

Markovian

MDP

T, R

$$V(s) \rightarrow E[G_t | s]$$

$$Q(s, a) \rightarrow E[G_t | s, a]$$

Bellman Eqn $V(s) \Leftrightarrow V(s')$

$$Q(s, a) \Leftrightarrow Q(s', a')$$

DP

Value Iteration
Policy Iter
 $\pi(a|s)$

Value Iter

$$V(s) = 0$$

for i

$$V(s) \leftarrow \max_a \left(R(s) + \gamma \sum_{s'} P(s' | s, a) V(s') \right)$$

\uparrow π^* \uparrow $\argmax(\)$

Policy Iter ✓

$V = 0, \pi \rightarrow \text{rand}$
evaluate π, V

$$V_\pi = \frac{R(s) + \gamma \sum_{s'} P(s' | s, \pi) V_\pi(s')}{1 - \gamma}$$

$V(a)$

T, R

RL: ~~T~~ T not known
 R not known

model-based \hat{T}, \hat{R}
model-free: $V \rightarrow$ Passive RL
 $Q \rightarrow$ Active RL

Passive RL
MC: \rightarrow Ignores Bellman
 $V(s) \downarrow \downarrow \downarrow \downarrow \downarrow \rightarrow T$
 \uparrow $\arg(R_i) \rightarrow i$

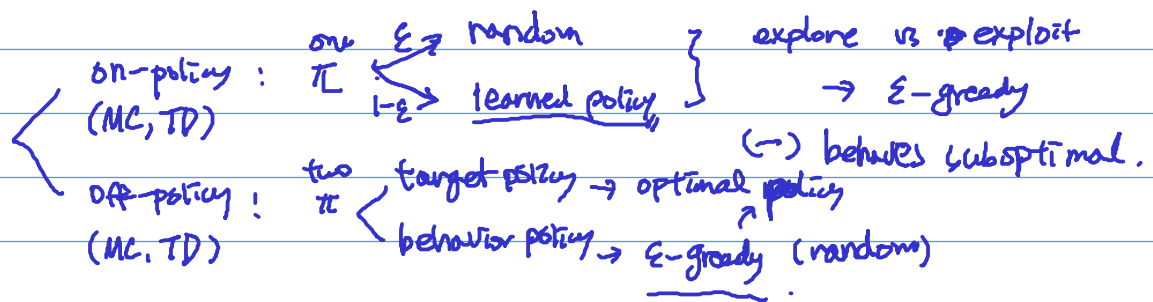
TD: Moving average $V \rightarrow$ Consistent, converges faster
one episode

$$s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_4 \rightarrow \dots \rightarrow s_T$$

$$V(s) \leftarrow V(s) + \alpha [R_k + \gamma V(s') - V_k(s)]$$

π_{fixed}

$\pi \rightarrow$ Improve Q -function



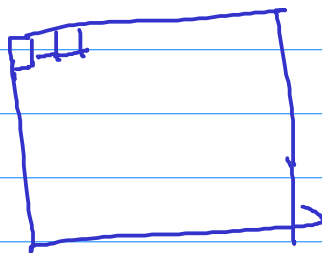
on-policy TD control : Sarsa

off-policy TD control : Q -learning

State \rightarrow Tables

$V(s)$

$Q(s,a)$



Approximation RL

function approximator \rightarrow "regression"

