

① Monte Carlo

Recap: - MDP [extensions: multi agent, imperfect information]

- Bellman equations
↳ VI / PI

Tabular assumption:

- X. S states is FINITE and TABULATED
- ✓. A actions is FINITE and TABULATED
- X. $\Delta: S \times A \rightarrow \text{Dist}(S \times \mathbb{R})$ is KNOWN

Now: Statistical framework

$$s, a \mapsto \text{Sample}(\Delta(s, a)) \mapsto (s', r)$$

$\sigma: S \rightarrow A$ (partial function)

Evaluation (σ) task

↳ $\text{val}_\sigma(s)$ for s all states σ is defined on

$$\text{val}_\sigma(s_0) = \mathbb{E}_{\sigma, s_0} \left[\underbrace{r_0 + \gamma r_1 + \gamma^2 r_2 + \dots}_{\text{total return}} \right]$$

↳ $q_\sigma(s_0, a)$ for s all states
 $a \in A$ action

Most natural idea: Monte Carlo

Sample a number N of trajectories from s_0 using σ

and $\widehat{val}_\sigma(s_0) = \frac{1}{N} \sum_{i=1}^N G_i$

i th trajectory:

$$\pi_i = s_0^i, a_0^i, r_0^i, s_1^i, a_1^i, r_1^i, s_2^i, \dots$$

$$G_i = r_0^i + \gamma r_1^i + \gamma^2 r_2^i + \dots$$

↑ total reward.

truncate
to length
 K

Works for s_0 . What about

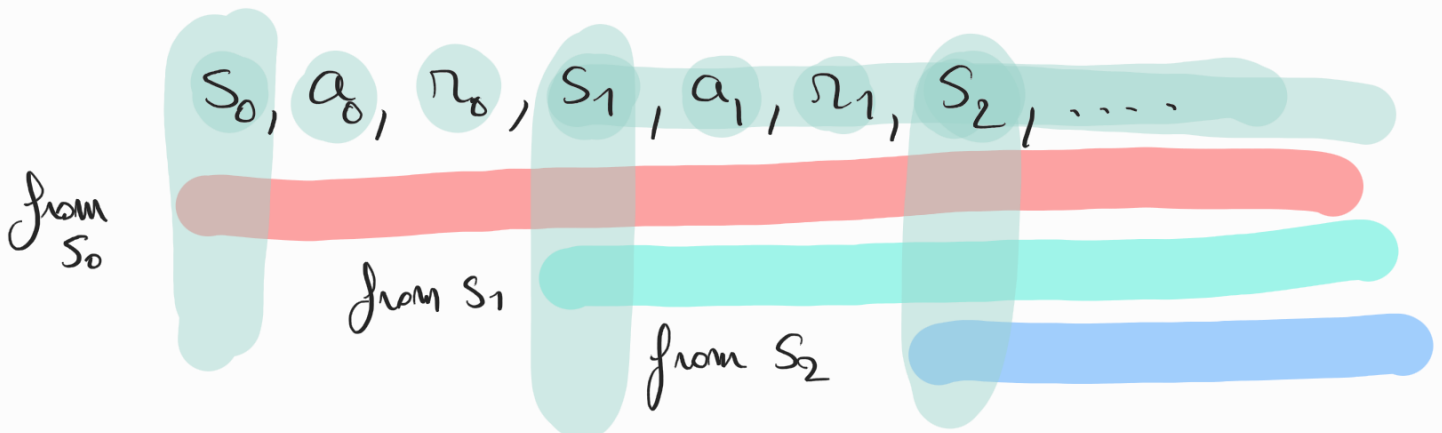
$val_\sigma(s)$ for all $s \in S$ that σ "knows"

$$val_\sigma(s) = \sum_{s', r} \Delta(s, \sigma(s))(s', r) (r + \gamma val_\sigma(s'))$$

↳ dynamic programming is not available
in our statistical setting!

idea:

one sample \rightarrow many trajectories!



How to use them?

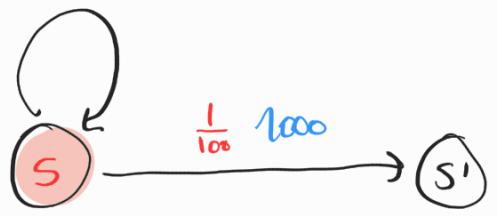
$s \in S$

$$val_{\pi}(s) = \frac{99}{100} \cdot (1 + val_{\pi}(s)) + \frac{1}{100} \cdot 1000$$

- ① average over all trajectories which start in the first occurrence of s
- ② average over all trajectories which start in any occurrence of s
- ③ average over all trajectories which start in the last occurrence of s

$$\frac{99}{100} \cdot 1$$

$$\gamma = 1$$



probability
reward

$$\hat{val}_{\pi}(s) = \begin{cases} \text{last visit} & 1000 \quad \times \\ \text{all visits} & 1099 \quad \checkmark \\ \text{first visit} & 1099 \quad \checkmark \end{cases}$$



actual value

$$val_{\pi}(s) = \frac{99}{100} \cdot (1 + val_{\pi}(s)) + \frac{1}{100} \cdot 1000$$

$$100v = 99(1+v) + 1000$$

$$v = 1099$$

last visit

all trajectories have total reward 1000

first visit

"Typical sample":

$s, 1, s, 1, s, 1, \dots, s, 1000, s'$

←-----→

~ 99

Value: $\underbrace{1 + 1 + \dots + 1}_{99} + 1000 \approx 1099$

all visits

calculations are hard to do by hand,
but it does converge to 1099...

CONCLUSION:

BOTH all visits and first visits converge to
the right value

last visit does NOT.