

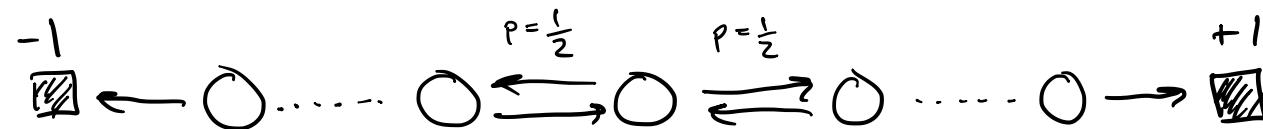
### Targets

$$TD(0) : R_t + \gamma V(s_{t+1})$$

$$n\text{-step TD} : R_t + \gamma R_{t+1} + \dots + \gamma^{n-1} R_{t+n-1} + \gamma^n V(s_{t+n})$$

$$\text{Monte Carlo} : R_t + \gamma R_{t+1} + \dots + \gamma^{T-t-1} R_{T-t}$$

# Chain Markov Reward Process (Length=19)



Mini 5-chain  $\alpha = 0.5 \gamma = 1$

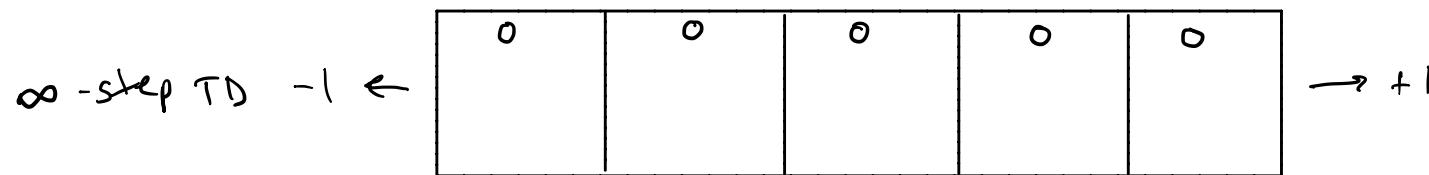
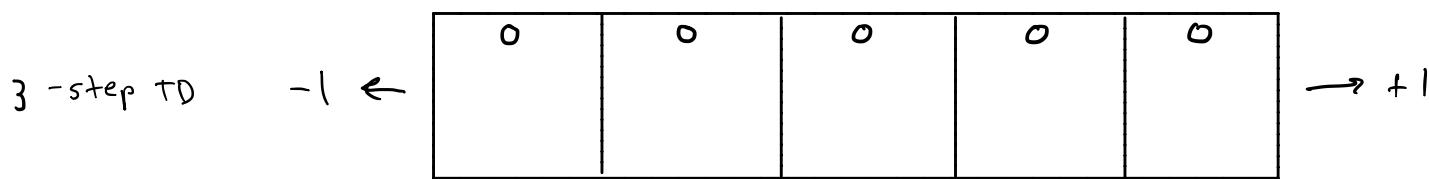
True  $V(s)$ :  $A: -\frac{2}{3}$     $B: -\frac{1}{3}$     $C: 0$     $D: \frac{1}{3}$     $E: \frac{2}{3}$

Experiences

$B, C, D, E, +1$

$C, B, A, -1$

$A, B, C, D, E, +1$



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

# Chain Markov Reward Process (Length=19)



Mini 5-chain  $\alpha = 0.5 \gamma = 1$

True  $V(S)$ :

$A$	$B$	$C$	$D$	$E$
$-\frac{2}{3}$	$-\frac{1}{3}$	0	$\frac{1}{3}$	$\frac{2}{3}$

Experiences

✓  $B, C, D, E, +1$

$C, B, A, -1$

$A, B, C, D, E, +1$

1-step TD

-1	$\leftarrow$	0	0	0	0	<del>0.5</del>	$\rightarrow +1$
----	--------------	---	---	---	---	----------------	------------------

3-step TD

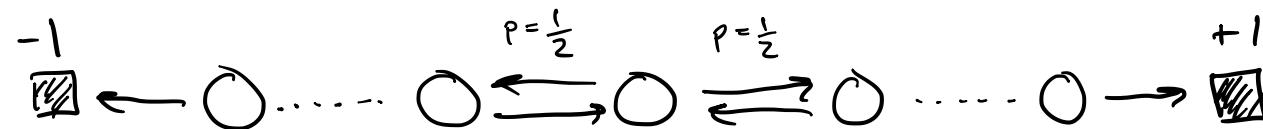
-1	$\leftarrow$	0	0	<del>0.5</del>	<del>0.5</del>	<del>0.5</del>	$\rightarrow +1$
----	--------------	---	---	----------------	----------------	----------------	------------------

$\infty$ -step TD

-1	$\leftarrow$	0	<del>0.5</del>	<del>0.5</del>	<del>0.5</del>	<del>0.5</del>	$\rightarrow +1$
----	--------------	---	----------------	----------------	----------------	----------------	------------------

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

# Chain Markov Reward Process (Length=19)



Mini 5-chain  $\alpha = 0.5 \gamma = 1$

$$\text{True } V(s) : \begin{matrix} -\frac{2}{3} \\ A \end{matrix} \quad \begin{matrix} -\frac{1}{3} \\ B \end{matrix} \quad \begin{matrix} 0 \\ C \end{matrix} \quad \begin{matrix} \frac{1}{3} \\ D \end{matrix} \quad \begin{matrix} \frac{2}{3} \\ E \end{matrix}$$

Experiences

✓ B, C, D, E, +1

✓ C, B, A, -1

A, B, C, D, E, +1

1-step TD	-1	$\leftarrow$	<table border="1" style="display: inline-table; vertical-align: middle;"> <tr><td><del>0</del></td><td>0</td><td>0</td><td>0</td><td><del>0</del></td></tr> <tr><td>-0.5</td><td></td><td></td><td></td><td>0.5</td></tr> </table>	<del>0</del>	0	0	0	<del>0</del>	-0.5				0.5	$\rightarrow +1$
<del>0</del>	0	0	0	<del>0</del>										
-0.5				0.5										

