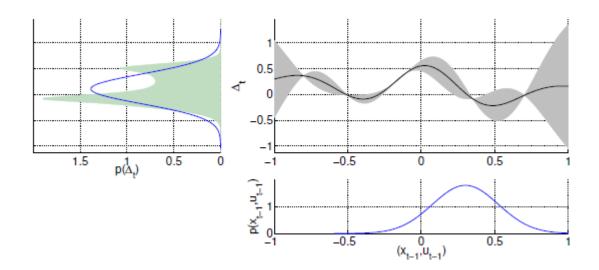
$$\mathcal{N}\left(\begin{bmatrix} \mu_{t-1} \\ A\mu_{t-1} + b \end{bmatrix}, \begin{bmatrix} \Sigma_{t-1} & \Sigma_{t-1}A^T \\ A\Sigma_{t-1} & A\Sigma_{t-1}A^T \end{bmatrix}\right)$$
 (۱۱-۳)معادله

که در این معادله، کوواریانس متقابل $\sum_{t=1}^T A^T$ است. برای خیلی از پارامتربندیهای که در این معادله، کوواریانس متقابل $\sum_{t=1}^T A^T$ است. برای خیلی از پارامتربندیهای کنترل کننده جالب دیگر، میانگین و کوواریانس میتوانند به صورت تحلیلی انجام شوند[118]، گرچه ممکن است $p(\mathbf{x}_{t-1}, \mathbf{u}_{t-1})$ دیگر دقیقا گوسی نباشد.

از الان به بعد، یک توزیع توام گوسی $p(\widetilde{\mathbf{x}}_{t-1}|\widetilde{\mu}_{t-1},\widetilde{\Sigma}_{t-1})$ در زمان t-1 فرض می کنیم $p(\widetilde{\mathbf{x}}_{t-1})=\mathcal{N}(\widetilde{\mathbf{x}}_{t-1}|\widetilde{\mu}_{t-1},\widetilde{\Sigma}_{t-1})$ در زمان $\mathbf{x}:=[\mathbf{x}^T \mathbf{u}^T]^T$ که $\mathbf{x}:=[\mathbf{x}^T \mathbf{u}^T]^T$ انتگرال گیری که توزیع نشان داده شده در معادله(۱۲-۳) را پیشبینی می کنیم، روی متغیر تصادفی $\widetilde{\mathbf{x}}_{t-1}$ انتگرال گیری می کنیم.

$$p(\Delta_t) = \int p(f(\tilde{x}_{t-1})|\tilde{x}_{t-1})p(\tilde{x}_{t-1})d\tilde{x}_{t-1}$$
 (17-۳) a solution

احتمال انتقال $p(f(\tilde{\mathbf{x}}_{t-1})|\tilde{\mathbf{x}}_{t-1})$ از توزیع پسین GP به دست میآید. محاسبه توزیع پیشبینی کننده دقیق $p(f(\tilde{\mathbf{x}}_{t-1})|\tilde{\mathbf{x}}_{t-1})$ به صورت تحلیلی لاینحل است. بنابراین، $p(\Delta_t)$ را با یک گوسی با میانگین و واریانس دقیق (تطبیق ممان) تقریب میزنیم. شکل ۲-۳ این سناریو را نشان می دهد. توجه داشته باشید که برای محاسبه میانگین می و واریانس σ_{Δ}^2 توزیع پیشبینی کننده، توزیع پیشبینی کننده (به معادلات معادلات معادله (۲-۳) و معادله (۱۰-۳) مراجعه شود) کفایت نمی کند چون $\tilde{\mathbf{x}}_{t-1}$ به طور قطعی داده نشده است.



شکل ۲-۳ پیشبینی \mathbf{GP} در یک ورودی نامطمئن. توزیع ورودی $\mathbf{F}(x_{t-1},u_{t-1})$ گوسی فرض شده است (پنل سمت راست پایین). انتشار آن از طریق مدل \mathbf{GP} (پنل سمت راست بالا)، توزیع سایه خورده است $\mathbf{p}(\Delta_t)$ در پنل سمت چپ بالا را نتیجه می دهد. $\mathbf{p}(\Delta_t)$ با یک گوسی با میانگین و واریانس دقیق تقریب زده می شود (پنل سمت چپ پایین). $\mathbf{p}(\Delta_t)$

فرض می کنیم که میانگین μ_{Δ} و کوواریانس Σ_{Δ} توزیع پیشبینی کننده $p(\Delta_t)$ معلوم هستند. آنگاه، یک تقریب گوسی به توزیع حالت مورد نظر $p(\mathbf{x}_t)$ به ترتیب میانگین و کوواریانس به صورت معادلات معادله $p(\mathbf{x}_t)$ دارد.

$$\mu_t = \mu_{t-1} + \mu_{\Delta}$$
 معادله(۱۳-۳)

$$\Sigma_t = \Sigma_{t-1} + \Sigma_{\Delta} + cov[x_{t-1}, \Delta_t] + cov[\Delta_t, x_{t-1}]$$
 (۱۴-۳)معادله

$$cov[x_{t-1}, \Delta_t] = cov[x_{t-1}, u_{t-1}] \Sigma_u^{-1} cov[u_{t-1}, \Delta_t]$$
 (۱۵-۳)معادله

محاسبه کوواریانسهای متقابل مورد نیاز در معادله(۳-۱۵) به پارامتربندی سیاست بستگی دارد، اما غالبا می تواند به صورت تحلیلی محاسبه شود.

در ادامه، میانگین μ_{Δ} و واریانس σ_{Δ}^2 توزیع پیشبینی کننده $p(\Delta_t)$ را محاسبه می کنیم (به معادله(۱۲-۳) در ادامه، میانگین μ_{Δ} و واریانس σ_{Δ}^2 توزیع پیشبینی کننده σ_{Δ}^2 توزیع پیشبینی کننده σ_{Δ}^2 را محاسبه می کنیم (به معادله(۱۲-۳) و را محاسبه (به معادله(۱۲

۱) میانگین: با پیروی از قانون امید ریاضیهای مکرر، به دست می آوریم:

$$\mu_{\Delta} = \mathbb{E}_{\tilde{\mathbf{x}}_{t-1}} \left[\mathbb{E}_f[f(\tilde{\mathbf{x}}_{t-1})|\tilde{\mathbf{x}}_{t-1}] \right] = \mathbb{E}_{\mathbf{x}_*} \left[m_f(\tilde{\mathbf{x}}_{t-1}) \right]$$

$$= \int m_f(\tilde{\mathbf{x}}_{t-1}) \mathcal{N} \left(\tilde{\mathbf{x}}_{t-1} \middle| \tilde{\mu}_{t-1}, \tilde{\Sigma}_{t-1} \right) d\tilde{\mathbf{x}}_{t-1} = \beta^T q$$
(۱۶-۳)معادله

که $\mathbf{q} \in \mathbb{R}^n$ به صورت زیر هستند: $\mathbf{q} = [q_1, \cdots, q_n]^{\mathrm{T}}$ که $\mathbf{g} = (K + \sigma_{\varepsilon}^2 \mathbf{I})^{-1} y$

$$\begin{split} q_i &= \int k(\mathbf{x}_i, \mathbf{x}_*) \mathcal{N} \big(\tilde{\mathbf{x}}_{t-1} \big| \tilde{\mu}_{t-1}, \tilde{\Sigma}_{t-1} \big) d\tilde{\mathbf{x}}_{t-1} \\ &= \frac{\alpha_a^2}{\sqrt{\left| \tilde{\Sigma}_{t-1} \Lambda^{-1} + \mathbf{I} \right|}} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \tilde{\mu}_{t-1})^T \big(\tilde{\Sigma}_{t-1} + \Lambda \big)^{-1} (\mathbf{x}_i - \tilde{\mu}_{t-1}) \right) \,. \end{split}$$

۲) واریانس: با استفاده از قانون واریانس کل، به دست می آوریم:

$$\begin{split} \sigma_{\Delta}^2 &= \mathbb{E}_{\tilde{\mathbf{x}}_{t-1}} \big[m_f(\tilde{\mathbf{x}}_{t-1})^2 \big] - \mathbb{E}_{\tilde{\mathbf{x}}_{t-1}} \big[m_f(\tilde{\mathbf{x}}_{t-1}) \big]^2 + \mathbb{E}_{\tilde{\mathbf{x}}_{t-1}} \big[\sigma_f^2(\tilde{\mathbf{x}}_{t-1}) \big] \\ &= \beta^T Q \beta + \alpha^2 - tr((K + \sigma_{\varepsilon}^2 I)^{-1} Q) - \mu_{\Delta}^2 + \sigma_{\varepsilon}^2 \end{split} \tag{1V-Y}$$
معادله

که منظور از tr(.)، مجموع عناصر روی قطر اصلی ماتریس است. عناصر $Q \in \mathbb{R}^{n imes n}$ به صورت زیر هستند:

$$\begin{split} Q_{ij} &= k(\mathbf{x}_i, \tilde{\mu}_{t-1}) k\big(\mathbf{x}_j, \tilde{\mu}_{t-1}\big) \big| 2\tilde{\Sigma}_{t-1} \Lambda^{-1} + \mathbf{I} \big|^{-\frac{1}{2}} \times \exp\Big(\frac{1}{2} z_{ij}^T \big(2\tilde{\Sigma}_{t-1} \Lambda^{-1} + \mathbf{I}\big) \tilde{\Sigma}_{t-1} z_{ij}\Big) \\ & . \\ \mathcal{Z}_{ij} := \Lambda^{-1} (\varsigma_i + \varsigma_j) \,\,_{\mathfrak{Z}} \,\,_{\mathfrak{Z}} := (\mathbf{x}_i - \tilde{\mu}_{t-1}) \end{split}$$
 که

توجه داشته باشید که هر دو μ_{Δ} و μ_{Δ} از لحاظ کاربردی بهترتیب به واسطه $\tilde{\mu}_{t-1}$ و $\tilde{\mu}_{t-1}$ به میانگین μ_{Δ} و عود داشته باشید که هر دو Σ_{t-1} و راز لحاظ کاربردی بهترتیب به واسطه Σ_{t-1} و معادله Σ_{t-1} رجوع شود). کوواریانس Σ_{t-1} سیگنال کنترلی وابسته هستند (به معادلات معادله Σ_{t-1}) و معادله (۱۲-۳) بفهمیم که عدم قطعیت در مورد تابع پنهان Σ_{t-1} انتگرال گیری می شود، چیزی که صریحا عدم قطعیت مدل را در نظر می گیرد.

94

¹ Trace

۳-۱-۳ یادگیری کنترل کننده از طریق جستجوی غیرمستقیم سیاست

از بخش 7-1-T می دانیم که چگونه پیش بینی های یک گامه را پشت سر هم انجام دهیم تا تقریبهای گوسی J^π در از توزیع های پیش بینی کننده $p(\mathbf{x}_1), \cdots, p(\mathbf{x}_T)$ را به دست آوریم. برای ارزیابی امیدریاضی حاصل $p(\mathbf{x}_1), \cdots, p(\mathbf{x}_T)$ معادله (Y-T)، محاسبه امید ریاضی مقادیر هزینه $p(\mathbf{x}_1)$ نسبت به توزیع های پیش بینی کننده حالت باقی مانده است.

$$\mathbb{E}[c(x_t)] = \int c(x_t) \mathcal{N}(x_t | \mu_t, \Sigma_t) dx_t$$
 , $t = 0, \dots, T$ (۱۸-۳)معادله

فرض می کنیم که هزینه c طوری انتخاب شده است (برای مثال چندجملهای) که معادله c می تواند به صورت تحلیلی حل شود.

