

# Reinforcement Learning with MATLAB

اسپرنی

نمایندگی یادگیری تقویتی (قسمت اول و دوم دوره)

نمایندگی یادگیری

(1) نمایندگی یادگیری

(2) Multi-Agent Reinforcement Learning

(3) نمایندگی تقویتی یادگیری (FMDP)

(4) Dynamic Programming

(5) Monte Carlo

(6) Temporal Difference

انواع یادگیری تقویتی:

✓ Supervised Learning

✓ Unsupervised Learning

✓ Reinforcement Learning

اجزای یادگیری تقویتی:

✓ Environment

✓ Agent

✓ Policy

✓ Reward

✓ Exploitation and Exploration

نمایندگی یادگیری تقویتی به دنبال یادگیری در محیطی است که در آن پاداش و تنبیه وجود دارد.



در یادگیری تقویتی، یادگیری از طریق پاداش و تنبیه انجام می‌گیرد.

اینکه شما در این دوره یاد بگیرید که چگونه در محیطی یاد بگیرید که در آن پاداش و تنبیه وجود دارد.

sam



دارنده است به تروه که نه تری بیست را بستم  
دارنده و تری می شوند.

میں اس سے پہلے /

شَرِّ الدَّائِمِ K-means

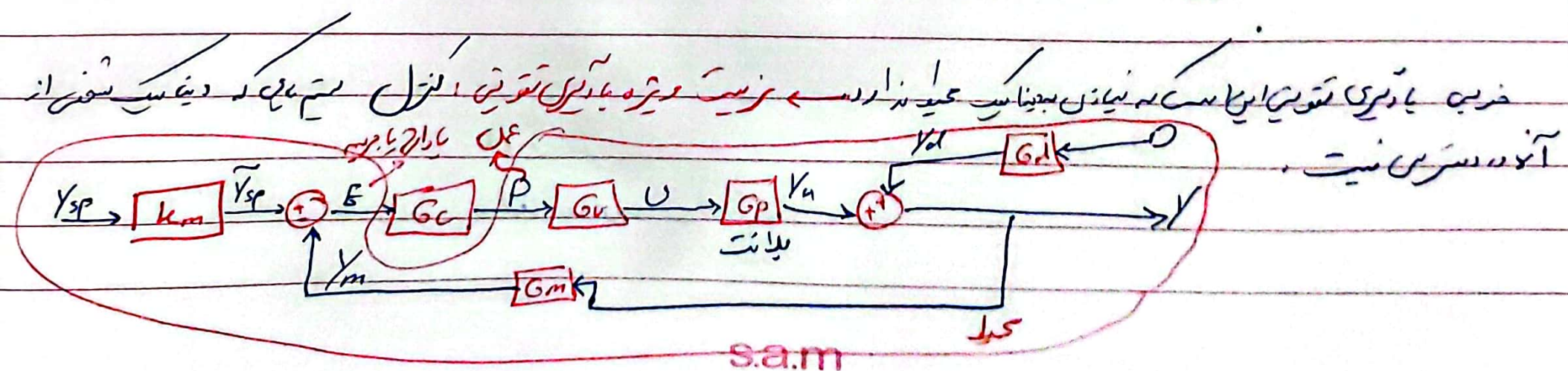
عالم در هر حقیقتی (یا حالتی) می‌تواند حرکت را (یا عملی را) تشخیص دهد.

Penalty      Reward

$$\max L(x) \text{ s.t. } \min -L(x)$$

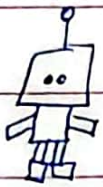
★ در یادگیری تفاوتی و غیر تفاوتی با دینی است. هر کسی که با دینی تفاوتی دارد، باید در یادگیری تفاوتی و غیر تفاوتی با دینی است. هر کسی که با دینی تفاوتی دارد، باید در یادگیری تفاوتی و غیر تفاوتی با دینی است.

خوبین یادگیری تقویتی این است که نیامدن به جایاس محسوس ندارد ← **زیست** و **تیره** یادگیری تقویتی کنترل مستقیم است که در اینجا یک شخص از آن یاد می‌گیرد.





در عرض کنترل کننده  $G$  در درس کنترل خطی باید بقیه دنیا که ما را به اینم و نحوه یادگیری تقویتی میانی به داشتن دنیا نیست.  
 و  $G$  یا کنترل کننده عامل خواسته شد که باید با تعامل با آنرا یاد بگیرد.

