

۱) فصل یک به انواع سیستم‌های یادگیری

۱) Supervised Learning: تابعی مطلوب بر آن داده می‌شود و مناسب برای آن آموزشی

صورت می‌گیرد.

۲) Unsupervised Learning: تابعی مطلوب را اختیار agent قرار نمی‌گیرد، agent با کشف ویژگی‌های مشترک بین آنها

طبقه بندی می‌کند.

۳) Reinforcement Learning (RL): agent در محیط به دنبال بر آوردن سودهای گوناگون است که در مرحله (step) زمانی مناسب

با state محیط یک action (عمل) را انتخاب می‌کند و به محیط اعمال می‌کند و سپس بر اساس یادگیری یا مجازات دریافتی از محیط وظیفه خود

را یاد می‌گیرد.

۴) Semi-supervised Learning: بدون نظارت است که تعداد محدودی از یادگیری‌ها نیز توسط یک معلم در اختیار agent قرار می‌گیرد

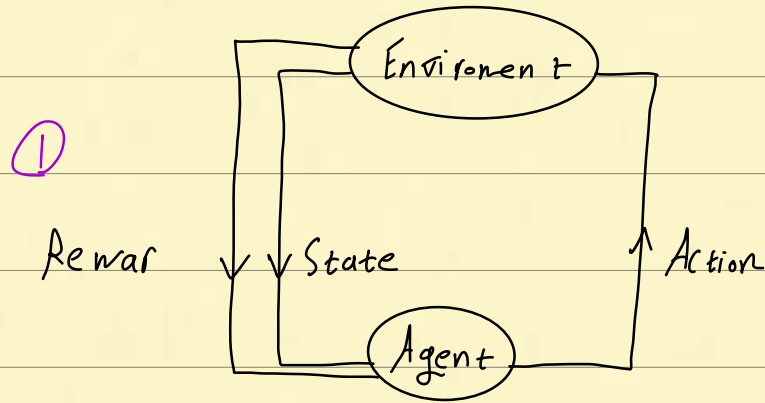
و agent از این یادگیری‌ها برای بهبود عملکرد یادگیری بدون نظارت استفاده می‌کند

۵) Self-supervised Learning: معلم وجود ندارد و خود agent خود به خود با توجه به تکلیف‌های خود را آموزش می‌دهد.

* AGI (Artificial general intelligent): هدف آن ترکیب هوشی‌های بالا و ربات انسان نما است.

اما بحث ما اینجا در مورد AL است که در آینده به صورت کلی‌نمای آن است:

نمای کلی سیستم RL به صورت مقابل:



فرآیند RL بدون مشخص کردن نحوه عمل برای عامل یا agent در صفا با تعیین پارامتر و نتیجه بیان آموزشی می‌باشد.

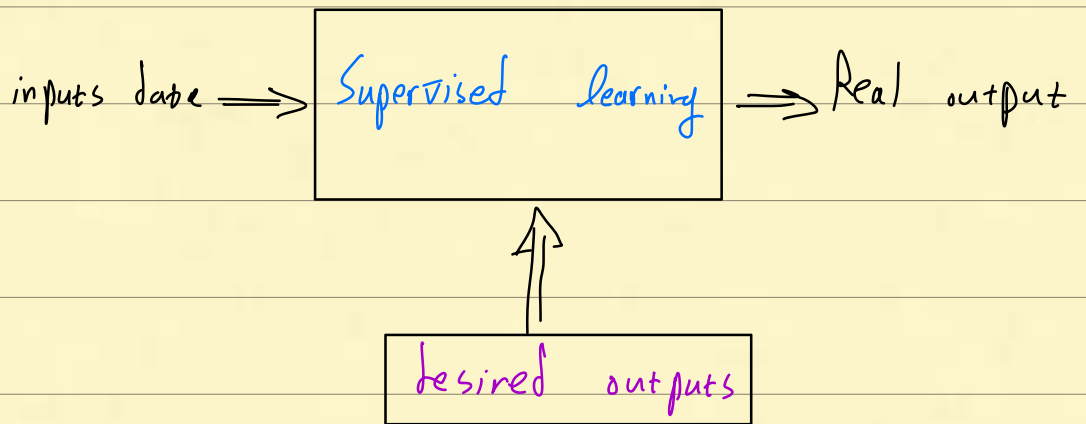
تفاوت RL با یادگیری با نظارت:

در با نظارت بر اساس ورودی خروجی آموزشی داده می‌شود اما در RL صرفاً بر اساس پارامترهای گذشته است و agent در جهت ماکسیم‌سازی

پارامترها در زمان آینه خواننده را انجام می‌دهد.

