

Lecture 5

- on policy learning → "Learn on the job"
You learn and ~~inf~~ about the policy which you are following. In simple words you ~~are~~ are doing the task in one specific way and you are learning about your method while doing it.
You sample from a policy, π , and at the same time you are evaluating that same policy.
- off Policy → "Look over someone's shoulder".
There are two policies, the one ~~is~~ you are using to ~~is~~ take actions 'u' and the other policy that you want to evaluate say ' π' , learning about ' π ' while drawing samples from 'u' is off policy learning
- Monte Carlo control → we can't use the previous monte carlo control method, i.e., first evaluating the $V(s)$ and then act greedily. Because to act greedily, we must know the $\alpha^a(s)$ for the state. To be able to estimate the α value for an action in a state, we need the state transition probability upon taking the action.

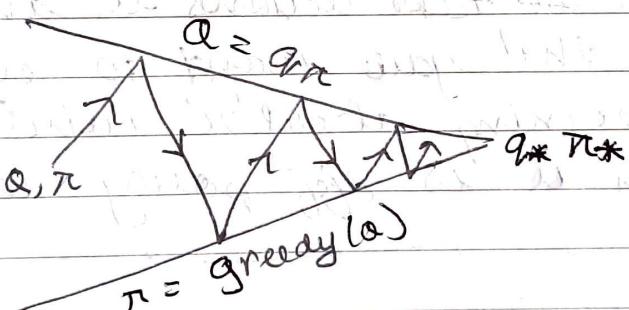
$$\alpha^a(s) = R_s^a + \sum_{s' \in S} P_{ss'}^a V(s')$$

As ~~the~~ we are studying model free control, we don't have P values and thus this method will break.

Instead we follow greedy policy improvement over $\alpha(s,a)$.

$$\pi'(s) = \underset{a \in A}{\operatorname{argmax}} \alpha(s,a)$$

~~for~~ π is policy evaluation, instead of generating ~~generating~~ V^π values for all states we generate $\alpha^\pi(s,a)$ values for all state action pair and directly follow greedy policy over all $\alpha(s,a)$ values for a state.



Policy evaluation - Monte carlo policy evaluation,
 $\alpha = q_\pi$

Policy improvement \rightarrow we can't take greedy policy. Because if we keep taking greedy right from the beginning, ~~the~~ there may be some state action pairs which may be unexplored.

Instead we use ϵ -greedy Exploration, with probability ϵ we choose randomly and with $1-\epsilon$ we choose greedily.

$$\pi(a|s) = \begin{cases} \epsilon / m + (1-\epsilon) \text{ if } a^* = \arg \max_{a \in A} Q(s, a) \\ \epsilon / m \quad \text{otherwise} \end{cases}$$

→ Another method for Exploration - Exploitation
 GLIE → Greedy in the limit with infinite exploration
 Two things must satisfy
 ↳ All state-action pairs are explored infinity many times

$$\lim_{k \rightarrow \infty} N_k(s, a) = \infty$$

It means that there is no state-action pair left which is not sufficiently explored. This is that is we must know the a^* for every state-action pair which is the expected return generated by taking that particular action in that state. This is possible by taking ∞ samples of the return generated by a^* and doing an empirical mean.



The policy eventually becomes greedy.

$$\lim_{K \rightarrow \infty} \pi_K(a|s) = 1 \quad (a \in \arg\max_{a \in A} Q_K(s, a))$$

We want to ~~eg~~ eventually act greedily. We eventually always want to take the best actions with more Q value and not act ~~so~~ stochastically.

