

Lecture 18 - March 21, 2023

- Temporal Difference Learning

- TD(0)

- SARSA

- Q-Learning

- On-Policy Vs. Off-Policy

↓

- Expected SARSA

- Double Q-Learning

- Multi-Step Bootstrapping

- SARSA-Lambda

- Actor-Critic Method

HW4 → Due March 31

Exam 2 → Tues, April 4

Project 3 → Due April 14

TA's office hour:

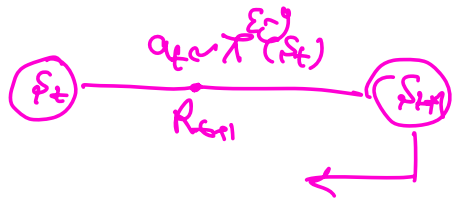
Wednesdays, 2pm-3pm (in-person)

Fridays, 2pm-3pm (virtual)

SARSA & Q-Learning

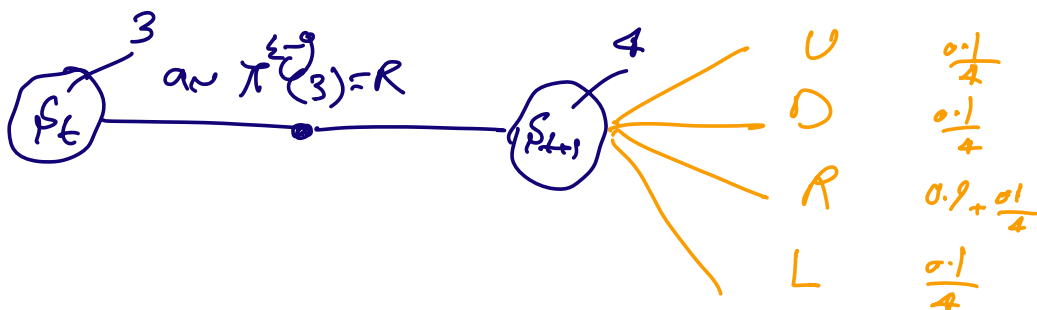
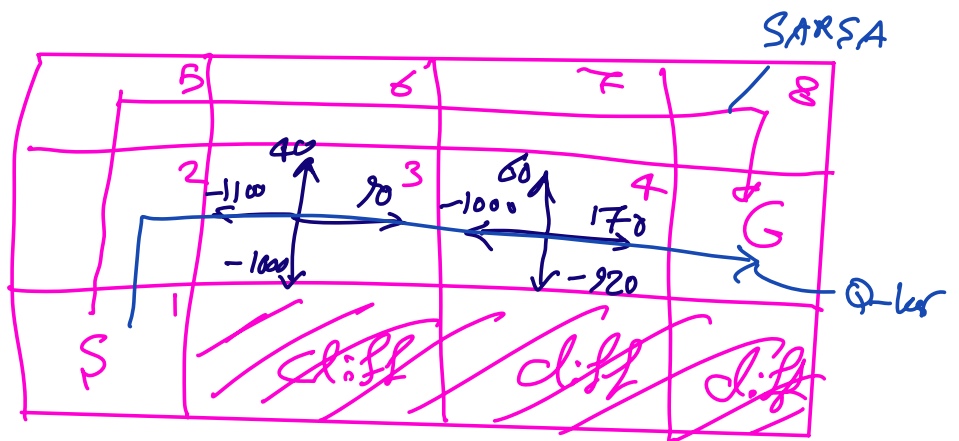
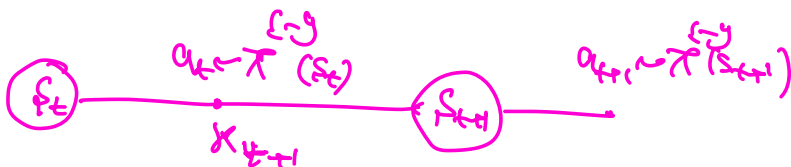
Q-Learning

$$Q(S_t, a_t) = Q(S_t, a_t) + \alpha [R_{t+1} + \gamma \max_{a \in A} Q(S_{t+1}, a) - Q(S_t, a_t)]$$

