$S_1$

$S_2$

$S_1$
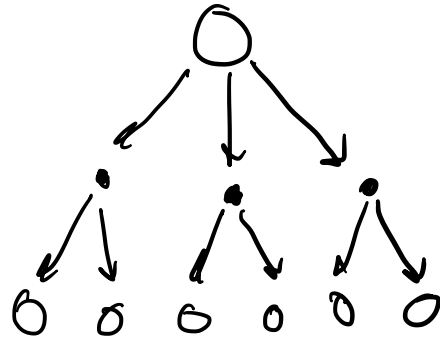
Bellman eqn. for $V_\pi$,
Dynamic programming

Monte Carlo
est. of $V_\pi$

Sample of return:

$$G_{s_t} = \sum_{i=t}^{T-1} r_i \qquad \Rightarrow \qquad V(s) = E[G_{s_t}]$$

Monte Carlo
est. of $V_\pi$

Bellman eqn. for $V_\pi$,
Dynamic programming

MC:
- only sampled transitions
- All the way to end of episode
- no bootstrapping

DP:
- All possible transitions
- only one step
- bootstrapping

Sample of return:
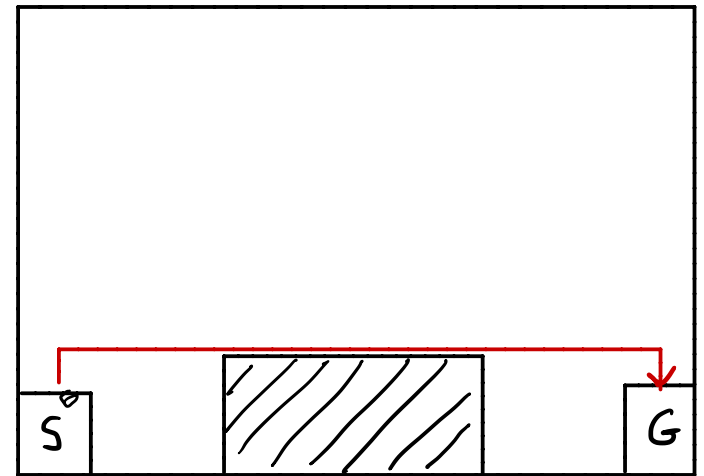
$$G_{s_t} = \sum_{i=t}^{T-1} r_i \quad \Rightarrow \quad V(s) = E[G_{s_t}]$$

$G_i^{s, \pi}$: $i^{th}$ return starting from state $s$, collected from policy $\pi$

On policy : $V_\pi(s) = \dfrac{1}{N} \sum\limits_{i=1}^{N} G_i^{s, \pi}$

Prediction

$\qquad\qquad\; Q_\pi(s,a) = \dfrac{1}{N} \sum\limits_{i=1}^{N} G_i^{s, a, \pi}$

$G_i^{s,\pi}$ : $i^{th}$ return starting from state $s$, collected from policy $\pi$



On policy : $V_\pi(s) = \frac{1}{N} \sum\limits_{i=1}^{N} G_i^{s,\pi}$
Prediction

$\qquad\qquad Q_\pi(s,a) = \frac{1}{N} \sum\limits_{i=1}^{N} G_i^{s,a,\pi}$

$G_i^{s,\pi}$: $i^{th}$ return starting from state $s$, collected from policy $\pi$



