

n-step Bootstrapping

- Bridge TD(0) & MC
- allow smooth transition based on task
- debate from America
- + Bootstrapping works best over large time periods / more change
- + TD(0) removes time-step boundary
- used for eligibility traces, but discussed later
- prediction problem first

n-step TD Prediction

- Again MC = all rewards to evaluate
- TD(0) = 1 step
- n-step = TD of n steps
- Consider update for S_t from sequence $S_t, R_{t+1}, S_{t+1}, \dots, R_T, S_T$
- + MC needs complete return

$$G_T = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T$$

- + MC, return is target of updates
- + TD(0) update reward + discounted est value of next state one-step return

$$G_{t:t+1} = R_{t+1} + \gamma V_t(S_{t+1})$$

- + So, now let go to $n=2, n=3$

$$G_{t:t+2} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V_t(S_{t+2})$$

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n})$$

- basically use discounted reward and "truncated" return at $V_{t+n-1}(S_{t+n})$
- + if $t+n > T$, $G_{t:t+n} = G_t$ (terms beyond $T=0$)
- can't really use though because R_{t+n} & V_{t+n-1} aren't known until seen at $t+n$
- + then

$$V_{t+n}(S_t) = V_{t+n-1}(S_t) + \alpha [G_{t:t+n} - V_{t+n-1}(S_t)]$$

and $V_{t+n}(S) = V_{t+n-1}(S) \forall S \neq S_t$

- x can't make updates till $n-1$ steps
- x make additional $n-1$ updates after termination

EXERCISE 7.1

$$G_{t:t+n} - V(S_t) = R_{t+1} + \gamma G_{t+1:t+n} - V(S_t) - \gamma V(S_{t+1}) + \gamma V(S_{t+1})$$

$$= \delta_t + \gamma (G_{t+1:t+n} - V(S_{t+1}))$$

$$= \delta_t + \gamma (\delta_{t+1} + \gamma (G_{t+2:t+n} - V(S_{t+2})))$$

$$= \delta_t + \gamma \delta_{t+1} + \gamma^2 (G_{t+2:t+n} - V(S_{t+2}))$$

$$= \delta_t + \gamma \delta_{t+1} + \dots + \gamma^{n-1} (R_{t+n} + \gamma V(S_{t+n}) - V(S_{t+n-1}))$$

$$= \sum_{k=0}^{n-1} \gamma^k \delta_{t+k}$$

EXAMPLE 7.1

- reasoning example for n-step TD
- $n=2, T=3$
- $t=0, T=3 \Rightarrow t-n+1 = -1$ do nothing
- $t=1, T=3 \Rightarrow t-n+1 = 0 \Rightarrow \gamma^0 R_1 + \gamma R_2 + \gamma^2 V(S_2)$
- $t=2, T=3 \Rightarrow t-n+1 = 1 \Rightarrow \gamma^1 R_2 + \gamma^2 R_3$
- $t=3, T=3 \Rightarrow t-n+1 = 2 \Rightarrow \gamma^0 R_3$
- how do determine a good n value?
- Empirical analysis
- + example here found intermediate value best

n-step SARSA

- combine it with control methods
- previous was one-step SARSA or SARSA(0)
- first redefine with action-values

$$Q_{t+n}(S_t, A_t) = Q_{t+n-1}(S_t, A_t) + \alpha [G_{t:t+n} - Q_{t+n-1}(S_t, A_t)]$$

- n-step expected SARSA:
- + single sequence till last state, then $V_{t+n-1}(S_{t+n})$

$$G_{t:t+n} = R_{t+1} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n \bar{V}_{t+n-1}(S_{t+n})$$

with $\bar{V}_t(S)$ as expected approximate value

$$\bar{V}_t(S) = \sum_a \pi(a|S) Q_t(S, a)$$

- * used throughout later algos
- * if S is terminal, $\bar{V}_t(S) = 0$

n-step Off Policy Learning

- n-step TD
- $V_{t+n}(S_t) = V_{t+n-1}(S_t) + \alpha p_{t:t+n-1} [G_{t:t+n} - V_{t+n-1}(S_t)] \quad 0 \leq t \leq T$

