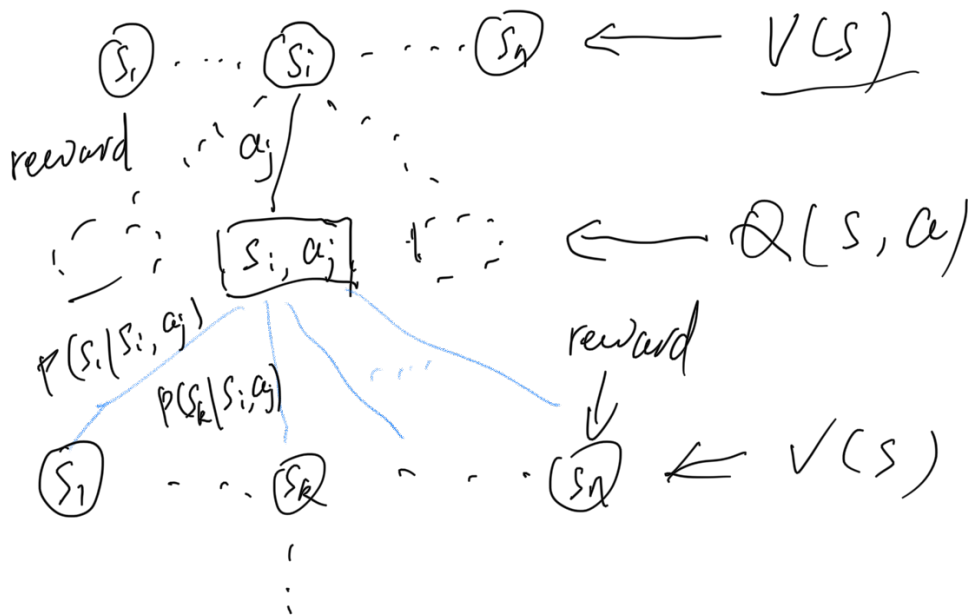


Lecture 12

$$V(s) = \max_{a_i} Q(s, a_i) \quad \text{i action}$$

$$Q(s, a) = \underbrace{r(a)} + \underbrace{\sum_{s'} P(s'|s, a) (r(s') + \lambda V(s'))}_{\text{stochastic}} \leftarrow$$



reward from action vs. reward at a state

Deterministic vs. stochastic

$$(s, a) \rightarrow s'$$

$$P(s'|s, a) \quad s' \in \{s_1, \dots, s_n\}$$

	s_1	s_2	\dots	s_n
a_1	s_5	s_{16}	\dots	\dots
a_2	\vdots	\vdots	\dots	\dots
\vdots	\vdots	\vdots	\dots	\dots
a_m	\vdots	\vdots	\dots	\dots

$$P(s, a, s') \text{ or } T(s, a, s')$$

Deterministic:

$$\dots r(s') + \lambda V(s') \leftarrow$$

$$Q(s, a) = r(a) + \gamma V(s')$$

stochastic:

$$Q(s, a) = r(a) + \sum p(s'|s, a) (r(s') + \gamma V(s')) \leftarrow$$

on policy vs. off policy

Q-learning — off policy

Q-table with some initial values.

Go through a number of episode
start an episode with a S

Go through the episode step by step
(Q table)

↑ $Q = Q_{\text{new}}$
(assign updated Q_{new} as the Q)

→ based on Q , choose an action a ,
calculate V .

$$V = \max_a (Q(s, a))$$

↑ decide with a to take

	s	s_1	s_2	...	s_n
a			10		
a_1			20		
a_2			15		
a_3			50		
a_4			-10		
a_m			8		

Q-table

so take action a_{a_3} .

→ Take action a , get reward, get to a state s'
if deterministic or real-time

is current Q

episode — a sequence of
action & state.

Goal state is the end of
an episode.

or use the length limit
to stop an episode.

compute current

$$Q_{\text{curr}} = \underline{r} + \lambda V(s')$$

if stochastic, in simulation after data collection (known transition probability $P(s'|sa)$)

$$Q_{\text{curr}} = V(a) + \sum_{s'} P(s'|sa) (r(s') + \lambda V(s'))$$

→ update Q table

$$Q_{\text{new}}(s,a) = Q(s,a) + \underbrace{\alpha}_{\text{learning rate}} (Q_{\text{curr}}(s,a) - Q(s,a))$$

$$s \leftarrow s'$$

continue until to the end of the episode