On policy: $V_\pi(s) = \frac{1}{N} \sum_{i=1}^{N} G_i^{s,\pi}$
Prediction

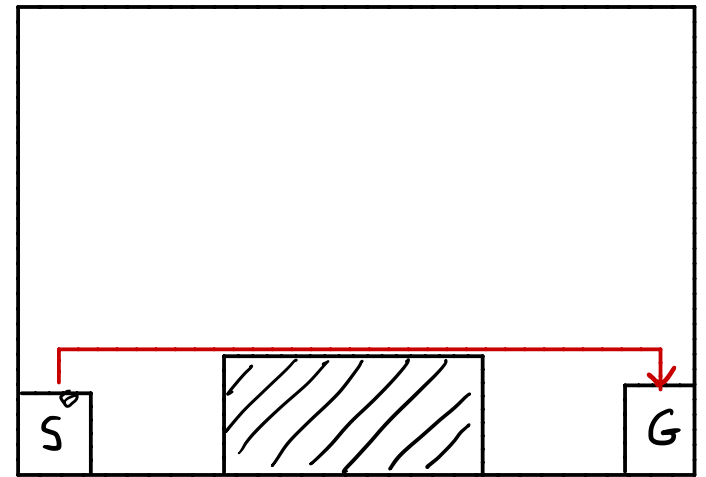$Q_\pi(s,a) = \frac{1}{N} \sum_{i=1}^{N} G_i^{s,a,\pi}$

Off policy: $V_{\pi'}(s) = \frac{1}{N} \sum_{i=1}^{N} G_i^{s,\pi} \cdot \rho_i$
Prediction

$Q_{\pi'}(s,a) = \frac{1}{N} \sum_{i=1}^{N} G_i^{s,a,\pi} \cdot \rho_i$

where $\rho_i = \prod_{k=t_i}^{T_i-1} \frac{\pi'(A_k|S_k)}{\pi(A_k|S_k)}$

$G_i^{s,\pi}$: $i^{th}$ return starting from state $s$, collected from policy $\pi$



On policy: $V_\pi(s) = \frac{1}{N} \sum_{i=1}^{N} G_i^{s,\pi}$
prediction

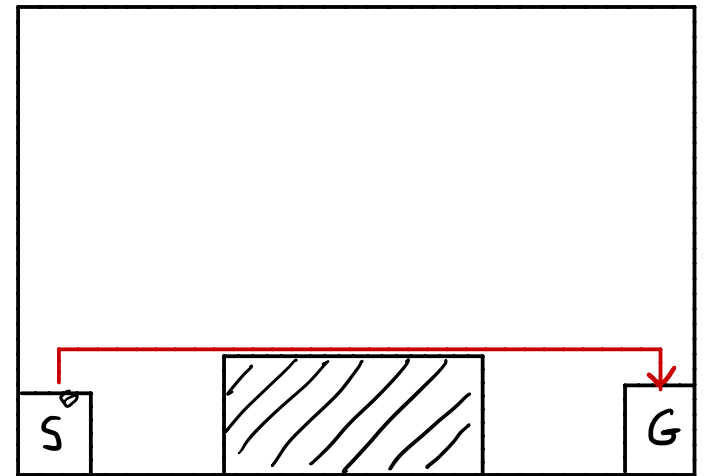$Q_\pi(s,a) = \frac{1}{N} \sum_{i=1}^{N} G_i^{s,a,\pi}$

Off policy: $V_{\pi'}(s) = \frac{1}{N} \sum_{i=1}^{N} G_i^{s,\pi} \cdot \rho_i$
prediction

$Q_{\pi'}(s,a) = \frac{1}{N} \sum_{i=1}^{N} G_i^{s,a,\pi} \cdot \rho_i$

where $\rho_i = \prod_{k=t_i}^{T_i-1} \frac{\pi'(A_k|S_k)}{\pi(A_k|S_k)}$

Consider: On-policy control vs. off-policy control w/ $\epsilon$-greedy exploration

What are $\underbrace{V_* \text{ and } \pi_*}_{\text{off-policy}}$ vs. $\underbrace{\tilde{V}_* \text{ and } \tilde{\pi}_*}_{\text{on-policy}}$ ?

Safe off policy evaluation:

Return probabilistic lower bound $V_\pi^{lb}$ such that:

$$V_\pi > V_\pi^{lb} \quad \text{with prob. } 1-\delta \qquad \text{Given: } \pi, \delta, \text{ data from } \pi_b$$

without ever running policy $\pi$!

