11 Apr 25

## Simplest "TD" method



$$V(S_t) \leftarrow V(S_t) + \alpha \left[ r_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right]$$

⎰ move a small step $\alpha$ in the
⎱ direction of error

estimate of $V(S_t)$ but at time $t+1$

$V(S_t)$ at time $t$ — my prediction as I land at $S_t$

TD error $\delta$

— TD like MC - do not require complete env only experience (sampling)

- can be fully incremental (bootstrapping)
- can learn even w/o the final outcome.

⟶  TD prediction

- $\pi \xrightarrow{\text{get}} v_\pi \text{ or } q_\pi$
- no knowledge of $p$ and $r$ but access to real system / sample model.

## TD(0)

$$V_{k+1}(S_t) \leftarrow V_k(S_t) + \alpha \left[ r_{t+1} + \gamma V_k(S_{t+1}) - V_k(S_t) \right]$$

take $v_k$ steps of updates

get sample $V_k \to$ best estimate of $v_\pi$ acc to policy $\pi$
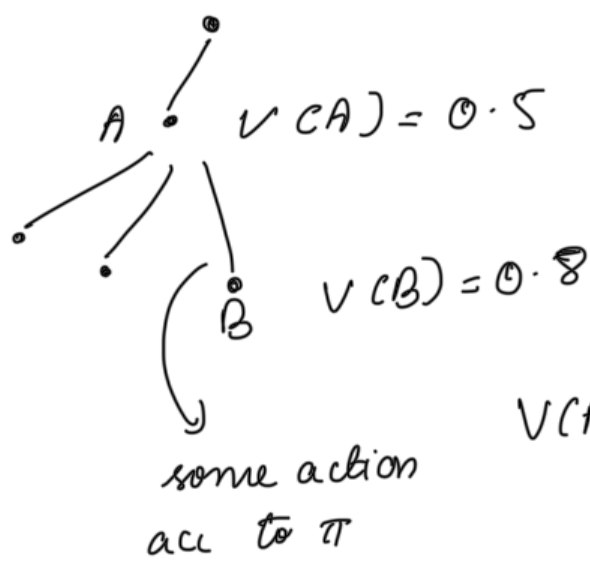
$$V_\pi = E_\pi \left\{ R_{t+1} + \gamma V_\pi(S_{t+1}) \right\}$$

Transition acc to MDP, act acc to best policy.

MC / Exhaustive Search: Go all the way to the end.
TD & DP : Stop after a step (bootstrapping)

eg. TD Update Example



r A → B : 0

$\alpha$ : 0.2     Step size

$\gamma = 0.9$     Discount factor

$$V(A) = V(A) + \alpha [R + \gamma V_B - V_A]$$
$$= 0.5 + 0.2 [0.9(0.8) - 0.5]$$
$$= 0.544$$

⟶ MC vs TD Updates

Same trajectory, value function by MC vs TD

Batch TD : Take a batch / set of episodes
loop TD
till convergence

fixed policy.

Consider :    A · B

AOBO (terminate)
B1
BO
A1 BO

B1
BO

MC
$$V(A) = \frac{0 + 1}{2} = 0.5$$     look at AOBO
                                          A1 BO

$$V(B) = \frac{0 + 1 + 0 + 0 + 1 + 0}{6} = \frac{2}{6} = \frac{1}{3} = 0.33$$
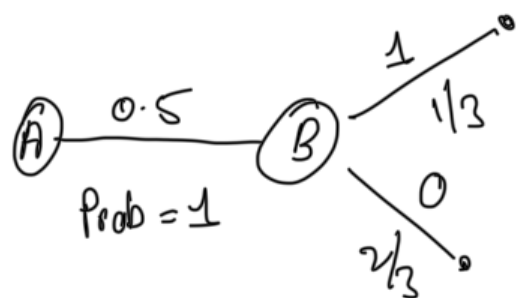
$\gamma = 1$

TD
$$V(A) = E[R + \gamma V(B)] = 1 + \gamma V(B) = 1 + \gamma \left(\frac{1}{3}\right)$$

every time
I go to B

$$= 0.833$$

$$V(B) = E[R] = \frac{1}{3}$$

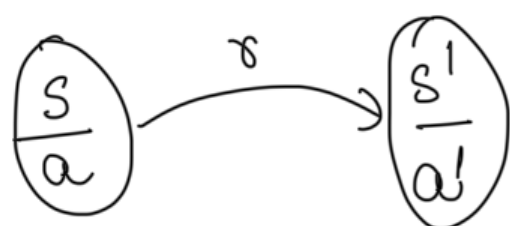"implicitly" forming an MDP even though I am only given samples of an MDP.



which is correct? happen because finite data.
Lot of data → both converge to same.