3-step TD	-1	$\leftarrow$	<table border="1" style="display: inline-table; vertical-align: middle;"> <tr><td><del>0</del></td><td><del>0</del></td><td><del>0</del></td><td><del>0</del></td><td><del>0</del></td></tr> <tr><td>-0.5</td><td>-0.5</td><td>-0.5</td><td>0.5</td><td>0.5</td></tr> </table>	<del>0</del>	<del>0</del>	<del>0</del>	<del>0</del>	<del>0</del>	-0.5	-0.5	-0.5	0.5	0.5	$\rightarrow +1$
<del>0</del>	<del>0</del>	<del>0</del>	<del>0</del>	<del>0</del>										
-0.5	-0.5	-0.5	0.5	0.5										

$\infty$ -step TD	-1	$\leftarrow$	<table border="1" style="display: inline-table; vertical-align: middle;"> <tr><td><del>0</del></td><td><del>0</del></td><td><del>0</del></td><td><del>0</del></td><td><del>0</del></td></tr> <tr><td>-0.5</td><td>-0.25</td><td>-0.25</td><td>0.5</td><td>0.5</td></tr> </table>	<del>0</del>	<del>0</del>	<del>0</del>	<del>0</del>	<del>0</del>	-0.5	-0.25	-0.25	0.5	0.5	$\rightarrow +1$
<del>0</del>	<del>0</del>	<del>0</del>	<del>0</del>	<del>0</del>										
-0.5	-0.25	-0.25	0.5	0.5										

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

# Chain Markov Reward Process (Length=19)



Mini 5-chain  $\alpha = 0.5 \gamma = 1$

$$\text{True } V(S) : \begin{matrix} -\frac{2}{3} \\ A \end{matrix} \quad \begin{matrix} -\frac{1}{3} \\ B \end{matrix} \quad \begin{matrix} 0 \\ C \end{matrix} \quad \begin{matrix} \frac{1}{3} \\ D \end{matrix} \quad \begin{matrix} \frac{2}{3} \\ E \end{matrix}$$

$1\text{-step TD}$	$-1 \leftarrow$ <table border="1" style="border-collapse: collapse; width: 100%;"> <tr> <td style="text-align: center; padding: 5px;"><del>0</del> <del>-0.5</del> -0.25</td><td style="text-align: center; padding: 5px;"><del>0</del> 0</td><td style="text-align: center; padding: 5px;"><del>0</del> 0</td><td style="text-align: center; padding: 5px;"><del>0</del> 0.25</td><td style="text-align: center; padding: 5px;"><del>0</del> <del>0.5</del> 0.75</td></tr> </table> $\rightarrow +1$	<del>0</del> <del>-0.5</del> -0.25	<del>0</del> 0	<del>0</del> 0	<del>0</del> 0.25	<del>0</del> <del>0.5</del> 0.75
<del>0</del> <del>-0.5</del> -0.25	<del>0</del> 0	<del>0</del> 0	<del>0</del> 0.25	<del>0</del> <del>0.5</del> 0.75		

Experiences

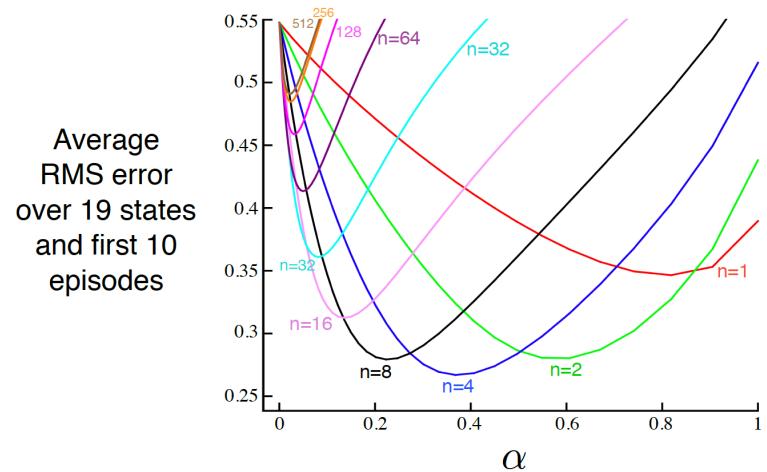
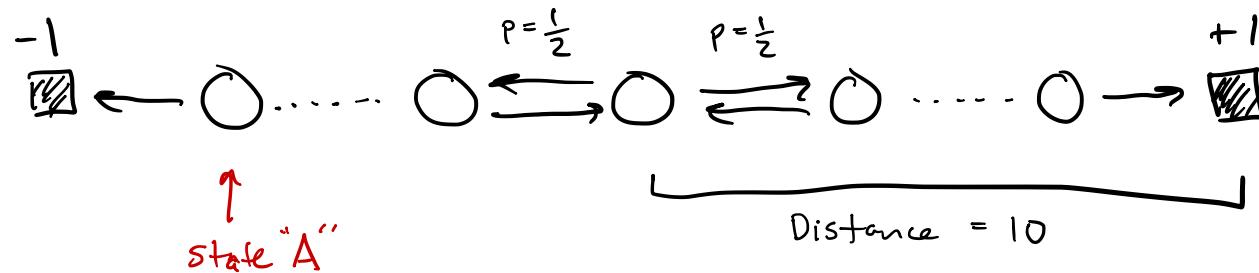
- ✓ B, C, D, E, +1
- ✓ C, B, A, -1
- ✓ A, B, C, D, E, +1

