



$$E_b \left[ (1 - \rho_t) V_{h-1}(s_t) \right]$$

$$E_b [\rho_t] = 1 \quad \& \quad \rho_t \text{ independent of } V_{h-1}(s_t)$$

behavior policy  $b$   
target policy  $\pi$

# Off-policy learning with importance sampling: State value

Without control variate:

$$\rho_t = \frac{\pi(a_t | s_t)}{b(a_t | s_t)}$$

$$\tilde{G}_{t:h} = \rho_t (R_{t+1} + \gamma \tilde{G}_{t+1:h}) + \boxed{?}$$

Current state value  $\times$   
 Target state value  $\circ$

With control variate:

$$\hat{G}_{t:h} = \rho_t (R_{t+1} + \gamma \hat{G}_{t+1:h}) + \boxed{(1 - \rho_t) V_{h-1}(s_t)}$$

Control variate

$$\rho_t \approx 0$$

$$0 (R_{t+1} + \gamma \hat{G}_{t+1:h}) + 1 V_{h-1}(s_t) = V_{h-1}(s_t)$$

$$E_b[\underbrace{\rho_{t+1}} \underbrace{Q_{h-1}(s_{t+1}, A_{t+1})}] = E_{\pi}[Q_{h-1}(s_{t+1}, A_{t+1})] = \boxed{V_{h-1}(s_{t+1})} \quad E_b[V_{h-1}(s_{t+1})] = V_{h-1}(s_{t+1})$$

# Off-policy learning with importance sampling: Action value

Without control variate:

$$E[E[R] - R] = E[E[R]] - E[R] = 0$$

$$\tilde{G}_{t:h} = \underline{R_{t+1}} + \gamma \underline{\rho_{t+1}} \tilde{G}_{t+1:h}$$

$$E[s] = s$$

$$Q_{\pi}(s_t, A_t)$$

$$\bar{V}(s) = \sum_a \pi(a|s) Q(s,a)$$

$\rho_t$  tells us the ratio of  $\pi(A_t|s_t)$  to  $b(A_t|s_t)$

With control variate:

$$E_b[\text{control variate}] = 0$$

$$E_b[E_b[\rho_{t+1} Q \dots] - \rho_{t+1} Q \dots]$$

$$\hat{G}_{t:h} = R_{t+1} + \gamma(\rho_{t+1} \hat{G}_{t+1:h} + \underbrace{E_b[\rho_{t+1} Q_{h-1}(s_{t+1}, A_{t+1})]} - \rho_{t+1} Q_{h-1}(s_{t+1}, A_{t+1}))$$

$$= R_{t+1} + \gamma(\rho_{t+1} \hat{G}_{t+1:h} + \underbrace{V_{h-1}(s_{t+1})} - \rho_{t+1} Q_{h-1}(s_{t+1}, A_{t+1}))$$

$$= R_{t+1} + \gamma \rho_{t+1} (\hat{G}_{t+1:h} - Q_{h-1}(s_{t+1}, A_{t+1})) + \gamma \bar{V}_{h-1}(s_{t+1})$$

if  $\rho_{t+1} = 0$

$$\bar{V}(s) = \sum_a \pi(a|s) Q(s, a)$$

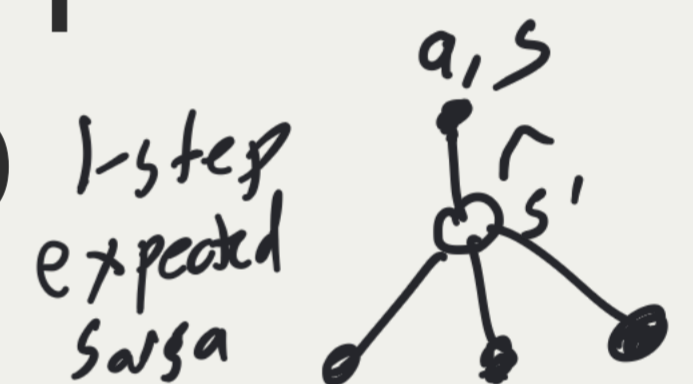
We are estimating action values  $Q(s, a)$

If we don't know  $b$ , can we do importance sampling? No. Requires  $p$  which requires  $b$

