

انواع یاد گیری ماشین

- یادگیری با نظارت
- يادگيري بدون نظارت
 - يادگيري تقويتي
- یادگیری نیمه نظارتی

1. يادگيري نظارتشده SUPERVISED LEARNING

در این روش، مدل از دادههای برچسبدار استفاده میکند. این بدان معناست که ورودیها همراه با خروجیهای مشخص (برچسب) به مدل داده میشوند.

هدف: یادگیری نگاشت میان ورودی و خروجی، بهطوری که مدل بتواند برای دادههای جدید پیشبینی دقیق انجام دهد.

کاربردها:

- طبقهبندی Classification : پیشبینی دسته ای که یک نمونه به آن تعلق دار د (مانند شناسایی ایمیلهای اسپم).
 - رگرسیون Regression : پیشبینی مقادیر عددی (مانند پیشبینی قیمت خانه).

مثال:

• پیش بینی نمره دانش آموزان با استفاده از داده هایی شامل ساعات مطالعه و نمرات قبلی. Linear Regression, SVM,

2. يادگيري بدون نظارتUNSUPERVISED LEARNING

در این روش، داده ها بدون برچسب هستند و مدل باید ساختار یا الگوهای پنهان موجود در داده ها را کشف کند. هدف: کشف کند. هدف: کشف الگوها و ارائه ی بینشی در مورد داده ها.

کاربردها:

- خوشهبندی Clustering: گروهبندی داده ها بر اساس شباهت ها (مانند تقسیم مشتریان به دسته های مشابه برای بازاریابی).
- کاهش ابعاد Dimensionality Reduction : سادهسازی داده ها با کاهش تعداد ویژگی ها (مانند PCAبرای فشر دهسازی داده).

مثال:

گروهبندی تصاویر مشابه از مجموعهای بدون برچسب.

K-Means, PCA,

3. يادگيري تقويتي REINFORCEMENT LEARNING

این روش بر اساس تعامل مدل عامل یا (Agent)با یک محیط پویا کار میکند. عامل بر اساس اقدامات خود، باز خور دی در قالب پاداش دریافت میکند و هدف آن یادگیری سیاستی است که پاداش تجمعی بلندمدت را به حداکثر برساند.

هدف: یافتن بهترین مجموعه اقدامات برای بهبود عملکرد در طول زمان.

کاربردها:

- بازی ها مانند شطرنج.
 - کنترل رباتها.
- مدیریت منابع در سیستمهای پیچیده.

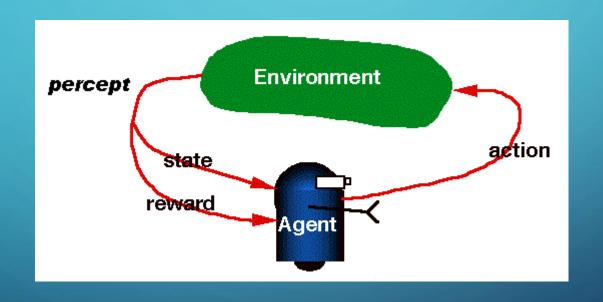
مثال:

• آموزش یک ربات برای یادگیری نحوه حرکت در محیط ناشناخته.

Q-Learning, PPO, DQN

ویژگی	یادگیری با نظارت	یادگیری بدون نظارت	يادگيرى تقويتى
دادهها	برچسبدار	بدون برچسب	تعامل با محیط
هدف	پیشبینی دقیق	كشف الكوها	بیشینهسازی پاداش
كاربردها	پیشبینی و دستهبندی	خوشهبندی و اکتشاف	تصمیمگیری و کنترل

اجزا اصلی یادگیری تقویتی



اجزا اصلی یادگیری تقویتی

[.عامل (ساده، هدفگر ا، سو دمندی، یادگیر نده)

عامل، موجودیت یا برنامهای است که تصمیم میگیرد چه کاری انجام دهد. هدف اصلی عامل، یادگیری از طریق تعامل با محیط است تا پاداش بیشتری به دست آورد.

مثال: یک ربات که در اتاق حرکت میکند و اشیای خاصی را جمعآوری میکند.

2.محيط (environment)

محیط شامل همه چیز هایی است که عامل با آن در تعامل است، از قوانین گرفته تا شرایط اطراف. وقتی عامل کاری انجام میدهد، محیط به آن بازخور د میدهد، معمولاً به شکل پاداش یا تغییر وضعیت.

مثال: تخته شطرنج برای یک عامل بازیکننده شطرنج.

(state) حالت.3

حالت، مشاهده ایی از وضعیت محیط است که عامل در هر لحظه در آن قرار دارد.

مثال: موقعیت فعلی مهرههای شطرنج روی تخته.

4.اقدام (action)

اقدام، کاری است که عامل میتواند در یک حالت خاص انجام دهد. همه اقدامات ممکن، فضایی به نام "فضای اقدام" تشکیل میدهند. مثال: حرکت مهره سرباز در شطرنج.