$3\text{-step TD}$	$-1 \leftarrow$ <table border="1" style="border-collapse: collapse; width: 100%;"> <tr> <td style="text-align: center; padding: 5px;"><del>0</del> <del>-0.5</del> 0</td><td style="text-align: center; padding: 5px;"><del>0</del> <del>-0.5</del> 0</td><td style="text-align: center; padding: 5px;"><del>0</del> <del>0.5</del> <del>-0.25</del> 0.375</td><td style="text-align: center; padding: 5px;"><del>0</del> <del>0.5</del> 0.75</td><td style="text-align: center; padding: 5px;"><del>0</del> <del>0.5</del> 0.75</td></tr> </table> $\rightarrow +1$	<del>0</del> <del>-0.5</del> 0	<del>0</del> <del>-0.5</del> 0	<del>0</del> <del>0.5</del> <del>-0.25</del> 0.375	<del>0</del> <del>0.5</del> 0.75	<del>0</del> <del>0.5</del> 0.75
<del>0</del> <del>-0.5</del> 0	<del>0</del> <del>-0.5</del> 0	<del>0</del> <del>0.5</del> <del>-0.25</del> 0.375	<del>0</del> <del>0.5</del> 0.75	<del>0</del> <del>0.5</del> 0.75		

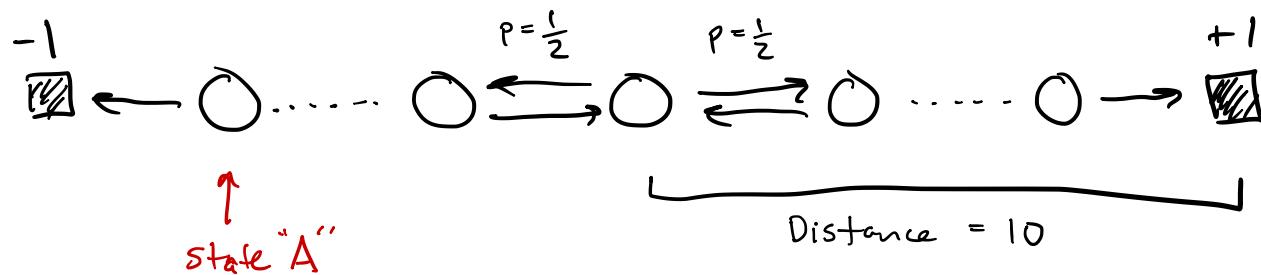
$\infty\text{-step TD}$	$-1 \leftarrow$ <table border="1" style="border-collapse: collapse; width: 100%;"> <tr> <td style="text-align: center; padding: 5px;"><del>0</del> <del>-0.5</del> 0.25</td><td style="text-align: center; padding: 5px;"><del>0</del> <del>-0.5</del> <del>-0.25</del> 0.375</td><td style="text-align: center; padding: 5px;"><del>0</del> <del>0.5</del> <del>-0.25</del> 0.375</td><td style="text-align: center; padding: 5px;"><del>0</del> <del>0.5</del> 0.75</td><td style="text-align: center; padding: 5px;"><del>0</del> <del>0.5</del> 0.75</td></tr> </table> $\rightarrow +1$	<del>0</del> <del>-0.5</del> 0.25	<del>0</del> <del>-0.5</del> <del>-0.25</del> 0.375	<del>0</del> <del>0.5</del> <del>-0.25</del> 0.375	<del>0</del> <del>0.5</del> 0.75	<del>0</del> <del>0.5</del> 0.75
<del>0</del> <del>-0.5</del> 0.25	<del>0</del> <del>-0.5</del> <del>-0.25</del> 0.375	<del>0</del> <del>0.5</del> <del>-0.25</del> 0.375	<del>0</del> <del>0.5</del> 0.75	<del>0</del> <del>0.5</del> 0.75		

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

# Chain Markov Reward Process (Length=19)



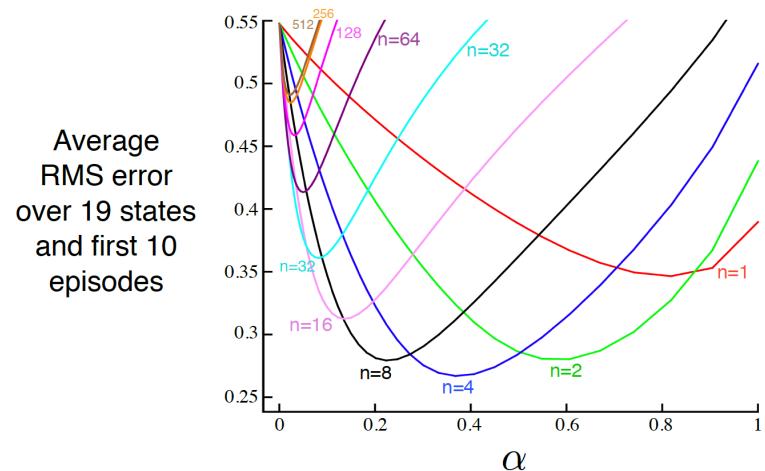
# Chain Markov Reward Process (Length=19)



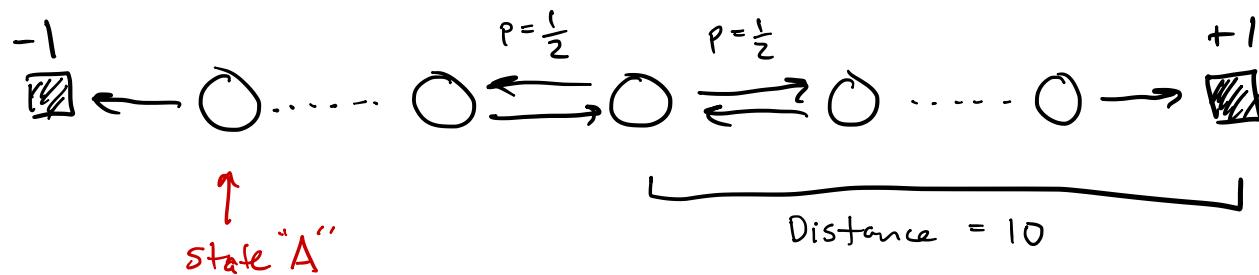
Influence state A value via a +1 reward :

MC: must get to  $+1$  goal at least once after visiting A

TD(0): must visit  $+1$  goal at least once anytime and then wait for value to propagate.



# Chain Markov Reward Process (Length=19)



Influence state A value via a +1 reward :

MC: must get to  $+1$  goal at least once after visiting A

TD( $\alpha$ ): must visit  $+1$  goal at least once anytime and then wait for value to propagate.

