

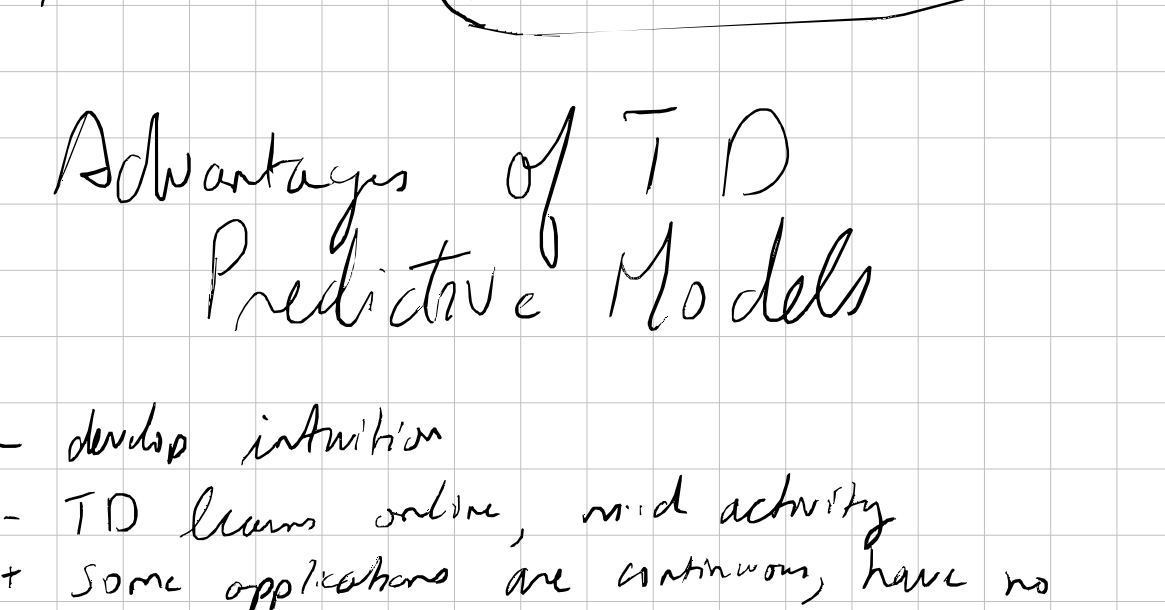
# Temporal-Difference Learning

- central & novel to RL
- MC wombo combo with DP
- + TD can learn from experience w/o dynamics like MC
- + TD can update estimates from other estimates w/o final outcome (bootstrap) like DP
- \* relationship between all is throughout rest of book
- once again, focus on evaluation or prediction problem first, then iteration or control problem of finding optimal policy
- + find  $v_{\pi}$  off  $\pi_0$ , then new  $\pi$  of  $v_{\pi}$

## TD Prediction

- MC basic form
- $$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$
- + [Constant  $\alpha$  MC] (step-size param)
  - MC waits till end of ep because it must know full future  $G_t$
  - TD only needs  $R_{t+1}$ , simplest form:
- $$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$
- MC update target  $G_t$  | TD update target  $R_{t+1} + \gamma V(S_{t+1})$
  - Called TD(0) / one-step TD
  - bootstraps on estimate
  - remember

- $$v_{\pi}(s) = \mathbb{E}_{\pi} [G_t | S_t] \quad (1)$$
- $$= \mathbb{E}_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t] \quad (2)$$
- $$= \mathbb{E}_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t] \quad (3)$$
- MC uses estimate of (1), DP uses estimate of (3)
  - + MC use sample return, hence estimate
  - + DP uses  $V(S_{t+1})$ , because  $v_{\pi}(S_{t+1})$  is not known.
  - x "estimate" difference
  - TD does both, sample expected return & new estimate  $V$ , not  $v_{\pi}$



- sample not expected because we don't look at the whole distribution of possible successors like in DP

## TD error

- $$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$
- between current est.  $S_t$  and better on  $R_t + \gamma V(S_{t+1})$
- + it can be available at  $t+1$
  - \* Suppose  $V$  isn't updated in the ep, like in MC methods, then error looks like:
- $$G_t - V(S_t) = R_{t+1} + \gamma(G_{t+1} - V(S_t) + \gamma V(S_{t+1}) - \gamma V(S_t))$$
- $$= \delta_t + \gamma(G_{t+1} - V(S_{t+1}))$$
- $$= \delta_t + \gamma(\delta_{t+1} + \gamma(G_{t+2} - V(S_{t+2})))$$
- $$= \sum_{k=0}^{T-t} \gamma^k \delta_{t+k}$$