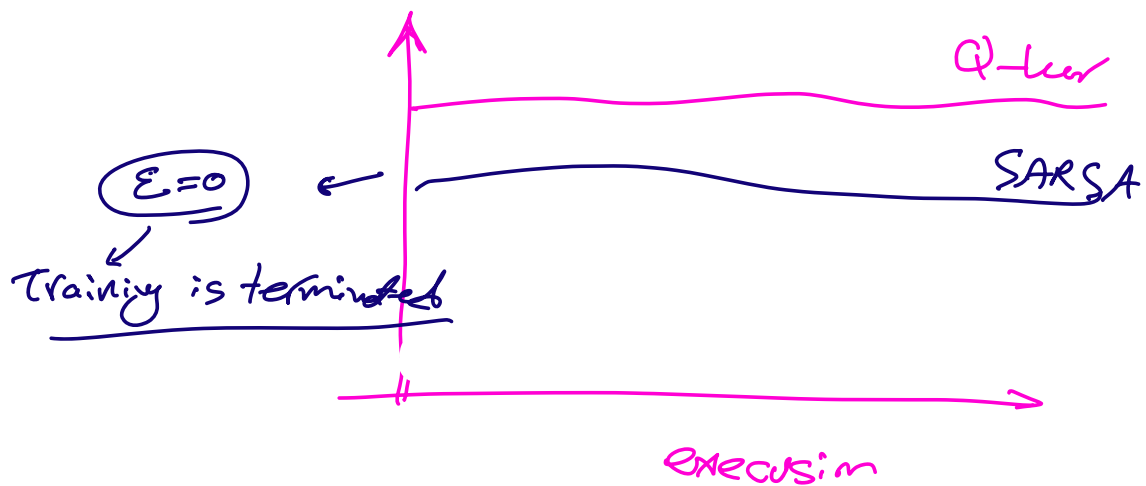
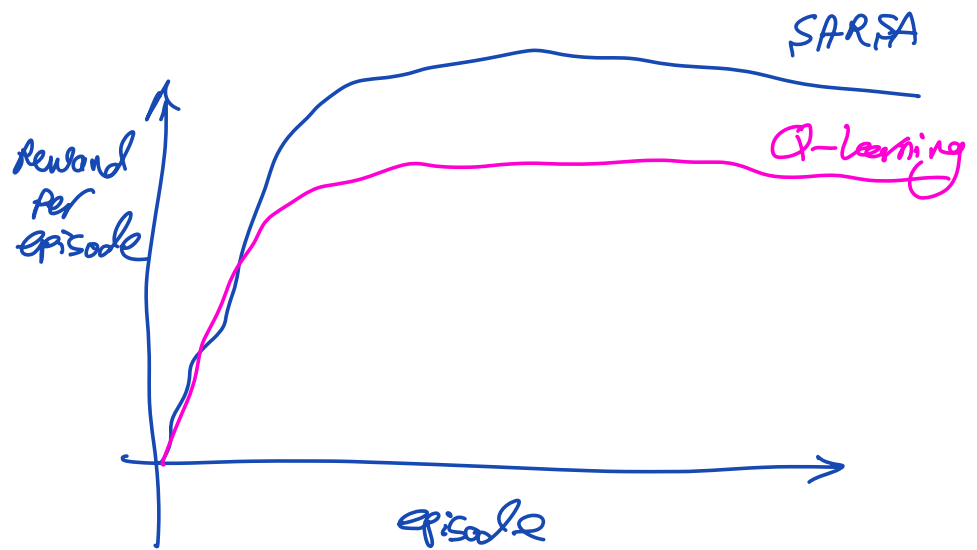


SARSA:

$$Q(S_t, a_t) = Q(S_t, a_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, a_{t+1}) - Q(S_t, a_t)]$$



$$Q(3, R) = Q(3, R) + \alpha \left[-1 + \gamma \frac{Q\text{-learning } Q(4, R) - 170}{\text{SARSA } Q(4, a)} - Q(3, R) \right]$$



Can we have replacement of ϵ -greedy?