→ GLIE Monte Carlo (continued)

Sample K-th episode using π^*

$$\{s_1, a_1, r_1, \dots, s_T\}$$

For each state and action -
we upto-update the $Q(s, a)$ using
incremental average
we calculate the number of times
~~we visited~~ visited the state and
took that action

$$N(s_t, a_t) \leftarrow N(s_t, a_t) + 1$$

and update $Q(s_t, a_t)$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{1}{N(s_t, a_t)} (G_t - Q(s_t, a_t))$$

Then implement policy improvement

$$\epsilon \leftarrow 1/k$$

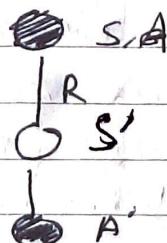
$$\pi \leftarrow \epsilon\text{-greedy } (\alpha)$$

as k increase ϵ decreases, we eventually get a greedy policy in the end as

$$\epsilon \rightarrow 0$$

~~Maximizing Q function~~
 Updating action value function with Sarsa

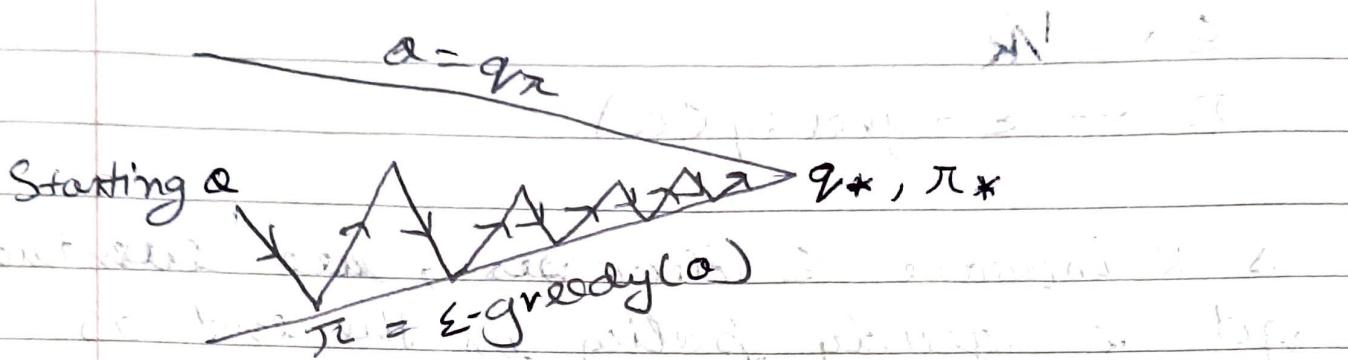
We start with some state and pick an action ~~according to ~~greedy~~~~ policy. where $\alpha(s,a)$ we need to evaluate. We get one reward and end up in a state. We then sample π and pick an action $\pi(A)$ and uses it $\alpha(s',a)$ to get the $\alpha(s,a)$



$$\alpha(s,a) \leftarrow \alpha(s,a) + \alpha(R + \gamma \alpha(s',a') - \alpha(s,a))$$

SARSA update

→ on policy control with SARSA.



we update our policy ~~in~~ after every time step.

We sample from our policy and takes an action, updates the policy $\pi(s, a)$, then we ~~will~~ update our policy π w.r.t. the new Q we have and then take another action w.r.t. the new policy π , correct our policy again and this continues until the end of episode (if it is terminating).

Pseudo code →



- Initialize $\alpha(s,a)$ & $s \in S, a \in A(s)$ arbitrarily and $\alpha(s, \text{terminal state}) = 0$
- Repeat (for each episode):
 - Initialize $s \in S$
 - Choose A from s using policy derived from α
 - Repeat (for each step of episode):
 - take action A , observe R, s'
 - Choose A' from s' using policy derived from α
 - $\alpha(s, A) \leftarrow \alpha(s, A) + \alpha(R + \gamma \alpha(s', A') - \alpha(s, A))$
 - $s \leftarrow s'; A \leftarrow A'$
 - until s is terminal
- Convergence of Sarsa.

Sarsa converges to the optimal action value function under the following conditions

- ↳ Gittin sequence of policy $\pi_t(a|s)$
- ↳ Robbins-Monro sequence of step size α_t

$$\sum_{t=1}^{\infty} \alpha_t = \infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