## EXERCISE 6.1 What if $V$ changed by a small amount?

- Let  $\Delta V(S_t) = V_{t+1}(S_t) - V_t(S_t)$  (1)
- $$G_t - V_t(S_t) = R_{t+1} + \gamma G_{t+1} - V_t(S_t)$$
- $$= R_{t+1} + \gamma G_{t+1} - V_t(S_t) + \gamma V_t(S_{t+1}) - \gamma V_t(S_{t+1})$$
- $$= \delta_t + \gamma(G_{t+1} - V_t(S_{t+1}))$$
- $$= \delta_t + \gamma(G_{t+1} - V_{t+1}(S_{t+1}) - \Delta V_t(S_{t+1}))$$
- $$= \delta_t + \gamma(G_{t+1} - V_{t+1}(S_{t+1})) - \gamma \Delta V_t(S_{t+1})$$
- $$= \delta_t + \gamma(\delta_{t+1} + \gamma(G_{t+2} - V_{t+2}(S_{t+2})) - \gamma \Delta V_{t+1}(S_{t+2})) - \gamma \Delta V_t(S_{t+1})$$
- $$= \sum_{k=0}^{T-t} \gamma^k \delta_{t+k} + \left( \sum_{k=0}^{T-t} \gamma^{k+1} \Delta V_{t+k}(S_{t+k+1}) \right)$$

## Advantages of TD Predictive Models

- develop intuition
- TD learns online, mid activity
- + some applications are continuous, have no episodes
- + MC must discount wrong experimentation, slowing learning
- TD methods converge
- Open question, which is faster? MC or TD?
- + in practice TD seems faster, but no formal proof
- x question itself is questionable

## EXERCISE 6.3

$$V(A) = V(A) + \alpha [R_T + \gamma V(T) - V(A)]$$

$$= 0.5 + 0.1(0 + -0.5)$$

$$= 0.5 - 0.05 = 0.45$$

Rest of states  $V(S) = 0.5$ , so updates to  $V(S) = 0$  except  $\epsilon$  terminal state, where value = 0

## EXERCISE 6.4

I don't think any larger or different values would help. It would only overshoot the true value so the error term would bring it back down. Most of these algos have to adjust the step size down over time to avoid that problem

## Optimality of TD(0)

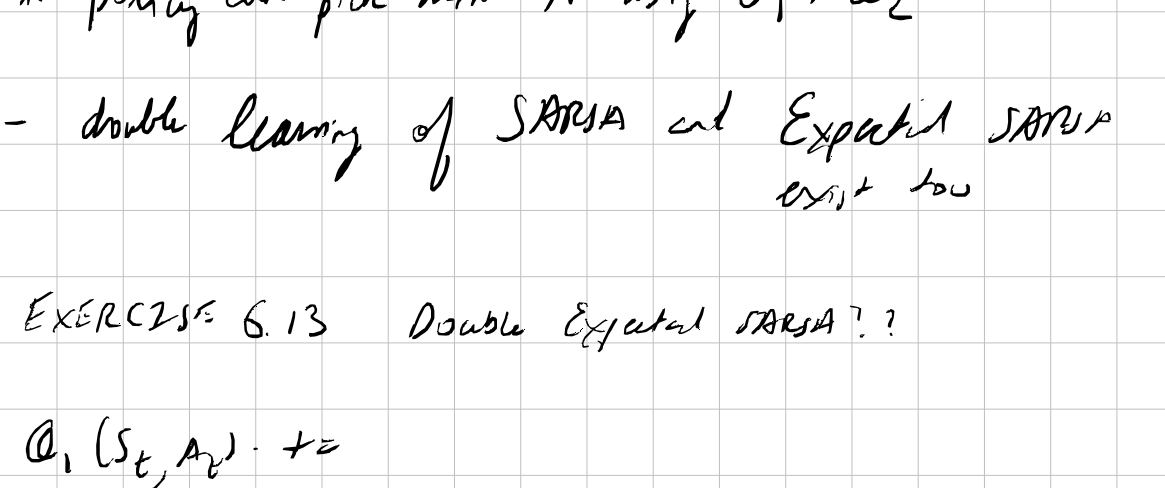
- batch updating  $\rightarrow$  compute increments, but don't update  $V$  after all increments are computed
- + keep representing same experience data till convergence
- in TD(0) converges to answer regardless of  $\alpha$  (if small)
- + in MC, converges to diff answer

## EXAMPLE 6.3: Random Walk

- batch updates of TD(0), MC
- MC reaches a  $V$  that is optimal for its experience
- TD(0) outperforms

+ MC is optimal in limited way, TD optimal more relevant to prediction

## EXAMPLE 6.4: 8 eps



What is  $V(A)$ ? Markov process  $V(A) = 3/4$

However! MC attempt to reduce RMS on training data with est.  $V(A) = 0$ ! TD(0) will update A on every visit from

- MC minimizes error on seen data
- TD(0) estimates value for maximum likelihood Markov model

+ certainty-equivalence estimate because it assumes that the est. of underlying process was known for certain

- x TD(0) converges to J
- \* non batch does not converge the way batch TD does.
- almost always impossible to compute certainty equivalence est due to space and time needs in states.

