

## Off-policy vs On-policy

- Let's introduce two types of policy
- 1.) Behaviour policy  
 $b(a|s)$
  - 2.) Target policy  
 $\pi(a|s)$

- We collect data with  $b\text{CA}(s)$  and we use data to estimate  $Q^\pi, V^\pi$  where  $\pi$  is the target,
- When  $b\text{CA}(s)$  is  $\pi(a|s)$  we say the algorithm is "on-policy"

	Pros	Cons
on-policy	<ul style="list-style-type: none"> <li>- generally simpler</li> <li>- Many on-policy algo are formulated for continuous action domains.</li> </ul>	<p>cannot do</p>
off-policy	<ul style="list-style-type: none"> <li>→ General &amp; powerful</li> <li>→ Better Exploration</li> <li>→ Can learn from observation</li> <li>→ Reuse Experience</li> <li>→ Decouple data-collection + learning</li> </ul>	<p>→ data has greater variance</p> <p>→ Harder to be formulated</p>

# Q-learning vs SARSA

Q-learning

estimate  $Q^*$

$$Q^*(s, a) \leftarrow Q(s, a) + \alpha [R + \max_a Q(s', a) - Q(s, a)]$$

$\rightarrow S, A, R, S'$

$\underbrace{b(A|s)}$

$\rightarrow \pi^*$  is

the target policy

## SARSA

$$Q(s, a) \leftarrow Q(s, a) + \alpha [R + \gamma Q(s', a') - Q(s, a)]$$

$A'$  is chosen from  $Q$

$\therefore \text{SARSA} = \text{arg}\max_a Q(s, a)$

# Off-policy prediction

with Importance

Sampling

- MC-control is off-policy

$$\rightarrow Q^{\pi}(s, a)$$

because we use  $\pi$  to collect data.

- Can we estimate  $Q^{\pi}(s, a)$  while collecting data from  $b(s|s)$

$$Q^\pi(s, a) = \mathbb{E}^\pi [G | s, a]$$

MC-prediction  $\Rightarrow \mathbb{E}^\pi [g_i]$

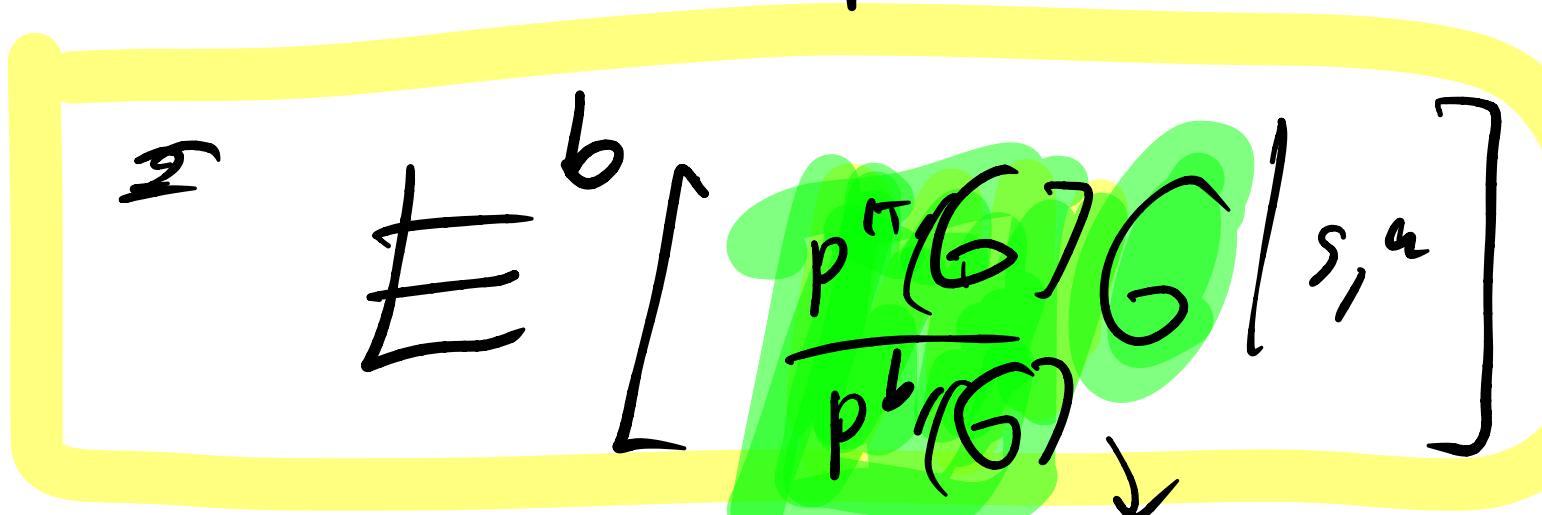
$$\mathbb{E}^\pi [g_i] = \frac{1}{n} \sum_n g_i \approx \frac{1}{n} \sum_n g_i$$

$$Q^\pi(s, a) = \sum_{g_i} p^\pi(g_i) g_i$$

$$= \sum_{g_i} p^\pi(g_i) \frac{p^b(g_i)}{p^b(g_i)}$$

$$= \sum_{g_i} p^b(g_i) \left[ \frac{p^\pi(g_i)}{p^b(g_i)} \right] g_i$$

$$= \sum_{g_i} p^b(g_i) \left[ \frac{p^\pi(g_i)}{p^b(g_i)} \cdot g_i \right]$$



$$G^\pi(s,a) = E^\pi G(s,a)$$

$$\rightarrow \beta = \frac{p^\pi}{p^b} \quad \text{Importance weight}$$

$$+ Q^\pi(s,a) = E^\pi G(s,a)$$

$$\rho = \frac{p^\pi(g)}{p^b(g)}$$

trajectory

$$= \frac{p(g(\tau)) p^\pi(\tau)}{p(g(\tau)) p^b(\tau)}$$

$$\mathcal{T} = \{S_t, A_t, S_{t+1}, A_{t+1}, \dots\}$$

$$p(\mathcal{T}) = p(S_t) \cdot \prod p(A_t | S_t) p(S_{t+1} | S_t, A_t) \\ \cdot p(A_{t+1} | S_{t+1}) p(S_{t+2} | S_{t+1}, A_{t+1})$$

