Model free control.

On-policy → learning on the job, observing self behaviour
off-policy → learning following someone else's behaviour.

~~Model~~
→ optimize the value function of an unknown MDP.

why?
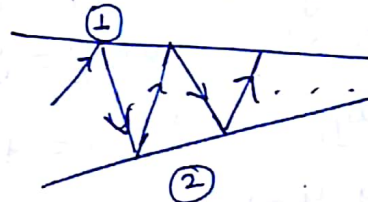→ most common problems are either.

i) model/MDP unknown, but experience can be sampled
ii) MDP known, but it is too expensive to compute even 1 step using the full model, except by sampling

On-policy learning → learning from job → while following policy π, learn that policy π.

Off-policy learning → observe someone do something and sampling trajectories from human's behavior (eg.), learn self behaviour (eg. robot).

In DP, alternating betn policy evaluatn & iteration converges to the optimal policy

we only used iterative method, we'll try to change ① & ②.



Eg. Monte Carlo control:

1-way → do monte carlo policy evaluatn & then do greedy policy improvement.

Problems:

① computing V(s) requires a model, but we want model free

$$V(s) =• \max_a \left( R_s^a + P_{ss'}^a V(s') \right)$$

But in model free, we want to learn a policy without any model. The $q$-function defined earlier → means that final score we would get at the end of a game (eg.) if we took action $a$, when in state $s$. → $q(s,a)$

So, no model dependence now. Just a state and possibly actions that can be taken.

→ greedy over $q(s,a)$ → action value function

$\pi'(s) = \text{argmax}_a \ q(s,a)$

For each $(s,a)$ pair taken m across all experience, to compute $q(s,a)$

- Proposal changes:
  → MC evaluate $q = q_\pi$
  → Greedy policy improvement!

2nd problem → choosing greedily → we never visit some states or take some actions which means we don't evaluate them correctly & hence, not selecting them.

eg. to illustrate that point. → suppose 2 states → $\hat{a}_1$ & $\hat{a}_2$
upon choosing both once, we found:
$$V(s_1) = -1, \quad V(s_2) = +1 \quad \text{we choose } s_2 \text{ greedily}$$
now say, $\quad V(s_2) = +3$ → mean $= 2$, again chooses $s_2$
and so on....

This shows that we haven't even considered state $s_1$ only after 1 seeing and hence, don't know what value it might have as we keep choosing $s_2$ infinitely.

$\varepsilon$- greedy exploration ↙ solution to ensure continual exploration

→ prob-$\varepsilon$ → choose random action
→ $1-\varepsilon$ → choose the greedy option.

Thus:
$$\pi(a/s) = \begin{cases} \varepsilon/m + (1-\varepsilon) & \rightarrow \text{ for greedy action} \\ \varepsilon/m & \rightarrow \text{ for any other action} \end{cases}$$

→ keep exploring

$\varepsilon$ greedy → policy improvement (ensures)
$$V_{\pi'}(s) \geqslant V_{\pi}(s).$$

$\pi, \pi' \rightarrow$ both $\varepsilon$-greedy policies

$$q_{\pi}(s, \pi'(s)) = \sum_a \pi'(a/s) \, q_{\pi}(s, a)$$

$$= (1-\varepsilon) \max_a q_{\pi}(s, a) + \frac{\varepsilon}{m} \sum_a q(s, a)$$

Now, $\max_a$
greater than
any weighted
$$\geq (1-\varepsilon) \sum_a \frac{\pi(s/a) - (\varepsilon/m)}{1-\varepsilon} \, q_{\pi}(s, a) + \frac{\varepsilon}{m} \sum_a q(s, a)$$
sum

$$= \sum_a \pi(s/a) \, q_{\pi}(s, a) = V_{\pi}(s)$$

$$\therefore V_{\pi'}(s) \geqslant V_{\pi}(s). \text{ (from earlier similar derivation)}$$

∴ Finally:      (MCPE)
① MC policy evaluate  $Q = q_{\pi}$   ── softening g the
② $\varepsilon$-greedy policy improvement ↙ ──  greedy policy

Idea → no need to completely ~~utili~~ evaluate the policy, when
we can get a better policy with only a few initial
steps (Similar to DP lecture).

.One extreme → do update every episode, basically on the
most fresh ~~da~~ information/ estimate g the value function.

after
every     ① MC PE → $Q = q_{\pi}$          ↗ which were visited in the
episode                                            episode.
          ② $\varepsilon$-greedy → For those states

How to ensure we get $\pi^*$?                    as when at $\pi^*$, we
→ Trade off b/w exploration & exploitation.      won't have
                                                  much
                                                  exploration.