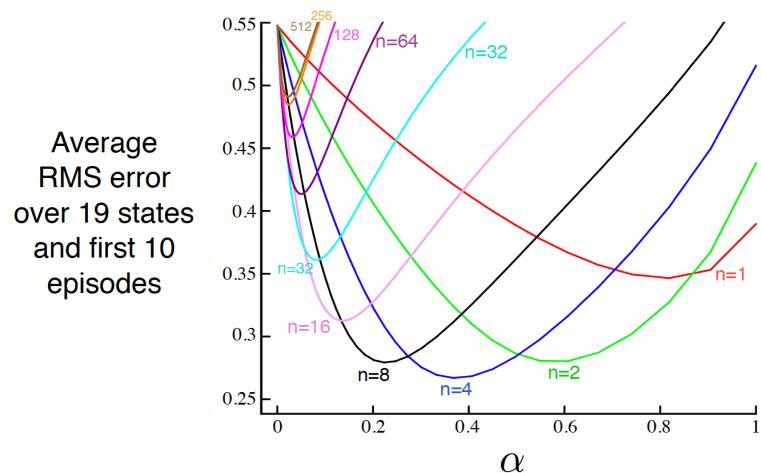
Statistical Properties :

MC: Unbiased, high variance  
Fast propagation of reward (all visited states)

TD( $\alpha$ ): Biased updates, low variance  
Slow propagation of reward (most recent state)

n-step TD : Medium variance  
Medium propagation of reward

Correct balance problem specific !