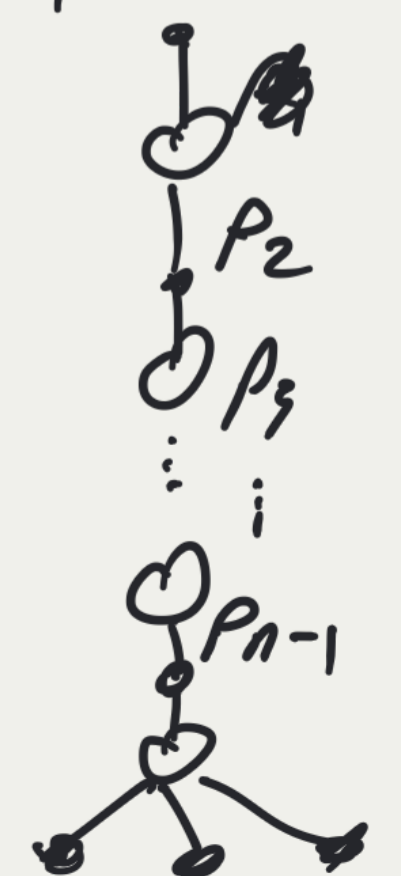
# Off-Policy learning **without** importance sampling: n-step Tree Backup

one-step:

$$G_{t:t+1} = \underline{R_t} + \underline{\gamma} \sum_a \pi(a|S_{t+1}) \underline{Q_t(S_{t+1}, a)}$$