برای اعمال جستجوی مبتنی بر گرادیان سیاست به منظور پیدا کردن پارامترهای ψ که J^π را کمینه می کنند (۷-۳) به معادله(۷-۳) رجوع شود)، ابتدا ترتیب مشتق گیری و جمع کردن در معادله(V-۳) را عوض می کنیم. با در نظر گرفتن $\mathcal{E}_t := \mathbb{E}_{\mathbf{x}_t}[c(\mathbf{x}_t)]$ به دست می آوریم:

$$\frac{d\mathcal{E}_t}{d\theta} = \frac{d\mathcal{E}_t}{d\mu_t} \frac{d\mu_t}{d\psi} + \frac{d\mathcal{E}_t}{d\Sigma_t} \frac{d\Sigma_t}{d\psi}$$
 (۱۹-۳)معادله

مشتقات کلی میانگین ψ و کوواریانس Σ_t مربوط به $p(\mathbf{x}_t)$ نسبت به پارامترهای سیاست ψ می توانند به مشتقات کلی میانگین ψ اعمال مکرر قاعده زنجیرهای به معادلات معادله(\mathbf{r} - \mathbf{r})، معادله(\mathbf{r} - \mathbf{r})، معادله(\mathbf{r} - \mathbf{r})، معادله(\mathbf{r} - \mathbf{r}) محاسبه شوند. این کار، محاسبه مشتقات جزئی \mathbf{r} و معادله(\mathbf{r} - \mathbf{r}) محاسبه شوند. این کار، محاسبه مشتقات در مقاله [118] به صورت تحلیلی دربرمی گیرد. در اینجا جزئیات بیش تر را حذف کردهایم، اما این مشتقات در مقاله [118] به صورت تحلیلی محاسبه شدهاند. این کار امکان استفاده از روشهای بهینه سازی غیرمحدب مبتنی بر گرادیان استاندارد از قبیل دو CG یا L-BFGS را فراهم می کند که بردار پارامتر بهینه شده ψ را برمی گردانند.

۱-۴-۳-۳- برنامهریزی با محدودیتهای فضای حالت

در ساز و کار یادگیری تقویتی کلاسیک، فرض بر این است که یادگیرنده از هیچکدام از محدودیتهای در فضای حالت باخبر نیست، بلکه باید دیوارها و غیره را با برخورد به آنها و متحمل جریمه سنگین شدن کشف کند. در ساز و کار رباتیک، این فرض کلی اما نه لازم، خیلی مطلوب نیست چون ربات می تواند صدمه ببیند.

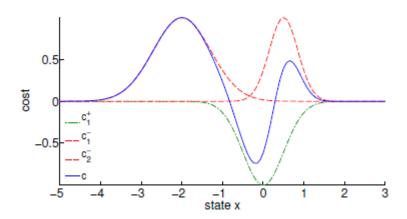
اگر محدودیتها (برای مثال موانع) در فضای حالت از قبل معلوم باشند، ترجیح بر این است که این دانش قبلی را مستقیما در برنامهریزی و یادگیری سیاست بگنجانیم. پیشنهاد میکنیم که موانع را به صورت مناطق "نامطلوب" تعریف کنیم، به عبارت دیگر مناطقی که ربات قرار است از آنها پرهیز کند. "نامطلوب بودن" را به صورت جریمه در تابع هزینه لحظهای c تعریف می کنیم. از همین رو، تابع هزینه c را به صورت زیر تعریف مىكنيم

$$c(x) = \sum_{k=1}^{K} c_k^+(x) - \sum_{j=1}^{J} \iota_j c_j^-(x)$$
 معادله(۲۰-۳)

که c_k^+ حالتهای مطلوب (برای مثال حالت هدف) هستند و c_i^- حالتهای نامطلوب (برای مثال موانع) هستند که با $t_i \geq 0$ وزن دهی شدهاند. مقادیر بزرگتر برای t_i باعث می شوند سیاست بیش تر از یک حالت نامطلوب خاص پرهیز کند. در این پایان نامه، همیشه $\iota_i=1$ را قرار می دهیم. برای c_k^+ و c_i^- ، نمایی های مجذور منفی شده را انتخاب می کنیم که در هنگام میانگین گیری مطابق با توزیع حالت، بین اکتشاف و بهرهبرداری توازن ایجاد می کند[118]. نماییهای مجذور با عرضهای بالقوه مختلف Σ_k^+ نرمالنشده هستند. عرضهای تکتک محدودیتها c_i^- معین می کنند که محدودیتها چقدر "نرم" هستند. محدودیتهای سفت و سخت با نماییهای مجذور c_i^- مجنور است که در آن است که در آن این ایده مربوط به مرجع [119] مجذور c_i^- مجنور است که در آن برنامه ریزی با دینامیک کاملا معلوم و کنترل کننده تکهای خطی انجام شده است.

شکل ۳-۳، معادله(۲۰-۳) را با دو جریمه c_i^- و یک یاداش c_k^+ نشان می c_i دهد. این شکل نشان می c_i دهد که اگر یک حالت نامطلوب و یک حالت مطلوب به هم نزدیک باشند، هزینه کلی c تا حدی بین دو هدف توازن ایجاد $x_*^+\in \operatorname{argmin} c^+(x)$ می کند. علاوه بر این، دیگر مثل گذشته حالت بهینه $x_*\in \operatorname{argmin} c(x)$ نیست: کمی دور شدن از حالت هدف (دور از حالت نامطلوب) بهینه است.

¹ Peaked



شکل ۳-۳ تابع هزینهای که محدودیتها (برای مثال موانع) را با "نامطلوب" کردن آنها در نظر می گیرد. منحنیهای با خطچین نشان داده شده اجزای منفرد c_j^+ و c_k^+ هستند، به معادله (۲۰-۳) رجوع می گیرد. منحنی با خط یکنواخت نشان داده شده جمع آنها c_j^- است. c_j^-

 $p(\mathbf{x}_t)$ امید ریاضیهای هزینه در معادله(۲۰-۳) و مشتقات نسبت به میانگین μ_t و کوواریانس Σ_t توزیع حالت (۱۹-۳) قاعده میتوانند برای هر c_j^+ و محاسبه شده و در آخر با هم جمع شوند. سپس، براساس معادله(۱۹-۳) قاعده زنجیرهای را برای جستجوی سیاست مبتنی بر گرادیان اعمال می کنیم.

بیان محدودیتها در قالب نامطلوب بودن در تابع هزینه در معادله(۳-۲۰)، هنوز اجازه برنامهریزی بلندمدت کاملا احتمالاتی را میدهد و امکان هدایت ربات در فضای حالت را فراهم می کند بدون اینکه با برخورد به موانع، آنها را "تجربه" کند.

برخوردهای درون چارچوب کاری استنتاج بیزی می توانند منع شوند، اما اکیدا از امید ریاضی حذف نمی شوند. این موضوع به این معنا نیست که میانگین گیری عدم قطعیتها اشتباه است – در عوض به ما می گوید که انتظار نمی رود که با یک اطمینان مشخص محدودیتها نقض شوند. توصیف صادقانه عدم قطعیت قابل پیش بینی غالبا باارزش تر از ادعا کردن با اطمینان کامل و گاه و بیگاه نقض غیر منتظره محدودیتهاست.

۴-۳- جمع بندی

در این فصل، روشهای پایه و پیشنهادی برای حل مسئله انتقال جسم به درون یک حفره توسط بازوی رباتیک شبیه سازی شده بیان شدند. برای شبیه سازی در محیطهای گسسته از الگوریتم Q-Learning استفاده شد و به عنوان روش پیشنهادی این پایان نامه، افزودن مسیرهای راهنما که از درون ویدئوهای ضبط شده استخراج شده اند شرح داده شد. برای شبیه سازی در محیطهای پیوسته از الگوریتم PILCO استفاده شد و قابلیت در نظر گرفتن محدودیتهای فضای حالت (مانند دیوار و مانع) به عنوان روش پیشنهادی به مدل PILCO افزوده شد.

فصل چهارم

۴- تحلیل و ارزیابی نتایج

تحلیل و ارزیابی نتایج

هدف از این فصل بررسی آزمایشهای صورت گرفته در زمینه ی این پایاننامه و ارزیابی عملکرد مدل پیشنهادی است. ابتدا در بخش $^{+}$ ۱ به معرفی مجموعه دادههای مورد استفاده در آزمایشها خواهیم پرداخت و ویژگیهای مجموعه دادههای مورد استفاده را بیان میکنیم. سپس در بخش $^{+}$ ۲ پیشپردازشهای صورت گرفته بر روی مجموعه داده شرح داده میشود. در بخش $^{+}$ ۳ معیارهای ارزیابی مورد استفاده توضیح داده میشوند و همچنین پارامترهای مختلف مدل به طور خلاصه بیان میشوند و در مورد نحوه ی تنظیم آنها بحث خواهد شد. در بخش $^{+}$ ۴ نیز آزمایشهای مختلف صورت گرفته برای بررسی پارامترهای مدل و همچنین مقایسه ی روش پیشنهادی با روشهای پیشین ارائه خواهند شد.

۱-۴- معرفی محیطهای شبیهسازی

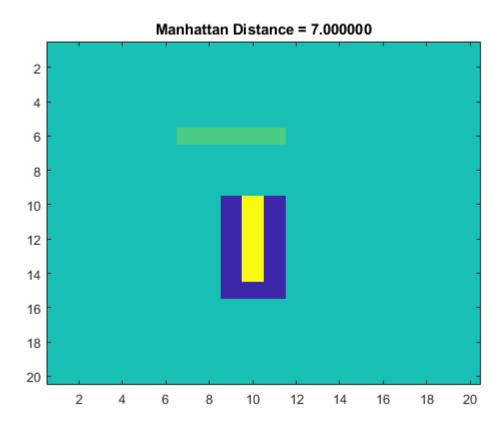
در این پایانامه، محیط شبیهسازی به دو صورت گسسته و پیوسته در نظر گرفته شده است. در محیط گسسته از الگوریتم Q-Learning و Q-Learning با راهنما برای یادگیری مسیر از موقعیت اولیه به موقعیت هدف استفاده شده است. همچنین، در محیط پیوسته از الگوریتم PILCO برای یادگیری مسیر از موقعیت اولیه به موقعیت هدف استفاده شده است. ویژگیهای محیط گسسته و محیط پیوسته به ترتیب در بخشهای ۲-۱-۱ و ۲-۱-۲ شرح داده شده اند.

۱-۱-۴- ویژگیهای محیط گسسته

محیط گسسته در نظر گرفته شده در این پایاننامه، یک محیط دوبعدی شبکهبندی شده است. این محیط دارای ابعاد $m \times n$ میباشد که m و n به ترتیب طول (جهت x) و عرض (جهت x) را مشخص میکنند. x_i میباشد که در آن x_i یک عدد طبیعی فرچک تر یا مساوی x_i است.

در این محیط شبیه سازی یک جسم وجود دارد که می تواند ابعاد مختلفی داشته باشد و با توجه به ابعاد آن و موقعیتش در محیط، سلول هایی از محیط را اشغال کند. ربات، جسم را از یک نقطه اولیه مشخص به یک نقطه هدف مشخص می رساند.

در این محیط شبیه سازی بعضی از نقاط به عنوان دیوار یا مانع در نظر گرفته شده اند که ربات امکان انتقال جسم به این خانه ها را ندارد. سه دیوار در محیط در نظر گرفته شده است که این سه دیوار در کنار یکدیگر تشکیل یک حفره را می دهند که موقعیت هدف می تواند در این حفره قرار گیرد. شکل ۲-۱ یک نمونه از محیط شبیه سازی گسسته با ابعاد ۲۰ در ۲۰ و جسم با ابعاد ۱ در ۵ را نشان می دهد.



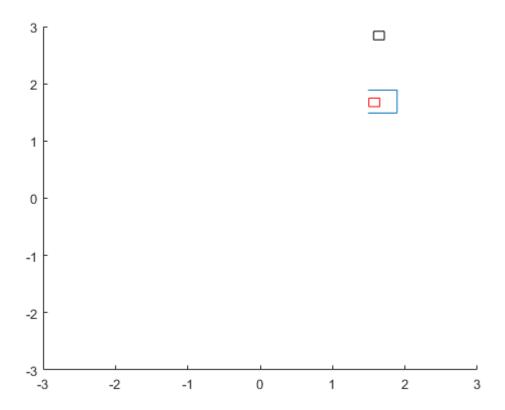
شکل ۱-۴ نمایی از یک محیط شبیهسازی گسسته. مستطیل سبز رنگ جسم گرفته شده توسط ربات را نشان می دهد که باید به موقعیت هدف که با مستطیل زرد رنگ مشخص شده است منتقل شود. موقعیت هدف درون حفرهای که توسط سه دیوار سرمهای رنگ مشخص شده قرار گرفته است.