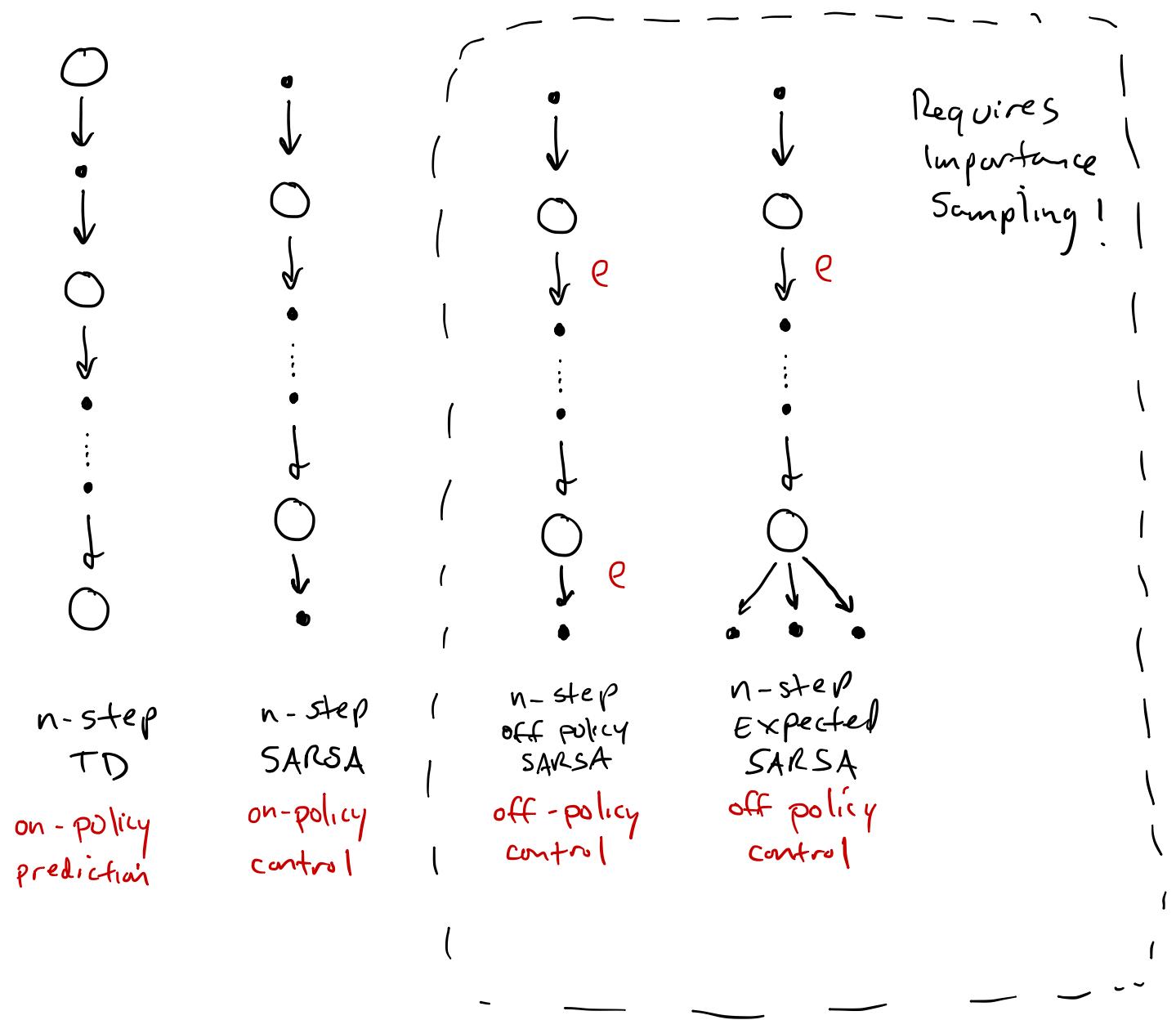


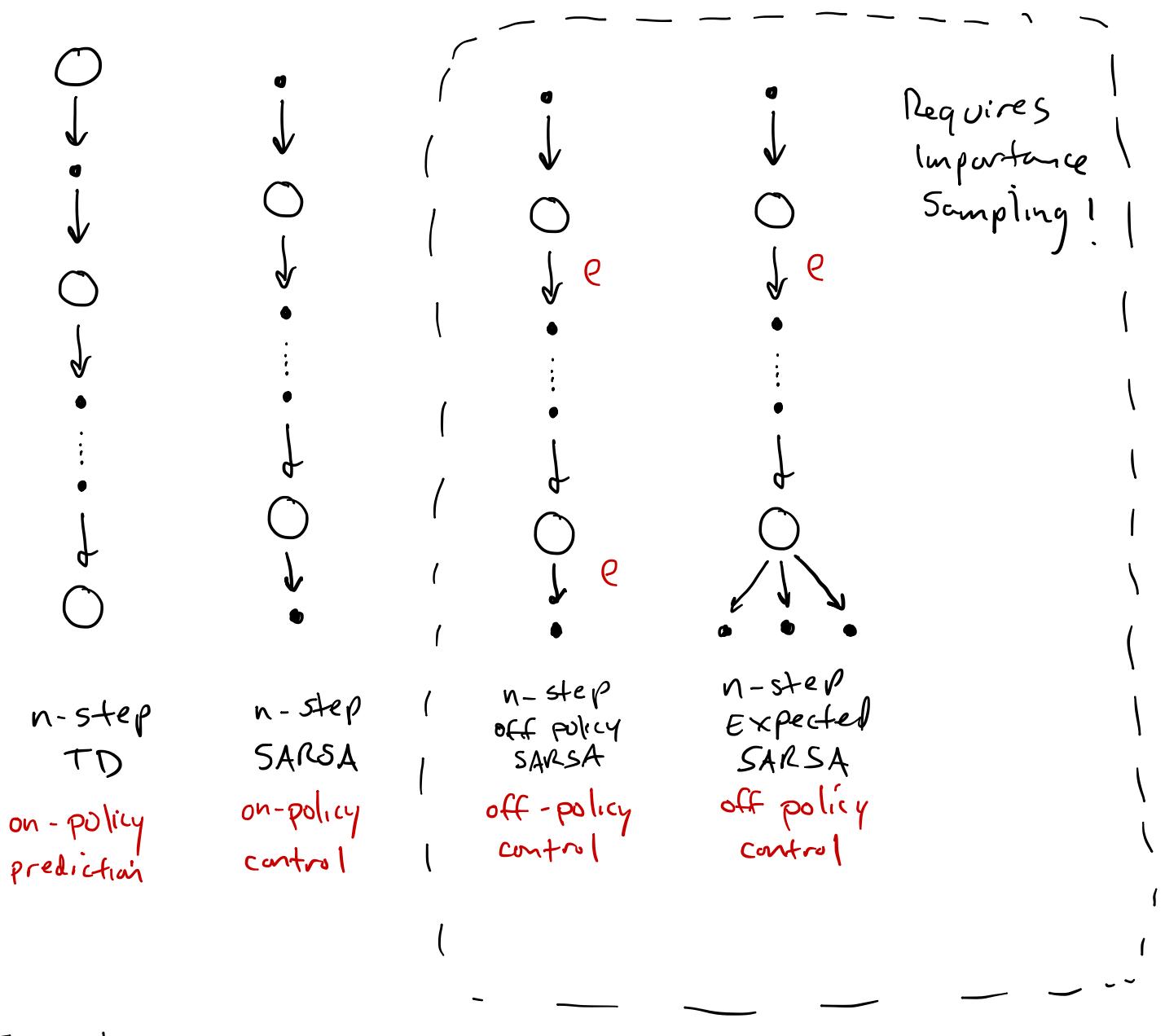


n-step  
TD  
on-policy  
prediction



n-step  
SARSA  
on-policy  
control

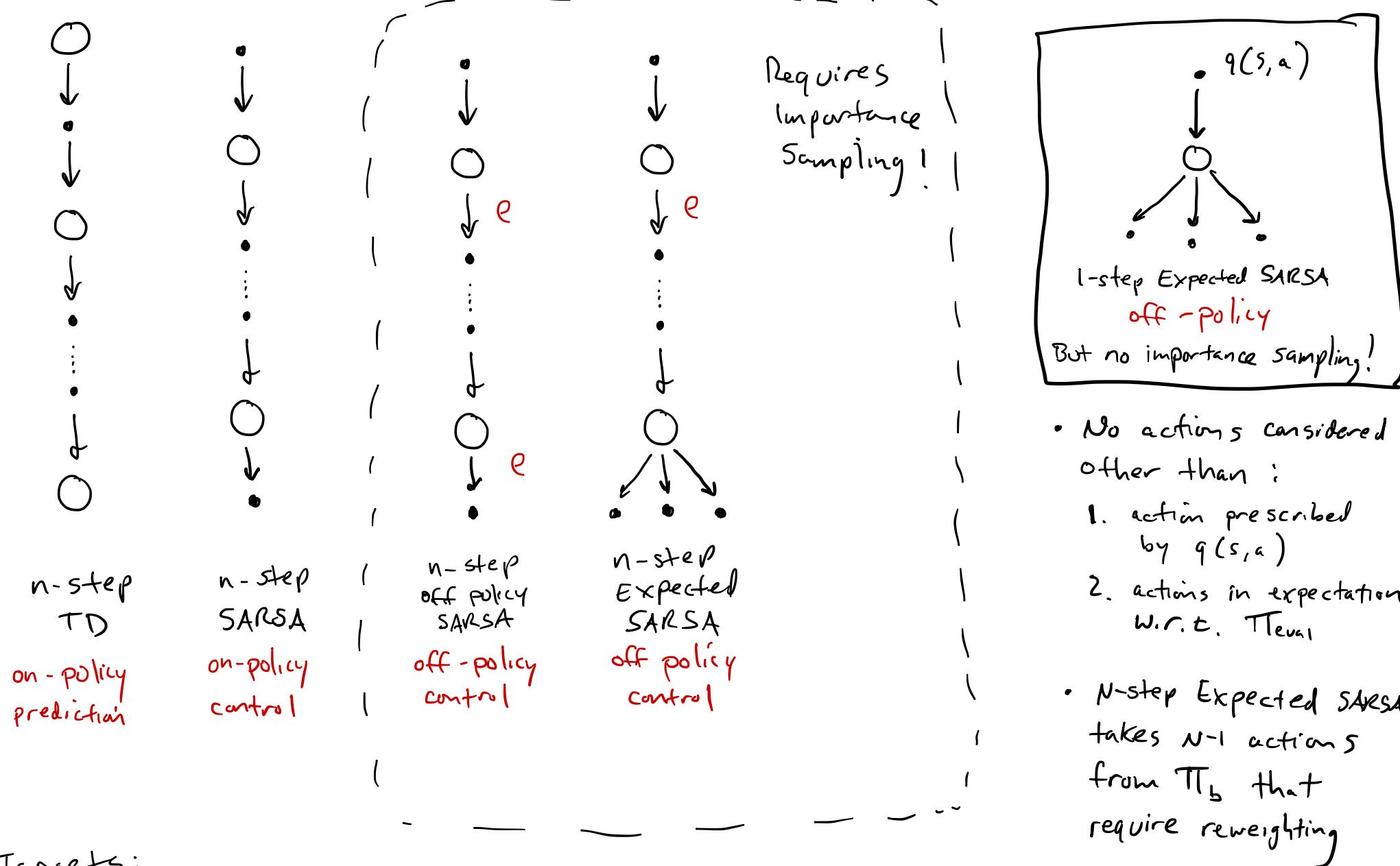




Targets:

$$n\text{-step SARSA} : R_t + \gamma R_{t+1} + \dots + \gamma^{n-1} R_{t+n-1} + \gamma^n Q(s_{t+n}, a_{t+n})$$

$$n\text{-step off-policy SARSA} : E_{t:t+n-1} \left[ R_t + \gamma R_{t+1} + \dots + \gamma^{n-1} R_{t+n-1} + \gamma^n Q(s_{t+n}, a_{t+n}) \right]$$

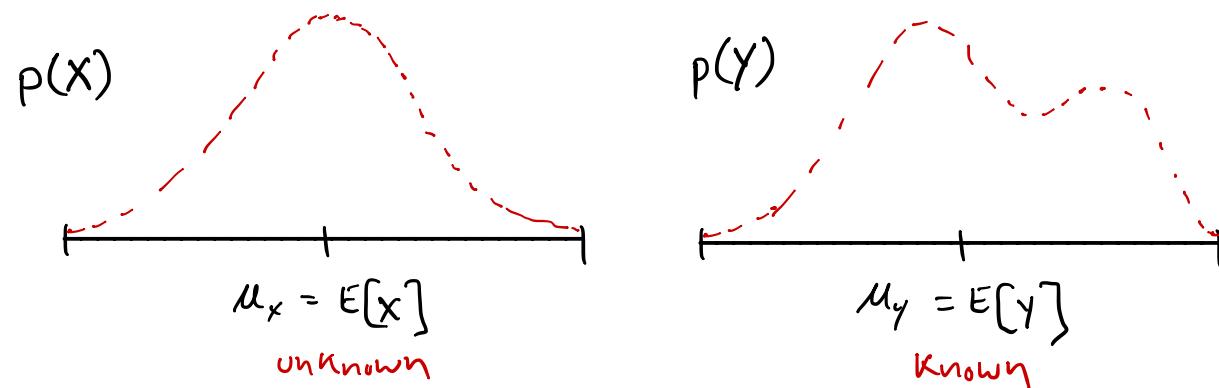


Targets:

$$\text{n-step SARSA} : R_t + \gamma R_{t+1} + \dots + \gamma^{n-1} R_{t+n-1} + \gamma^n Q(s_{t+n}, a_{t+n})$$

$$\text{n-step off-policy SARSA} : E_{t:t+n-1} \left[ R_t + \gamma R_{t+1} + \dots + \gamma^{n-1} R_{t+n-1} + \gamma^n Q(s_{t+n}, a_{t+n}) \right]$$

## Control Variates

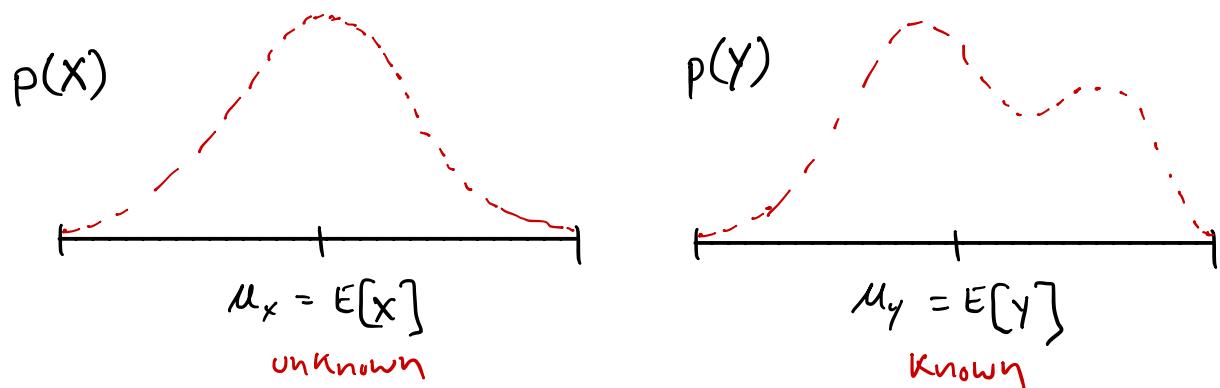


Assumption:  $X$  is positively correlated with  $Y$

i.e. when  $X > E[X]$ , it is likely that  $y > E[y]$

ex: rain and traffic

## Control Variates



Assumption:  $X$  is positively correlated with  $Y$   
i.e. when  $X > E[X]$ , it is likely that  $y > E[y]$   
ex: rain and traffic