Boltzman policy = Softmax policy

$$\pi^{\epsilon\text{-g}}(s_t) = \begin{cases} \arg\max Q(s_t, a) & 1 - \epsilon \\ \text{Random} & \epsilon \end{cases}$$

$$\pi^{\text{Bolt}}(a | s_t) = \frac{e^{\frac{Q(s_t, a)}{\tau}}}{\sum_{a' \in A} e^{\frac{Q(s_t, a')}{\tau}}} \quad : \text{Differentiability}$$

$$a \in \{a^1, a^2\}$$

$$Q(s_t, a^1) = 10 \rightarrow \pi(a^1 | s_t) = \frac{e^{10}}{e^{10} + e^{100}} = 0.001$$

$$Q(s_t, a^2) = 100 \rightarrow \pi(a^2 | s_t) = \frac{e^{100}}{e^{100} + e^{10}} = 0.999$$

Consider Stochastic & Reward is Stochastic

Q-Learning { Works very well in high stochasticity
" " " large state spaces
→ Batch Data ($\epsilon \rightarrow$ as you want)

Issues: { Overestimation of Q-Values
(Double Q-Learning)
Risky in training

SARSA: { - Conservative (both training & Execution)
- Converge Fast

Issues → { Performs poorly in high rewards with high variance
on-policy → cannot arbitrary ϵ X

Expected SARSA:

(Robust but slower)

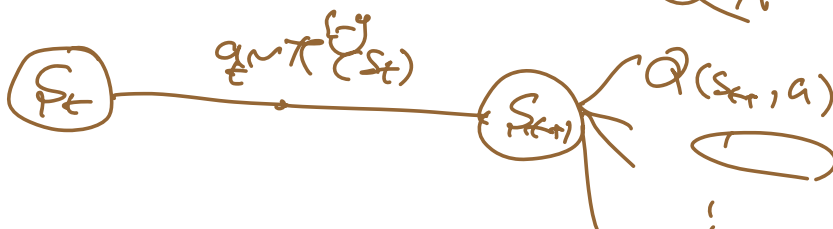
$$\begin{cases} Q(4, R) = 170 & 0.925 \\ Q(4, D) = -920 & 0.025 \\ Q(4, L) = -1600 & 0.025 \\ Q(4, U) = 60 & 0.025 \end{cases}$$

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha [R_{t+1} + \gamma \underbrace{0.925 \frac{Q(4, R)}{170} + 0.025 Q(4, D) + 0.025 Q(4, L) + 0.025 Q(4, U)}_{-Q(s_t, a_t)} - Q(s_t, a_t)]$$

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha [R_{t+1} + \gamma \sum_{a \in A} \pi(a | s_{t+1}) Q(s_{t+1}, a) - Q(s_t, a_t)]$$

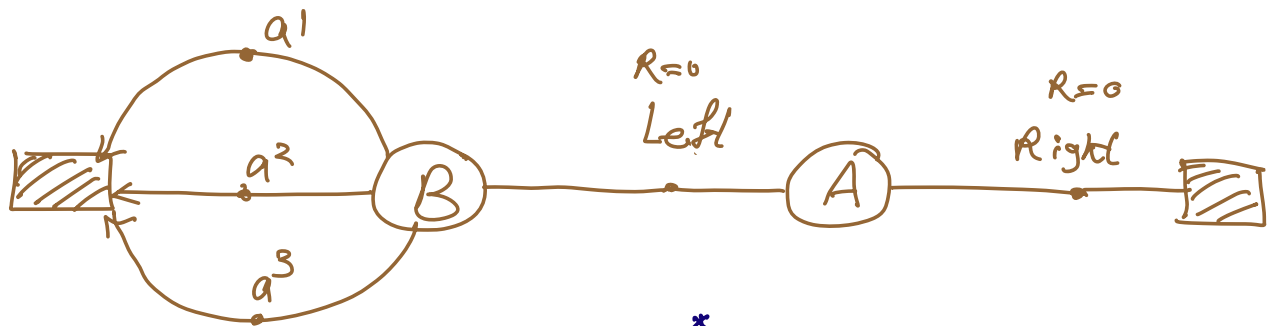
on-policy \Rightarrow you can not set ϵ arbitrary

$$Q \sim \pi^{\epsilon, \gamma}$$



$$\frac{\epsilon}{4} \\ 1 - \epsilon + \frac{\epsilon}{4} \\ \frac{\epsilon}{4} \\ \vdots$$