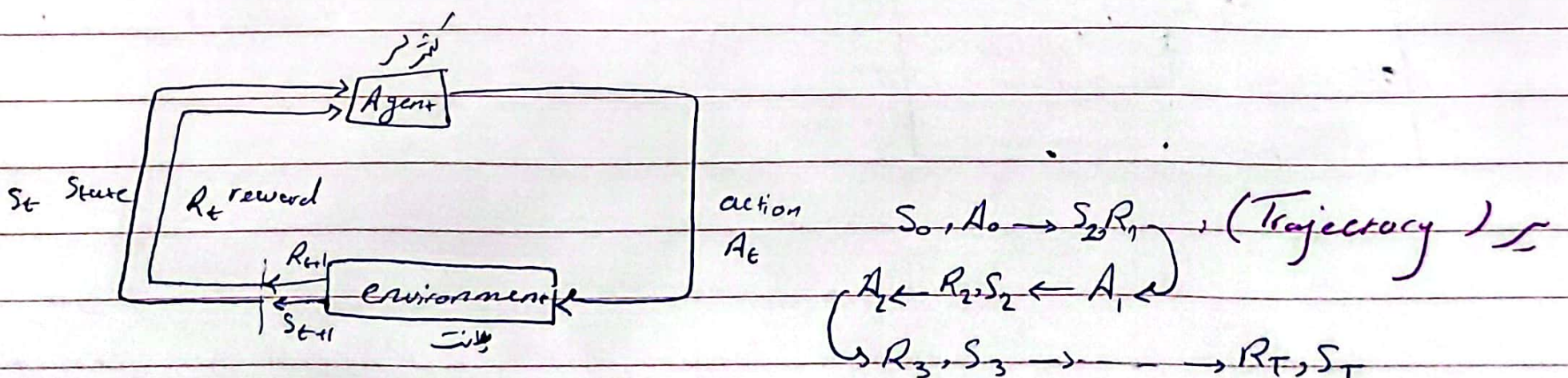


نقل دتیرا کنترل تمام شدن دبات  
 به حد تمام شدن  
 دفع اعتنا 2  
 به عبور از موانع

در بازی ما هم عامل یاد می گیرد که به یاد ریاضت بهترین یاداری یا کسی می جریه بازی را تمام کند.  
 در سال 1997 سوپر کامپیوتر Deep Blue در 19 حرکت بهترین شلر خنچ بازی را شکست برد.

اجزای یادگیری تقویتی  
 اجزای اصلی: عامل و محیط  
 Environment Agent  
 Reward action

تفصیل بی اصلی: حالت، State، پاداش: Reward، عمل، Action



اصل کار یادگیری تقویتی

عامل از حالت فعلی اطلاع دارد،  $S_t$   
 عامل در حالت فعلی عمل را انجام می دهد،  $A_t$   
 محیط به عامل جدیدی برود و پاداش را به عامل خواهد داد،  $R_{t+1}$  و  $S_{t+1}$   
 عامل با دانستن دریافت را بر دانه می کند.  
 عامل به حالت جدید هم عمل را انجام می دهد و ...

s.a.m



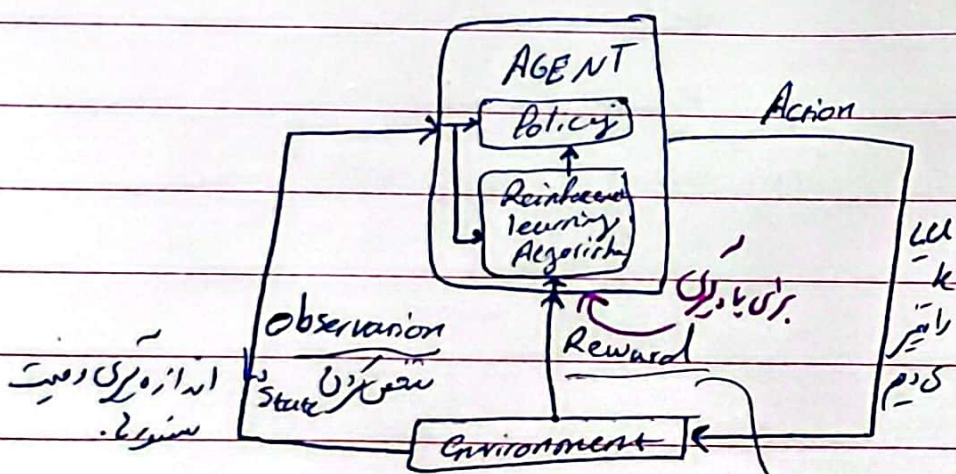
یادگیری تقویتی در حقیقت یک تجارت از حالت به عمل است.  
 تابعی که در یک حالت خاص یک عمل خاص را انتخاب می کند، سیاست نام دارد.  
Policy