رreward) پاداش.5

پُاداش، یک عدد است که عامل پس از انجام یک اقدام در یک حالت دریافت میکند. این عدد به عامل میگوید که آیا عملکردش خوب بوده یا نه. هدف عامل: جمعآوری بیشترین پاداش در طول زمان. مثالها: +10 برای رسیدن به هدف, -1 برای برخورد با مانع.

6.سیاست (policy)

سیاست، یک برنامه یا راهنما است که مشخص میکند عامل در هر حالت چه کاری انجام دهد. انواع سیاست:

- قطعی: برای هر حالت، دقیقاً یک اقدام مشخص می شود.
- احتمالی: احتمال انجام هر اقدام تعریف می شود. مثال: اگر در گوشه زمین باشی، به سمت مرکز حرکت کن.

(value function) تابع ارزش.

تابع ارزش، به عامل کمک میکند تخمین بزند که یک حالت یا اقدام چقدر خوب است و در بلندمدت چه پاداشی به همراه دارد. دو نوع اصلی:

- تابع ارزش حالت: مشخص میکند که بودن در یک حالت چقدر ارزشمند است.
- تابع ارزش اقدام: نشان میدهد که انجام یک اقدام در یک حالت خاص چقدر خوب است. مثال: ارزش حرکت یک مهره شطرنج ممکن است به احتمال برد یا بهبود وضعیت تخته بستگی داشته باشد.







الگوريتمهای معروف يادگيری تقويتی

Q - Learning vs SARSA (State Action Reward State Action) Algorithm

Q - Learning (Off policy)

Updated Q Value Current Q Value Target Q Value Current Q Value

$$Q(s,a) = Q(s,a) + \alpha \left[r + \max_{a'} \gamma Q(s',a') - Q(s,a) \right]$$

 α = Learning Rate

Target policy is always Greedy Policy

Behaviour Policy

SARSA (State Action Reward State Action) Algorithm (On policy)

Updated Q Value Current Q Value Target Q Value Current Q Value $Q(s,a) = Q(s,a) + \alpha \left[r + \gamma Q(s',a') - Q(s,a) \right]$ $\alpha = \text{Learning Rate}$ Target Policy is always same as

by Dr. Pankaj Kumar Porwal (BTech - IIT Mumbai, PhD - Cornell University): Principal, Techno India NJR Institute of Technology, Udaipur

- الف) الكوريتمهاى كلاسيك (Classic RL)
 - 1. مبتنی بر ارزش، بدون مدل (Q-Learning)

یادگیری ارزش حالت-اقدام (Q(s,a)

ویژگی: یادگیری خارج از سیاست

State, Action, Reward, State', Action')SARSA .2) مبتنی بر ارزش، بدون مدل

> مشابه Q-Learningهاما یادگیری داخل سیاست (On-Policy)

ویژگ <i>ی</i>	Q-Learning	SARSA
سیاست یادگیری	(Off-Policy) یادگیری خارج از سیاست	(On-Policy) یادگیری داخل سیاست
تصمیمگیری	بر اساس بهترین اقدام ممکن در حالت بعدی	بر اساس اقدامی که عامل طبق سیاست فعلیاش انجام میدهد
کارپرد	مناسب برای کاوش و یادگیری کلی	مناسب برای بهبود تدریجی یک سیاست مشخص

• ب) الگوريتمهای يادگيری عميق تقويتی (Deep Reinforcement Learning)

تعریف کلی:

الگوریتمهای یادگیری عمیق تقویتی ترکیبی از **یادگیری تقویتی RL و یادگیری عمیق DL**هستند. در این روشها، از **شبکههای عصبی عمیق** برای تقریب توابع پیچیده مانند تابع ارزش یا سیاست استفاده میشود. این رویکرد برای حل مسائل پیچیدهتر با فضای حالت و اقدام بزرگ طراحی شده است.

(DQN) Deep Q-Network:1

ترکیب Q-Learning با شبکههای عصبی برای حل مسائل با فضای حالت بسیار بزرگ.

به جای ذخیره مقدار (s,a) در یک جدول، از یک شبکه عصبی برای تخمین آن استفاده میکنیم.

دو ابزار کلیدی در DQN

Replay Buffer : تجربیات گذشته عامل را ذخیره میکند و به صورت تصادفی از آنها استفاده می شود تا داده ها همبستگی کمتری داشته باشند. Target Network: یک کپی از شبکه اصلی برای پایدار کردن یادگیری استفاده می شود.