Double Q-Learning (overestimation of Q-values)



$$R(B, a^1) \sim \mathcal{N}(0, 1)$$

$$R(B, a^2) \sim \mathcal{N}(-0.2, 1)$$

$$R(B, a^3) \sim \mathcal{N}(0.3, 1)$$

$$Q^*(A, R) = 0$$

$$\pi^*(A) = L$$

$$\pi^*(B) = a^3$$

$$Q^*(A, L) = 0.3$$

Q-learning

Episode 1

Random State: B $\xrightarrow{a^2 \sim \pi^*(B) = \begin{cases} 0.33 & a^1 \\ 0.33 & a^2 \\ 0.33 & a^3 \end{cases}}$ T

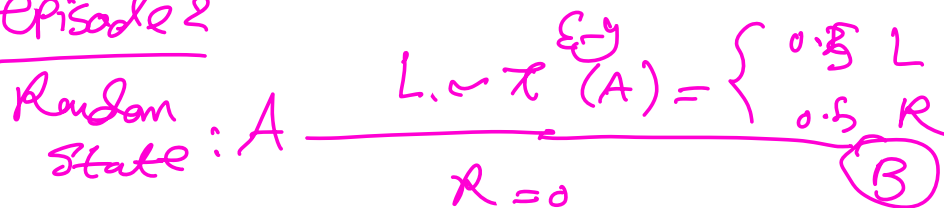
$R \sim \mathcal{N}(-0.2, 1) = 0.8$

$$Q(B, a^2) = Q(B, a^2) + \alpha [0.8 + \gamma \max_q Q(T, q) - \underbrace{Q(B, a^2)}_{=0.4}]$$

$$Q(B, a^2) = E[\underbrace{R_{t+1} + \gamma R_{t+2} + \dots}_0 \mid S_t = B, a_t = a^2] = -0.2$$

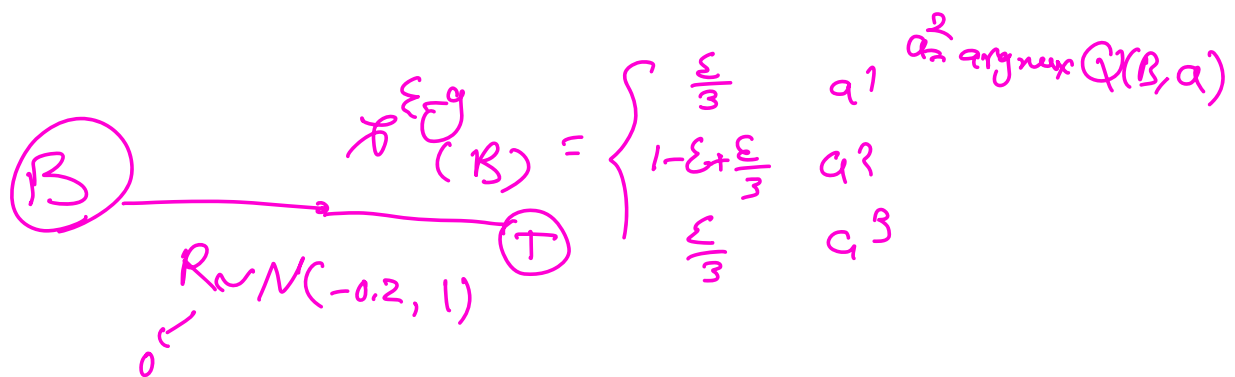
$\rightarrow \mathcal{N}(-0.2, 1)$

episode 2



$$Q(A, L) = Q(A, L) + \alpha \left[\underbrace{\mathcal{R}}_0 + \underbrace{\gamma \max_a Q(B, a)}_{0.4} - Q(A, L) \right]$$

$$= 0.2$$



$$Q(B, a^2) = \underbrace{Q(B, a^2)}_{0.4} + \alpha \left[\underbrace{\mathcal{R}}_0 + \underbrace{\gamma \max_a Q(B, a)}_{0.4} - \underbrace{Q(B, a^2)}_{0.4} \right] = 0.4$$