۲-۱-۲ ویژگیهای محیط پیوسته

محیط پیوسته در نظر گرفته شده در این پایاننامه، یک محیط دوبعدی شده است. این محیط دارای ابعاد i محیط پیوسته در نظر گرفته شده در این پایاننامه، یک محیط دوبعدی شده است. این محیط $m \times n$ میباشد که m و m به ترتیب طول (جهت m) و عرض (جهت m) را مشخص می نقطه m در محیط شبیه سازی با زوج مرتب m (m) مشخص می شود که در آن m یک عدد حقیقی بین m و m است.

در این محیط شبیه سازی یک جسم وجود دارد که می تواند ابعاد مختلفی داشته باشد و با توجه به ابعاد آن و موقعیتش در محیط، بخشی از محیط را اشغال می کند. ربات، جسم را از یک نقطه اولیه مشخص به یک نقطه هدف مشخص می رساند.

در این محیط شبیه سازی بعضی از نقاط به عنوان دیوار یا مانع در نظر گرفته شده اند که امکان عبور ربات از این خانه ها وجود ندارد. سه دیوار در محیط در نظر گرفته شده است که این سه دیوار در کنار یکدیگر تشکیل یک حفره را می دهند که موقعیت هدف می تواند در این حفره قرار گیرد. این سه دیوار با سه پاره خط که دارای طول و معادله خط مشخص هستند معین می شود. معادله این خطوط به صورت x=a یا است که a و عدد ثابت را مشخص می کنند. شکل ۲-۴ یک نمونه از محیط شبیه سازی با ابعاد ۶ در ۶ و جسم با ابعاد a در ۱ در ۱ در ۱ در ۱ در ۱ نشان می دهد.



شکل ۴-۲ نمایی از یک محیط شبیهسازی پیوسته. مربع به رنگ مشکی، جسم گرفته شده توسط ربات را نشان می دهد که باید به موقعیت هدف که با مربع قرمز رنگ نشان داده شده است منتقل شود. موقعیت هدف درون حفرهای که توسط سه دیوار آبی رنگ مشخص شده قرار گرفته است.

۲-۲- نحوه آمادهسازی ویدئوهای مسیر راهنما

یکی از الگوریتمهای استفاده شده در این پایاننامه، Q-Learning با مسیر راهنما است. در این الگوریتم لازم است که مسیر های راهنما به عنوان ورودی به الگوریتم داده شوند. ما در این پایاننامه، مسیرهای راهنما را از ویدئوهایی که توسط خودمان ضبط شده است استخراج کردهایم. در این پروژه از چهار ویدئوی تهیه شده استفاده گردید که در هر ویدئو جسم از نقطه اولیه مشخصی توسط یک فرد به یک نقطه هدف که داخل یک حفره قرار گرفته است منتقل می شود. برای اینکه بتوانیم از مسیرهای نشان داده شده در ویدئوها در الگوریتم مورد نظر استفاده کنیم لازم است تا مختصات مربوط به جسم در ویدئو تشخیص داده شود. برای این منظور از الگوریتم ردیابی شی آ استفاده شده است. با اعمال این الگوریتم آبر روی ویدئوهای ضبط شده، مختصات مربوط به جسم تشخیص داده شده و مسیر راهنما مشخص می گردد. این مسیر، برای راهنمایی مدل به عنوان ورودی به الگوریتم داده می شود. در شکل - نمایی از یکی از این ویدئوها نشان داده شده است.



شکل ۴-۳ نمایی از ویدئوهای مسیر راهنما.

¹ Trajectory

² Object tracking

³ https://github.com/bikz05/object-tracker

۴-۳- تنظیم پارامترها و معیار ارزیابی

در این بخش قصد داریم به معرفی پارامترهای مربوط به الگوریتمهای مورد استفاده در این پایاننامه بپردازیم. در بخش ۴-۳-۱ به معرفی پارامترها و معیار ارزیابی الگوریتم PILCO میپردازیم. ادامه در بخش ۲-۳-۴ به معرفی پارامترها و معیار ارزیابی الگوریتم PILCO میپردازیم.

۹-۳-۱- پارامترها و معیار ارزیابی الگوریتم Q-Learning

پارامترهای مختلفی در الگوریتم Q-Learning وجود دارد که در این بخش به معرفی آنها و شرایط تنظیم آنها میپردازیم. یکی از ابتداییترین پارامترهای موجود در این الگوریتم، ابعاد فضای حالت است که به طور دلخواه تعیین میشود. در اکثر آزمایشها، فضای حالت به صورت مربعی در اندازههای 9×9 ، 81×18 و 36×36 در نظر گرفته شده است.

در ابتدای همه آزمایشها، موقعیت اولیه جسم نیز به طور دلخواه مشخص می شود. در تمامی آزمایشهای صورت گرفته در این پایان نامه، موقعیت اولیه جسم در ناحیه بالا سمت راست فضای حالت در نظر گرفته شده است. هم چنین، موقعیت نهایی برای جسم به عنوان یکی از پارامترهای آزمایش مطرح است که در ابتدای آزمایش باید مشخص شود. در اکثر آزمایشها، این پارامتر وسط فضای حالت و درون حفره در نظر گرفته شده است.

مشابه بسیاری از الگوریتمهای یادگیری، در الگوریتم Q-Learning نیز پارامتری تحت عنوان ضریب یادگیری الگوریتم که با آلفا (α) نشان داده می شود مطرح است. این پارامتر، نسبت حرکت الگوریتم در جهت اکتشاف موقعیتهای جدید و بهرهبرداری آز دانش پیشین را مشخص می کند. مقدار این پارامتر عددی بین صفر و یک است که هر چه به یک نزدیک تر باشد، الگوریتم بیش تر اکتشاف می کند و هر چه به صفر نزدیک تر باشد، بیش تر باشد، بیش تر به بهرهبرداری از دانش پیشین روی می آورد. با توجه به آزمایشهای صورت گرفته، در اکثر آزمایشات مقدار (α) برای این پارامتر در نظر گرفته شده است.

یکی دیگر از پارامترهای مطرح در این الگوریتم، پارامتر گاما (γ) میباشد که ضریب کاهش نام دارد. بدیهی است پاداش دریافتی در دورههای قبلی دارد.

¹ Exploration

² Exploitation

³ Discount factor

پارامتر ضریب کاهش وظیفه تنظیم تاثیر تدریجی پاداشها در تکرارهای مختلف را دارد، به این صورت که هر چه این عدد بزرگتر باشد اثر پاداشهای قبلی دیرتر از بین میرود. مقدار این پارامتر عددی بین صفر و یک است که با توجه به بررسیهای انجام شده در آزمایشهای مختلف، در اکثر آزمایشات مقدار ۵/۰ به عنوان مقدار مناسب برای این پارامتر در نظر گرفته شده است.

با توجه به وجود روند تکراری در الگوریتم Q-Learning یک شرط خاتمه برای آن باید در نظر گرفته شود که از پارامتر تعداد دور تکرار به این منظور استفاده شده است. الگوریتم بعد از رسیدن به تعداد تکرار لازم متوقف می شود و این نقطه از اجرا به عنوان نقطه همگرایی در نظر گرفته می شود. در تمام آزمایشات این پارامتر ۱۰۰ در نظر گرفته شده است چرا که با بررسی آزمایشات مختلف به این نتیجه رسیده ایم که بعد از این تکرار تغییر قابل توجهی در نتایج الگوریتم دیده نمی شود.

در الگوریتم Q-Learning یک ماتریس Q وجود دارد که تخمین مناسب بودن هر عمل در هر حالت را نشان میدهد. این ماتریس در ابتدا مقداردهی اولیه شده و در طی اجرای الگوریتم به تدریج به مقادیر مناسب هر عمل در هر حالت همگرا میشود. تعیین مقدار اولیه مناسب برای این ماتریس میتواند باعث همگرایی سریع تر و بهتر الگوریتم شود. لذا، این مقدار اولیه میتواند به عنوان یکی از پارامترهای الگوریتم در نظر گرفته شود. این ماتریس در الگوریتم ولیه میشود در حالی که در الگوریتم سیر راهنما، مقدار اولیه این ماتریس با توجه به دانش استخراج شده از ویدئوهای راهنما تعیین می گردد.

برای مقایسه آزمایشات مختلف که با پارامترهای مختلف و در شرایط مختلف صورت می گیرند، نیاز به یک معیار ارزیابی مقایسهای وجود دارد. معیارهای ارزیابی مختلفی در کارهای پیشین در نظر گرفته شده است که ما در این پایاننامه از معیار میانگین تعداد عملهای لازم برای رسیدن به هدف استفاده کردهایم. در هر تکرار الگوریتم جسم بعد از انجام تعدادی عمل به هدف می رسد. اگر تعداد عملهای انجام شده توسط ربات در تکرار i ام را با i نشان دهیم در الگوریتمی که دارای k تکرار می باشد، معیار میانگین تعداد عملهای لازم برای رسیدن به هدف با توجه به معادله i محاسبه می شود:

$$avgActCount = \frac{1}{k} \times \sum_{i=1}^{k} n_i$$
 (۱-۴)معادله

۲-۳-۲- پارامترها و معيار ارزيابي الگوريتم PILCO

در الگوریتم PILCO، پارامترهای بسیار زیادی وجود دارد. این پارامترها، در دستههای «تنظیمات کلی»، «تنظیمات تابع سیاست»، «تنظیمات تابع هزینه»، «تنظیمات مدل دینامیکی GP» و «تنظیمات بهینه سازی (پارامترهای سیاست)» قرار می گیرند. در این بخش، به تفکیک به معرفی اجمالی پارامترهای هر یک از این دستهها می پردازیم.

• تنظيمات كلى

پارامتر th بیانگر زمان نمونهبرداری است. یا به عبارتی دیگر 1/dt فرکانس نمونهبرداری را مشخص می کند. یکی دیگر از پارامترهای کلی، پارامتر T است که طول افق پیشبینی در واحد ثانیه را نشان می دهد. پارامتر μ 0 میانگین توزیع حالت آغازی مرکز جسم، یعنی μ 1 را مشخص می کند. پارامتر μ 2 نشان دهنده ماتریس کوواریانس توزیع μ 2 است. μ 3 است. μ 4 تعداد دفعاتی است که حلقه اصلی الگوریتم PILCO اجرا می شود. μ 5 تعداد مالی الولیه با یک سیاست تصادفی است. این rollout ها برای جمع آوری یک مجموعه داده اولیه برای آموزش اولین مدل دینامیکی فرآیند گوسی استفاده می شوند. معمولا این پارامتر ۲ در نظر گرفته می شود. پارامتر μ 4 تعداد حالات اولیه ای است که برای آنها سیاستی آموخته می شود. در PILCO توزیعهای حالت آغازی مرکز جسم به صورت معادله (۲-۴) بیان می شود که به توزیعهای با میانگینهای μ 4 متفاوت و اما با ماتریسهای کوواریانس مشتر ک μ 5 مربوط است.

$$p(x_0) = \sum_{i=1}^K \mathcal{N}(\mu_0^{(i)}, \Sigma_0)$$
 (۲-۴) معادله

• تنظيمات Plant

dynamics پارامتری است که تابع 'ODE شبیه سازی سیستم را مشخص می کند. در حالتی که از ربات واقعی استفاده شود، نیازی به مشخص کردن معادله مدل دینامیکی نیست. اما در حالت شبیه سازی لازم است: است معادله مدل دینامیکی مسئله شبیه سازی شده را مشخص کنیم که در معادله (۳-۴) آمده است:

_

¹ Ordinary Differential Equation

$$rac{dz}{dt} = egin{bmatrix} Z_1 \\ Z_2 \\ rac{f_1(t) - b imes z_3}{m} \\ rac{f_2(t) - b imes z_4}{m} \end{bmatrix}$$
 (٣-۴) معادله

در معادله $f_1(t)$ نیروی وارد بر جسم در حالت $z=[x,y,v_x,v_y]$ نیروی وارد بر جسم در معادله $z=[x,y,v_x,v_y]$ نیروی وارد بر جسم در جهت $z=[x,y,v_x,v_y]$

پارامتر noise ماتریس کوواریانس نویز اندازه گیری در نظر گرفته شده در فرآیند شبیهسازی است. در آزمایشهای صورت گرفته در این پایاننامه، فرض بر این است که از یک نویز گوسی با میانگین صفر استفاده شده است. این نویز، در طول rollout تراژکتوری به حالت (پنهان) اضافه می شود.