Safe off policy evaluation:

Return probabilistic lower bound $V_\pi^{lb}$ such that:

$$V_\pi > V_\pi^{lb} \quad \text{with prob. } 1-\delta \qquad \text{Given: } \pi, \delta, \text{ data from } \pi_b$$

without ever running policy $\pi$!

Confidence bounds: Chernoff - Hoeffding inequality

with probability at least $1-\delta$:

$$\mu \geq \frac{1}{n} \sum_{i=1}^{n} X_i - b \sqrt{\frac{\log(1/\delta)}{2n}} \qquad \text{for } 0 \leq X_i \leq b$$

Safe off policy evaluation:

Return probabilistic lower bound $V_\pi^{lb}$ such that:

$$V_\pi > V_\pi^{lb} \quad \text{with prob. } 1-\delta \qquad \text{Given: } \pi, \delta, \text{ data from } \pi_b$$

without ever running policy $\pi$!

Confidence bounds: Chernoff - Hoeffding inequality

with probability at least $1-\delta$:

$$\mu \geq \frac{1}{n} \sum_{i=1}^{n} X_i - b \sqrt{\frac{\log(1/\delta)}{2n}} \qquad \text{for } 0 \leq X_i \leq b$$

$$V_\pi \geq \frac{1}{n} \sum_{i=1}^{n} G_i^{\pi_b} \cdot \rho_i^{\pi, \pi_b} - G_{max} \sqrt{\frac{\log(1/\delta)}{2n}} \qquad \text{for } 0 \leq G_i \leq G_{max}$$

Given returns $G_1 \cdots G_n$ from policy $\pi_b$ AND $\rho_i = \prod\limits_{k=t_i}^{T_i-1} \dfrac{\pi(A_k|S_k)}{\pi_b(A_k|S_k)}$ THEN: $V_\pi(s) =$

Given returns $G_1 \cdots G_n$ from policy $\Pi_b$ AND $\rho_i = \prod_{k=t_i}^{T_i - 1} \frac{\Pi(A_k | S_k)}{\Pi_b(A_k | S_k)}$ THEN: $V_\Pi(s) =$

OIS:

$$\frac{1}{n} \sum_{i=1}^{n} G_i \rho_i$$

Given returns $G_1 \cdots G_n$ from policy $\pi_b$  <span style="color:red">AND</span>  $\rho_i = \prod_{k=t_i}^{T_i-1} \frac{\pi(A_k|S_k)}{\pi_b(A_k|S_k)}$  <span style="color:red">THEN:</span>  $V_\pi(s) =$

OIS:

$$\frac{1}{n} \sum_{i=1}^{n} G_i \rho_i$$

WIS:

$$\frac{1}{n} \sum_{i=1}^{n} \frac{G_i \rho_i}{\sum_{j=1}^{n} \rho_j}$$

Given returns $G_1 \cdots G_n$ from policy $\pi_b$ **AND** $\rho_i = \prod\limits_{k=t_i}^{T_i-1} \dfrac{\pi(A_k|S_k)}{\pi_b(A_k|S_k)}$ **THEN:** $V_\pi(s) =$

OIS:

$$\frac{1}{n} \sum_{i=1}^{n} G_i \rho_i$$

WIS:

$$\sum_{i=1}^{n} \frac{G_i \rho_i}{\sum_{j=1}^{n} \rho_j}$$

PDIS:

$$\frac{1}{n} \sum \tilde{G}_i, \text{ where:}$$

$$\tilde{G}_i = \rho_{1:1} R_1 + \gamma \rho_{1:2} R_2 + \ldots + \gamma^{n-1} \rho_{1:t} R_t$$

and

$$\rho_{a:b} = \prod_{k=a}^{b} \frac{\pi(A_k|S_k)}{\pi_b(A_k|S_k)}$$