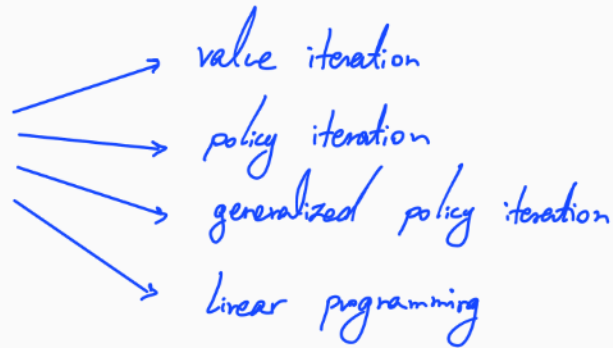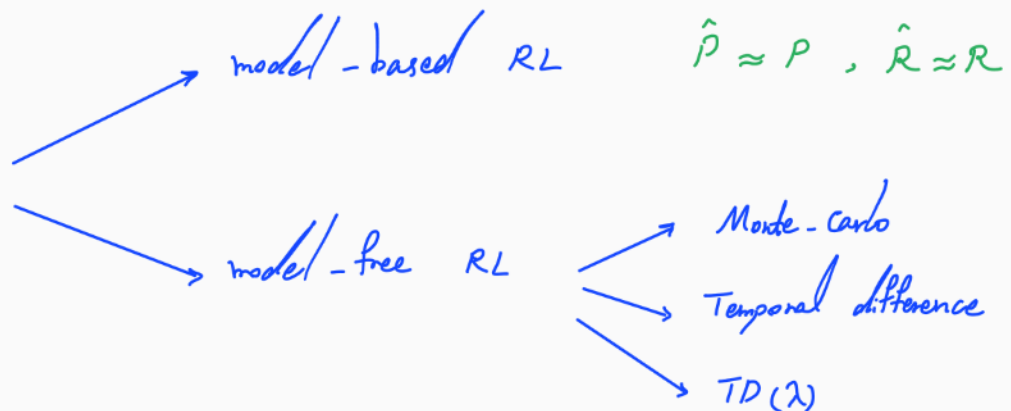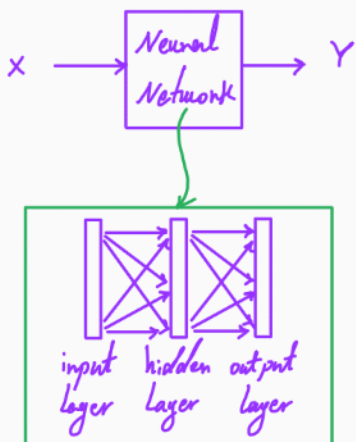So far, we have been dealing with planning settings where the full environment model ( transition function $P$ and reward function $R$ ) is known.

→ value iteration
→ policy iteration
→ generalized policy iteration
→ Linear programming

Now, we move to learning settings where the environment model is not fully known.

→ model-based RL        $\hat{P} \approx P$ , $\hat{R} \approx R$

→ model-free RL → Monte-Carlo
→ Temporal difference
→ TD ($\lambda$)

Function Approximation

$X \longrightarrow$ Neural Network $\longrightarrow Y$

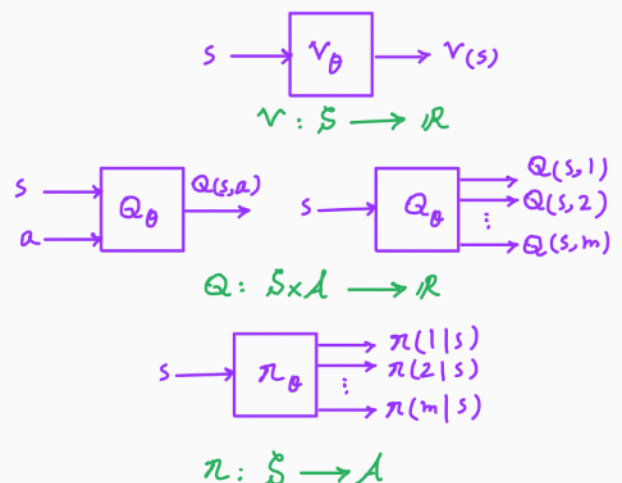input hidden output
layer layer layer

Functions of interest in RL

value functions $V$, $Q$
$V^*$, $Q^*$

Policy $\pi$ / $\pi^*$

Environment model $\hat{P}$, $\hat{R}$
(in model-based RL)

RL with NN approximators (Deep RL)

$s \longrightarrow V_\theta \longrightarrow V(s)$
$V : S \longrightarrow \mathbb{R}$

