

در FMDP هم انتقال در حالت نگه داشتن و دارای چندین مقدار عدد هستند.

$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, S_3$

$$P: S \times R \times S \times R \rightarrow [0, 1]$$

تایید مایل FMDP

بازه فراموش (پوشش) و روی د در فریب (تاریک) تراداند

$$P(s', r | s, a) \equiv P\{S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a\} \rightarrow \{ارو\}$$

همه مجموعه مشخصی کند، اگر در حالت s با a و a را انجام بدهیم، مقدار انتقال در دو بار در حالت جدید s' شود و پاداش را دریافت کنیم.

چون در RL اغلب موارد تعدادی هستند با هم و یک رابطه شکل تابع انتقال تفریق کنیم ولی اگر $deterministic$ بود این انتقال با اکتساب میسر.

	A^2		B
			$+5$
			B'
			19
	A'		
	22		

$$P(22, 9 \mid 2, L) = 0$$

اگر در 2 با a و a را انجام بدهیم تا به 22 می رویم و پاداش 9 خواهد بود
تبدیل و در انتقال میسر شود، ولی اگر مثل زیر باشد:

$$P(22, 10 \mid 2, L) = 1$$

خاصیت هم بداندی MDP اینکه برای عملیاتی تعدادی می تواند امتداد به عملیاتی که R در تمام آنها باشد: $R \in \{ \frac{3}{10}, \frac{4}{6}, \frac{2}{3}, \frac{1}{1} \}$

اگر عملیات باشد

$$P(22, 10 \mid 2, L) = 3$$

$$P(20, 1 \mid 19, R) = 0$$

نالی

چون بداندی منتها: A به A' و B به B' بگردیم پاداشی نمیگیریم انتقال بلا ضرر است هر چند اگر با a و a در 19 میگردیم به 2

می رویم

باید به A و A' بگردیم و در A' بگردیم

نکته: دینامیک سیستمی MDP که تابع منفرد احتمال است.

به ازای انجام عملی در هر حالت خاص، حالت بعدی و پاداش در داخل محیط MDP هستند.

$$f(n) \rightarrow \int_{-\infty}^{\infty} f(x) dx$$

$$\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1, \text{ for all } s \in S, a \in A(s)$$

این معادله برای هر حالت s و هر اقدام a برقرار باشد.
 تغییر در احتمال s' و r است. s' و r اثر به احتمال را جمع کنیم باید بود به دست بیاید.

ماست مارکوف (حالت بعدی و پاداش تنها به حالت قبل و عمل انجام شده در آن وابسته است).

نتیجه: حالت بعدی به تاریخچه وابسته نیست و تنها به حالت قبل و عمل آن قبل دارد. این حالت در t و $t+1$ به یکدیگر وابسته است و به ازای هر t به ازای $t+1$ وابسته است.
 آخرین حالت قبل و اولین انجام شده.
 گزافه واضح نیست و در واقع حالت دی قبل قبل به حالت بعدی تأثیری ندارند.

← تابع P که نشان دهنده عمل را مشخص می کند.

① $p(s', r | s, a) \rightarrow$ این را داریم و دست می یابیم
 تابع زیر است:

② $p(s' | s, a) = \Pr\{s_t = s' | s_{t-1} = s, A_{t-1} = a\}$ و $\sum_{r \in R} p(s', r | s, a)$
 (توابع تویلی حایه ای)

یادآوری:
 $p(x, y)$
 $p_x(x) = \int_{-\infty}^{\infty} p(x, y) dy$

reward را به ریاضت کنیم به شکل s_t و a_{t-1}
 آن a را انجام دهیم.

③ $r(s, a) = E[R_t | s_{t-1} = s, A_{t-1} = a] = \sum_{r \in R} r \sum_{s' \in S} p(s', r | s, a)$
 اگر r deterministic باشد این expected value و با عددی ثابت بازمی آید.