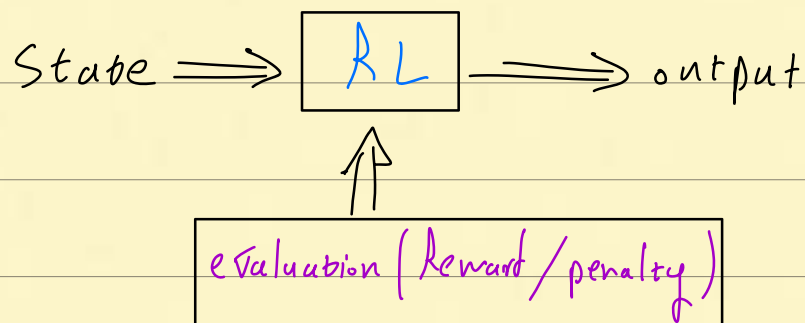
در تفاوت به صورت Offline است و در RL یادگیری به صورت Online است.

a) Supervised learning:



هدف؟ حداقل کردن خطای بین Real و desired است.

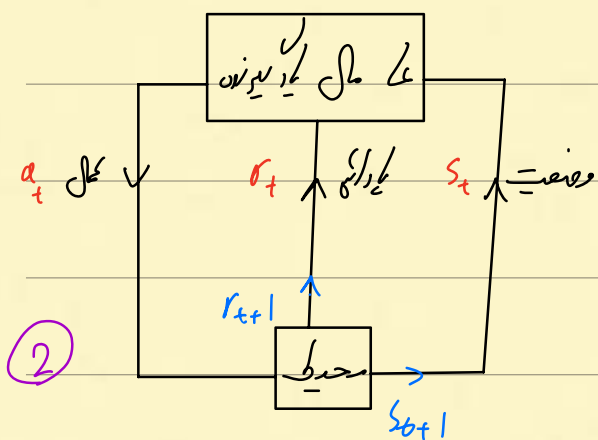
b) Reinforcement learning:



هدف؟ بهجوریک حداکثر پارامترها را در زمان معین است.

* در نهایت پارامتری های آینده و انتهای بازه زمانی مهم تر است

فصل دوم ← بخشی های اصلی یک سیستم RL



در شکل کنونی، روبرو RL عامل یادگیرنده (agent) با توجه به عامل S_t و پارامتری a_t

در یک بازه زمانی متناهی اقدامی را انجام می دهد و محیط از آن اقدام یک state جدید S_{t+1}

و پارامتری جدید به ترتیب S_{t+1} و r_{t+1} تولید می کند.

بخشی های سیستم RL ← بلوک های اصلی در این بخش است:

1 محیط (environment) 2 عامل یادگیرنده (Agent) 3 سیاست (Policy)

4 تابع پاداش (Reward function) 5 تابع ارزش (Value function)

* در برخی سیستم های RL یک مدل از محیط نیز داریم.

Policy (سیاست) ← نحوه رفتار کردن عامل یادگیرنده را مشخص می کند، بنابراین از حالت های محیط به عملی است که در آن حالت های توانمند

انتخاب می شوند و به حالت کلی داریم:

1) یک تابع رفتاری داریم است. 2) یک جدول جستجو است.

3) تابع احتمالاتی از اعمال

«Reward (پاداش)» ← این تابع هر حالت یا state از محیط را به یک عدد اسکالر و یک پاداشی نشان می‌دهد.

این عدد میزان رضایت محیط از وارد شدن به آن state را نشان می‌دهد.

«Value (ارزنی)» ← برخلاف Reward، یا پاداشی تعیین می‌کند در بلند مدت چه چیزی سودمند است.

به عبارت دیگر ارزشی هر state که مقدار پاداشی را معین می‌کند که یک عامل می‌تواند انتظار داشته باشد با شروع از آن در آینده جمع آوری کند.

* امکان دارد محلی پاداشی لحظه‌ای (V) می‌دانسته باشد اما ارزشی آن بالا باشد چرا که در انتها برای agent پاداشی های بیشتری دارد.

* Value یا ارزشی وابسته به سیاست است که از آینده داریم و تحت یک سیاست بیان می‌کنیم.

* به عنوان Value (ارزنی) بسیار سخت تر از محاسبه Reward (پاداشی) است چون پاداشی مستقیم از محیط است.

محیط اول ارزشی را با به تخصیص زرد.

«روشهای حل مسائل RL» v و Q جستجوی سیاست: تمامی سیاست‌های ممکن جستجوی و آموزش و الگوریتم‌های ML بدون آنکه

از این تابع ارزشی استفاده کنیم و یاد بگیریم درستی از محیط دانسته باشیم بهترین سیاست را می‌یابیم.

* هر چه مقدار بزرگتر حجم محاسبات نشانگر

مثال) جستجوی تمام وکالت آتی در یک بازی شطرنج.

2) تعامل با محیط: سعی در پیدا کردن سیاست بهینه با تعامل با محیط است.

روشهایی داریم قبل از آن: - برنامه ریزی دینامیک به معنی برنامه ریزی، آفون وی محاسبات بالا هستند.