- + $p_{t:t+n-1}$ because last step is only the state visited
- x $V=0$ if $p=0$ because behavior policy takes action π would never take
- for action values

$$Q_{t+n}(S_t, A_t) = Q_{t+n-1}(S_t, A_t) + \alpha p_{t:t+n-1} [G_{t:t+n} - Q_{t+n-1}(S_t, A_t)]$$

- $p_{t:t+n-1}$ because A_t already taken, and A_{t+n} has been chosen
- these generalize because if on-policy $p=1$
- n-step expected SARSA was form $t+n-1$ because last value is V , across all a

On Decision Methods with Control Variables

- previous algos are not efficient with off policy sampling ratios

- + use notation $G_{t:h} = R_{t+1} + \dots + G_{t+h-1:h} \quad t \leq h \leq T$
- $\pi(a|s) = 0$ for a chosen by b
- $V(s) = 0$, resulting in high variance
- + instead:

$$G_{t:h} = p_t (R_{t+1} + \gamma G_{t+1:h}) + (1-p_t) V_{h-1}(S_t)$$

and $G_{h:h} = V_{h-1}(S_h)$

- when behavior chooses action that target would never choose, Do NOT update V , just ignore sample

- + second term called control variate
- + use with conventional n-step TD learning equation
- action values, a little different

$$G_{t:h} = R_{t+1} + \gamma (p_{t+1} (G_{t+1:h} + \bar{V}_{h-1}(S_{t+1}) - p_{t+1} Q_{h-1}(S_{t+1}, A_{t+1})))$$

$$= R_{t+1} + \gamma p_{t+1} (G_{t+1:h} - Q_{h-1}(S_{t+1}, A_{t+1})) + \gamma \bar{V}_{h-1}(S_t)$$

- ends with $G_{h:h} = Q_{h-1}(S_h, A_h)$
- if $h > T$, $G_{h-1:h} = R_t$

- analogous to expected SARSA when combined with learning step
- off-policy learning with importance sampling is good, but slow and must be controlled for high variance
- + can be improved using control variates, average step-size in accordance with variances, invariant updates

Off Policy Learning Without Importance Sampling: The n-step Tree Backup Algorithm

- Backup example on right

- for ref central spine are sampled S, A, R
- target previously was made up of only sampled data
- + now target will include bootstrapped (estimated) values of actions NOT taken
- x only considers leaf nodes
- x actions actually taken are not considered
- x weight actions by likelihood down the chain
- x $\pi(a|S_{t+1})$ for 1st level
- $\pi(A_{t+1}|S_{t+1}) \pi(a'|S_{t+2})$ for 2nd level
- $\pi(A_{t+1}|S_{t+1}) \pi(A_{t+2}|S_{t+2}) \pi(a'|S_{t+3})$

- 3-step tree backup is more like 6 "half" steps
- + sample action to state & from state to all possible actions in policy

- n-step tree backup equation
- + start with 1-step return = expected SARSA

$$G_{t:t+n} = R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q_{t+n-1}(S_{t+1}, a)$$

- + 2-step:

$$G_{t:t+2} = R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1}) Q_{t+1}(S_{t+1}, a) + \gamma \pi(A_{t+1}|S_{t+1}) (R_{t+2} + \gamma \sum_a \pi(a|S_{t+2}) Q_{t+2}(S_{t+2}, a))$$

- + general recursion

$$G_{t:t+n} = R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1}) Q_{t+n-1}(S_{t+1}, a) + \gamma \pi(A_{t+1}|S_{t+1}) G_{t+1:t+n}$$