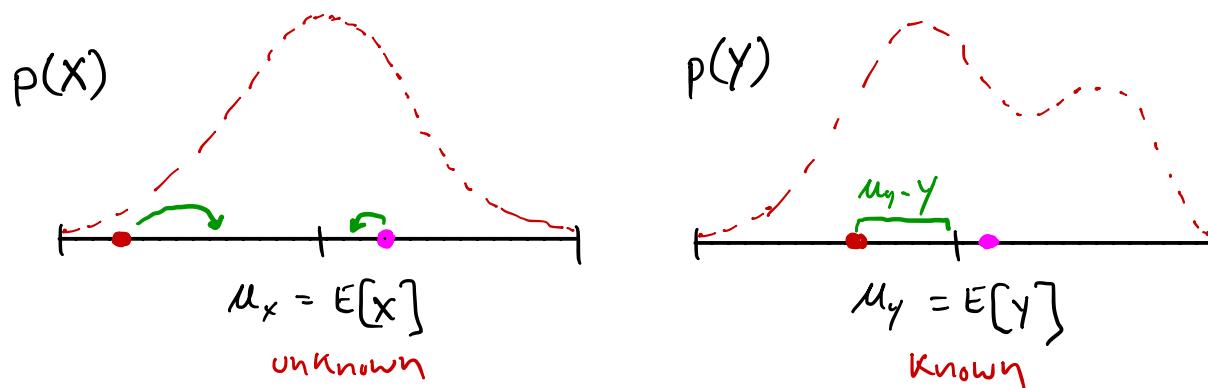
Sample  $x \sim p(x)$  and  $y \sim p(y)$

Consider two different estimators:

$$\hat{\mu}_x = x$$

$$\hat{\mu}_x = x + [\mu_y - y]$$

## Control Variates



Assumption:  $X$  is positively correlated with  $Y$   
i.e. when  $X > E[X]$ , it is likely that  $y > E[Y]$   
ex: rain and traffic

Sample  $x \sim p(x)$  and  $y \sim p(y)$

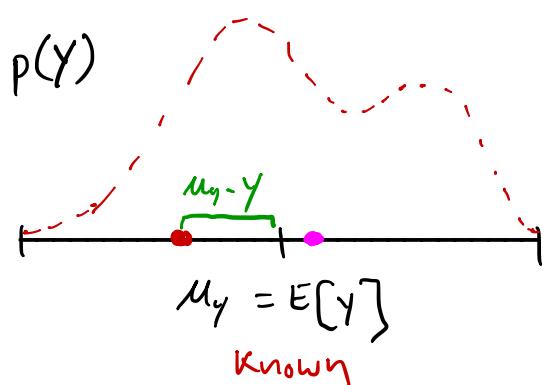
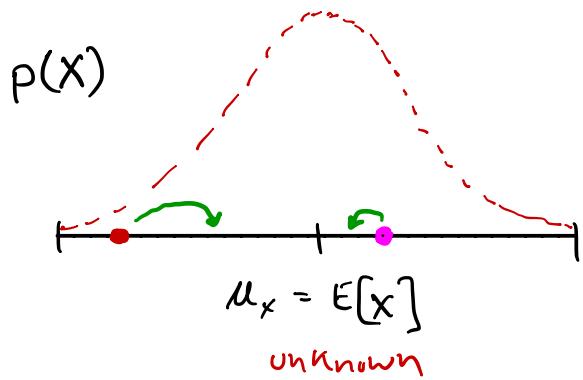
Consider two different estimators:

$$\hat{\mu}_x = X$$

$$\hat{\mu}_x = X + [\mu_y - Y]$$

Expectation zero!

## Control Variates



Assumption:  $X$  is positively correlated with  $Y$   
 i.e. when  $X > E[X]$ , it is likely that  $y > E[Y]$   
 ex: rain and traffic

Sample  $x \sim p(x)$  and  $y \sim p(y)$

Consider two different estimators:

$$\hat{\mu}_x = x$$

$$\hat{\mu}_x = x + [\mu_y - y]$$

Expectation zero!

## n-Step PDIS with CV

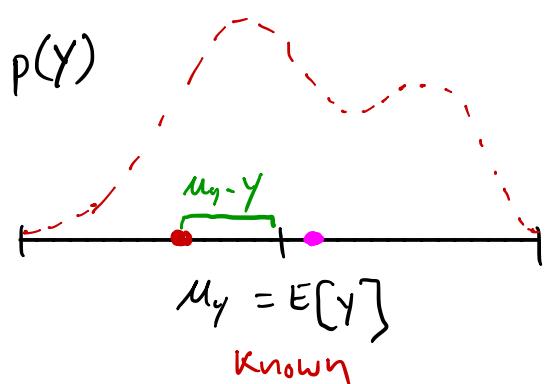
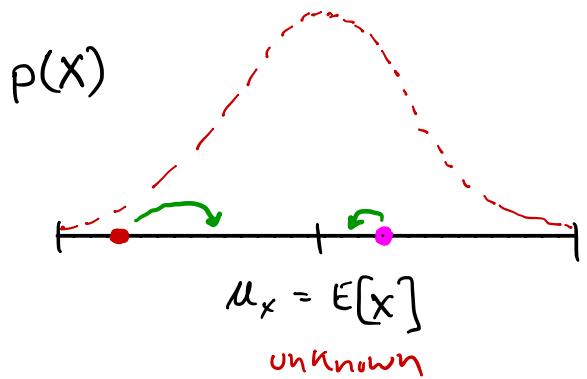
$$G_{t:h} = \ell_t (R_{t+1} + \gamma G_{t+1:h}) + ((1-\ell_t) V_{h-1}(s_t)) \quad CV$$

$$E[\ell_t] = 1, \text{ thus:}$$

$$E[1-\ell_t] = 0$$

$$\rightarrow E[(1-\ell_t) V_{h-1}(s_t)] = 0$$

## Control Variates



Assumption:  $X$  is positively correlated with  $Y$   
i.e. when  $X > E[X]$ , it is likely that  $y > E[Y]$

ex: rain and traffic

sample  $x \sim p(x)$  and  $y \sim p(y)$

Consider two different estimators:

$$\hat{\mu}_x = x$$

$$\hat{\mu}_x = x + [\mu_y - Y]$$

Expectation zero!