- منت کارلو ← معنی برنامه ریزی با agent است (غیر وابسته به مدل)، آفون پیشه ولی شروع هستند.

- تقاضی ← ترکیب مدل بالا است که خود شامل Q-learning, SARSA, ... است.

«مدل محیط» ← یک مدل محیط خوب می تواند رفتار محیط را پیش بینی کند که در برنامه ریزی دینامیک جهت و مسیر.

«فرم مسئله Tic-Tac-Toe» ← action: انتخاب هر خانه، state: تعداد و وضعیت 0, X

Reward: تعداد سطر و ستون و ردیف است

1) فصل سوم ← انواع ساختار های هوشمند RL

agent ← هر چیزی که قادر است محیط پیرامون خود را از طریق حسگرهایی رد کند و بر محیط توسط عملگرهایی عمل کند.

با به صورت یک نوع agent داریم:

1) agent های که در قبل مدل محیط به آنها داده شده است اما تابع ارزشی (value func) ندارند.

* نوع نشان دادن معروف داریم: $R_{ss'}^a$ یا $r(s_t, a_t) \leftarrow$ مقدار پاداشی گرفته شده از انجام action به محیط در حالت s

ورفتن به حالت s است.

- P_{ss}^a یا $\Pr(S_{t+1} | S_t, a_t)$ ← احتمال رفتن محیط به حالت S در صورت عمل a در حالت S است.

حال در مدل اول agent ما P_{ss}^a و R_{ss}^a را داریم.

این مدلها بر اساسی مدلی که از محیط دارند تابع ارزشی را به طور مستمر تخمین میزنند و از آنجا بابت مدل محیط همه مشخصات مورد نیاز از محیط

را در بردارند، لذا برای تخمین تابع ارزشی (value func) نیاز تعامل مستمر با محیط نداریم.

پس تخمین تابع ارزشی بصورت **offline** (جوابی از محیط) است.

(2) **agent** های که هیچکدام از مولفههای محیط یا تابع ارزشی به آنها داده نشده است. پس بنابراین نیاز دارند

تا بصورت مستمر با محیط تعامل کنند و بتوانند تابع ارزشی را تخمین بزنند.

این کار یعنی تخمین تابع ارزشی به صورت **online** است.

- مبتنی بر مدل به اول agent با تعامل همراهِ محیط یک مدل از محیط را از ارزشی میبیند و سپس به کمک آن مدل

تابع ارزشی را بصورت آنلاین تخمین میزنند که روش **Adaptive programming** یا **ADP** است.

- **model-free** از مدل ← بدون اینکه مدلی از محیط به دست آید agent بدون واسطه و از طریق تعامل مستمر با محیط

به یادگیری تابع ارزشی میپردازد که روشی که **Mont Carlo** یا **Tempored difference** هستند.

انواع محیط های سیم های RL ←

1) محیط مارکوف (Markov) محیطی که در آن پیشینی وضعیت آینده صرفاً به آخرین وضعیت (state)

ما وابسته است مثال: بازی شطرنج، کنترل ماشین، کنترل هواپیما

* یعنی از کارها در محیط غیر مارکوف است مثال پیش بینی رفتار انسان، پیش بینی آب و هوا و ...

2) محیط مسئله / پیوسته محیطی که در آن فضای حالت عمل (action-state) مسئله با مسئله پیوسته

مثال) در هزارتو یا راهروی مارپیچ (Maze) محیطی با ابعاد 4x7 داریم که حاوی یک منبع در یکی از سلولها و یک agent

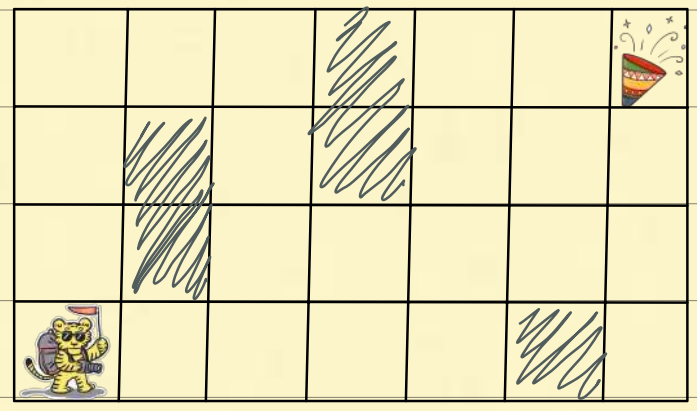
جستجو کننده در یکی از خانه ها دارد همچنین تعدادی مانع نیز در مسیر قرار داده شده است

امکالی که باعث (فردی agent از مسیر می شود یا باعث برخورد به مانع ها می شود را پاداش -1) و (رسیدن به منبع پاداش $+1$)

و سایر پاداش ها $r=0$ است

همانطور که مشاهده می شود این امکان مسئله هستند و به بالا، پایین، چپ و راست خلاصه می شوند

3



* این یک محیط مارکوف است