→ n-step Sarsa →

$$n=1 \quad q_t^{(1)} = R_{t+1} + \gamma Q(s_{t+1})$$

$$n=2 \quad q_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 Q(s_{t+2})$$

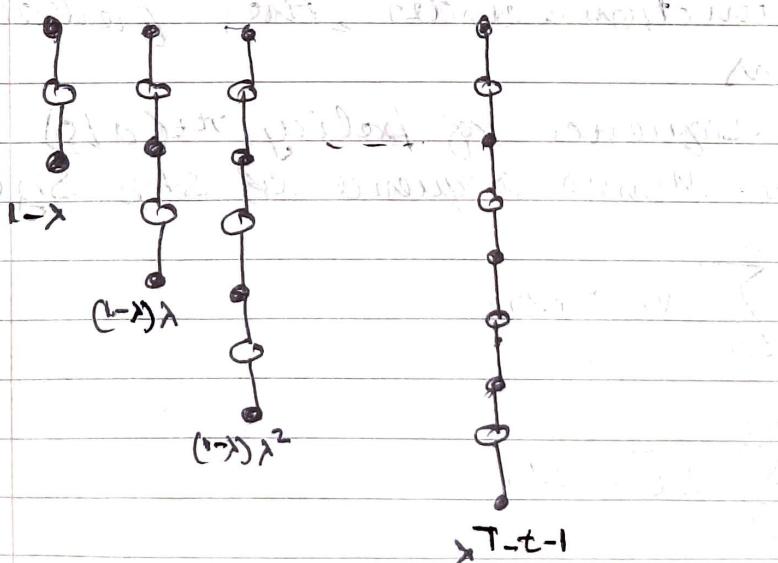
$$n=\infty \quad q_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

$$(mc) \quad q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(s_{t+n})$$

$$q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(s_{t+n})$$

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \alpha [q_t^{(n)} - Q(s_t, A_t)]$$

→ Forward view Sarsa (λ)



$$q^\lambda := (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} q_t^{(n)}$$



$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (q^* - Q(s_t, a_t))$$

→ Backward view SARSA (λ)

$$E(s, a) = 0$$

$$E_t(s, a) = \gamma \lambda E_{t-1}(s, a) + I(s_t = s, a_t = a)$$

This is done for all action-value pair

Similar to model free prediction.

$Q(s, a)$ is updated for every state s and action a .

We update in proportion to the eligibility trace

$$\delta_t = R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha \delta_t E_t(s, a)$$

Pseudo code →

- Initialize $\hat{Q}(s,a)$ arbitrarily for all $s \in S$ and $a \in A(s)$.
- Repeat (for each ~~stack~~ episode):
 - $E(s,a) = 0 \forall s \in S, a \in A$
 - Initialize s, A
 - Repeat for each step in episode
 - Take action A , observe R, s'
 - Choose A' from s' using policy
 - $s \leftarrow R + \gamma \hat{Q}(s', A') - \hat{Q}(s, A)$
 - $E(s, A) \leftarrow E(s, A) + 1$
 - for all $s \in S, a \in A(s)$
 - $\hat{Q}(s, a) \leftarrow \hat{Q}(s, a) + \frac{1}{E(s, a)} \hat{Q}(s, a)$
 - $E(s, a) \leftarrow \gamma \lambda E(s, a)$
 - until s is terminal

- Off-policy learning \rightarrow

Evaluate some other policy while following one policy.

That is draw samples and rewards from one policy and use that to evaluate another policy.

→ Evaluate target policy $\pi(a|s)$ to compute $V^\pi(s)$ or $q^\pi(s, a)$ while following behaviour policy ~~$\pi(a|s)$~~ $\mu(a|s)$

→ Importance Sampling →

$$\mathbb{E}_{x \sim \mu} [f(x)] = \sum p(x) f(x)$$

$$= \sum \alpha(x) \frac{p(x)}{\alpha(x)} f(x)$$

~~$\mathbb{E}_{\alpha} [p(x) f(x)]$~~

→ Monte Carlo learning doesn't work for off-policy learning. ~~(Explained later)~~

~~→ Off-policy TD learning.~~

Given a starting state s_t , the probability of the subsequent state-action trajectory; $A_t, s_{t+1}, A_{t+1}, \dots, s_T$ occurring under policy π is

$$\prod_{k=t}^{T-1} \pi(A_k | s_k) p(s_{k+1} | s_k, A_k)$$

p is the state transition probability and π is the policy.