→ MC — converges to min least squares estimate of return.
TD — certainty equivalence estimate.



But Q learning we don't use next action

→ SARSA

Look at TD evaluation now
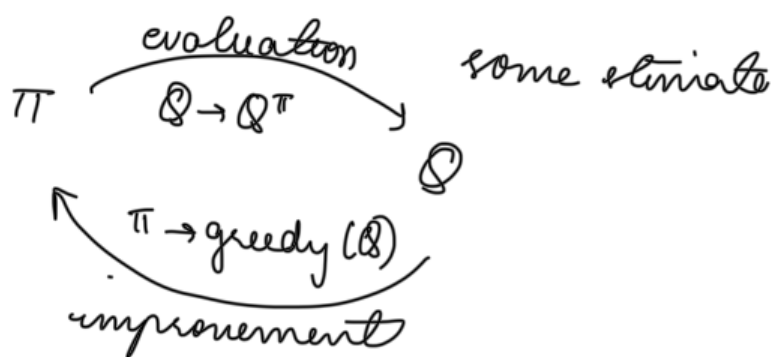
TD control
find the optimal policy.

GPI:



$\pi \xrightarrow{\text{evaluation}} Q \to Q^\pi$ some estimate $Q$

$\pi \to$ greedy $(Q)$
improvement

Policy Evaluation : use $TD(0)$

Policy Improvement : make greedy wrt current value function

**Note:**
We estimate action values rather than state values in the absence of model.

**ε- Greedy Policies:**

$$a^* \leftarrow \text{argmax}_a \, Q(s, a)$$

$$\forall \, a \in A(s):$$

$$\pi(a|s) = \begin{cases} 1 - \varepsilon + \dfrac{\varepsilon}{|A(s)|} & \text{if } a = a^* \\ \varepsilon / |A(s)| & \text{if } a \neq a^* \end{cases}$$

→ any ε greedy policy wrt $Q$ following $\pi$ is an impr over
any ε - soft policy is assured by the policy improvement
$\pi(a|s)$ is atleast     theorem.
ε for every a

**(1) SARSA : On - Policy TD control**

samples are
to policy we
are trying to
evaluate & improve

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \, Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \right]$$

TD(0) for $Q$

$$Q(s_{t+1}, a_{t+1}) = 0 \quad \text{if } s_{t+1} \text{ is terminal}$$

**Sarsa Algorithm :**
    initialize $Q(s, a)$ arbitrarily
    For each episode                    $\pi :$ derived from $Q$
        $S$
        $a \leftarrow \pi(s)$  (eg. ε greedy)

For each step in episode

$$S \xrightarrow[r]{a} S'$$

(improvement)  $\qquad a' = \pi(s')$

(evaluation)  $\qquad Q(s,a) \leftarrow Q(s,a) + \alpha\left[r + \gamma Q(s',a') - Q(s,a)\right]$

$$s \leftarrow s' \qquad a \leftarrow a'$$
until $s$ is terminal

$Q = q_\pi$  (Policy eval)



$q_*$    $\pi_*$

$\pi = \varepsilon\text{-greedy}(Q)$

(Policy imp)

## Convergence

- all $(s,a)$ vis $\infty$
- converges in the limit to the greedy policy $\varepsilon \to 0$

$$(\text{GLIE})$$

$\longrightarrow Q\text{-learning}$

One-step $Q$-learning

Temporal Difference

$$\hat{q}(S_t, a_t) \leftarrow \hat{q}(S_t, a_t) + \alpha\left[r_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a) - \hat{q}(S_t, a_t)\right]$$

in SARSA $\quad Q(S_{t+1}, a_{t+1})$

Bellman Optimality Equation:

$$Q^*(s,a) = E\left\{ r_{t+1} + \gamma \max_{a'} Q^*(S_{t+1}, a') \mid S_t = S, a_t = a \right\}$$
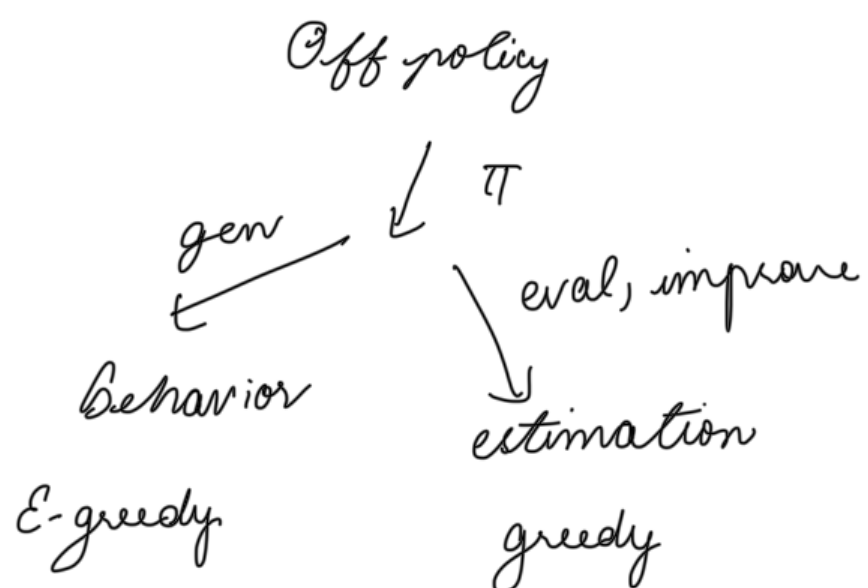
## Stochastic Averaging Rule:

$$G(x) \approx \frac{1}{n} \sum_{i=1}^{n} x_i$$

$$\overline{x_{n+1}} = \frac{1}{n+1}\left( x_{n+1} + \overline{x_n} \cdot n \right)$$

$$= \frac{1}{n+1}\left( x_{n+1} + (n+1)\overline{x_n} - \overline{x_n} \right)$$

$$= \overline{x_n} + \frac{1}{n+1}\left( x_{n+1} - \overline{x_n} \right)$$

$$= \overline{x_n} + \alpha \left( x_{n+1} - \overline{x_n} \right)$$

new est = old est + $\alpha$ (new sample - old est)

TD(0) $\rightarrow$ expectations in Bellman Eq to an average

Q learning $\rightarrow$ optimality

$\rightarrow$ Q learning : Off policy TD control

Off policy

gen $\swarrow$ $\pi$

eval, improve

Behavior

estimation

$\varepsilon$-greedy

greedy

## Q - learning Algorithm:

initialize $Q(s,a)$ arbitrarily

for each episode

initialize $s$

for each step of ep

can be
uniform random or
anything you like ~~~~~~

$u = \pi(s)$ (policy derived from $Q$)
eg $\varepsilon$-greedy

$$(s) \xrightarrow[r]{a} (s')$$

estimation policy ~~~
is always
greedy.

