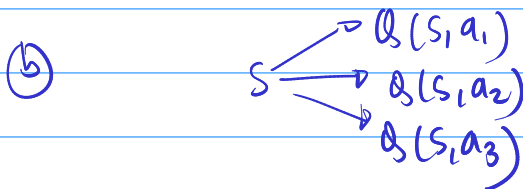
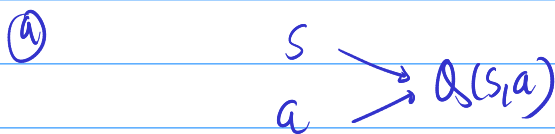


Real-time reference notes / Lecture 01

- ① RL basics: States, actions, rewards, returns
- ② Intuition behind value-based methods, tabular version
- ③ Filling up the Q-table using ϵ -greedy algorithms
- ④ Scalability of tabular RL
 - Ⓐ size of table
 - Ⓑ continuous states & actions
 - Ⓒ visiting each (s, a) pair
- ⑤ Approximation of $Q(s, a)$ using neural networks



- ⑥ Terminology
 - Ⓐ Model-free & Model-based
 - Ⓑ online & offline
 - Ⓒ on-policy & off-policy

⑦ Bellman equation

↳ ϵ -greedy using Bellman equation

⑧ MKS paper (ATARI)

⑨ Issues with basic DQN

↳ Bootstrapping
↳ Stabilisation
↳ Last layer discrimination
↳ Sampling memory



$$\dot{x} = Ax + Bu$$

State

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

$$x_{k+1} = A^* x_k + B^* u_k$$

action

transition

$$J = \sum_0^{\infty} [x_k^T Q x_k + u_k^T R u_k]$$

cost (t_k)

$$G_t = \sum_{\tau=t}^{\infty} \gamma^{\tau-t} r_{\tau}$$

s_t, a_t

$$|r_t| \leq R_0$$

$$\max G_t = \frac{R_0}{1-\gamma}$$

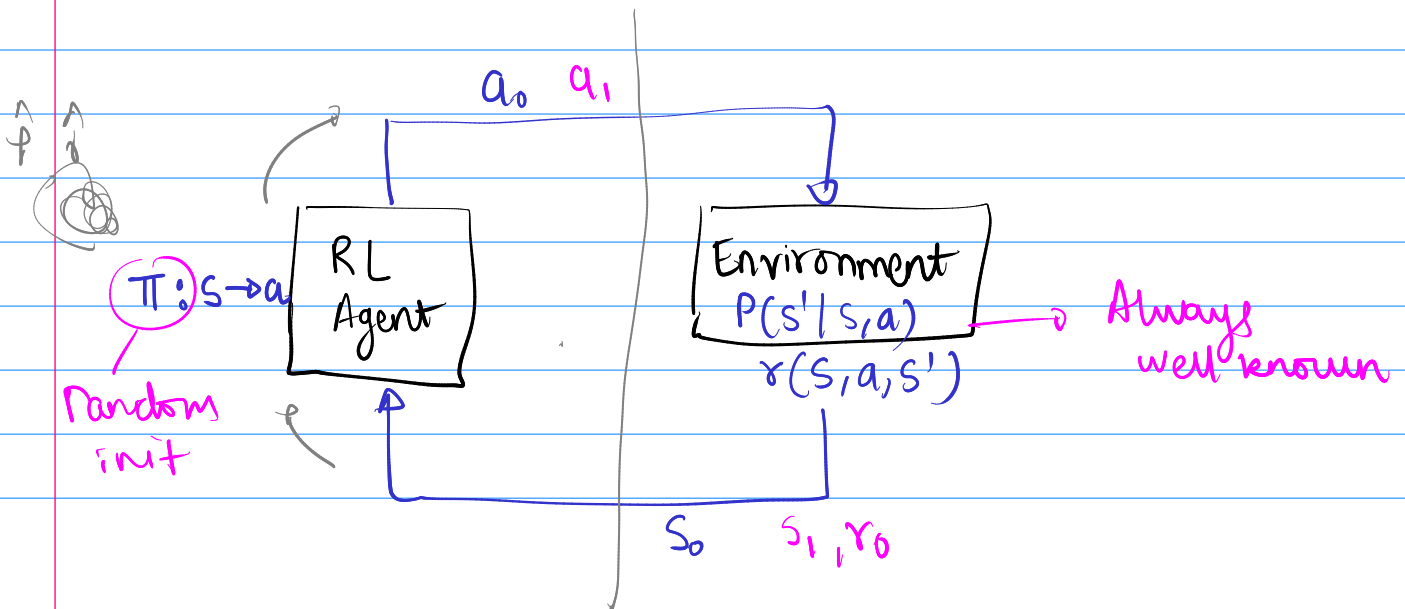
$$P(s_{t+1} | s_t, a_t)$$

transition prob.

Markov decision process (MDP)

$s_{t-1}, a_{t-1}, s_{t+1}, \dots, s_0, a_0$
independent

$(s, a, r, P, \gamma) \rightarrow$ defines MDP

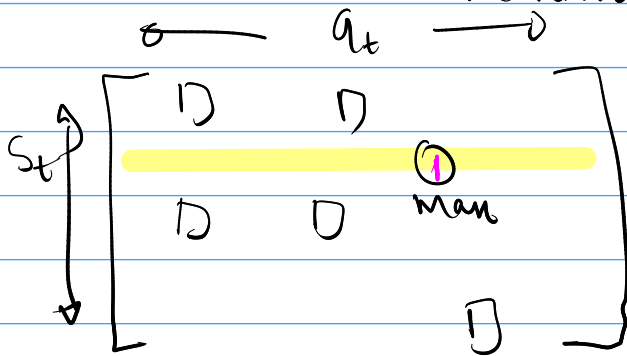


$$\max_{a_t} \boxed{G_t = r_t + \gamma^1 r_{t+1} + \gamma^2 r_{t+2} + \dots} = r_t + \gamma G_{t+1}$$

$\pi: s_t \rightarrow a_t$
(r_t)

$\hat{G}_t(a_t=0)?$ $\hat{G}_t(a_t=1)?$

maximum



Tabular RL

ϵ -greedy exploration

\rightarrow w.p. ϵ choose randomly (exploration)
 \rightarrow w.p. $(1-\epsilon)$ choose argmax (exploitation)

$$\hat{G}_t \checkmark \quad G_t = r_t + \gamma G_{t+1}$$

$$Q(s_t, a_t) = r_t + \gamma Q(s_{t+1}, a_{t+1})$$

rewrite

$$\boxed{Q(s_t, a_t) = r_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})}$$

Bellman equation

neural net
for this regression problem

$$L = [Q(s_t, a_t) - r_t - \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})]^2$$

$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \Rightarrow \text{circle} \quad 0 \rightarrow Q$

$$\begin{matrix} 0 \\ s_t & 0 \\ & 0 \end{matrix} \longrightarrow Q(s_t, a_t)$$

$$\begin{matrix} 0 \\ a_t & 0 \\ & 0 \end{matrix}$$

$$\begin{matrix} 0 \\ 0 \\ 0 \end{matrix} \longrightarrow \begin{matrix} 0 \\ 0 \end{matrix} \begin{matrix} \rightarrow Q(s_t, a_1) - [r_t + \gamma Q(s_{t+1})] \\ \rightarrow Q(s_t, a_2) - Q(s_t, a_2) \end{matrix}$$