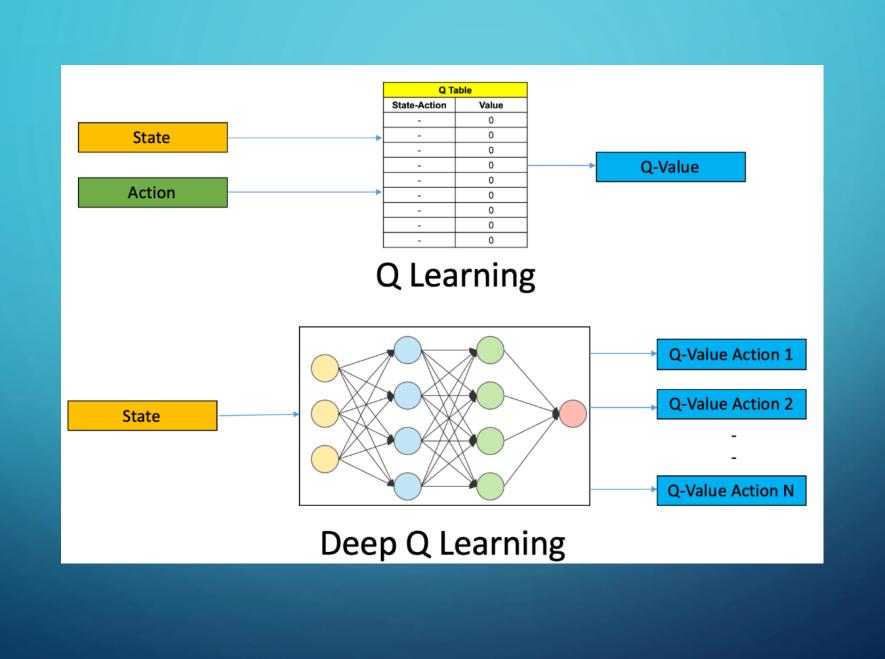
مثال ساده:

فرض کنید عامل در حال بازی کردن بازی آتاری Breakout است عامل تصویر صفحه بازی را میبیند و شبکه عصبی تخمین میزند کدام حرکت (چپ یا راست) بیشترین پاداش را دارد و عامل حرکت میکند، پاداش دریافت میکند و شبکه Qرا بهروزرسانی میکند.

ویژگیها:

مناسب برای مسائل گسسته Discrete

نقطه عطف در یادگیری تقویتی، اولین بار توسط DeepMindبرای حل بازی های آتاری معرفی شد.



A2C, A3C:2

ایده اصلی:

ترکیبی از دو بخش:

:Actorیاد می گیرد چه اقدامی در هر حالت انجام دهد (سیاست). (Critic: می کند که این اقدام چقدر خوب بوده است (ارزش).

نحوه کار:

Actor هر دو توسط شبکههای عصبی جداگانه مدلسازی میشوند.

Actorسیاست را تولید میکند، و Criticبازخورد میدهد تا Actorبهتر شود.

در A3Cچند عامل Agents به صورت موازی اجرا میشوند و تجربیات خود را به اشتراک میگذارند. مثال ساده:

تصور کنید عامل یک بازوی رباتیک است که باید یک جسم را بلند کند:

- 1. Actor تصمیم میگیرد که چگونه بازو را حرکت دهد.
- 2. Critic بازخور د می دهد که این حرکت چقدر به دستیابی به هدف کمک کرده است.
 - 3. Actor بر اساس این بازخورد سیاست خود را بهبود میدهد.

ویژگیها:

مناسب برای مسائل پیچیده با فضای حالت و اقدام بزرگ.

از اجرای موازی برای افزایش سرعت یادگیری استفاده میکند.

نوع	ویژگی بارز	مزایا	معايب	الگوريتم
ارزش	یادگیری Off-Policy	ساده و کارا	مناسب برای فضای گسسته	Q-Learning
ارزش	یادگیری On-Policy	سياست پايدارتر	کندتر از Q-Learning	SARSA
ارزش	استفاده از شبکه عصبی	مناسب فضای حالت بزرگ	پیچیدگی محاسباتی بیشتر	DQN
Actor-Critic	یادگیری موازی	سر عت بیشتر	نیازمند منابع بیشتر	A3C

كاربردهاى يادگيرى تقويتى

رباتیک:

- یادگیری حرکات پیچیده برای مسیریابی و دستکاری اشیاء.
- مثال: رباتهایی که در کارخانهها کار میکنند یا رباتهای جراحی.

بازیها:

- موفقیت چشمگیر در بازی هایی مانند شطرنج و .Go
- مثال: AlphaGoاز یادگیری نقویتی برای شکست قهرمان جهانی استفاده کرد.

وسایل نقلیه خودمختار:

- تصمیمگیری در زمان واقعی برای رانندگی ایمن و کارآمد.
 - مثال: پارک کردن ماشین با آزمون وخطا

مدیریت منابع:

- بهینهسازی تخصیص منابع و زمانبندی.
 - مثال: مدیریت انرژی یا ترافیک.

محدودیتهای یادگیری تقویتی

کارایی نمونه Sample Efficiency

- · نیاز به تعاملهای فراوان برای یادگیری.
 - بسیار محاسباتی و زمانبر.

تریدآف اکتشاف و بهرهبرداری Exploration vs. Exploitation

• تعادل میان کاوش اقدامات جدید و بهر هبر داری از دانش موجود دشوار است.

مهندسی پاداش Reward Engineering

- طراحی پاداش مناسب بسیار پیچیده است.
- پاداش نادر ست ممکن است به رفتار های غیر منتظره منجر شود.

ملاحظات اخلاقی Ethical Considerations

• ممکن است عامل رفتارهای مضریا غیراخلاقی بیاموزد

Book https://www.amazon.com/Deep-Reinforcement-Learning-Hands-Qnetworks/dp/1788834240 • Code https://github.com/PacktPublishing/Deep-Reinforcement-Learning-Hands-On