پارامتر ctrl، بیانگر کنترلکنندهای است که باید اعمال شود. در آزمایشهای صورت گرفته در این پایاننامه، first-order-hold و -first و -first و -first مستند.

prop، تابعی را مشخص می کند که توزیع $p(x_{t+1})$ را با استفاده از توزیع $p(x_t)$ (پیشبینی احتمالاتی یک گامه) محاسبه می کند. در آزمایشهای صورت گرفته، از تابعی به نام propagated استفاده شده که توزیع حالت پیشبینی کننده $p(x_{t+1})$ و مشتقات جزئی که برای جستجوی سیاست مبتنی بر گرادیان نیاز هستند را محاسبه می کند.

• تنظیمات تابع سیاست

b و w تابع سیاستی که در آزمایشهای این پایاننامه استفاده شده، تابع خطی است که پارامترهای آن w و w هستند. این دو پارامتر به ترتیب، وزنهای خطی (ماتریس w w w و بایاس (ماتریس w w و بایاس (ماتریس w w و بایاس (ماتریس w و نظر گرفته شده است. توضیحات بیش تر در مورد این سیاست، در بخش روش پیشنهادی آورده شده است.

_

¹ State

سمی الگوریتم u را معین می کند. به این صورت که الگوریتم u را معین می کند. به این صورت که الگوریتم $[-u_{max},u_{max}]$ در بازه ی

• تنظیمات تابع هزینه

در این بخش، پارامترهای تابع هزینه آنی را معرفی می کنیم. پارامتر α ، بیانگر تابع هزینه مورد استفاده است. در آزمایشهای صورت گرفته، از تابع هزینه اشباع شده (که با تفریق یک گوسی نرمال نشده از عدد است. در آزمایشهای صورت گرفته، از تابع هزینه اشباع شده (که با تفریق یک گوسی نرمال نشده از معادله، است. در این معادله، احاصل می شود) با پراکندگی σ_c استفاده شده است که در معادله (۴-۴) بیان شده است. در این معادله PILCO حالت هدف را مشخص می کند که یکی دیگر از پارامترهای مورد استفاده در الگوریتم است که باید از قبل معین شود.

$$c(x) = 1 - exp\left(-\frac{1}{2\sigma_c^2} \|x - x_{target}\|^2\right) \in [0,1]$$
 هعادله (۴-۴) معادله

gamma، پارامتر ضریب کاهش است و از آنجایی که مسئله مورد نظر ما، یک مسئله با افق متناهی است، مقدار این پارامتر را ۱ در نظر گرفتهایم.

پارامتر σ_c ، بیانگر پراکندگی یا عرض تابع هزینه است. بر حسب تجربه، بهتر است این پارامتر $\frac{\|\mu_0 - x_{target}\|}{10}$ در نظر گرفته شود.

expl، پارامتر اکتشاف از نوع 'UCB است. مقادیر منفی، باعث اکتشاف بیشتر میشوند و مقادیر مثبت، باعث ماندن سیاست در مناطق با عملکرد پیشبینی خوب میشوند. در آزمایشهای صورت گرفته، مقدار این پارامتر را صفر در نظر گرفته ایم تا از هر نوع اکتشاف یا بهرهبرداری اضافی جلوگیری کنیم.

• تنظیمات مدل دینامیکی GP

در این بخش، پارامترهای مدل دینامیکی GP را معرفی می کنیم.

فی پراکنده را مشخص GP یک پارامتر اختیاری است و زمان تعویض از $a \times b \times c$ های پراکنده را مشخص induce $a \times b \times c$ است که در آن $a \times b \times c$ است. عداد ورودیهای القایی است.

¹ UCB-type

Inducing iputs $^{ec{v}}$ اگر اندازه مجموعه آموزشی از a بزرگ تر شود، G کامل به طور خودکار به تقریب پراکندهاش تغییر میcند.

برابر با ۰ مشخص می کند که هیچ ورودی القایی وجود ندارد. c یا برابر با ۱ است یا برابر با تعداد ابعاد قابل پیشبینی است که به تعداد خروجی های مدل دینامیکی مربوط می شود. برای c=1 ورودی های القایی بین همه c=1 ها به اشتراک گذاشته می شوند. در غیر اینصورت، مجموعه های ورودی های القایی، برای هر بعد قابل پیشبینی به طور جداگانه یاد گرفته می شوند.

پارامتر trainOpt، تعداد جستجوهای خطی برای آموزش فراپارامتر GP را مشخص می کند که (عدد اول) توسط GP کامل و (عدد دوم) توسط GP پراکنده استفاده می شود.

• تنظیمات بهینهسازی (یادگیری سیاست)

در این بخش، پارامترهای اختیاری برای یادگیری سیاست را تعریف می کنیم. به طور کلی، از یک بهینه ساز مبتنی بر گرادیان غیرمحدب استفاده شده است.

پارامتر length، حداکثر تعداد جستجوهای خطی را مشخص میکند که بهینهساز بعد از آن، بهترین مجموعه پارامتر تاکنون را برمی گرداند.

پارامتر MFEPLS، حداکثر تعداد ارزیابیهای تابع در هر جستجوی خطی است. یا جستجوی خطی، موفق MFEPLS به پیدا کردن یک مجموعه پارامتر با گرادیان نزدیک به صفر میشود یا موفق نمیشود و بعد از MFEPLS بار ارزیابی تابع (و گرادیان) متوقف میشود.

برای مقایسه آزمایشات مختلف که با پارامترهای مختلف و در شرایط مختلف صورت می گیرند، نیاز به یک معیار ارزیابی مقایسهای وجود دارد. معیارهای ارزیابی مختلفی در کارهای پیشین در نظر گرفته شده است که ما در این پایاننامه از معیار درصد موفقیت ربات در رساندن جسم به هدف استفاده کردهایم. هر آزمایش شامل n بار تلاش ربات برای رساندن جسم از موقعیت اولیه به موقعیت هدف است، اگر در این n بار تکرار ربات موفق شود که جسم را در n تکرار به هدف برساند، معیار درصد موفقیت ربات در رساندن جسم به هدف با استفاده از معادله n محاسبه می شود:

$$successRate = \frac{k}{n} \times 100$$
 (4-4)

۴-۴- نتایج آزمایش

در این بخش آزمایشهای مختلف صورت گرفته در زمینه ی بررسی پارامترهای مختلف بر نتایج انتقال جسم به درون حفره در محیطهای گسسته و پیوسته گزارش شدهاند. در ادامه، به ترتیب بخشهای «بررسی تاثیر اندازه جسم گرفته شده توسط ربات»، «بررسی تاثیر قابلیت چرخاندن جسم گرفته شده توسط ربات»، «بررسی تاثیر اندازه فضای حالت» و «بررسی تاثیر در نظر گرفتن محدودیتهای فضای حالت در محیط پیوسته» خواهند آمد.

۱-۴-۴ بررسی تاثیر اندازه جسم گرفته شده توسط ربات

اندازه جسم گرفته شده توسط ربات، یکی از عوامل موثر در عملکرد الگوریتم به کار گرفته شده توسط ربات است. با توجه به این موضوع در این آزمایش قصد داریم تا به بررسی عملکرد ربات در حالات مختلف اندازه جسم گرفته شده توسط ربات بپردازیم. آزمایش بر اساس دو اندازه یک سلولی و دو سلولی انجام شده است و از معیار میانگین تعداد عملهای انجام شده توسط ربات برای رسیدن به مقصد، به منظور مقایسه حالتها استفاده شده است. لیست پارامترهای مورد استفاده در این آزمایش، در جدول ۲-۱ آمده است.

جدول ۴-۱ تنظیمات پارامترهای مورد استفاده در آزمایش بررسی تاثیر اندازه جسم گرفته شده توسط ربات

توضیح	مقدار	پارامتر
فضای حالت مربعی در نظر گرفته شده است	۹،۱۸،۳۶	فضای حالت
	Q-Learning	مدل مورد بررسی
جسم در نقطه راست بالای فضای حالت قرار می گیرد	راست بالا	موقعيت اوليه جسم
جسم باید درون حفره قرار گیرد	درون حفره	موقعیت نهایی جسم
ضریب یادگیری	٠/٨	آلفا (α)
ضریب عامل کاهش	٠/۵	گاما (۲ٖ)
	1	تعداد دور تکرار
		الگوريتم
	ندارد	قابلیت چرخاندن جسم
پارامتر مورد آزمایش	1,7	اندازه جسم

جدول ۲-۴ نتایج حاصل از آزمایش بررسی تاثیر اندازه جسم گرفته شده توسط ربات

میانگین تعداد عملها	اندازه جسم برحسب تعداد سلول	فضاي حالت
V8/84	١	
۸۹/۸۶	٢	9 × 9
WY - / W9	١	
44./84	۲	18 × 18
Y1A9/Y	١	
٣٣. 8/A	۲	36 × 36

همانطور که از قبل قابل پیشبینی بود، انتقال جسم بزرگتر به هدف پیچیدگی بیشتری دارد لذا به تعداد عملهای بیشتری نیاز دارد. همانطور که از نتایج گزارش شده در جدول ۴-۲ نیز مشخص است، در تمامی فضاهای حالت آزمایش شده تعداد عملهای لازم برای انتقال جسم با اندازه ۲ سلول به هدف بیشتر از تعداد عملهای لازم برای انتقال جسم با اندازه ۱ سلول است.

۲-۴-۴ بررسی تاثیر قابلیت چرخاندن جسم گرفته شده توسط ربات

جهت جسمی که در اختیار ربات قرار گرفته است، می تواند در نحوه انتقال جسم به هدف موثر باشد. در بسیاری از موارد، چرخاندن جسم می تواند باعث سهولت در انتقال آن به هدف شود. از این نقطه نظر، وجود توانایی چرخاندن جسم توسط ربات می تواند عامل موثری در عملکرد ربات باشد. با توجه به این موضوع در این آزمایش قصد داریم تا به بررسی عملکرد ربات در دو حالت مختلف بپردازیم. در حالت اول، ربات توانایی چرخاندن جسم را ندارد ولی در حالت دوم، این توانایی برای ربات در نظر گرفته شده است. آزمایش بر اساس این دو حالت انجام شده است و از معیار میانگین تعداد عملهای انجام شده توسط ربات برای رسیدن به مقصد، به منظور مقایسه حالتها استفاده شده است. لیست پارامترهای مورد استفاده در این آزمایش، در جدول ۴-۳ به طور خلاصه آمده است.

جدول ۴-۳ تنظیمات پارامترهای مورد استفاده در آزمایش بررسی تاثیر قابلیت چرخاندن جسم گرفته شده توسط ربات

توضیح	مقدار	پارامتر
فضای حالت مربعی در نظر گرفته شده است	۹،۱۸،۳۶	فضای حالت
	Q-Learning	مدل مورد بررسی
جسم در نقطه راست بالای فضای حالت قرار می گیرد	راست بالا	موقعيت اوليه جسم
جسم باید درون حفره قرار گیرد	درون حفره	موقعیت نهایی جسم
ضریب یادگیری	٠/٨	(α) آلفا
ضریب عامل کاهش	٠/۵	گاما (γ)
	١	تعداد دور تکرار
		الگوريتم
پارامتر مورد آزمایش	دارد، ندارد	قابلیت چرخاندن جسم
	٢	اندازه جسم

جدول ۴-۴ نتایج حاصل از بررسی تاثیر قابلیت چرخاندن جسم گرفته شده توسط ربات

میانگین تعداد عملها	قابلیت چرخاندن جسم	فضاي حالت
184/04	دارد	
ለዓ/ለ۶	ندارد	9 × 9
۹۰۰/۷۵	دارد	
44.184	ندارد	18 × 18
	دارد	
WW + 8/A	ندارد	36 × 36

در حالتی که ربات قابلیت چرخاندن جسم را نداشته باشد، در صورتی که در ابتدای کار جهت جسم معادل با جهت هدف جسم نباشد امکان رسیدن جسم به هدف وجود ندارد. با توجه به این موضوع، در تمام آزمایشهای صورت گرفته در حالتی که ربات قابلیت چرخاندن جسم را ندارد، جهت اولیه جسم معادل با جهت نهایی جسم در نظر گرفته شده است. بنابراین، پیچیدگی مسئله در حالت وجود قابلیت چرخاندن جسم بیش تر از حالت دیگر است چرا که علاوه بر ۴ عمل حرکت در چهار جهت اصلی، باید عمل چرخاندن جسم نیز توسط ربات صورت گیرد. با توجه به پیچیدگی بیش تر موجود در این حالت، تاخیر در رساندن جسم به هدف قابل پیشبینی است. نتایج گزارش شده در جدول ۴-۴ نیز موید این موضوع در تمامی فضاهای حالت می باشد.