Thus the relative probability of the trajectory under the target and behaviour policy

$$P_t^T = \frac{\prod_{k=t}^{T-1} \pi(A_k | s_k) p(s_{k+1} | s_k, A_k)}{\prod_{k=t}^{T-1} \mu(A_k | s_k) p(s_{k+1} | s_k, A_k)}$$

$$\Rightarrow \frac{\prod_{k=t}^{T-1} \pi(A_k | s_k)}{\prod_{k=t}^{T-1} \mu(A_k | s_k)} \text{ which only depends}$$

upon μ and π and not on rest of the MDP.

Let t denote timestep which doesn't end or terminate after the end of the episode but rather continues to increase in the next episode.

Let $T(s)$ denote ~~all the~~ the set of all the time steps t in which state s is visited. This can be first visit or every visit.

Let $T(t)$ denote the first time of termination following time t and b_t denote

The return after t up through $\bar{T}(t)$. Then $\{g_t\}_{t \in \mathcal{T}(s)}$ are the returns that pertain to state s , and $\{p_t^{\bar{T}(t)}\}_{t \in \mathcal{T}(s)}$ are corresponding importance sampling ratios. To estimate $v_{n(s)}$, we simply scale the returns by ratios and average the returns.

Averaging can be

- 1) Simple average by dividing by the number of times the state is visited.

$$v(s) = \frac{\sum_{t \in \mathcal{T}(s)} p_t^{\bar{T}(t)} g_t}{|\mathcal{T}(s)|}$$

- 2) weighted average

$$v(s) = \frac{\sum_{t \in \mathcal{T}(s)} p_t^{\bar{T}(t)} g_t}{\sum_{t \in \mathcal{T}(s)} p_t^{\bar{T}(t)}}$$

or zero if denominator is 0.

- 1) has high variance and 2) has high bias.

→ Incremental updates → for ordinary mean, the incremental update would follow the

~~Mean~~ $\leftarrow m$

$$v(s) \leftarrow v(s) + \frac{1}{N(s)} (\cancel{g_t} - v(s))$$

However for weighted average things would

differ:

Suppose we have a ~~the~~ series of returns $c_1, c_2, \dots, c_{n-1}, c_n$, all starting in the same state and each with a corresponding random weight w_i . We need

$$V_{CS}^* = \frac{\sum_{k=1}^{n-1} w_k c_k}{\sum_{k=1}^n w_k}; n \geq 2$$

we can keep it up-to-date as we generate more ~~c_n~~, c_n .

to For this we must maintain the cumulative sum c_n of the weights given to the first n returns. Therefore the update rule becomes

$$V_{n+1} = V_{n+1} w_n [c_n - v_n]; n \geq 1$$

$$\text{and } c_n \leftarrow c_n + w_{n+1}$$

\rightarrow is not memory if one uses



→ off Policy T-D learning

It is in simple terms defined as

$$\alpha(s_t, A_t) \leftarrow \alpha(s_t, A_t) +$$

$$\alpha [R_{t+1} + \gamma \max_{a \in A} \alpha(s_{t+1}, a)]$$

$$- \alpha(s_t, A_t)]$$

In this case the learned action value function, α , directly approximates q^* , the ~~per~~ optimal action value function independent of the policy being followed.

However if we want to learn a policy π (target policy) while following the behaviour policy μ we take one step according to our behaviour policy μ and record the reward. Then we update our $V(s_t)$ as

$$V(s_t) \leftarrow V(s_t) + \alpha \left(\frac{\pi(A_t | s_t)}{\mu(A_t | s_t)} (R_{t+1} - V(s_{t+1})) \right)$$

$V(s_t)$ is the expected value of all the returns that we can generate. The expectation depends upon the probability of getting a certain a_t . Higher the p_a , a higher would be the weightage of a_t in the $\pi(a_t)$.

This probability is determined by two distributions, the policy and the transition probability.

So to know what would be the probability of getting some return or when followed a target policy while generating or samples from behavior policy, we use importance sampling.