## EXERCISE 6.7 off policy TD(0) with $\pi$ / b & $p_{\pi,b}$

Incremental Update becomes:

$$V(S_t) = V(S_t) + \alpha p_{\pi,b} [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

## Soma: On-Policy TD Control

- explore vs exploit in control, again.
- first, use state-action values:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

- $Q(S_{t+1}, A_{t+1}) = 0$  if  $t+1 = T$  (terminal)
- quadruple in var:  $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$



- convergence depends on policy's relationship with  $Q$
- +  $\epsilon$ -soft or  $\epsilon$ -greedy with  $\epsilon = 1/b$
- x state-action pairs must be visited infinitely
- x policy must converge to limit of greedy policy

## SARSA (on policy TD control) for $Q \approx q_{\pi}$

- $\alpha \in (0, 1)$
- $\epsilon > 0$  (small)
- $Q(s, a) \forall s \in S^+, a \in \mathcal{A}(s)$  arbitrarily except  $Q(\text{terminal}, \cdot) = 0$

for each ep

init  $S$

choose  $A$  from  $S$  based on  $Q$  ( $\epsilon$ -greedy)

for each step in ep:

$S' \leftarrow \text{step}(S, A)$

choose  $A'$  from  $S'$

$$Q(S, A) = Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$$

$S = S'$

$A = A'$

until termination

- important to remember that if a policy chosen for MC gets stuck, ep never turns, MC never updates  $Q$
- + SARSA online doesn't suffer, will put a new policy

## Q-learning: Off-Policy TD Control

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

- $Q$  directly approximates  $q_{\pi^*}$ , regardless of policy
- + makes it "off policy" on update

## EXAMPLE: Cliff Walking



- SARSA takes action into account, so stays away from cliff
- Q-learning will learn optimal faster, but will die a lot because it doesn't care about action taken

## EXPECTED SARSA

- what if you liked Q learning, but also considered the likelihood of each action, not just the max value
- + hence, obtain expected value

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \mathbb{E}_{\pi} [Q(S_{t+1}, A_{t+1}) | S_t] - Q(S_t, A_t)]$$

$$= Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \sum_a \pi(a | S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t)]$$

- deterministic movement as SARSA in expectation

## Expected SARSA

- outperforms SARSA, does better on live than Q-learning
- + if done with off-policy strategy, it is the same as Q-learning
- x Expected SARSA generalizes Q-learning
- \* dominates most well known TD methods



## Maximization Bias & Double Learning

- current algos use maximum of estimate for true value
- leads to significant positive bias

## EXAMPLE



- Expected value of left trajectory is -0.1
- right 0
- left always worse
- current control prefers left because of maximization bias
- + positive value over optimal

- fix by varying problem as split between estimating value and picking best action
- + don't use the same samples for both
- +  $Q_1, Q_2 \approx q_{\pi^*}$
- +  $A^* = \arg \max_a Q(a)$  use  $Q_1$  to pick

+  $Q_2(A^*) = Q_2(\arg \max_a Q_1(a))$  now use  $Q_2$  to estimate value

\* then flip flop every timestep to evenly distribute bias

## \* Double Q-learning

$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha [R_{t+1} + \gamma Q_2(S_{t+1}, \arg \max_a Q_1(S_{t+1}, a)) - Q_1(S_t, A_t)]$$

- \* policy can pick action  $A$  using  $Q_1$  +  $Q_2$
- double learning of SARSA and Expected SARSA exist too

## EXERCISE 6.13 Double Expected SARSA?

$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha [R_{t+1} + \gamma \sum_a \pi(a | S_{t+1}) Q_2(S_{t+1}, a) - Q_1(S_t, A_t)]$$

## Games, Afterstates, & Other Special Cases

- what happens when different state actions lead to the same afterstate?

- + games often have this ability
- + some states must have same value
- x speeds up learning
- occur in many tasks impossible to specify to all the specializations
- + can still introduce appropriate value function and policy iteration in GPIT

## Summary

- these TD methods are simple & powerful
- they are one step, tabular, model free
- + soon, expect to non-step (MC link) & model forms (planning & DP)
- x after, further approximation instead of tabular (neural nets)