آر جیز، فاصله باشد $R \in [5/3, 5/6]$
 اگر $5/3 \times 3 + 5/6 \times 6 = 4.5$
 اگر قرار بداریم Action را انجام دهیم 4.5

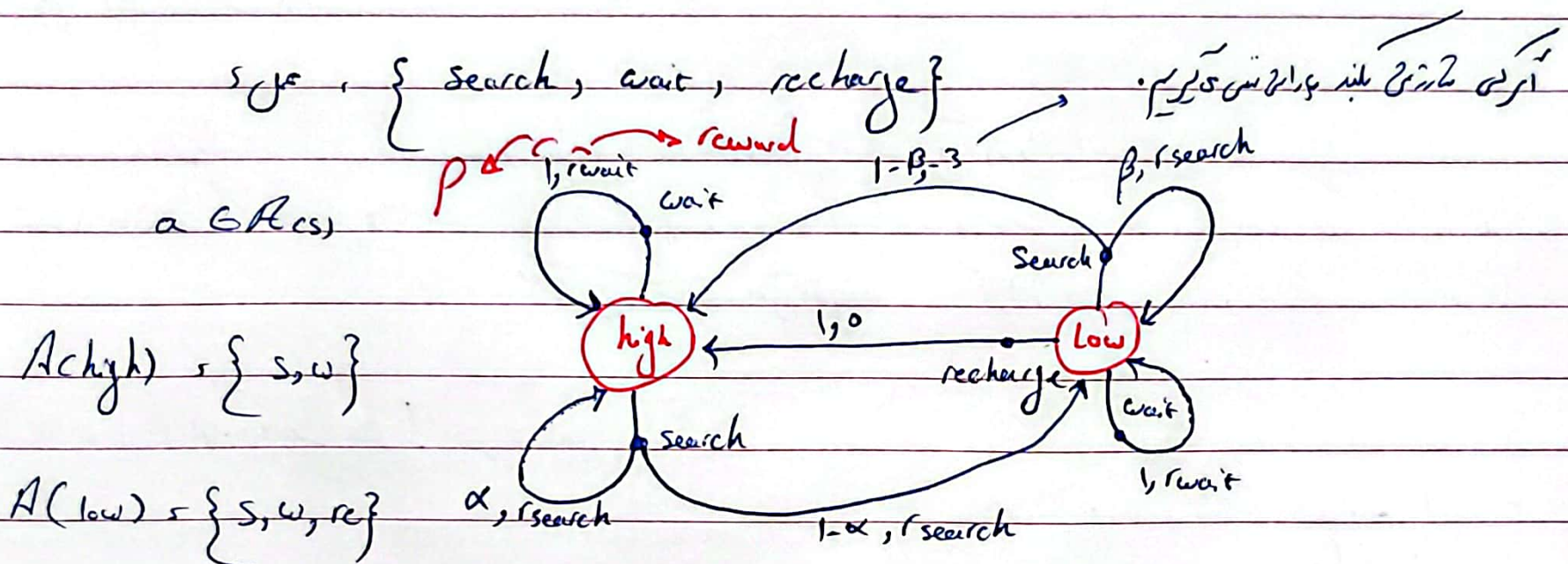
یادآوری:
 $E[x], \sum x p_x(x)$

sam

$$Q(s, a, s') \in E [R_t | s_{t-1} = s, A_{t-1} = a, S_t = s'] = \sum_{r \in R} r \frac{P(s', r | s, a)}{P(s' | s, a)} \quad \left. \begin{array}{l} \text{جائزہ} \\ \text{نشان دہی} \end{array} \right\}$$

جائزہ، اندازہ، جمع بہتر ① کے طریقہ میں RL داخل ہیں۔

رہت جمع اس کی مثال: $\{low, high\}$ حالت (برائے وقت بھری)



$$P(high, r_s | high, s) = \alpha$$

$$P(low, r_s | high, s) = 1 - \alpha$$

تکلیف بہتر

قبل

s	a	s'	$P(s' s, a)$	$r(s, a, s')$
high	search	high	α	r_{search}
"	"	low	$1 - \alpha$	"
low	"	high	$1 - \beta$	-3
"	"	low	β	r_{search}
high	wait	high	1	r_{wait}
"	"	low	0	-
low	"	high	0	-
"	"	low	1	r_{wait}
"	recharge	high	1	0
"	"	low	0	-