Idea for balancing the 2 ideas → GLIE.

→ come up with schedule such that 2 conditions are met.

① all state-action pairs explored. infinitely many times

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

② converges to the optimal greedy policy

$$\lim_{k \to \infty} \pi_k(a|s) = I\left(a^* = \arg\max_{a' \in A} q_k(s, a')\right)$$

One idea → schedule / decay $\varepsilon$ for $\varepsilon$-greedy.

→ hyperbolic schedule → $\varepsilon_k = 1/k$.

## GLIE MC control

→ sample $k^{th}$ episode following $\pi$

$$\downarrow \{s_1, A_1, R_1 \ldots R_T\}$$

→ update mean qvalue & $N(s_t, A_t)$

$$N(s_t, A_t) = N(s_t, A_t) + 1$$

$$q(s_t, A_t) = q(s_t, A_t) + \frac{1}{N(s_t, A_t)}\left(G_T - q(s_T, A_t)\right)$$

↳ this mean is not the actual statistical mean as we are accumulating values as we improve the policy

↓
changing over time

→ Improve policy

$$\varepsilon \to 1/k$$
$$\pi \to \varepsilon\text{-greedy}(q)$$

GLIE makes sure that these collected statistics converge to the actual ~~mean~~ over time.
value

$P(s, a) \rightarrow q^*(s, a)$

Iterate over the entire process.

→ considerably more efficient<del>ly</del> that updating after a batch

Initialization → In this case, for the 1st time we observe a state $N(s^t, a^t) = 1$, and $Q(s^t, a^t) = G_T$, so init doesn't matter. But for weighing other than $\frac{1}{N(s^t, a^t)}$, it will affect much more.
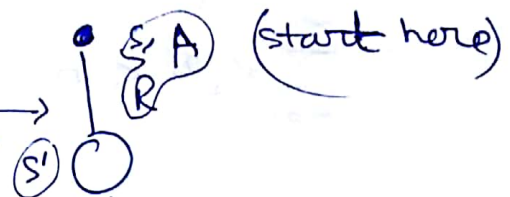
Now, we'll show TD.

TD vs MC

⊕ Variance low   → online   → incomplete sequences.

- look at $P(s, a)$ for model-free
- slot in TD learning for policy evaluation
- ε-greedy

} can update after 1 step., no need to wait till episode end.

So, update policy after each step.
Called SARSA.

$P(s, A) \Leftarrow P(s, A) + \alpha \left( R + \gamma Q(s', A') - P(s, A) \right)$.

sarsa update.
In every time step, update value function.

→ started with s. took action A (from current policy) and reached s'. → update only $P(s, A)$ → only for this pair

sample from environment



(S, A) (start here)
R
(S')
(A')
→ SARSA

sample from policy

SARSA algo for on-policy

① Initialize $Q(s, a) \forall s, \forall a$. $Q(terminal, \cdot) = 0$.

② Repeat for each episode:

    Initialize s      using policy derived

    Choose A → sampled from $Q$.

    Repeat (for each step):

        Take action A, got reward R, landed in $S'$.

        choose A' from $S'$, using policy derived from $Q$

$$Q(S, A) \leftarrow Q(S, A) + \alpha ( R + \gamma Q(S', A') - Q(S, A))$$

        $S \leftarrow S', A' \leftarrow A$.

Sarsa → on-policy algorithm.

↳ converges to the optimal policy, just need a GLIE ①

policy → eg. $\epsilon_k = \frac{1}{k}$

② $\alpha_t$ step sizes follows.

$$\sum_{t=1}^{\infty} \alpha_{t'} = \infty \quad \text{(sufficiently large to move } Q \text{ values)}$$

$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty \quad \left(\text{step size becomes small)} \atop \text{changes to } \gamma \quad \text{or else} \atop Q \text{ values} \quad \text{noise}\right)$$

In practise, sometimes don't depend on both ① & ②. ③

Again apply similar concept of TD(λ) → n-step SARSA.

to get the best of both worlds.

    $n = 1$ . . . (SARSA) ——

    ⋮

    $n = \infty$      (MC) . .

n-step $Q$ return:

$$Q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \cdots \gamma^n Q(S_{t+n})$$

n-step Q return:

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

Do this for all. n.c.
average over "n";

Sarsa $(\lambda) \rightarrow Q(S_t, A_t) = Q(S_t, A_t) + \alpha\,(q_t^{\lambda} - Q(S_t, A_t))$

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

↓
Forward view.