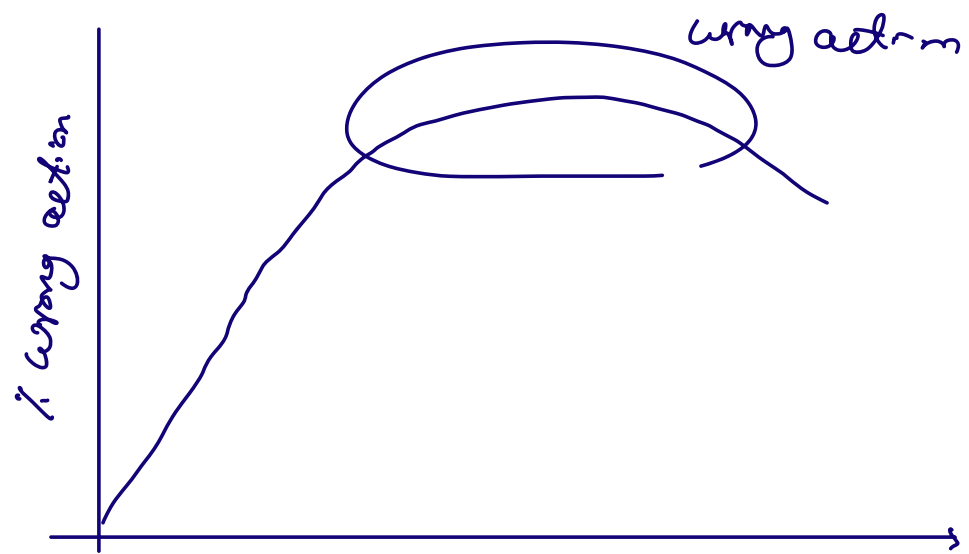
$$R^{a^1} \sim \mathcal{N}(-0.5, 5)$$

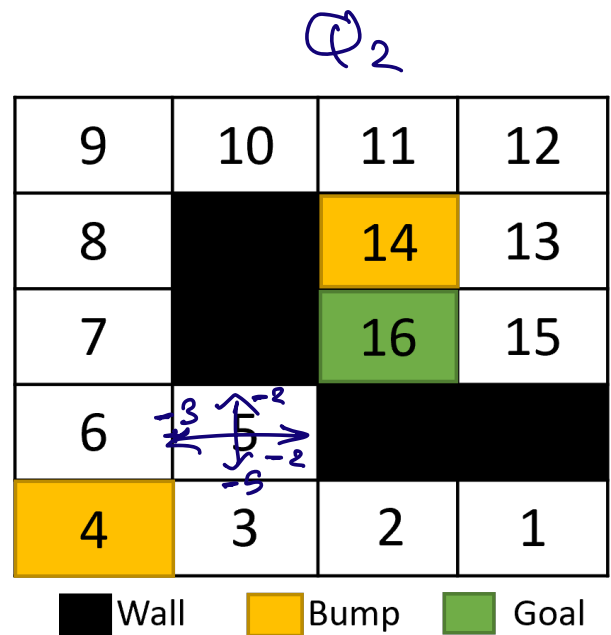
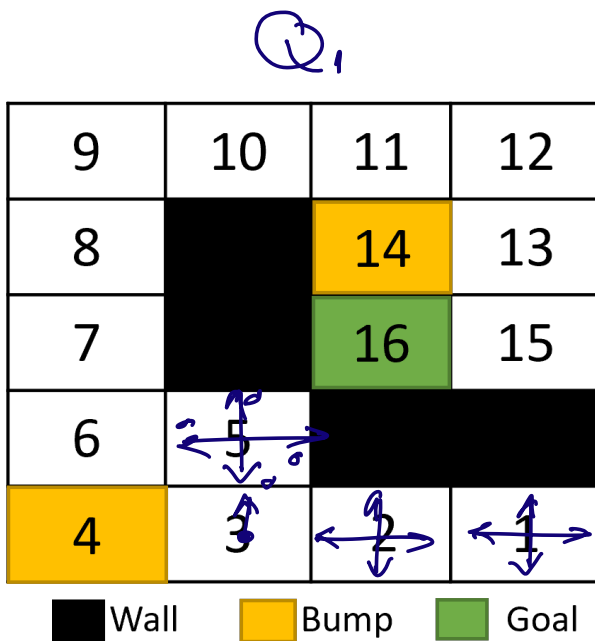
$$R^{a^2} \sim \mathcal{N}(-0.5, 5)$$