We make the Gt obtained for that particular action and acknowledging the fact that the contribution of that G_t^{π} would not be that same as $G_t^{\pi'}$ or in other words, if we sample the episodes from that state many times we will get ~~a that pert~~ to ~~the~~ one particular state as different number of times while following the behavior policy than following target policy, hence ~~the~~ to get the $V(s)$ we need multiply the ~~target~~ $R_{t+1} + \gamma V(s_{t+1})$ by the ratio of ~~in a row~~

Probability of it occurring in R_{t+1}
Probability of it occurring in u



$$\pi(A_t | s_t) \cdot P(s_{t+1} | s_t, A_t) \\ \pi(A_t | s_t) \cdot P(s_{t+1} | s_t, A_t)$$

$$\Rightarrow \frac{\pi(A_t | s_t)}{\pi(A_t | s_t)}$$

this would deliver draws or concentrate all the weightage of one particular act according to its relative probability of occurring while following the target policy

∴ the T-D equation becomes (SARSA)

$$V(s_t) \leftarrow V(s_t) + \alpha \left(\frac{\pi(A_t | s_t)}{\pi(A_t | s_t)} (R_t + \gamma V(s_{t+1})) - V(s_t) \right)$$

This has much lower variance than multi carlo off policy learning

→ Q-learning → with Q-values

- ~~Q-Learning~~
- Q-Learning →

We follow a ~~normal~~ normal process
 SARSA \Rightarrow update policy.

We take first action using our behaviour $\neq \pi$ and get a reward while calculating the $q(s_{t+1}, a)$.
 We use target policy.
 \rightarrow Let $Q(s_t, a^*)$ denote Q-value at state s and following target policy π , and $q(s_t, a)$ be the action value for state s and following target π .

\Rightarrow (Eqn -

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (R_{t+1} + \gamma q(s_{t+1}, a')) - Q(s_t, a_t)$$

Think of it, in simple terms we are using π as just to generate episodes of data and ultimately updating q using those \neq episodes.

→ π_θ policy control with α learning

In this case the target policy is greedy policy w.r.t $\alpha(s,a)$

$$\hat{\pi}(s_{t+1}) = \underset{a'}{\operatorname{argmax}} \alpha(s_{t+1}, a')$$

$$\pi(s_{t+1}) = \underset{a}{\operatorname{argmax}} \alpha(s_{t+1}, a)$$

The behavior and the target policy can improve.

The behaviour policy will be ϵ -greedy w.r.t $\alpha(s,a)$ so then we have sufficient explanation.

$$\alpha(s, a_t) \leftarrow \alpha(s, a_t) + \alpha(R_{t+1} + \gamma \alpha(s_{t+1}, a'))$$

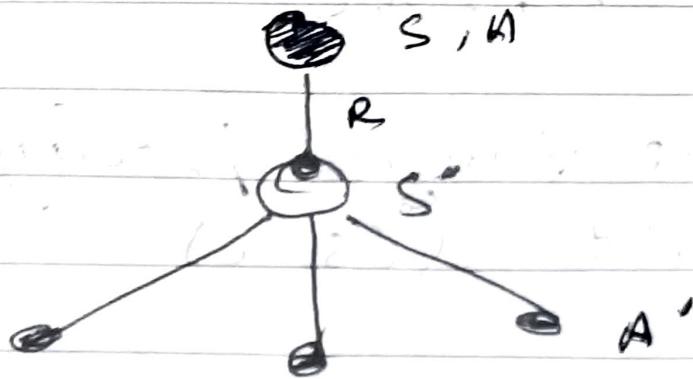
$$R_{t+1} + \gamma \alpha(s_{t+1}, a) \rightarrow \text{target}$$

$$\Rightarrow R_{t+1} + \gamma \alpha(s_{t+1}, \underset{a''}{\operatorname{argmax}} \alpha(s_{t+1}, a'))$$

$$\Rightarrow R_{t+1} + \gamma \alpha$$

$$\Rightarrow R_{t+1} + \underset{a''}{\operatorname{max}} \gamma \alpha(s_{t+1}, a')$$

Basically go ϵ -greedy while sampling and completely greedy while bootstraping or calculating the target.



$$Q(s, a) \leftarrow Q(s, a) + \alpha (R + \gamma \max_{a'} Q(s', a'))$$

α'

$\max_{a'} Q(s', a')$

It converges to optimal action value function.