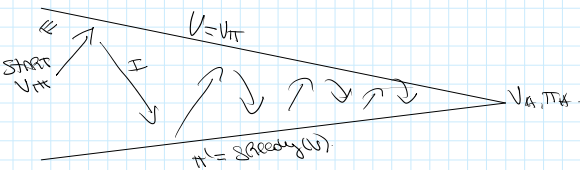


GOAL: Find  $V^*$  /  $\pi^*$  w/o knowledge of MDP dynamics

# Recap: Generalised Policy Iteration (GPI)



- We want to make this framework model-free.

## PROBLEM #1:

- We can't do model-free policy evaluation

For  $V_{\pi}$ , because acting greedy would require knowledge of the MDP

$$\pi' = \text{greedy}(\pi) = \underset{a}{\text{argmax}} Q_{\pi}(s, a) = \underset{a}{\text{argmax}} R_0 + \sum_s P(s'|s, a) V_{\pi}(s')$$

- Instead, we can do policy evaluation on  $q_{\pi}(s, a)$ !

↳ Greedy Policy Improvement over  $q$  is model-free

## PROBLEM #2:

- Because we are doing Q-policy evaluation with

sampling  $\Rightarrow$  if our policy  $\pi$  is deterministic, we won't have estimates for  $Q_{\pi}(s, a) \forall a \neq \pi(s)$ .

$\Rightarrow$  Policy improvement won't select best actions

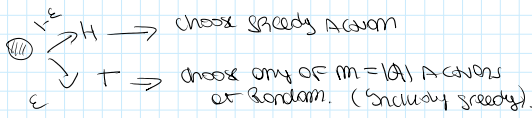
## $\epsilon$ -Greedy Exploration.

- IDEA: Modify Policy Improvement step to make sure we don't get stuck.

$$\pi'(s) = \begin{cases} \underset{a}{\text{argmax}} Q_{\pi}(s, a) & \text{if } a = \underset{a}{\text{argmax}} Q_{\pi}(s, a) \\ \epsilon & \text{otherwise} \end{cases}$$

$$(m-1) \frac{\epsilon}{m} + \frac{\epsilon}{m} + 1 - \epsilon = \epsilon - \frac{\epsilon}{m} + \frac{\epsilon}{m} - \epsilon + 1 = 1 \quad \square$$

Intuition: We flip a coin, H comes up w/ prob  $1-\epsilon$  and  $T$  w/ prob  $\epsilon$ .



## IS $\epsilon$ -Greedy A Policy Improvement?

According to last class, this is enough to

Verify  $Q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s) \forall s \in S$ .

Because  $\Rightarrow V_{\pi}(s) \geq V_{\pi}(s) \forall s \in S$ .

Proof:

$$\begin{aligned} Q_{\pi}(s, \pi'(s)) &= \sum_{a \in A} \pi'(s, a) Q_{\pi}(s, a) = \\ &= \sum_m \sum_{a \in A} q_{\pi}(s, a) + (1-\epsilon) \max_a q_{\pi}(s, a) \leq \\ &= \sum_m \sum_{a \in A} q_{\pi}(s, a) + (1-\epsilon) \sum_{a \in A} \frac{\pi(s, a) - \epsilon/m}{1-\epsilon} q_{\pi}(s, a) \\ &\stackrel{\text{So Because}}{\leq} \sum_{a \in A} \frac{\pi(s, a) - \epsilon/m}{1-\epsilon} q_{\pi}(s, a) \leq \max_a q_{\pi}(s, a) \left[ \sum_{a \in A} \frac{\pi(s, a) - \epsilon/m}{1-\epsilon} \right] \\ &= \max_a q_{\pi}(s, a) \left[ \frac{1-\epsilon}{1-\epsilon} \right] \\ &= \sum_a q_{\pi}(s, a) [\pi(s, a)] = q_{\pi}(s, \pi(s)) \\ &= V_{\pi}(s) \quad \square \end{aligned}$$

## MONTE-CARLO POLICY ITERATION / MONTE-CARLO CONTROL.

POLICY EVALUATION: MONTE-CARLO POLICY EVALUATION FOR  $Q$ .

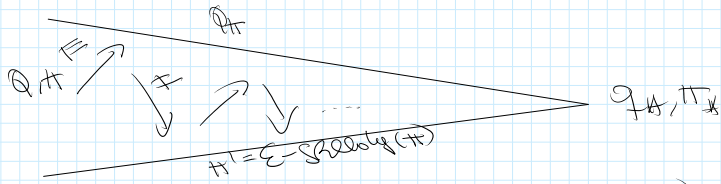
the theorem does not apply in this case because we only have an estimate

## POLICY EVALUATION: MONTE-CARLO POLICY EVALUATION FOR Q.

POLICY IMPROVEMENT:  $\epsilon$ -GREEDY POLICY IMPROVEMENT.

How can we make this more Efficient?

- IMPROVE your POLICY AFTER Every EPISODE.