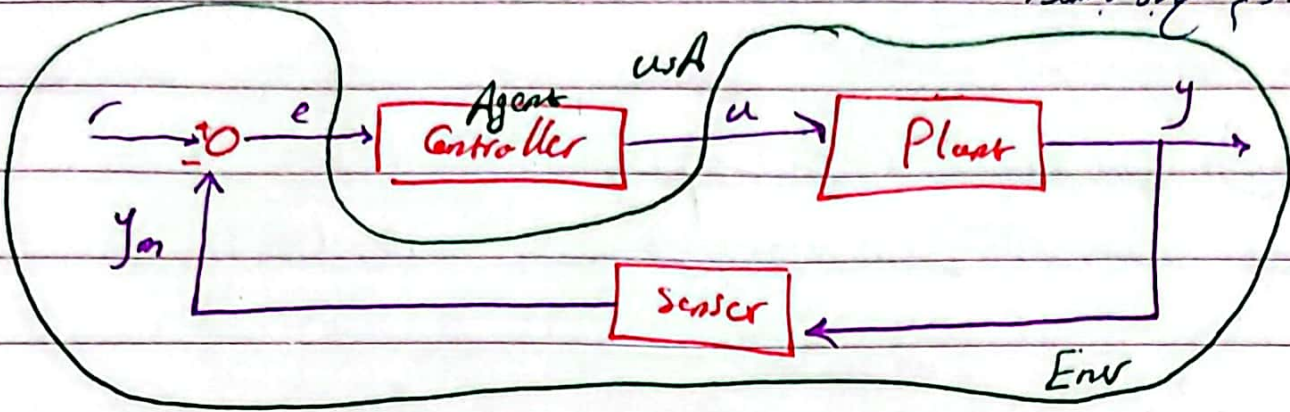
sam

اداره ترانژیدای تقسیم بار کلاف عدد

★ در MDP لازم نیست متغیر و پیافور دان داشت باشد و می تواند متغیر باشد و این عمل را نشان دهد.

time step $t \rightarrow 0 \rightarrow 1 \rightarrow 2$

عمل در MDP هم می تواند متغیر باشد و هم سلسله باشد.



یک حالت دیگر است که می توانیم بایسیم سیستم حالتی به راب عنوان Env در نظر بگیریم. در این A در ی ترسیم هلاک RL یادای س درین ی کند.

★ در یادگیری تقویتی اهداف عامل در یاداری نهفته می شود. هدف عامل انتخاب عملی است که بانی یاداری در مدالانی دست باز می شود. انتخاب یاداری باید به دست انجام شود تا به نتیجه برسد و نتواند به نتیجه حاصل یاداری است. نمایه دانشی بهر از طریق یاداری به عامل داده شوند و یاداری برای بیان هدف دست نهاده رسیدن به هدف.

Return \rightarrow Cumulative Reward

به یاداری مجموع بازگشتی داریم هدف عامل بازگشتی آن است.
Return \rightarrow جمع یاداری های دریا تر به راب

$S_0, A_0, R_1, S_1, A_1, \dots$

G_0, R_1, \dots, R_T

G_0, R_1, R_2, \dots

★ عبادی ده ای (Episode Task): تعداد به بازی تکرار و از راز قرار تر \rightarrow دارای حالت ای دست Terminal State
★ عبادی بیو ته (Continued Task): تعداد به \rightarrow تکرار ترانژیدای متفر T_{∞} Terminal state نامند.

در عبادی ده ای و عبادی بیو ته حالت اولیه در ده بیو به صورت لا شعل انتخاب می شود.

طری عبادی بیو ته اغلب بازگشتی می شود به با تریب نرخ عفت Discount Rate به یاداری ای این اصیت بیو ته داده می شود.
Discount Factor sam

$$G_t \equiv R_{t+1} + \delta R_{t+2} + \delta^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \delta^k R_{t+k+1} = \sum_{k=t+1}^{\infty} \delta^{k-t-1} R_k$$

عالم سری کند مریقر به اسد ببرد و اسن را بیزمان ای اصلان مکرر کند.

حالت صاف

$\delta = 0$ → عدم توجه به آینده و تنها یاداری فعلی کان-سمی شود.

$\delta = 1$ → تمامی یاداری های دریاثر در آینده اصصا کلبان دارند (دانشی)

$$G_t \equiv R_{t+1} + \delta R_{t+2} + \delta^2 R_{t+3} + \dots \\ = R_{t+1} + \delta (R_{t+2} + \delta R_{t+3} + \dots)$$

$$G_t \equiv R_{t+1} + \delta G_{t+1}$$

$$R_{t+1} \rightarrow G_t = \sum_{k=0}^{\infty} \delta^k, \quad \frac{1}{1-\delta}$$

$$\text{اثر از } k_0 \rightarrow \frac{\delta^{k_0}}{1-\delta}$$

برای ص از انو باید شروع کنیم و رسم به G_t مالا.

