

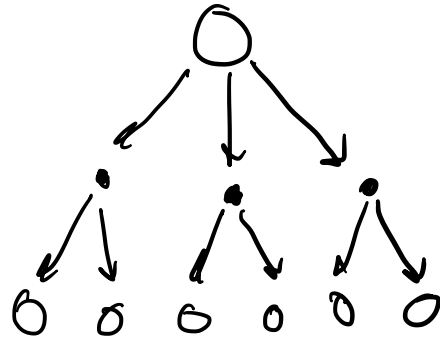
Bellman eqn. for V_π ,
Dynamic programming

Monte Carlo
est. of V_π

Sample of return:

$$G_{s_t} = \sum_{i=t}^{T-1} r_i$$

$$\Rightarrow v(s) = E[G_{s_t}]$$



Bellman eqn. for V_π ,
Dynamic programming

MC:

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

DP:

- All possible transitions
- only one step
- bootstrapping

Monte Carlo
est. of V_π

Sample of return:

$$G_{s_t} = \sum_{i=t}^{T-1} r_i$$

$$\Rightarrow v(s) = E[G_{s_t}]$$

First-visit MC prediction, for estimating $V \approx v_\pi$

Input: a policy π to be evaluated

Initialize:

$V(s) \in \mathbb{R}$, arbitrarily, for all $s \in \mathcal{S}$

$Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

Loop forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless S_t appears in S_0, S_1, \dots, S_{t-1} :

Append G to $Returns(S_t)$

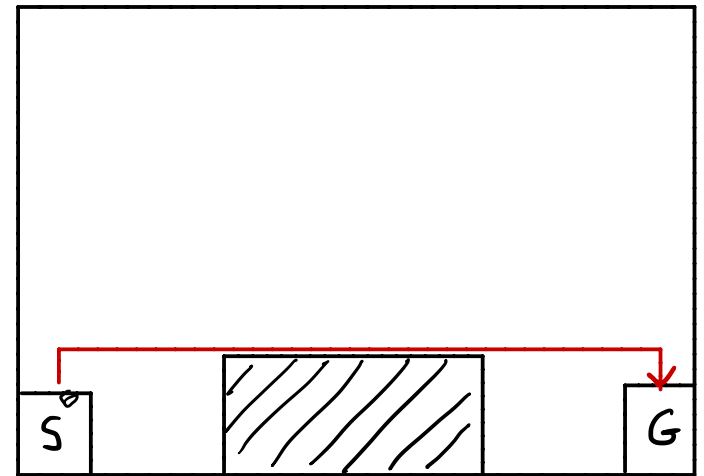
$V(S_t) \leftarrow \text{average}(Returns(S_t))$

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

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

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

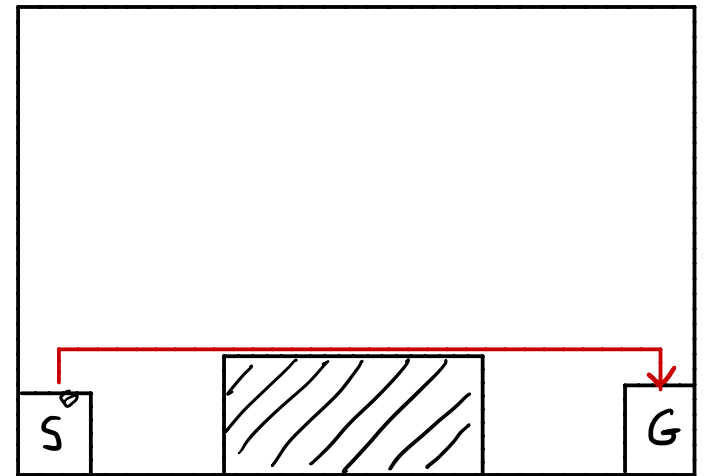
$G_i^{s,\pi}$: i^{th} return starting from state s , collected from policy π