R^{a^1}	Q^{a^1} \downarrow μ^{a^1}	R^{a^2}	μ^{a^2}	$\max_a \mu^a$
3	3	-8	-8	3
-2	0.5	-5	-6.5	0.5
4	1.67	-8	-7	1.67
\vdots	\vdots	\vdots	\vdots	\vdots
	<u>-0.5</u>		<u>-0.5</u>	<u>1.72</u> $E[\max_a \mu]$

$$Q = Q + \alpha [R + \gamma \max_a \underbrace{Q(s, a)}_{\mu} - Q]$$

$$Q\text{-learning: } Q_{\pi}(s, a) = E_{\pi} [R + \gamma \max_{a'} Q(s, a') \mid s_t = s, a_t = a]$$





$Q \sim \pi(s) = \begin{cases} \text{argmax}_{a \in A} Q_1(s, a) + Q_2(s, a) \\ \text{Random} \end{cases}$
 $Q_1(3, a) + Q_2(3, a)$
 $1 - \epsilon$
 ϵ

w.p. 0.5

$$Q_1(3, u) = Q_1(3, u) + \alpha [R + \gamma Q_2(5, \text{argmax}_{a \in A} Q_1(5, a)) - Q_1(3, u)]$$

w.p. 0.5

$$Q_2(3, u) = Q_2(3, u) + \alpha [R + \gamma Q_1(5, \text{argmax}_{a \in A} Q_2(5, a)) - Q_2(3, u)]$$

Double Q-Learning

Initialize $Q_1(s, a)$ and $Q_2(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$, arbitrarily

Initialize $Q_1(\text{terminal-state}, \cdot) = Q_2(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize S

Repeat (for each step of episode):

Choose A from S using policy derived from Q_1 and Q_2 (e.g., ϵ -greedy in $Q_1 + Q_2$)

Take action A , observe R, S'

With 0.5 probability:

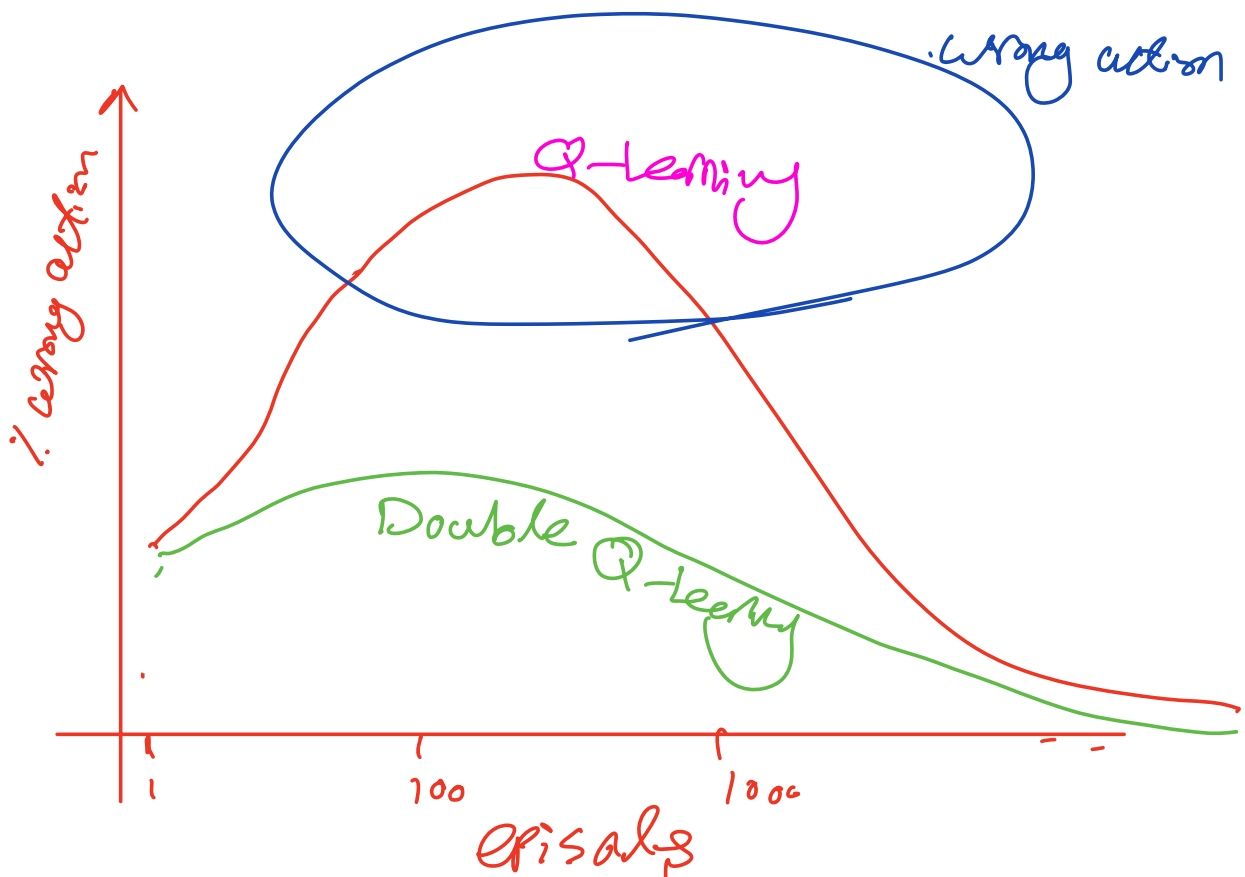
$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha \left(R + \gamma Q_2(S', \arg \max_a Q_1(S', a)) - Q_1(S, A) \right)$$