GLIE (Greedy in Limit with Infinite Exploration)

come up w/ a schedule for exploration. st

- (1) Every state and Action is visited  $\infty$  times
- (2) the policy we obtain in GPE eventually becomes greedy w/ respect to  $Q$ .

$\epsilon$ -greedy is GLIE if  $\epsilon \rightarrow 0$ . ( $\epsilon = \frac{1}{k}$ ).

THEOREM: GLIE MONTE-CARLO CONTROL CONVERGES TO  $q_*$

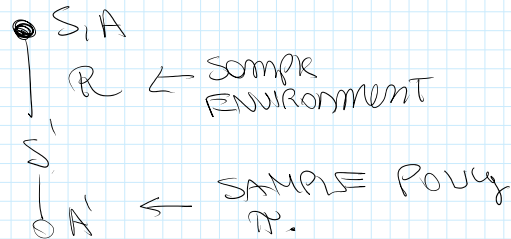
↳ First Algorithm to fully solve on Map  
W/ unknown Dynamics.

## TEMPORAL DIFFERENCE Learning Control.

2 ways  $\rightarrow$  ON-POLY.  
 $\downarrow$  OFF-POLY.

## SARSA / ON-POLICY TO CONTROL.

POLICY EVALUATION step = EVALUATE  $Q^{\pi}(s, a)$  using  
TD learning.



$$Q(S, A) \leftarrow Q(S, A) + \lambda (R_{t+1} + \gamma Q(S', A') - Q(S, A)).$$

POLICY IMPROVEMENT:  $\epsilon$ -GREEDY POLICY IMPROVEMENT.

↳ CAN BE DONE AFTER EVERY STEP.

## OFF POLICY Learning.

- GOAL: EVALUATE  $\pi(s|a)$  to compute  $V_{\pi}$  or  $Q_{\pi}(s,a)$  while following / sampling from  $\mu(s|a)$ .

- N: target Police.
- U: Behavioural Police

Apply in this case

Because we only have an estimate of  $Q_{\pi}$ . And then

/ Assumes  $Q_{it}$  has been computed exactly!  
we might not have monotonic improvement.

## Q-LEARNING (EVALUATE $q_{\pi^*}$ )

- Goal: OFF-POLICY Learning FOR  $q_{\pi^*}(s, a)$ .
- Next Action is chosen using behaviour policy  $A_{\pi} \sim \mu(s)$
- But we consider alternative successor  $A'$  from  $A' \sim \pi(s)$ . sampled from target policy.
- UPDATE  $Q(s, a)$  towards value of  $A'$ .

$$Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma Q(s', A') - Q(s, A)]$$

SPECIAL CASE :  $\mu = \epsilon$ -greedy( $Q$ )  $\leftarrow$  BEHAVIOURAL Policy

$\pi =$  greedy( $Q$ )  $\leftarrow$  TARGET Policy

$$Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma \max_{A'} Q(s', A') - Q(s, A)]$$

- In this case,  $Q$  DIRECTLY APPROXIMATES  $Q^*$  Independent OF the policy we follow  $\rightarrow$  OFF-POLICY (Follow any but target is  $\pi^*$ ).
- For convergence we require  $\mu$  to keep on visiting  $(s, a)$  ( $\epsilon$ -greedy is enough).

NOTE: WE COULD USE A COMPLETELY RANDOM  $\mu$ , AND THE ALGORITHM WOULD STILL CONVERGE TO  $q^*$ .

THEOREM: Q-Learning Converges to  $q^*$ .

QUESTION: IS IT TRUE THIS UPDATE DIRECTLY EVALUATES  $Q^*$  USING TD-LEARNING BECAUSE  $\pi^* = \text{greedy}(Q^*)$ ?

CAN WE THINK OF THE ALGO. AS SAMPLING FROM  $Q^*$  THE BELMAN OPTIMIZATION EQUATION??

$\pi^*$  satisfies  $\pi^* = \text{argmax}_\pi Q_{\pi^*}(s, a)$ . therefore,  $Q(s, a)$  EVALUATION USING TD-LEARNING WOULD LOOK SOMETHING LIKE THIS:

$$\begin{aligned} Q(s, A) &= Q(s, A) + \alpha [R + \gamma Q(s', \pi^*(s')) - Q(s, A)] \\ &= Q(s, A) + \alpha [R + \gamma Q(s', \text{argmax}_{A'} Q_{\pi^*}(s', a')) - Q(s, A)] \\ &= Q(s, A) + \alpha [R + \gamma Q(s', \text{argmax}_{A'} Q(s', a)) - Q(s, A)] \\ &= Q(s, A) + \alpha [R + \gamma Q(s', \text{max}_{A'}(s, a)) - Q(s, A)] \end{aligned}$$

PROBLEM: CAN WE FOLLOW  $\pi^*$ ? depends on  $Q_{\pi^*}$ .

$\hookrightarrow$  has been sampled randomly.

so maybe USE OFF-POLICY. Follow ANY EXPLICITLY policy that visits ALL  $(s, a)$  forever.