مرکزی که انجمنی دیگر در ارزیابی Policy Evaluation نامیده می شود به سیاست بهینه برسیم.  
 در ربات ای مثل سگ

حالت: زاویه، سرعت، شتاب، فعلی و هدف  
 یادمان: حالت جدید، مقدار به حالت مطلوب نزدیک است.  
عمل: حرکت یا گشت در اعمال به صورت فعلی

سیاست: انتخاب عمل مناسب در حالت های مختلف  $\pi(a|s)$   
 نیکبخت  $u-k$

کار یادگیری تقویتی تعیین بهترین سیاست ممکن است.  
تعیین سیاست بهینه از طریق تعالی با عمل



$k$  تعداد  $(u-k)$

وضعیت جدید، مقدار به مطلوب نزدیک است.

تغییر بین ربات و عمل باعث می شود سیاست بهینه از حالت به عمل فراهم شود.  
 الگوریتم های یادگیری تقویتی می توانند مداره در حال یادگیری باشند.

سیاست تابعی از حالت به عمل با پارامترهای قابل تنظیم  
 $\pi(a|s)$   
 نیکبخت حالت

تعیین بازی با دنبال  $look up$  سیستم

یادگیری: فرایندی سیستم برای تنظیم پارامترهای سیاست به منظور رسیدن به سیاست بهینه  
یادگیری تقویتی: به روز رسانی پارامترهای سیاست از روی بازمانده و هزینه

sam



**عمل و عامل** با توجه به اصدات تشریفاتی شوند.  
 در کل هر چیزی به غیر از عامل غیر از عیار را تکرار می دهد.  
 عیار می تواند هم **بسیار** و هم **نیز** باشد.

که دست بالاسر، سادس، فردی  
 که برت یادگیری بالاسر و تیر سازی حالات مختلف است

بعد از تعیین هدف عیار باید **اصدا** تعیین شوند.  
 با تعیین **یاداری** **باسب** برای عمل عامل که عامل انجام می دهد، اصدات تحت سزا خواهند کرد.  
 یاداری را باید به شکل باشد که **بهره سیاست** **انتخاب عمل** **بسیار** را نشان می دهند.

در یادگیری تقویتی **یاداری** اغلب اصدات و تابع از حالت و عمل هستند.

$$\text{Reward} = \text{Function}(\text{State}, \text{Action})$$

یاداری 3- سه مرحله 3 وادی

goodness → یاداری میزان خوب بودن انجام یک عمل در یک حالت خاص را نشان می دهد.

تقل  $LQR$  ،  $x \in A x + B u$

$$J = \frac{1}{2} \int_0^{\infty} (x^T Q x + u^T R u) dt$$

**یاداری** هدف از این است که این یاداری بنیم شود.  
 هدف این است که این  $a$  را به گونه ای انتخاب شوند که تغییراتی حالت بنیم به صفر صفر شوند.  $x \rightarrow 0$   
 صریح می نمایم  $x \rightarrow 0$  باید عبارت از بنیم شود.  $R$  نشان دهنده میزان صرف انرژی است و  $Q$  هم نشان دهنده صرف انرژی به صفر صفر است. می نمایم با کمترین صرف انرژی به صفر صفر شوند.

در کل تیزی در فرمولاسیون یاداری وجود ندارد. (**زیر بنیم یاداری** **تقدیر**) **sparse** تیزی کمتری به نونج **look up roll**

$$⑤ \xrightarrow{a} ⑤'$$

یاداری می تواند برای هر عمل، برای یک انتقال حالت خاص و یا برای هر یک از مودین باشد.

تقل **نیم سادس** :

در  $t$  و  $t+1$  به مودین دشاری داریم، **Action** انجام می دهیم و برای هر  $t$  که اکنون در حالت انتقال باشد، یاداری + اصدات می یابیم.