else:

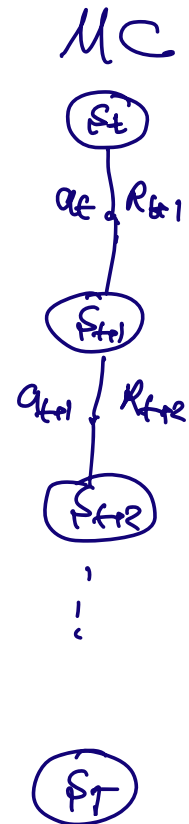
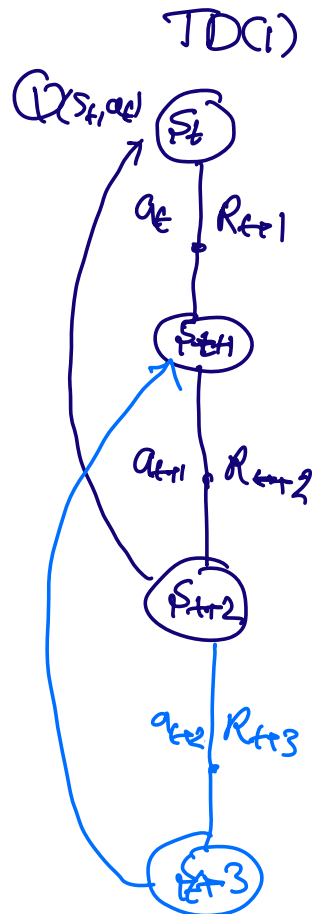
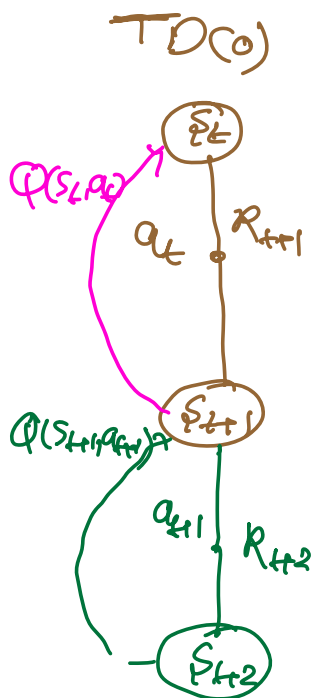
$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha \left(R + \gamma Q_1(S', \arg \max_a Q_2(S', a)) - Q_2(S, A) \right)$$

$S \leftarrow S'$;

until S is terminal



Multi-Step Bootstrapping



TD(0):

$$Q(S_t, a_t) = Q(S_t, a_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, a_t)]$$

TD(1)

$$Q(S_t, a_t) = Q(S_t, a_t) + \alpha [R_{t+1} + \gamma R_{t+2} + \gamma^2 \max_a Q(S_{t+2}, a) - Q(S_t, a_t)]$$

TD(n) larger variance \leftarrow Variance

TD(n) less Biased \leftarrow Biased
 here $n > 0$