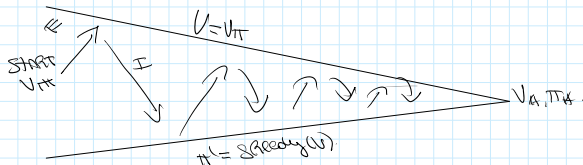


MODEL FREE CONTROL

GOAL: Find V^* / π^* 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_{π} , because acting greedily would require knowledge of the MDP

$$\pi' = \text{greedy}(\pi) = \underset{a}{\text{argmax}} Q_{\pi}(s, a) = \underset{a}{\text{argmax}} \left(R(s) + \sum_{s'} P(s'|s, a) V_{\pi}(s') \right)$$

- 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 π 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

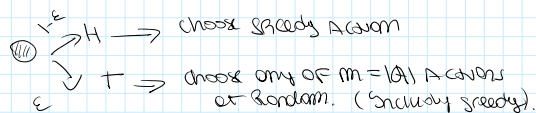
ϵ -Greedy EXPLORATION.

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

$$\pi'(a|s) = \begin{cases} \epsilon_m + 1 - \epsilon & \text{if } a = \underset{a}{\text{argmax}} Q_{\pi}(s, a) \\ \epsilon_m & \text{otherwise} \end{cases}$$

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

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



IS ϵ -GREEDY A POLICY IMPROVEMENT?

According to last class thm, it's 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:

$$Q_{\pi}(s, \pi'(s)) = \sum_{a \in A} \pi'(a|s) q_{\pi}(s, a) =$$

$$\epsilon \sum_{a \in A} q_{\pi}(s, a) + (1 - \epsilon) \max_a q_{\pi}(s, a) \leq$$

$$\epsilon \sum_{a \in A} q_{\pi}(s, a) + (1 - \epsilon) \sum_{a \in A} \frac{\pi(a|s) - \epsilon}{1 - \epsilon} q_{\pi}(s, a)$$

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

$$= \sum_a q_{\pi}(s, a) [r(a|s)] = q_{\pi}(s, \pi(s)) = V_{\pi}(s)$$

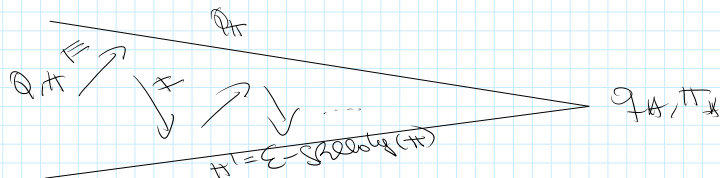
Monte-Carlo Policy Iteration / Monte-Carlo Control

POLICY EVALUATION: Monte-Carlo Policy Evaluation for Q .

POLICY IMPROVEMENT: ϵ -greedy policy improvement.

How can we make this more efficient?

- IMPROVE YOUR POLICY AFTER EVERY EPISODE.



GLIE (Greedy in Limit w/ Infinite Exploration)

Come up w/ a schedule for exploration. s.t

- (1) Every state and action is visited ∞ times
- (2) the policy we obtain in GLIE eventually becomes greedy w/ respect to Q .

ϵ -greedy is GLIE if $\epsilon \rightarrow 0$. ($\epsilon = 1/k$).

THEOREM: GLIE Monte-Carlo control converges to Q^*

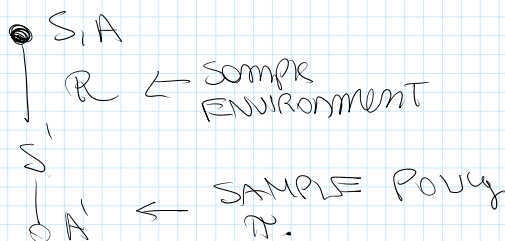
- FIRST ALGORITHM TO FULLY SOLVE ON MDP w/ unknown dynamics.

TEMPORAL DIFFERENCE Learning Control

2 ways \rightarrow ON-POLICY -
 \downarrow OFF-POLICY.

SARSA / ON-POLICY TO CONTROL

POLICY EVALUATION STEP: EVALUATE $Q^{\pi}(s, a)$ USING TO LEARNING.



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

POLICY IMPROVEMENT: ϵ -greedy policy improvement.

\rightarrow CAN BE DONE AFTER EVERY STEP.

the theorem does not apply in this case
 Because we only have an estimate of Q^{π} . and then
 Assumes Q^{π} has been computed exactly!
 we might not have monotonic improvement.

OFF POLICY Learning.

- GOAL: EVALUATE $\pi(s|a)$ to compute V_π or $q_\pi(s, a)$ while following / sampling from $\mu(s|a)$.

- π : target policy. • μ : Behavioural policy.

Q-LEARNING (EVALUATE q_π).

- Goal: OFF POLICY Learning for $q_\pi(s, a)$.
- Next Action is chosen using behaviour policy $A_{t+1} \sim \mu(s)$.
- 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

$\pi = \text{greedy}(Q) \leftarrow$ 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 μ but target is π^*).
- For convergence we require μ to keep on visiting (s, a) (ϵ -greedy is enough).

NOTE: WE COULD USE A completely random μ , And the Algorithm would still converge to q^* .

THEOREM: Q-Learning Converges to q^* .

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

CAN WE THINK OF THE ALGO. AS SAMPLING FROM Q_π THE BELMAN OPTIMIVITY EQUATION??

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

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

Problem: CAN WE FOLLOW π^* ? depends on Q_π .

↳ has been Unobserved Randomly.

so maybe use OFF-policy. Follow Any Exploratory Policy that visits ALL (S, a) forever.