نسبت و تابع ارزشی:

اختیار الان و بعدی یا دیگری تدریس **تابع ارزشی** را تجزیه می کنند.

تابع ارزشی حالت اول: نشان دهنده میزان ارزشمندی حالت s

تابع ارزشی حالت دوم: نشان دهنده میزان ارزشمندی انجام عمل a در حالت s

ارزشی حالت s نسبت به سایر حالت ها: $V_{\pi}(s) \equiv E_{\pi}[G_t | S_t = s] = E_{\pi}[\sum_{k=0}^{\infty} \delta^k R_{t+k+1} | S_t = s], \text{ for all } s \in S$

$q_{\pi}(s, a) \equiv E_{\pi}[G_t | S_t = s, A_t = a] = E_{\pi}[\sum_{k=0}^{\infty} \delta^k R_{t+k+1} | S_t = s, A_t = a]$

که بین ارزشی a برای حالت s میزنند.

$V_*(S) = \max_{a \in A(S)} q_*(S, a)$

معادلات بهینگی

$= \max_a E_{\pi_*} [G_t | S_t = S, A_t = a]$

$= \max_a E_{\pi_*} [R_{t+1} + \gamma V_{\pi_*}(S_{t+1}) | S_t = S, A_t = a]$

$= \max_a E [R_{t+1} + \gamma V_*(S_{t+1}) | S_t = S, A_t = a]$

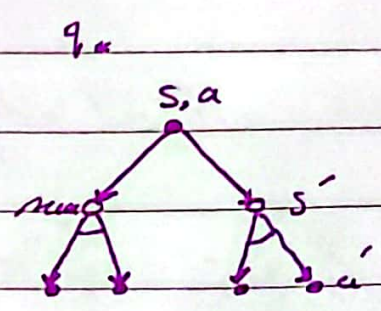
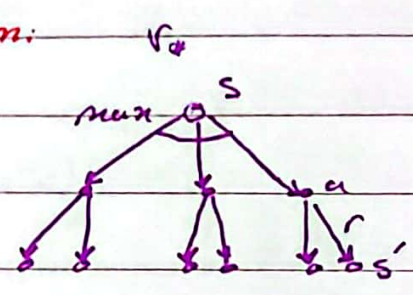
اینجا $E_{\pi_*}[V_{\pi_*}(S_{t+1})]$ هم می تونه این ترم

$= \max_a \sum_{S', r} P(S', r | S, a) [r + \gamma V_*(S')]$

$q_*(S, a) = E [R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') | S_t = S, A_t = a]$
 $= \sum_{S', r} P(S', r | S, a) [r + \gamma \max_{a'} q_*(S', a')]$

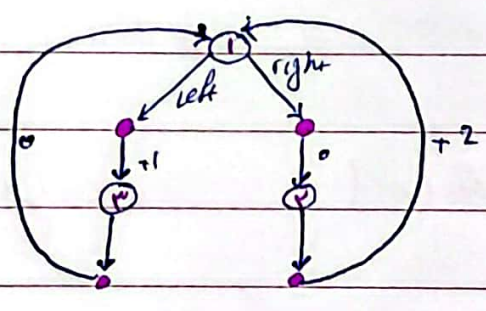
به صورت بهینه داریم:

backup diagram:



این ترم $\max_{a'} q_*(S', a')$ تکرار می شه

$\pi_*(s) = \arg \max_a q_*(S, a)$



دفعه اول: برای فرخ بخشیدن به هر کار و در نهایت بهینه کردن.
 چون می تونه بهینه از حالت return می تونه بهینه نیازیست.

$G_0 = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$

Right c

$= 0 + \gamma \cdot 2 + \gamma^2 \cdot 0 + \dots$
 $= 2\gamma + 2\gamma^3$

$V_*(1) = G_0 = 2\gamma(1 + \gamma^2 + \dots) = \frac{2\gamma}{1 - \gamma^2}$

$G_0 = 1 + \gamma \cdot 0 + \gamma^2 \cdot 1 + \dots$

left c

$= 1 + \gamma^2 + \gamma^4 + \dots = \frac{1}{1 - \gamma^2}$

s.a.m

y	Right, $\frac{2y}{1-y^2}$	Left, $\frac{1}{1-y^2}$
0	0	1
0.5	$\frac{4}{3}$	$\frac{4}{3}$
0.9	$\frac{1.8}{1-0.9^2}$	$\frac{1}{1-0.9^2}$

s.a.m