On policy prediction

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

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



On policy prediction

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

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

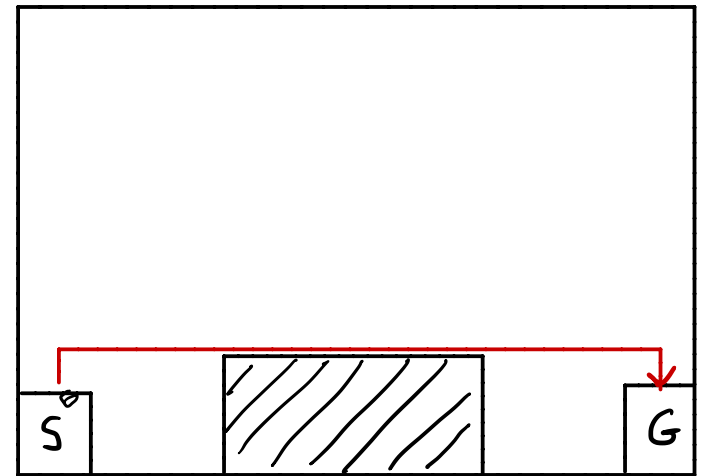
Off policy prediction

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

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

where $\rho_i = \frac{\prod_{k=t_i}^{T-1} \pi'(A_k | S_k)}{\prod_{k=t_i}^{T-1} \pi(A_k | S_k)}$

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



On policy prediction:

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

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

Off policy prediction:

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

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

where $\rho_i = \prod_{k=t_i}^{T-1} \frac{\pi'(A_k | S_k)}{\pi(A_k | S_k)}$

Consider: On-policy control vs. off-policy control w/ ϵ -greedy exploration
 What are $\underbrace{V_{\star} \text{ and } \pi_{\star}}_{\text{off-policy}}$ vs. $\underbrace{\tilde{V}_{\star} \text{ and } \tilde{\pi}_{\star}}_{\text{on-policy}}$?

Safe off policy evaluation:

Return probabilistic lower bound V_{π}^{lb} such that:

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

without ever running policy π !

Safe off policy evaluation:

Return probabilistic lower bound V_{π}^{lb} such that:

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

without ever running policy π !

Confidence bounds: Chernoff-Hoeffding inequality

with probability at least $1 - \delta$:

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

Safe off policy evaluation:

Return probabilistic lower bound V_{π}^{lb} such that:

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

without ever running policy π !

Confidence bounds: Chernoff-Hoeffding inequality

with probability at least $1 - \delta$:

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



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

Given returns $G_1 \dots G_n$ AND
from policy π_b

$$p_i = \prod_{k=t_i}^{T_i-1} \frac{\pi(A_k | S_k)}{\pi_b(A_k | S_k)} \quad \text{THEN: } V_{\pi}(s) =$$

Given returns $G_1 \dots G_n$ AND $P_i = \prod_{k=t_i}^{T_i-1} \frac{\pi(A_k | S_k)}{\pi_b(A_k | S_k)}$ THEN: $V_{\pi}(s) =$

OIS :

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

Given returns $G_1 \dots G_n$
from policy π_b

AND

$$p_i = \prod_{k=t_i}^{T_i-1} \frac{\pi(A_k | S_k)}{\pi_b(A_k | S_k)} \quad \text{THEN: } V_{\pi}(s) =$$

OIS:

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

WIS:

$$\sum_{i=1}^n \frac{G_i p_i}{\sum_{j=1}^n p_j}$$

Given returns $G_1 \dots G_n$ from policy π_b AND $p_i = \prod_{k=t_i}^{T_i-1} \frac{\pi(A_k|S_k)}{\pi_b(A_k|S_k)}$ THEN: $V_{\pi}(s) =$

OIS:

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

WIS:

$$\sum_{i=1}^n \frac{G_i p_i}{\sum_{j=1}^n p_j}$$

POIS:

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

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

and

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