## n-Step PDIS with CV

$$G_{t:h} = \ell_t (R_{t+1} + \gamma G_{t+1:h}) + ((1-\ell_t) V_{h-1}(s_t)) \quad CV$$

$$E[\ell_t] = 1, \text{ thus:}$$

$$E[1-\ell_t] = 0$$

$$\rightarrow E[(1-\ell_t) V_{h-1}(s_t)] = 0$$

Compare to without CV:

$$G_{t:h} = \ell_t [R_{t+1} + \gamma G_{t+1:h}]$$

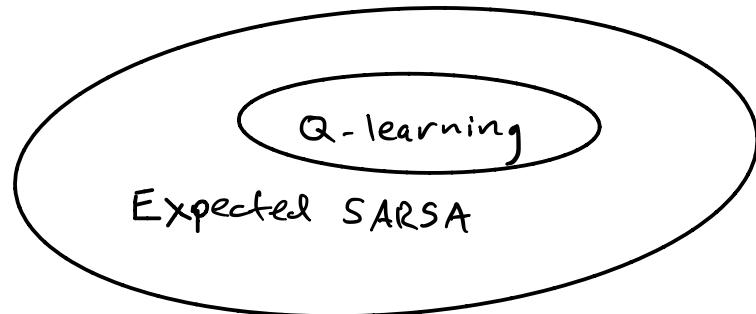
if  $\ell_t = 0$ :

$$CV: G_{t:h} = V_{h-1}(s_t)$$

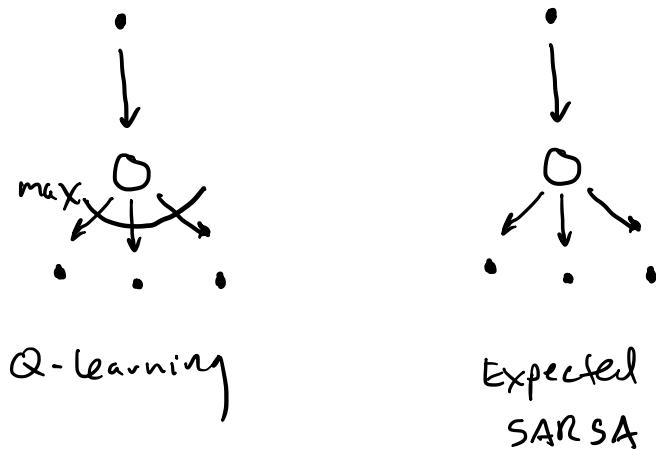
$$No\ CV: G_{t:h} = 0$$

"how much better or worse than average is this sample?"

## one step off policy

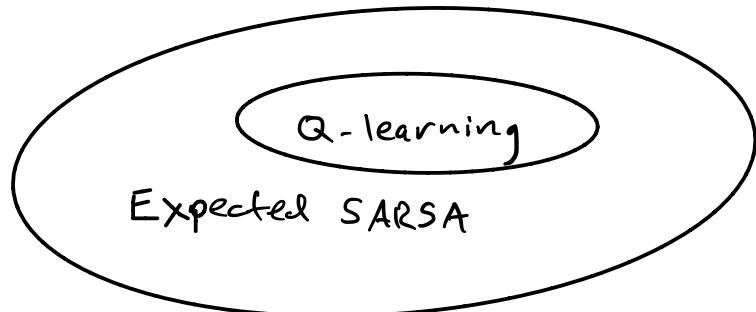


Same when  $\pi_{\text{eval}} = \pi_\star$

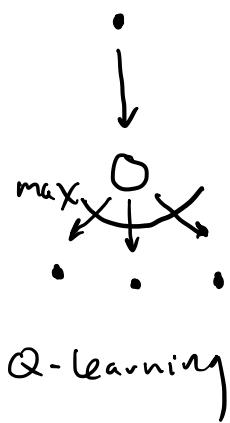


- Both off policy
- No importance sampling!

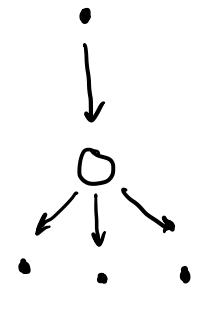
## one step off policy



Same when  $\pi_{\text{eval}} = \pi_\star$



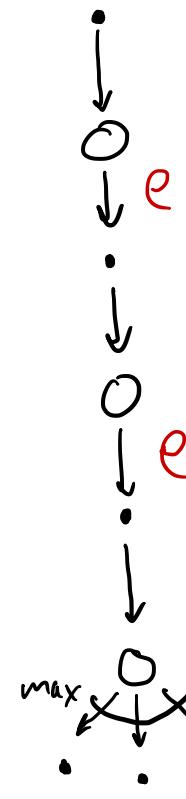
Q-learning



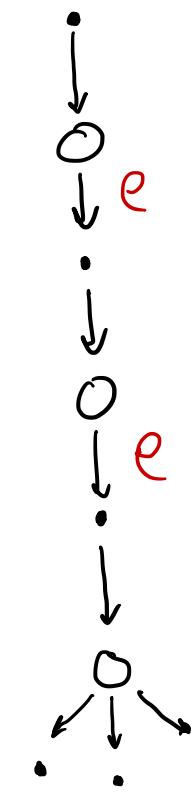
Expected  
SARSA

- Both off policy
- No importance sampling!

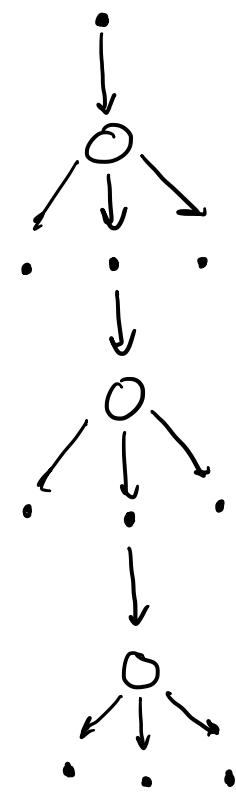
## n-step off policy



n-step  
Q-learning



n-step  
Expected  
SARSA



Tree  
Backups

Tree Backups Target:

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