۳-۴-۴- بررسی تاثیر تعداد مسیرهای راهنما

در حالت معمول الگوریتم Q-Learning که ماتریس Q با صفر مقداردهی اولیه می شود، در فضای حالت با ببعاد بزرگتر از 20×20 ربات در زمان معقول نمی تواند جسم را با شروع از موقعیت اولیهای که بیرون از حفره قرار دارد به موقعیت هدف درون حفره برساند. با توجه به این موضوع، یک راه حل برای کاهش زمان اجرا می تواند استفاده از مسیرهای راهنما باشد. این مسیرهای راهنما در قالب مقدار اولیه ماتریس Q در اختیار الگوریتم Q-Learning قرار می گیرند. مسیرهای راهنما می توانند از دو طریق در اختیار الگوریتم قرار بگیرند. در حالت اول، مسیرهای رساندن جسم به هدف به صورت دستی توسط عامل انسانی طراحی شده و در اختیار الگوریتم قرار گرفتهاند. استفاده از این روش، در فضاهای حالت با ابعاد بزرگ دارای پیچیدگی زیادی است. در حالت دوم، این مسیرهای راهنما از درون ویدئوهای ضبط شده و با استفاده از الگوریتمهای ردیابی شی استخراج می شوند. در هر دو حالت، تعداد مسیرهای راهنما که در اختیار الگوریتم قرار می گیرند می توانند در عملکرد الگوریتم تاثیر گذار باشند. هدف از این آزمایش بررسی تاثیر تعداد مسیرهای راهنما است. لیست پارامترهای مورد استفاده در این آزمایش، در جدول 4-8 به طور خلاصه آمده است. نتایج آزمایشهای صورت گرفته در سه حالت جسم ۱ سلولی، ۲ سلولی، ۲ سلولی با چرخش به ترتیب در جدول 4-8 ، جدول 4-8 و جدول 4-8 آمده است.

جدول ۴-۵ تنظیمات پارامترهای مورد استفاده در آزمایش بررسی تاثیر تعداد مسیرهای راهنما

توضيح	مقدار	پارامتر
فضای حالت مربعی در نظر گرفته شده است	9.11.75.144	فضای حالت
مسیرهای راهنما به صورت دستی در اختیار الگوریتم قرار گرفته است	Q- Learning with guidance	مدل مورد بررسی
جسم در نقطه راست بالای فضای حالت قرار می گیرد	راست بالا	موقعيت اوليه جسم
جسم باید درون حفره قرار گیرد	درون	موقعیت نهایی جسم
ضریب یادگیری	٠/٨	آلفا (α)
ضریب عامل کاهش	٠/۵	گاما (γ)
	1	تعداد دور تكرار الگوريتم
	دارد، ندارد	قابلیت چرخاندن جسم
	1.7	اندازه جسم

جدول ۴-۶ نتایج حاصل از بررسی تاثیر تعداد مسیرهای راهنما در انتقال جسم به اندازه ۱ سلول

میانگین تعداد عملها	تعداد مسیرهای راهنما	فضاي حالت
WW/10	١	
Y1/89	٢	9 × 9
۲۷/۵۱	٣	
71/87	۴	
194/77	1	
147	٢	18 × 18
114/08	٣	19 17 18
79/92	۴	
77.5	١	
1.78/4	٢	36 × 36
1 • 8 ٧/٧	٣	351150
۵۵۳/۲۲	۴	

جدول ۴-۷ نتایج حاصل از بررسی تاثیر تعداد مسیرهای راهنما در انتقال جسم به اندازه ۲ سلول بدون قابلیت چرخاندن آن

میانگین تعداد عملها	تعداد مسیرهای راهنما	فضاي حالت
۳۹/۸۵	1.	
7./71	۲	9 × 9
17/44	٣	
11/19	۴	
744/18	١	
744/71	۲	18 × 18
۸٩/٧۴	٣	
۵۲/۳	۴	
1078/4	١	
1.79/	۲	36 × 36
1777/4	٣	
900/90	۴	

جدول ۴-۸ نتایج حاصل از بررسی تاثیر تعداد مسیرهای راهنما در انتقال جسم به اندازه ۲ سلول با قابلیت چرخاندن آن

میانگین تعداد عملها	تعداد مسیرهای راهنما	فضاي حالت
74277/94	٣	
1246/14	۴	144 × 144

همان طور که از نتایج ارائه شده در جداول بالا نیز مشخص است، افزایش تعداد مسیرهای راهنما توانسته است میانگین تعداد عملهای لازم برای رسیدن به هدف را کاهش دهد. البته توجه به این موضوع حائز اهمیت است که افزودن مسیر راهنما باید دارای ارزش افزودهای برای الگوریتم باشد به این معنی که اگر چندین مسیر همگی در یک بخش از فضای حالت طراحی شده باشند و پراکندگی لازم را در فضای حالت نداشته باشند، لزوما با کاهش معیار ارزیابی روبرو نیستیم. یعنی تنها افزودن مسیر راهنما در بخشهایی از فضای حالت که در آنها مسیری طراحی نشده است، می تواند باعث بهبود عملکرد ربات گردد.

۴-۴-۴ بررسی تاثیر اندازه فضای حالت

در فضاهای حالت بزرگ، اجرای الگوریتم Q-Learning بدون استفاده از مسیر راهنما عملا امکانپذیر نیست و حتما باید از مسیر راهنما استفاده شود. طراحی دستی مسیر راهنما در فضاهای بزرگ با پیچیدگی زیادی همراه است و در نتیجه در چنین حالتی از مسیرهای استخراج شده از ویدئو استفاده می شود. در این آزمایش، هدف استفاده از مسیر استخراج شده از ویدئو در ابعاد مختلف فضای حالت می باشد. لیست پارامترهای مورد استفاده در این آزمایش، در جدول +- به طور خلاصه آمده است. در همه آزمایشهای این بخش از چهار مسیر راهنما استفاده شده است.

جدول ۴-۹ تنظیمات پارامترهای مورد استفاده در آزمایش بررسی تاثیر اندازه فضای حالت

توضیح	مقدار	پارامتر
پارامتر مورد بررسی (در جدول نتایج آورده شده است)		فضای حالت
	Q- Learning with guidance	مدل مورد بررسی
جسم در نقطه راست بالای فضای حالت قرار می گیرد	راست بالا	موقعیت اولیه جسم
جسم باید درون حفره قرار گیرد	درون حفره	موقعیت نهایی جسم
ضریب یادگیری	٠/٨	(α) آلفا
ضریب عامل کاهش	٠/۵	گاما (<i>ץ</i>)
	1	تعداد دور تكرار الگوريتم
	دارد	قابلیت چرخاندن جسم
در جدول نتایج آورده شده است		اندازه جسم

جدول ۴-۱۰ نتایج حاصل از بررسی تاثیر اندازه فضای حالت

میانگین تعداد عملها	اندازه جسم	فضاي حالت
VW 1 W/9	4 × 2	69 × 39
78987	5 × 3	86 × 49
۶۰۸۷۴/۸۸	9 × 5	138 × 80
180746/14	11 × 6	172 × 99
Y4187/81	15 × 9	231 × 132

۵-۴-۴- بررسی تاثیر در نظر گرفتن محدودیتهای فضای حالت در محیط پیوسته

در محیطهای پیوسته، از الگوریتم PILCO برای انتقال جسم به هدف استفاده کردهایم. این الگوریتم، از یک تابع هزینه برای ارزیابی سیاست استفاده میکند. در تابع هزینههای معمول هیچگونه اطلاعاتی در مورد محدودیتهای فضای حالت (مانند دیوارها و موانع موجود در فضای حالت) در نظر گرفته نمی شود. تغییر تابع هزینه به نحوی که بتوان محدودیتهای فضای حالت را به آن تزریق کرد می تواند نقش بسیار موثری در عملکرد ربات داشته باشد. در تابع هزینههای معمول از فاصله اقلیدسی جسم تا هدف به عنوان یک پارامتر جهتدهنده استفاده می شود. این در حالی است که در فضای حالتی که در آن مانع یا دیوار وجود داشته باشد، نزدیکی اقلیدسی به هدف لزوما به معنی مناسب بودن سیاست یاد گرفته شده نیست چرا که ممکن است جسم در فاصلهای بسیار نزدیک به هدف قرار داشته باشد اما دیوار مانع رسیدن به هدف شود. برای رویارویی با این چالش لازم است که میزان نزدیکی جسم به دیوار نیز به عنوان پارامتری جهتدهنده مورد توجه قرار بگیرد. در این پایان نامه، تابع هزینه به گونهای ساخته شده است که در آن هم مقایسه تابع هزینه پیشین که هیچ توجهی به محدودیتهای فضای حالت ندارد با تابع هزینه ارائه شده که معدودیتهای فضای حالت ندارد با تابع هزینه ارائه شده که محدودیتهای فضای حالت زیر نیز در نظر می گیرد است. مجموعه پارامترهای مورد استفاده در این آزمایش در جدول ۴-۱۱ به طور کامل آورده شده است. نتایج آزمایش نیز برای موقعیتهای اولیه مختلف جسم نیز در جدول ۴-۱۲ گزارش شده است.

جدول ۴-۱۱ تنظیمات پارامترهای مورد استفاده در آزمایش بررسی تاثیر در نظر گرفتن محدودیتهای فضای حالت

توضيح	مقدار	پارامتر
	PILCO	مدل مورد بررسی
زمان نمونهبرداری	•/1	dt
طول افق پیشبینی	١٠	T

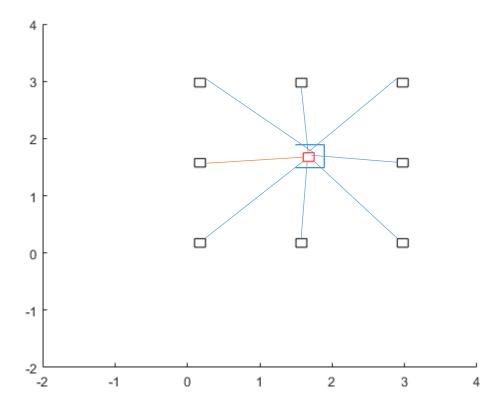
$\left[x,y,v_{x},v_{y} ight]$ بردار به شکل	پارامتر مورد بررسی	μ_0
	ماتریس با درایههای قطری	Σ_0
	۰/۰۱ و بقیه درایههای صفر	
جسم باید در نقطهای درون	(۰،۰،۱/۶،۱/۵)	حالت نهایی جسم
حفره و به طور ساکن قرار گیرد		
تعداد تکرارها برای بهینهسازی	۵۰	N
	١	J
۸ نقطه شروع اولیه مختلف در	٨	K
نظر گرفته شده است		
از مدل دینامیکی در نظر گرفته	dynamics_mini	dynamics
شده برای این شبیهسازی که در		
بخش روشهای پیشنهادی		
توضیح داده شده استفاده شده		
است		
نویز اندازهگیری هر بعد از حالت	ماتریس با درایههای قطری	noise
۰/۰۰۱ در نظر گرفته شده است	۰/۰۰۱ و بقیه درایههای صفر	
	zero-order-hold	ctrt
	propagated	prop
حداکثر نیروی وارد بر جسم در	(۵۰،۵۰)	maxU
هر جهت ۵۰ نیوتن در نظر گرفته	(/	
شده است		

c c	تابع هزينه اشباع شده	به معادله (۴-۴) رجوع شود
گاما (γ)	١	ضریب کاهش
σ_c		به فاصله بین موقعیت اولیه و
		نهایی جسم بستگی دارد
expl	•	عرض تابع هزینه که از نوع UCB است
طول و عرض جسم	(•/١۵، •/١۵)	در آزمایشات این بخش جسم مربعی در نظر گرفته شده است
	ماتریس با درایههای صفر و با ابعاد (۱٬۰٬۳۰۰)	ورودیهای القایی مشترک
trainOpt	(۵۰۰،۳۰۰)	حداکثر تعداد جستجوهای خطی به هنگام آموزش مدل دینامیکی GP
length	10.	حداکثر تعداد جستجوهای خطی
MFEPLS	٣٠	حداکثر تعداد ارزیابیها به ازای هر جستجوی خطی