- for $t \leq T-1, n \geq 2, G_{T-1:t+n} = R_T$

- combine with general update step

$$Q_{t+n}(S_t, A_t) = Q_{t+n-1}(S_t, A_t) + \alpha [G_{t:t+n} - Q_{t+n-1}(S_t, A_t)]$$

- for $0 \leq t \leq T$, & other values are unchanged

$$Q_{t+n}(S, a) = Q_{t+n-1}(S, a) \quad S \neq S_t, a \neq A_t$$

EXERCISE 7.11

Show that

$$G_{t:t+n} = Q(S_t, A_t) + \sum_{k=t}^n \delta_k \prod_{i=t+1}^k \gamma \pi(A_i|S_i)$$

$$\delta_t = R_{t+1} + \gamma \bar{V}_t(S_{t+1}) - Q(S_t, A_t)$$

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

$$G_{t:t+n} = R_{t+1} + \gamma (\bar{V}_t(S_{t+1}) - Q(A_{t+1}, S_{t+1})) + \gamma \pi(A_{t+1}|S_{t+1}) (R_{t+2} + \gamma (\bar{V}_{t+1}(S_{t+2}) - Q(A_{t+2}, S_{t+2})) + \dots)$$

- $R_{t+1} + \gamma \bar{V}_t(S_{t+1}) - \gamma Q(A_{t+1}, S_{t+1}) + \gamma \pi(A_{t+1}|S_{t+1}) R_{t+2} + \gamma \pi(A_{t+1}|S_{t+1}) \gamma \bar{V}_{t+1}(S_{t+2}) - \gamma \pi(A_{t+1}|S_{t+1}) Q(A_{t+2}, S_{t+2}) + \dots$
- $\textcircled{1} + Q(S_t, A_t) - Q(S_t, A_t)$
- $(R_{t+1} + \gamma \bar{V}_t(S_{t+1}) - Q(S_t, A_t))(1) + (R_{t+2} + \gamma \bar{V}_{t+1}(S_{t+2}) - Q(S_{t+1}, A_{t+1}))(\gamma \pi(A_{t+1}|S_{t+1}))$
- $Q(S_t, A_t) + \sum_{k=t}^n \delta_k \prod_{i=t+1}^k \gamma \pi(A_i|S_i)$

* A Unifying Algorithm

n-step $Q(\sigma)$

- how do we take n-step SARSA, expected SARSA, & tree backup, all together?

- $\sigma = \sqrt{\text{degree of sampling}}$
- + $\sigma=1$ full sample
- $\sigma=0$ full expectation, ignore sample

- first, tree backup with $t+n=h$

$$G_{t:h} = R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1}) Q_{t+h-1}(S_{t+1}, a) + \gamma \pi(A_{t+1}|S_{t+1}) G_{t+1:h}$$

$$= R_{t+1} + \gamma \bar{V}_{h-1}(S_{t+1}) - \gamma \pi(A_{t+1}|S_{t+1}) Q_{t+h-1}(S_{t+1}, A_{t+1}) + \gamma \pi(A_{t+1}|S_{t+1}) G_{t+1:h}$$

$$= R_{t+1} + \gamma \pi(A_{t+1}|S_{t+1}) (G_{t+1:h} - Q_{h-1}(S_{t+1}, A_{t+1})) + \gamma \bar{V}_{h-1}(S_{t+1})$$

- then, slide linearly between σ

$$G_{t:h} = R_{t+1} + \gamma (G_{t+1:h} + (1-\sigma) \pi(A_{t+1}|S_{t+1}) (G_{t+1:h} - Q_{h-1}(S_{t+1}, A_{t+1})) + \gamma \bar{V}_{h-1}(S_{t+1}))$$

- if $h \leq t$, or $G_{T-1:T} = R_T$ if $h=T$

- combine with n-step SARSA, but leave out importance sampling, since it is in return itself.

Summary

- n-step: more computation, more memory
- + can be reduced with eligibility traces later