$$Q(s,a) \mathrel{+}= \alpha \left[ r + \gamma \max_{a'} Q(s',a') - Q(s,a) \right]$$
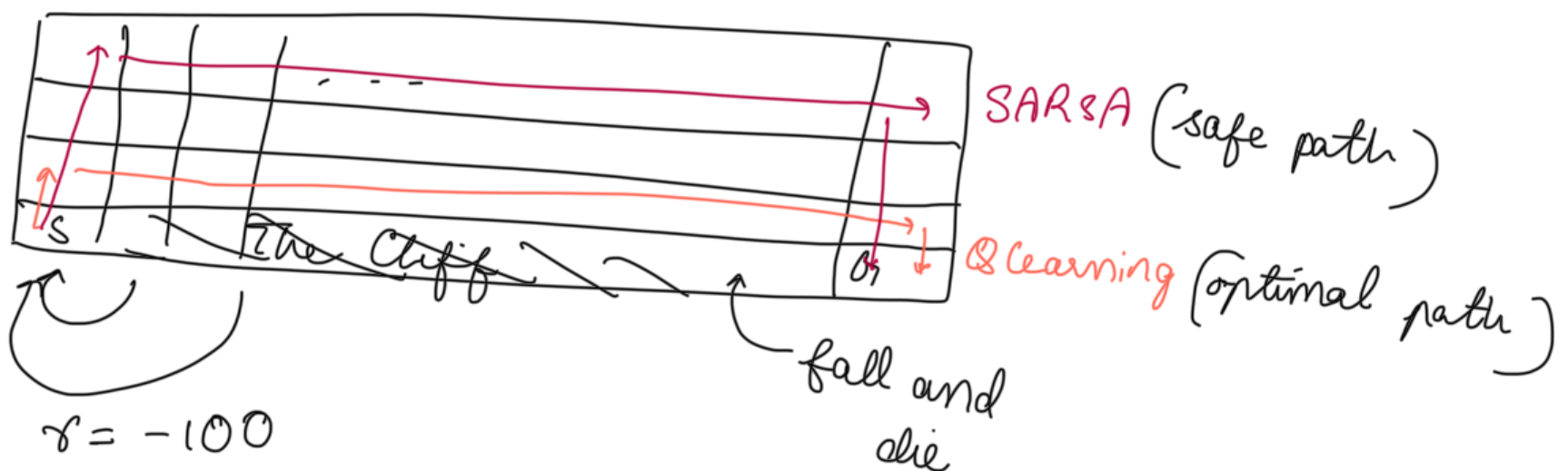
$s \leftarrow s'$

until $s$ is terminal

→ if $\varepsilon$-greedy

only $\varepsilon - \frac{\varepsilon}{|\mathcal{A}(s)|}$ times they differ

SARSA acc to exploratory action

$Q$      greedy

$\delta = -1$



SARSA (safe path)

Q Learning (optimal path)

→ fall and
die

$r = -100$

say 10% times I decide to go to cliff, I fall and update
my previous state in SARSA

Q-Learning: still update acc to greedy policy.
     assumes you will behave greedily in future
Sarsa: $\varepsilon$ greedy execution in future.
     (exploratory)
     incorporate cost of exploration

## Off Policy Learning

- Target $\pi(a|s) \xrightarrow{\text{eval}} V_\pi(s), q_{,\pi}(s,a)$

policy

- Behaviour policy $\mu(a|s)$  eg. if $\pi$ deterministic
observing somebody else (agent) but want to learn about $\pi$

- assumption of coverage:
  - $\pi(a|s) > 0 \rightarrow \mu(a|s) > 0$

- eg.  $\pi$ - greedy
        $\mu$ - $\varepsilon$ - greedy

$\rightarrow$ _Importance Sampling_:

$$E_{x \sim p}\left[f(x)\right] = \sum P(x) f(x)$$

$X \sim P$
but can only sample acc to $Q$

$$= \sum Q(x) \frac{P(x)}{Q(x)} f(x)$$

$$= E_{x \sim Q}\left[\underbrace{\frac{P(x)}{Q(x)}}_{w(x)} f(x)\right]$$

if I see it often in $P$ but not $Q$ it is importance to give it more weight
and if $P(x) \neq 0$ imp $Q(x) \neq 0$

$\rightarrow$ _Importance Sampling Ratio_

$$\rho_t^T = \frac{\prod_{k=t}^{T-1} \pi\left(\frac{A_k}{S_k}\right) P\left(\frac{S_{k+1}}{S_k A_k}\right)}{\prod_{k=t}^{T-1} \mu\left(\frac{A_k}{S_k}\right) P\left(\frac{S_{k+1}}{S_k, A_k}\right)}$$

$$= \frac{\prod_{k=t}^{T-1} \pi\left(A_k|S_k\right)}{\prod_{k=t}^{T-1} \mu\left(A_k|S_k\right)} = \log \sum \frac{\pi}{\mu}$$

$$k = t \quad | \quad \overbrace{\phantom{-14L}}$$

weighted average return for $V$

$$V(s) = \frac{\displaystyle\sum_{t=\tau(s)} p_t^{T(t)} \, G(t)}{\displaystyle\sum_{t \in C(s)} p_t^{T(t)}}$$

Random Variable: $G(t)$

trying to sample & average

$$\frac{N}{D}$$

$$G(t) = r_{t+1} + \gamma \, r_{t+2} + \cdots + \gamma^{T-t-1}$$

$$S_t \; A_t \; S_{t+1} \; \cdots \cdots$$