to make algo
online

Backward view → using Eligibility Traces ↑

$E_0(s, a) = 0$

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

$Q(s, a) \rightarrow$ updated for every $(s, a)$ pair

$S_t = R_{t+1} + \gamma\, Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)$

$Q(s, a) = Q(s, a) + \alpha\, E(s, a)\, S_t.$

• Init $Q(s, a) = $ arbitrarily
• Repeat - (episode)
   $E(s, a) = 0 \quad \forall s, a.$
   Init $S, A.$
    Repeat-(step):
      Take action $A$, observe $R, S'$
      Take $A'$ → from. $Q.$
      $S_t = R_t + \gamma\,(Q(S', A')) - Q(S, A).$
      $E(S, A) = E(S, A) + 1$
      $\forall a, s$
        $Q(s, a) = Q(s, a) + \alpha\, S_t\, E_t(s, a)$
        $E(s, a) = \lambda \gamma\, E(s, a)$
      $A \leftarrow A', S \leftarrow S'.$

until termination.

SARSA $(\lambda)$

↳ faster flow of info. back through time

SARSA(0)
↓
needs many steps to flow the info. back.

# Off-policy learning

→ evaluate target policy → to compute $V_\pi / q_\pi$ following
   behaviour policy → $\mu(a/s)$

Why?

→ Learn from obsv.g other agents ( eg. human) ⊥look
   at traces g behaviour, not just supervised learning)

→ Re-use experience, from old policies $\pi_1, \pi_2, \ldots \pi_{t-1}$
   generated. → using off-policy learning
                           can do that.

→ Learn about optimal policy
   While following exploratory policy

   We want policy to be deterministic but also want to
   explore. optimal off policy learning allows that as
   we can learn a deterministic optimal policy while
   following an exploratory policy.

→ Learn about multiple policies while following one policy.

Importance Sampling → estimate expectation g a different dist⁴

$$E_{x \sim p}\left[f(x)\right] = \sum_{x} P(x)\, F(x) = \sum_{x} Q(x)\, \frac{P(x)}{Q(x)}\, F(x)$$

eg. reward
    funct⁴       $= E_{x \sim Q}\left[\frac{P(x)}{Q(x)} F(x)\right]$

## use IS for off policy MC.

→ Sample returns from $\mu$ < $G_t$ to evaluate $\pi$.

→ weigh $G_T$ as per similarity betⁿ policies

→ multiply importance sampling correction

$$G_t^{\pi/\mu} = \frac{\pi(A_t/s_t)}{\mu(A_t/s_t)} \cdot \frac{\pi(A_{t+1}/s_{t+1})}{\mu(A_{t+1}/s_{t+1})} \cdots \frac{\pi(A_T/s_T)}{\mu(A_T/s_T)} G_t$$

→ Update value towards corrected return.
$$N(S_t) = V(S_t) + \alpha (G_t^{\pi/\mu} - V(S_t))$$

→ Very high variance → as multiplying so many ratios diminishes value.

→ MC → very bad off-policy

Thus, only TD learning for off-policy

Now, IS only after/upto 1 step.
$$N(S_t) = V(S_t) + \alpha \left( \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} (R_{t+1} + \gamma V(S_{t+1})) - V(S_t) \right)$$

→ policies need to be similar only over a single step.

→ much lower variance than MC (and can still bias off).

§-learning → Best with off-policy

→ make use of Q values.
→ no IS reqd.
→ next action from $\pi$  $A_t \sim \mu(\cdot|S_t)$ ← what actually took
→ also an alternative successor action $\pi$
(what we could have taken following target-policy in future   $A' \sim \pi(\cdot|S_t)$

→ update $Q(S_t, A_t)$ towards value of successive action
$$Q(S_t, A_t) \Leftarrow Q(S_t, A_t) + \alpha ( R_{t+1} + \gamma Q(S_{t+1}, A') - Q(S_t, A_t))$$

main
→ allow both behaviour & target-policies to improve
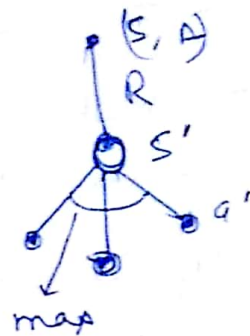→ learn greedy policy from exploratory policy.

$\pi$ → greedy wrt. $Q(s, a)$
$$\pi(S_{t+1}) = \arg\max_{a'} Q(S_{t+1}, a').$$

- $\mu \to \epsilon$-greedy w.r.t $Q(s,a)$.     Q-learning control
  or SARSA

Q-learning target:
$$R_{t+1} + Q(S_{t+1}, A')\gamma$$
$$= R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a')$$



map

$$Q(S,A) \leftarrow Q(S,A) + \alpha \left( R + \gamma \max_{a'} Q(s', a') - Q(s,a) \right)$$

$\to$ converges to $q^*(s,a)$