جدول ۲-۴ نتایج حاصل از بررسی تاثیر در نظر گرفتن محدودیتهای فضای حالت

درصد موفقیت	دانش در مورد محدودیتهای فضای حالت	μ_0
۵۸	دارد	
45	ندارد	[0.1,0.1,0,0]
۵۶	دارد	
۴۲	ندارد	[1.5,0.1,0,0]
۵٠	دارد	
٣٠	ندارد	[2.9,0.1,0,0]
44	دارد	
۲٠	ندارد	[2.9,1.5,0,0]
۵۴	دارد	
774	ندارد	[2.9,2.9,0,0]
۵۴	دارد	[1.5,2.9,0,0]
۴.	ندارد	
۵۴	دارد	[0.1,2.9,0,0]
45	ندارد	
99	دارد	[0.1,1.5,0,0]
٨٠	ندارد	

همانطور که نتایج گزارششده در جدول ۱۲-۴ نشان می دهد، در اکثر موارد (به جز حالت = μ_0 [0.1,1.5,0,0] که موقعیت ابتدایی جسم دقیقا مقابل حفره قرار دارد) تابع هزینه ای که محدودیتهای فضای حالت را نیز در نظر گرفته توانسته است براساس معیار ارزیابی درصد موفقیت رساندن جسم به هدف، نسبت به تابع هزینه ای که توجهی به محدودیتهای فضای حالت ندارد عملکرد بهتری داشته باشد. برخلاف تمامی موقعیتهای اولیه، در حالت [0.1,1.5,0,0] = μ با کاهش عملکرد براساس معیار ارزیابی درصد موفقیت رساندن جسم به هدف روبرو هستیم. تفاوت اصلی که بین این نقطه شروع و سایر نقاط اولیه وجود دارد این است که در نزدیک ترین مسیر از نظر فاصله اقلیدسی از این نقطه به هدف مانعی مشاهده نمی شود، این در حالی است که کوتاه ترین مسیر از نظر فاصله اقلیدسی از سایر نقاط اولیه به نقطه مشاهده نمی شود، این در حالی است که کوتاه ترین مسیر از نظر فاصله اقلیدسی از سایر نقاط اولیه به نقطه هدف توسط مانع بسته می شود. این موضوع در شکل ۴-۴ مشاهده می گردد.



شکل ۴-۴ کوتاه ترین مسیر از نظر فاصله اقلیدسی از موقعیتهای اولیه مختلف به هدف.

به طور کلی، در مواقعی که در کوتاه ترین مسیر از نظر فاصله اقلیدسی از نقطه اولیه جسم به هدف با مانع روبرو نیستیم، در نظر گرفتن پارامتر جهت دهی مرتبط با موانع در تابع هزینه مورد استفاده می تواند باعث ایجاد پیچیدگی برای الگوریتم گردد. برای مثال، در مواقعی که جسم بدون برخورد با دیوار، اما با فاصلهای

کم از آن وارد حفره می شود پارامتر جهت دهی مرتبط با موانع باعث انحراف جسم از مسیر معمول شده و باعث ایجاد اختلال در عملکرد الگوریتم می گردد. نتایج گزارش شده در جدول ۲-۱۲ برای نقاط اولیه دیگر نیز این طور نشان می دهد که هر چه مسیرهای کوتاه رساندن جسم از نقطه اولیه به هدف بیش تر با موانع همراه باشد، افزودن پارامتر جهت دهی مرتبط با موانع کمک به بیش تری به راهنمایی ربات می نماید.

۵-۴- جمع بندی

در این فصل ابتدا به معرفی مجموعه محیطهای شبیهسازی گسسته و پیوسته مورد استفاده در آزمایشهای انجام گرفته پرداختیم. سپس پارامترهای مختلف الگوریتم Q-Learning و الگوریتم PILCO و روش پیشنهادی بیان شدند و در مورد نحوه ی تنظیم هر کدام از آنها بحث شد. معیارهای ارزیابی مورد استفاده برای مقایسه ی مدلها نیز معرفی شدند. در بخش انتهایی این فصل، آزمایشهای مختلفی که در طول انجام این پروژه برای بررسی مدل پایه و همچنین مقایسه ی روش پیشنهادی با روشهای پیشین مورد استفاده قرار گرفتهاند آورده شده است. آزمایشها نشان میدهند که روش پیشنهادی ارائه شده در این پایان نامه که از تابع هزینه برای تزریق دانش در مورد فضای حالت در الگوریتمهای یادگیری تقویتی که به مساله ی انتقال جسم به درون حفره می پردازند استفاده می کند، توانسته است نتایج قابل قبولی را در مقابل روشهای پیشین که از دانش در مورد فضای حالت استفاده نمی کردند کسب کند.

فصل پنجم ۵- نتیجه گیری و پیشنهادها

نتیجه گیری و پیشنهادها

۱-۵- نتیجهگیری

در این پژوهش به مساله ی انتقال جسم به درون یک حفره که یکی از مسائل معمول در آموزش رباتها به حساب می آید پرداخته شد. در این مساله، هدف انتقال جسم از هر نقطه اولیهای به موقعیت هدفی است که درون یک حفره قرار گرفته است. این وظیفه جزء مطرح ترین مسائلی است که در صنایع مختلف بر عهده ی رباتها قرار داده می شود. الگوریتمهای مختلفی برای آموزش رباتها برای انجام وظایف ارائه شده است که در این پژوهش به بررسی دو دسته ی مهم از الگوریتمها یعنی الگوریتمهای یادگیری تقویتی و الگوریتمهای یادگیری تقلیدی پرداخته شد.

آزمایشها بر روی الگوریتمهای یادگیری تقویتی از قبیل Q-Learning نشان داده است که این گونه از الگوریتمها بدون در نظر داشتن دانش پیشین قادر به یادگیری مسئله مورد نظر در زمان معقول نخواهند بود. به همین دلیل به عنوان روش پیشنهادی در این پایاننامه، دو چارچوب برای تزریق دانش پیشین به دو الگوریتم Q-Learning و PILCO معرفی شد. در روش پیشنهادی ارائه شده برای الگوریتم تقویتی Q-Learning بیشین به دو الگوریتم و کاهش زمان یادگیری استفاده تقویتی Q-Learning از مسیرهای راهنما از درون ویدئوهای ضبط شده و با استفاده از مدلهای ردیابی جسم استخراج شد. این مسیرهای راهنما از درون ویدئوهای ضبط شده و با استفاده از مدلهای ردیابی جسم استخراج شدند. با افزودن مسیرهای راهنما به الگوریتم یادگیری تقویتی Q-Learning این الگوریتم به سمت یک الگوریتم یادگیری تقلیدی سوق داده شد. همچنین در روش پیشنهادی دیگری که برای تزریق دانش پیشین به الگوریتم PILCO ارائه شد، تابع هزینه الگوریتم PILCO به نحوی تغییر داده شد که بتواند محدودیتهای فضای حالت از جمله دیوارها را نیز در نظر بگیرد. آزمایشهای مختلفی برای بررسی عملکرد روشهای پیشنهادی انجام شد که همگی نشان از عملکرد بهتر این روشها به نسبت روشهای پایهی PLCO و Q-Learning دادد.

۲-۵- کارهای پیشنهادی

روش پیشنهادی ارائه شده در این پایاننامه مشابه خیلی از کارهای یادگیری تقویتی در مسائل کنترل پیوسته، برای اینکه اطلاعات هدف مورد نظر را در اختیار عامل قرار داده و همچنین برای ساده کردن مسئله اکتشاف به ساختن تابع پاداش مناسب متکی است. در حالی که در کار رباتیکی که در این پایاننامه

مورد بررسی قرار گرفت و در کارهای مشابه آن که در آنها به آسانی می توان موقعیت هدف نهایی را معین کرد، تنظیم یک تابع پاداش مناسب که منجر به راه حلهای خوب می شود دشوار است. می توان این کار سخت تنظیم توابع پاداش را با نمایشهای انسان از کار رباتیکی مورد نظر جایگزین کرد. این کار مسئله اکتشاف را آسان می کند بدون اینکه نیازی به تنظیم دقیق تابع پاداش باشد. همچنین، از آنجایی که در روش PILCO از فرآیندهای گوسی برای آموزش مدل دینامیکی سیستم مورد بررسی استفاده می کند و با توجه به اینکه پیچیدگی محاسباتی استنتاج تقریبی استفاده شده در فرآیندهای گوسی n بیش تر است و در مسئله بررسی شده در این پایان نامه، هرچه تعداد دادههای آموزشی فرآیند گوسی یا n بیش تر شود مدل بهتری از سیستم به دست می آید، می توان از روشهای جدیدتر از قبیل [120] استفاده کرد که هم از لحاظ داده – کارایی قابل مقایسه با روش PILCO هستند و هم پیچیدگی محاسباتیشان نسبت به تعداد دادههای آموزش مدل خطی است. در مقاله [120]، یک رویکرد یادگیری تقویتی مبتنی بر مدل ارائه شده است که از شبکههای عصبی بیزی برای تخمین مدل دینامیکی و تابع سیاست استفاده می کند.

مراجع

- [1] M. Večerík, T. Hester, J. Scholz, F. Wang, O. Pietquin, B. Piot, N. Heess, T. Rothörl, T. Lampe, and M. Riedmiller, "Leveraging Demonstrations for Deep Reinforcement Learning on Robotics Problems with Sparse Rewards," 2017.
- [2] C. J. C. H. W. and P. Dayan, "Q-learning," *Mach. Learn. J.*, vol. 8, 1992.
- [3] M. P. D. and C. E. Rasmussen, "PILCO: A Model-Based and Data-Efficient Approach to Policy Search," *Proc. Int. Conf. Mach. Learn.*, pp. 465–472, 2011.
- [4] A. Billard, S. Calinon, R. Dillmann, and S. Schaal, "Robot Programming by Demonstration," *Springer Handb. Robot.*, no. February, pp. 1371–1394, 2008.
- [5] S. Raza, S. Haider, and M. A. Williams, "Teaching coordinated strategies to soccer robots via imitation," 2012 IEEE Int. Conf. Robot. Biomimetics, ROBIO 2012 Conf. Dig., pp. 1434–1439, 2012.
- [6] S. Schaal, "Is imitation learning the route to humanoid robots?," *Trends Cogn. Sci.*, vol. 3, no. 6, pp. 233–242, 1999.
- [7] and others Claude Sammut, Scott Hurst, Dana Kedzier, Donald Michie, "Learning to fly," *Proc. 9th Int. Work. Mach. Learn.*, pp. 385–393, 2014.
- [8] D. Silver, J. A. Bagnell, and A. Stentz, "High Performance Outdoor Navigation from Overhead Data using Imitation Learning," *Robot. Sci. Syst.*, 2008.
- [9] M. J. Matarić, "Getting humanoids to move and imitate," *IEEE Intell. Syst. Their Appl.*, vol. 15, no. 4, pp. 18–24, 2000.
- [10] J. K. and J. R. Peters, "Policy search for motor primitives in robotics," *Adv. Neural Inf. Process. Syst.*, pp. 849–856.
- [11] E. Berger, H. Ben Amor, D. Vogt, and B. Jung, "Towards a Simulator for Imitation Learning with Kinesthetic Bootstrapping," *Work.* "The Universe Rob. Simulators", Int. Conf. Simulation, Model. Program. Auton. Robot., pp. 167–173, 2008.
- [12] and G. S. Christian Thurau, Christian Bauckhage, "Imitation learning at all levels of Game-AI," *Proc. Int. Conf. Comput. Games, Artif. Intell. Des. Educ.*, vol. 5, pp. 402–408, 2004.
- [13] J. Kober and J. Peters, "Imitation and reinforcement learning," *IEEE Robot. Autom. Mag.*, vol. 17, no. 2, pp. 55–62, 2010.
- [14] and J. P. Jens Kober, J. Andrew Bagnell, "Reinforcement Learning in Robotics: A Survey," *Int. J. Rob. Res.*, vol. 32, no. 11, pp. 1238–1274, 2013.
- [15] Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jayne, "Imitation Learning: A Survey of Learning Methods," *ACM Comput. Surv.*, vol. April, pp. 0–35, 2017.
- [16] B. D. Argall, S. Chernova, M. Veloso, and B. Browning, "A survey of robot learning from demonstration," *Rob. Auton. Syst.*, vol. 57, no. 5, pp. 469–483, 2009.
- [17] C. Nehaniv and K. Dautenhahn, "The Correspondence Problem," *Imitation Anim. Artifacts*, pp. 1–40, 2002.