$$P(\tau) = P(S_t) \prod_{k=t}^{T-1} \pi(A_k | S_k) \cdot$$

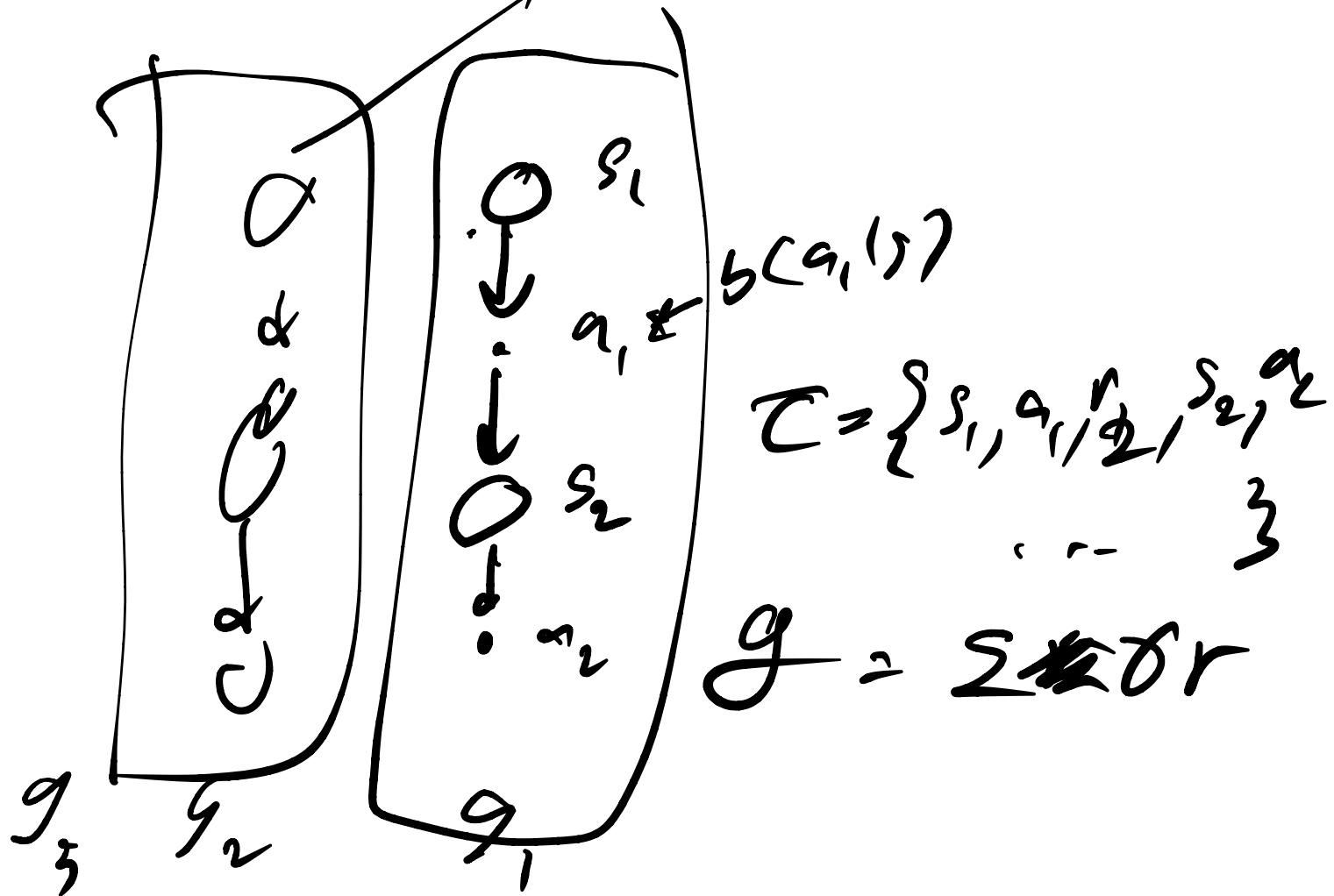
$k=t$

$$P(S_{k+1} | S_k, A_k)$$

$$\rho = \frac{P(\tau)}{P^b(\tau)} = \frac{P(S_t) \prod_{k=t}^{T-1} \pi(A_k | S_k) P(S_{k+1} | S_k, A_k)}{P(S_t) \prod_{k=t}^{T-1} b(A_k | S_k) P(S_{k+1} | S_k, A_k)}$$

$$\boxed{\rho = \prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{b(A_k | S_k)}}$$

$$Q(s,a) = \mathbb{E}^{\pi}[G \cdot \left[ \prod_{n=0}^{T-1} \frac{\pi(a_n | s_n)}{b(a_n | s_n)} \right]]$$



\* We can use the same idea  
for n-step SARSA

$$\begin{aligned}
 Q^\pi(s, a) &= E^\pi \left[ R_t + \gamma R_{t+1} + \dots + \gamma^{n-1} R_{t+n} \right] \\
 &= E^b \left[ \rho \left( \dots \right) \right] \\
 \rho &= \prod_{k=t}^{t+n} \frac{p(A_k | S_k)}{b(A_k | S_k)}
 \end{aligned}$$