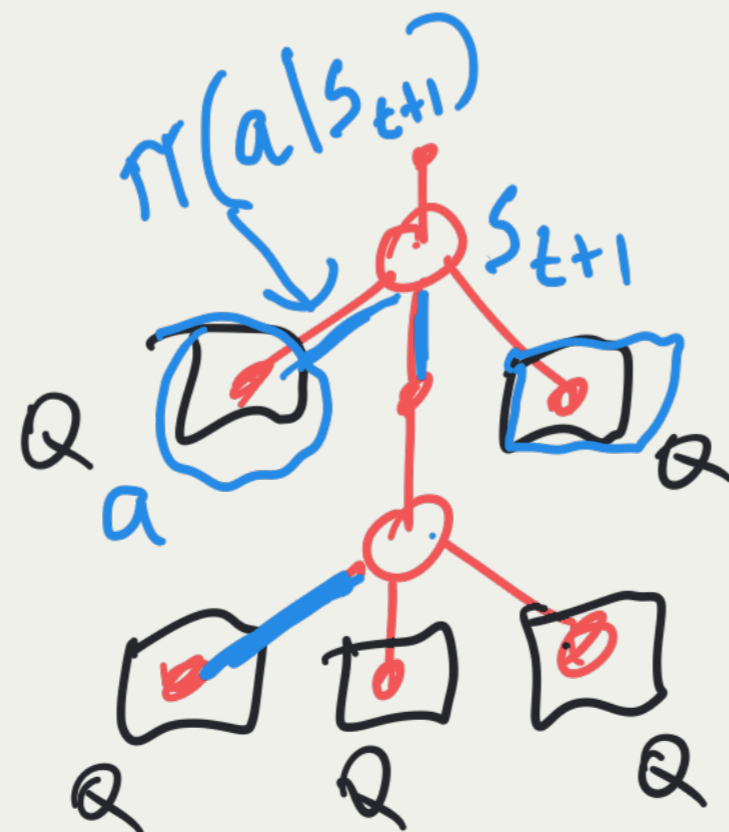
n-step expected  $S_{t+n} a$

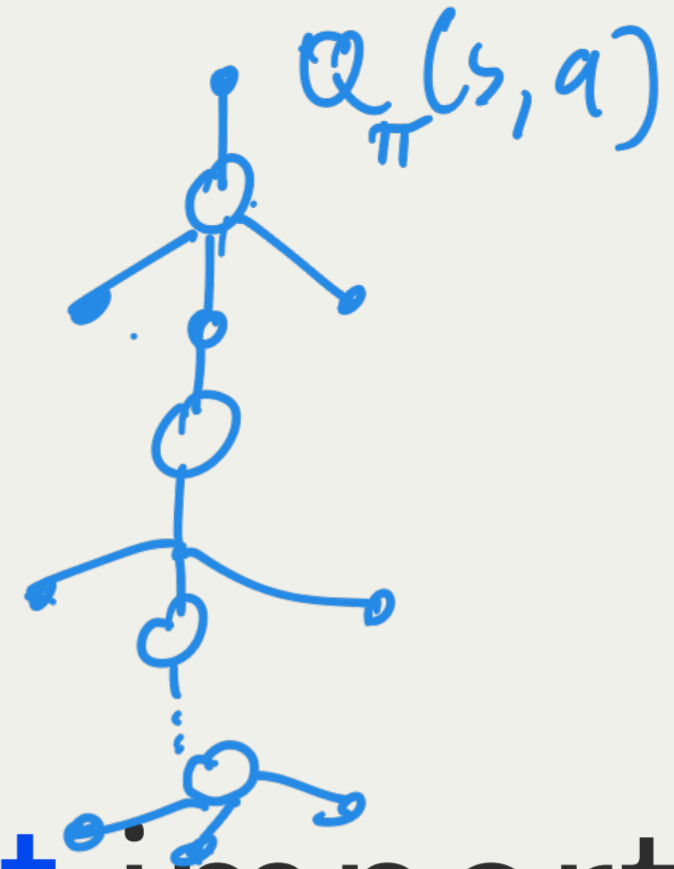


# Off-Policy learning **without** importance sampling: n-step Tree Backup

two-step:

$$\begin{aligned} G_{t:t+2} &= \underline{R_{t+1}} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1}) Q_t(S_{t+1}, a) \\ &+ \gamma \pi(A_{t+1}|S_{t+1}) \left( \underline{R_{t+2}} + \gamma \sum_a \pi(a|S_{t+2}) Q_{t+1}(S_{t+2}, a) \right) \\ &= R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1}) Q_t(S_{t+1}, a) + \gamma \pi(A_{t+1}, S_{t+1}) G_{t+1:t+2} \end{aligned}$$





# Off-Policy learning without importance sampling: n-step Tree Backup

n-step:

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

# n-step $Q(\sigma)$

$$\sum_{a \neq A_{t+1}} \text{value} = \underbrace{\sum_a \text{value}}_{\bar{V}(s)} - \text{value}@a = Q(s, A_{t+1})$$

tree-backup return:

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

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

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

per-decision importance sampling return:

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

combined:

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

$$\sigma_t \in \text{Range} [0, 1]$$

$\sigma_i$ 's can be different

For one paper

descending  $\sigma$

$$\sigma = 0$$

$$\sigma = 1$$

$$\sigma = .2$$

$$\sigma = .6$$

change over time from

start reduce sampling "

$\sigma @ 1$   
to  $0.95$   
"  $(0.95)^2$

⋮  
,

⋮  
 $(0.95)^n \approx 0$

to

19-state

random-walk

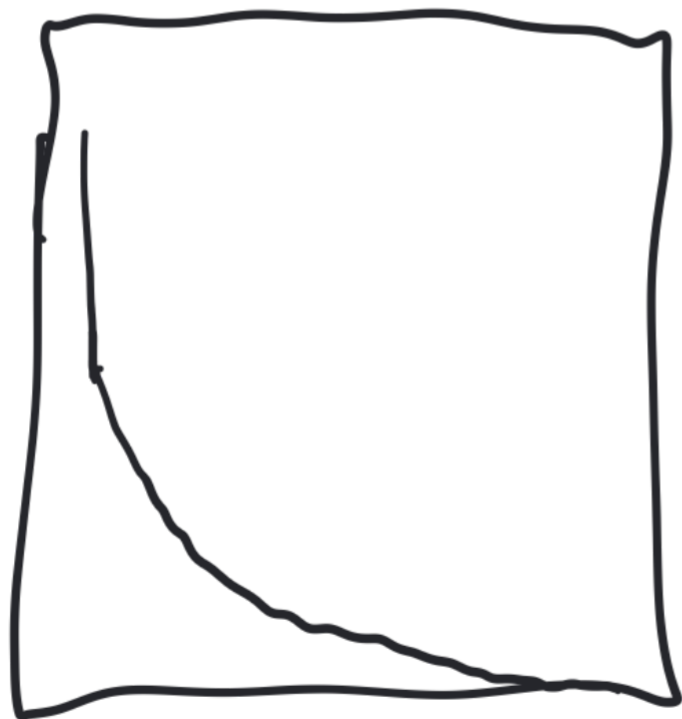
tree

backup



num

RMS



importance sampling

Ordinary importance sampling

weighted importance sampling

per-decision importance sampling

1-step



importance sampling  
descending  $\sigma$

$V(s)$