- [18] A. Schaal, Stefan, Ijspeert, Auke, Billard, "Computational approaches to motor learning by imitation," *Philos. Trans. R. Soc. London B Biol. Sci.*, vol. 358, no. 1431, 2003.
- [19] A. De Santis, B. Siciliano, A. De Luca, and A. Bicchi, "An atlas of physical human-robot interaction," *Mech. Mach. Theory*, vol. 43, no. 3, pp. 253–270, 2008.
- [20] S. Chernova and M. Veloso, "Confidence-based policy learning from demonstration using gaussian mixture models," *Proc. 6th Int. Jt. Conf. Auton. agents multiagent Syst.*, p. 233, 2007.
- [21] J. A. Bagnell and S. Ross, "Efficient Reductions for Imitation Learning," *Proc. 13th Int. Conf. Artif. Intell. Stat. 2010*, vol. 9, pp. 661–668, 2010.
- [22] N. D. Ratliff, D. Bradley, J. A. Bagnell, and J. Chestnutt, "Boosting Structured Prediction for Imitation Learning," *Adv. Neural Inf. Process. Syst.*, pp. 1153–1160, 2006.
- [23] S. Ross, G. Gordon, and J. A. (Drew) Bagnell, "A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning," *Proc. 14th Int. Conf. Artifical Intell. Stat.*, 2011.
- [24] R. Rahmatizadeh, P. Abolghasemi, and L. Bölöni, "Learning Manipulation Trajectories Using Recurrent Neural Networks," *arXiv*, 2016.
- [25] J. Nakanishi, J. Morimoto, G. Endo, G. Cheng, S. Schaal, and M. Kawato, "Learning from demonstration and adaptation of biped locomotion," *Rob. Auton. Syst.*, vol. 47, no. 2–3, pp. 79–91, 2004.
- [26] K. Mülling, J. Kober, O. Kroemer, and J. Peters, "Learning to select and generalize striking movements in robot table tennis," *Int. J. Rob. Res.*, vol. 32, no. 3, pp. 263–279, 2013.
- [27] a J. Ijspeert, J. Nakanishi, and S. Schaal, "Learning Attractor Landscapes for Learning Motor Primitives," *Adv. Neural Inf. Process. Syst.* 15, pp. 1547–1554, 2002.
- [28] J. Kober and J. Peters, "Learning motor primitives for robotics," 2009 IEEE Int. Conf. Robot. Autom., pp. 2112–2118, 2009.
- [29] S. Schaal, P. Mohajerian, and A. Ijspeert, "Dynamics systems vs. optimal control a unifying view," *Prog. Brain Res.*, vol. 165, pp. 425–445, 2007.
- [30] S. Calinon, Z. Li, T. Alizadeh, N. G. Tsagarakis, and D. G. Caldwell, "Statistical dynamical systems for skills acquisition in humanoids," *IEEE-RAS Int. Conf. Humanoid Robot.*, pp. 323–329, 2012.
- [31] L. Rozo, D. Bruno, S. Calinon, and D. G. Caldwell, "Learning optimal controllers in human-robot cooperative transportation tasks with position and force constraints," *IEEE Int. Conf. Intell. Robot. Syst.*, vol. 2015–Decem, pp. 1024–1030, 2015.
- [32] N. Chen, J. Bayer, S. Urban, and P. Van Der Smagt, "Efficient movement representation by embedding Dynamic Movement Primitives in deep autoencoders," *IEEE-RAS Int. Conf. Humanoid Robot.*, vol. 2015–Decem, pp. 434–440, 2015.

- [33] and G. C. Darrin C. Bentivegna, Christopher G. Atkeson, "Learning tasks from observation and practice," *Rob. Auton. Syst.*, vol. 47, no. 2, pp. 163–169, 2004.
- [34] S. Chernova and M. Veloso, "Teaching collaborative multi-robot tasks through demonstration," 2008 8th IEEE-RAS Int. Conf. Humanoid Robot. Humanoids 2008, pp. 385–390, 2008.
- [35] J. Togelius, R. De Nardi, and S. M. Lucas, "Towards automatic personalised content creation for racing games," *Proc.* 2007 IEEE Symp. Comput. Intell. Games, CIG 2007, pp. 252–259, 2007.
- [36] S. Calinon and A. Billard, "A framework integrating statistical and social cues to teach a humanoid robot new skills," *Proc IEEE Int. Conf. Robot. Autom. ICRA Work. Soc. Interact. with Intell. Indoor Robot.*, 2008.
- [37] L.-J. Lin, "Self-Improving Reactive Agents Based On Reinforcement Learning, Planning and Teaching," *Mach. Learn.*, vol. 8, pp. 293–321, 1992.
- [38] L. Lin, "Programming Robots Using Reinforcement Learning and Teaching.," *Aaai*, vol. 2, pp. 781–786, 1991.
- [39] F. Guenter, M. Hersch, S. Calinon, and A. Billard, "Reinforcement learning for imitating constrained reaching movements," *Adv. Robot.*, vol. 21, no. 13, pp. 1521– 1544, 2007.
- [40] K. S. et al. D Silver J Schrittwieser, "Mastering the game of Go with deep neural networks and tree search.," *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [41] T. Brys, A. Harutyunyan, V. U. Brussel, and M. E. Taylor, "Reinforcement Learning from Demonstration through Shaping," *Proc. Twenty-Fourth Int. Jt. Conf. Artif. Intell.*, no. Ijcai, pp. 3352–3358, 2015.
- [42] A. Y. Ng, D. Harada, and S. Russell, "Policy invariance under reward transformations: Theory and application to reward shaping," *Sixt. Int. Conf. Mach. Learn.*, vol. 99, pp. 278–287, 1999.
- [43] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," *Twenty-first Int. Conf. Mach. Learn. ICML '04*, p. 1, 2004.
- [44] N. Kohl and P. Stone, "Policy gradient reinforcement learning for fast quadrupedal locomotion," *IEEE Int. Conf. Robot. Autom. 2004. Proceedings. ICRA '04. 2004*, p. 2619–2624 Vol.3, 2004.
- [45] J. Peters and S. Schaal, "Reinforcement learning of motor skills with policy gradients," *Neural Networks*, vol. 21, no. 4, pp. 682–697, 2008.
- [46] A. J. I. A. J. Ijspeort, J. N. J. Nakanishi, and S. Schaal, "Learning rhythmic movements by demonstration using nonlinear oscillators," *IEEE/RSJ Int. Conf. Intell. Robot. Syst.*, vol. 1, no. October, pp. 958–963, 2002.
- [47] S. Levine and V. Koltun, "Guided Policy Search," *Proc. 30th Int. Conf. Mach. Learn.*, vol. 28, pp. 1–9, 2013.
- [48] M. Zhang, Z. McCarthy, C. Finn, S. Levine, and P. Abbeel, "Learning deep neural

- network policies with continuous memory states," *Proc. IEEE Int. Conf. Robot. Autom.*, vol. 2016–June, pp. 520–527, 2016.
- [49] X. Guo, S. Singh, H. Lee, R. Lewis, and X. Wang, "Deep Learning for Real-Time Atari Game Play Using Offline Monte-Carlo Tree Search Planning," *Adv. Neural Inf. Process. Syst.* 27, vol. 2600, pp. 3338–3346, 2014.
- [50] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-End Training of Deep Visuomotor Policies," 2015.
- [51] S. N. and D. Floreano, "Evolutionary Robotics: The Biology, Intelligence, and Technology," 2000.
- [52] N. Rokbani, A. Zaidi, and A. M. Alimi, "Prototyping a biped robot using an educational robotics kit," 2012 Int. Conf. Educ. e-Learning Innov. ICEELI 2012, 2012.
- [53] Y. Zhang, S. Wang, and G. Ji, "A Comprehensive Survey on Particle Swarm Optimization Algorithm and Its Applications," *Math. Probl. Eng.*, vol. 2015, 2015.
- [54] C. Zhang, Z. Zhen, D. Wang, and M. Li, "UAV path planning method based on ant colony optimization," *Control Decis. Conf. (CCDC)*, 2010 Chinese, pp. 3790–3792, 2010.
- [55] R. Aler, O. Garcia, and J. M. Valls, "Correcting and improving imitation models of humans for Robosoccer agents," *Proc. 2005 IEEE Congr. Evol. Comput.*, vol. 3, p. 2402–2409 Vol. 3, 2005.
- [56] J. Ortega, N. Shaker, J. Togelius, and G. N. Yannakakis, "Imitating human playing styles in Super Mario Bros," *Entertain. Comput.*, vol. 4, no. 2, pp. 93–104, 2013.
- [57] T. Sun, C. Huo, S. Tsai, and C. Liu, "Optimal UAV Flight Path Planning Using Skeletonization and Particle Swarm Optimizer," *Proc. 2008 IEEE Congr. Evol. Comput. (IEEE World Congr. Comput. Intell.*, pp. 1183–1188, 2008.
- [58] R. Cheng and Y. Jin, "A Social Learning Particle Swarm Optimization Algorithm for Scalable Optimization," *Inf. Sci.* (*Ny*)., vol. 291, no. 2015, pp. 43–60, 2015.
- [59] F. Gruau and K. Quatramaran, "Cellular Encoding for Interactive Evolutionary Robotics," *Proc. Eur. Conf. Artif. Life*, pp. 368–377, 1997.
- [60] J. C. Bongard and G. S. Hornby, "Combining fitness-based search and user modeling in evolutionary robotics," *Proceeding fifteenth Annu. Conf. Genet. Evol. Comput. Conf. GECCO '13*, pp. 159–166, 2013.
- [61] H.-I. Institute of Electrical and Electronics Engineers., Y.-C. IEEE Computer Society., C.-L. Chinese Automatic Control Society., IEEE Control Systems Society., Guo li jiao tong da xue (Taiwan), and Asian Control Association., "Evaluation of Human-Robot Arm Movement Imitation," *Control Conf. (ASCC)*, 2011 8th Asian, pp. 287–292, 2011.
- [62] H. El-Hussieny, S. F. M. Assal, A. A. Abouelsoud, S. M. Megahed, and T. Ogasawara, "Incremental learning of reach-to-grasp behavior: A PSO-based Inverse optimal control approach," *Proc.* 2015 7th Int. Conf. Soft Comput. Pattern

- Recognition, SoCPaR 2015, pp. 129–135, 2016.
- [63] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Trans. Knowl. Data Eng.*, vol. 22, no. 10, pp. 1345–1359, 2010.
- [64] L. Torrey, T. Walker, J. Shavlik, and R. Maclin, "Using advice to transfer knowledge acquired in one reinforcement learning task to another," *Lect. Notes Comput. Sci.* (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics), vol. 3720 LNAI, pp. 412–424, 2005.
- [65] L. Torrey and J. Shavlik, "Transfer Learning," *Handb. Res. Mach. Learn. Appl.*, pp. 2–4, 2009.
- [66] G. Kuhlmann and P. Stone, "Graph-based domain mapping for transfer learning in general games," *Mach. Learn. ECML* 2007, no. September, pp. 188–200, 2007.
- [67] T. Brys, A. Harutyunyan, M. E. Taylor, and A. Nowé, "Policy Transfer using Reward Shaping," *Auton. Agents Multiagent Syst.*, pp. 181–188, 2015.
- [68] M. Wulfmeier, P. Ondruska, and I. Posner, "Maximum Entropy Deep Inverse Reinforcement Learning," *AAAI Conf. Artif. Intell.*, pp. 1433–1438, 2015.
- [69] G. Lee, M. Luo, F. Zambetta, and X. Li, "Learning a Super Mario controller from examples of human play," *Proc. 2014 IEEE Congr. Evol. Comput. CEC 2014*, pp. 1–8, 2014.
- [70] K. Judah, A. P. Fern, and T. G. Dietterich, "Active Imitation Learning via Reduction to I.I.D. Active Learning," *Uai*, pp. 428–437, 2012.
- [71] S. Ikemoto, H. Amor, T. Minato, B. Jung, and H. Ishiguro, "Physical Human-Robot Interaction: Mutual Learning and Adaptation," *IEEE Robot. Autom. Mag.*, vol. 19, no. 4, pp. 24–35, 2012.
- [72] S. Calinon and A. G. Billard, "What is the Teacher's Role in Robot Programming by Demonstration? Toward Benchmarks for Improved Learning," *Interact. Stud.*, vol. 8, no. 3, pp. 441–464, 2007.
- [73] H. Daumé, J. Langford, and D. Marcu, "Search-based structured prediction," *Mach. Learn.*, vol. 75, no. 3, pp. 297–325, 2009.
- [74] H. M. Le, A. Kang, Y. Yue, and P. Carr, "Smooth Imitation Learning for Online Sequence Prediction," 2016.
- [75] A. Droniou, S. Ivaldi, and O. Sigaud, "Learning a repertoire of actions with deep neural networks," *IEEE ICDL-EPIROB 2014 4th Jt. IEEE Int. Conf. Dev. Learn. Epigenetic Robot.*, pp. 229–234, 2014.
- [76] H. Mayer, F. Gomez, D. Wierstra, I. Nagy, A. Knoll, and J. Schmidhuber, "A system for robotic heart surgery that learns to tie knots using recurrent neural networks," *Adv. Robot.*, vol. 22, no. 13–14, pp. 1521–1537, 2008.
- [77] M. L. Puterman, "Markov Decision Processes--Discrete Stochastic Dynamic Programming," 1994.

- [78] R. S. Sutton and A. G. Barto, "Reinforcement Learning: An Introduction," 1988.
- [79] R. Howard, "Dynamic Programming and Markov Processes," p. 42, 1960.
- [80] R. E. Bellman, Dynamic Programming. 1957.
- [81] M. L. Puterman and M. C. Shin, "Modified Policy Iteration Algorithms for Discounted Markov Decision Problems," *Manage. Sci.*, vol. 24, no. 11, pp. 1127– 1137, 1978.
- [82] A. G. Barto, S. J. Bradtke, and S. P. Singh, "Learning to act using real-time dynamic programming," *Artif. Intell.*, vol. 72, no. 1–2, pp. 81–138, 1995.
- [83] R. R. S. Sutton, "Generalization in Reinforcement Learning: Successful Examples Using Sparse Coarse Coding," *Adv. Neural Inf. Process. Syst.*, vol. 8, pp. 1038–1044, 1996.
- [84] R. S. Sutton, "Dyna, an integrated architecture for learning, planning, and reacting," *ACM SIGART Bull.*, vol. 2, no. 4, pp. 160–163, 1991.
- [85] A. Moore and C. Atkeson, "Prioritized sweeping: Reinforcement learning with less data and less time," *Mach. Learn.*, vol. 13, pp. 103–130, 1993.
- [86] C. G. Atkeson and J. C. Santamaria, "A comparison of direct and model-based reinforcement learning," *Proc. Int. Conf. Robot. Autom.*, vol. 4, no. April, pp. 3557–3564, 1997.
- [87] E. W. Aboaf, S. Drucker, and C. G. Atkeson, "Task-level robot learning: juggling a tennis ball more accurately," *Robot. Autom. 1989. Proceedings.*, 1989 IEEE Int. Conf., pp. 1290–1295, 1989.
- [88] J. Peters, K. Mülling, and Y. Altun, "Relative Entropy Policy Search.," *Proc. 24th Natl. Conf. Arti cial Intell.*, 2010.
- [89] J. Kober and J. Peter, "Policy search for motor primitives in robotics," *Mach. Learn.*, pp. 1–33, 2010.
- [90] M. P. Deisenroth, "A Survey on Policy Search for Robotics," *Found. Trends Robot.*, vol. 2, no. 1–2, pp. 1–142, 2011.
- [91] R. J. Willia, "Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning," *Mach. Learn.*, vol. 8, no. 3, pp. 229–256, 1992.
- [92] C. Daniel, G. Neumann, and J. Peters, "Hierarchical Relative Entropy Policy Search," *Proc. Int. Conf. Artif. Intell. Stat. (AISTATS 2012)*, vol. XX, pp. 273–281, 2012.
- [93] J. Schneider, "Exploiting Model Uncertainty Estimates for Safe Dynamic Control Learning," *Adv. Neural Inf. Process. Syst.*, 1997.
- [94] J. M. and C. G. Atkeson, "Minimax Differential Dynamic Programming: An Application to Robust Biped Walking," *Adv. Neural Inf. Process. Syst.*
- [95] S. Schaal and C. G. Atkeson, "Constructive incremental learning from only local information," *Neural Comput.*, vol. 10, no. 8, pp. 2047–2084, 1998.

- [96] D. C. M. and K. Glover, "Robust Controller Design Using Normalized Coprime Factor Plant Descriptions," *Control Inf. Sci.*, vol. 138, 1990.
- [97] K. J. Åström and B. Wittænmark, "Adaptive Control," *Dover Publ.*, 2008.
- [98] B. Wittenmark, "Adaptive dual control methods: An overview," *IFAC Symp. Adapt. Syst. Control Signal Process.*, pp. 67–72, 1995.
- [99] A. A. Fel'dbaum, "Dual-Control Theory, Parts I and II," *Autom. Remote Control*, vol. 21, no. 11, pp. 874–880, 1961.
- [100] S. G. Fabri and V. Kadirkamanathan, "Dual adaptive control of nonlinear stochastic systems using neural networks," *Automatica*, vol. 34, no. 2, pp. 245–253, 1998.
- [101] P. Abbeel, M. Quigley, and A. Y. Ng, "Using inaccurate models in reinforcement learning," *Proc. 23rd Int. Conf. Mach. Learn. ICML '06*, pp. 1–8, 2006.
- [102] C. E. Rasmussen and C. K. I. Williams, "Gaussian Processes for Machine Learning," *Massachusetts Inst. Technol.*, p. 267, 2006.
- [103] A. Y. Ng and M. Jordan, "PEGASUS: A Policy Search Method for Large MDPs and POMDPs," *Conf. Uncertain. Artif. Intell.*, vol. 94720, pp. 406–415, 2000.
- [104] and E. L. A. Y. Ng, A. Coates, M. Diel, V. Ganapathi, J. Schulte, B. Tse, E. Berger, "Autonomous helicopter flight via Reinforcement Learning," *Adv. Neural Inf. Process. Syst. 16*, vol. 16, pp. 363–372, 2004.
- [105] J. Ko, D. J. Klein, D. Fox, and D. Haehnel, "Gaussian Processes and Reinforcement Learning for Identi cation and Control of an Autonomous Blimp," *Proc. Int. Conf. Robot. Autom.*, pp. 742–747, 2007.
- [106] and E. Andrew Y. Ng, Adam Coates, Mark Diel, Varun Ganapathi, Jamie Schulte, Ben Tse, Eric Berger and Liang, "Autonomous Inverted helicopter flight via Reinforcement Learning," *Exp. Robot. IX. Springer*, pp. 363–372, 2004.
- [107] J. Bagnell and J. Schneider, "Autonomous helicopter control using reinforcement learning policy search methods," *Int. Conf. Robot. Autom.*, pp. 1615–1620, 2001.
- [108] M. Deisenroth, C. Rasmussen, and D. Fox, "Learning to control a low-cost manipulator using data-efficient reinforcement learning," *Robot. Sci. Syst. VII*, pp. 57–64, 2011.
- [109] W. S. Cleveland and S. J. Devlin, "Locally-Weighted Regression: An Approach to Regression Analysis by Local Fitting," *J. Am. Stat. Assoc.*, vol. 83, pp. 596–610, 1988.
- [110] B. D. O. A. and J. B. Moore, "Optimal filtering," Dover Publ., 2005.
- [111] S. J. Julier and J. K. Uhlmann, "Unscented Filtering and Nonlinear Estimation," *Comput. Eng.*, vol. 92, no. 3, pp. 401–422, 2004.
- [112] K. Xiong, H. Y. Zhang, and C. W. Chan, "Performance evaluation of UKF-based nonlinear filtering," *Automatica*, vol. 42, no. 2, pp. 261–270, 2006.
- [113] C. M. Bishop, "Pattern Recognition and Machine Learning," Inf. Sci. Stat., p. 407,

2006.

- [114] J. a. A. Nelder, R. Mead, B. J. a Nelder, and R. Mead, "A Simplex Method for Function Minimization," *Comput. J.*, vol. 7, no. 4, pp. 308–313, 1964.
- [115] S. J. Russell, P. Norvig, J. F. Canny, J. M. Malik, and D. D. Edwards, *Artificial Intelligence A Modern Approach*, vol. 27. Pearson Education, 2003.
- [116] J. Quiñonero-candela, C. E. Rasmussen, and R. Herbrich, "A unifying view of sparse approximate Gaussian process regression," *J. Mach. Learn. Res.*, vol. 6, pp. 1935–1959, 2005.
- [117] D. Nguyen-Tuong, M. Seeger, and J. Peters, "Model learning with local gaussian process regression," *Adv. Robot.*, vol. 23, no. 15, pp. 2015–2034, 2009.
- [118] M. Deisenroth, Efficient Reinforcement Learning Using Gaussian Processes. KIT Scientific Publishing, 2010.
- [119] M. Toussaint and C. Goerick, "A bayesian view on motor control and planning," *Stud. Comput. Intell.*, vol. 264, pp. 227–252, 2010.
- [120] Higuera, Juan Camilo Gamboa, David Meger, and Gregory Dudek. "Synthesizing Neural Network Controllers with Probabilistic Model-Based Reinforcement Learning." In 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 2538-2544. IEEE, 2018.

Abstract

Today, the use of robots in various industries and different tasks is rapidly increasing. Using robots decreases human resources' costs and increases accuracy of performance. One of the common tasks in various industries in which robots are used is moving an object from an initial position to a target position. Various algorithms are provided for training different tasks to robots. Two important categories of these algorithms are reinforcement learning and imitation learning algorithms. In imitation learning, an agent is taught to accomplish a task by learning a mapping between observations and actions from demonstrations. In reinforcement learning, the robot is trained to perform a task by considering rewards and penalties for each action.

In this thesis, based on imitation learning methods guiding paths have been used to improve the performance of the Q-Learning algorithm. To extract guiding paths, we have used demonstrations in which an object is moved into a hole by a human.

Recently, the PILCO algorithm is used to solve problems where the design of the dynamic model of the robot is very complicated or the dynamic model is uncertain. This algorithm, which uses probabilistic inference and models for learning control, has been used in various problems in the past, but has not been used to solve the problem of moving an object into a hole by simulated manipulator. Given the uncertainty inherent in this problem, PILCO algorithm can be used as an appropriate solution to solve this problem in continuous environments. Hence, in this thesis, an approach based on the PILCO algorithm is presented. This algorithm is modified in such a way that it takes into account state space constraints in its computational process. In order to evaluate the proposed methods, simulations are performed in discrete and continuous environments. It is worth mentioning that a simulated robot is used to move an object from an initial position to a target position. Simulation results in these two discrete and continuous environments show that the proposed algorithms perform better than the original algorithms, Q-learning algorithms and PILCO algorithms.

Key Words: Imitation Learning, Guiding Paths, Reinforcement Learning, Probabilistic Models, Robotics



Amirkabir University of Technology (Tehran Polytechnic)

Department of Robotics Engineering

M.Sc. Thesis

Probabilistic Model based Imitation Learning in Robotic Applications

By Ali Javadi

Supervisors
Dr. Saeed Shiry Ghidary
Dr. Ahmad Nickabadi