$s \longrightarrow Q_\theta \xrightarrow{Q(s,a)}$
$a \longrightarrow$

$s \longrightarrow Q_\theta \begin{array}{l} \rightarrow Q(s,1) \\ \rightarrow Q(s,2) \\ \vdots \\ \rightarrow Q(s,m) \end{array}$

$Q : S \times A \longrightarrow \mathbb{R}$

$s \longrightarrow \pi_\theta \begin{array}{l} \rightarrow \pi(1|s) \\ \rightarrow \pi(2|s) \\ \vdots \\ \rightarrow \pi(m|s) \end{array}$

$\pi : S \longrightarrow A$

# Monte carlo Method

Rather than having the full environment model ( transition function $P$ and reward function $R$ ), Monte Carlo methods only require experiences, i.e., sample trajectories (sequence of states, actions, rewards).
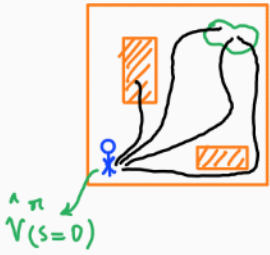
Model_free method

Monte carlo: Using random sampling to perform a computational task, such as estimation and optimization.

Assumption: The setting is episodic, i.e., the interactions happen in episodes of finite length.

[ The updates of MC-based methods happen after each episode ⟶ not fully online ]

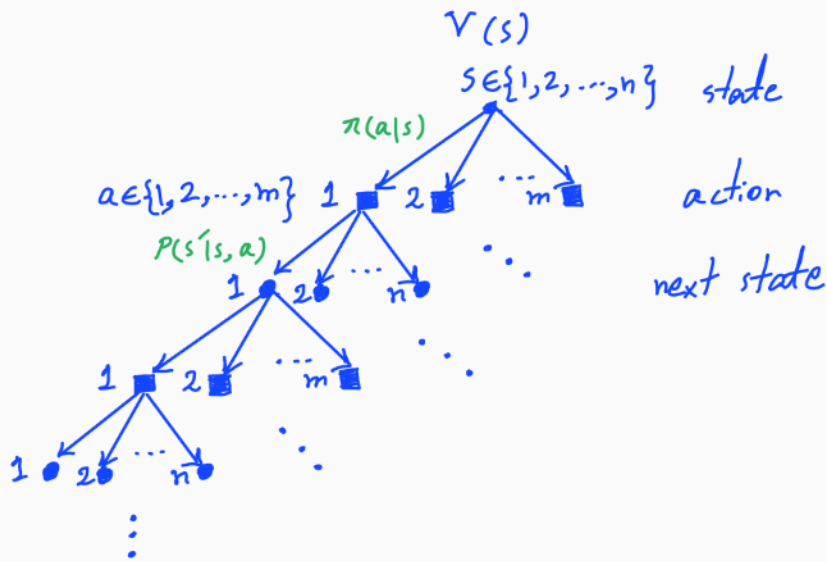MC for policy evaluation: Given a stationary policy $\pi$, we want to compute

$$V^{\pi}(s) = \mathbb{E}_{\substack{A_t \sim \pi \\ S_{t+1} \sim P(\cdot | S_t, A_t)}} \left[ \overbrace{\sum_{t=0}^{T} \gamma^t R(S_t, A_t)}^{\text{return } G} \,\Big|\, S_0 = s \right]$$

$\hat{V}^{\pi}(s=0)$

**Idea**: Estimate the expected return using the empirical mean of the returns of sample trajectories.

$$\boxed{\mathbb{E}_{P_x}[x] \approx \hat{\mathbb{E}}[X] = \frac{1}{N} \sum_{i=1}^{N} x_i}$$

sampled from $P_x$

x can be return G

$V(s)$

$s \in \{1, 2, \ldots, n\}$  state

$\pi(a|s)$

$a \in \{1, 2, \ldots, m\}$  1  2  $\cdots$  m   action

$P(s'|s, a)$

1  2  $\cdots$  n   next state

1  2  $\cdots$  m

1  2  $\cdots$  n

until the episode terminates

Let $T = \{z_1, z_2, \ldots, z_{|T|}\}$ denote a set of trajectories sampled from $s$

$$\longrightarrow V^{\pi}(s) = \mathbb{E}_{\substack{A_t \sim \pi \\ S_{t+1} \sim P(\cdot|S_t, A_t)}} \left[ \sum_{t=0}^{T} \gamma^t R(S_t, A_t) | S_0 = s \right] \approx \frac{1}{|T|} \left( \sum_{z \in T} \underbrace{\sum_{t=0}^{|z_i|} \gamma^t R(s_t^i, a_t^i)}_{\text{return for trajectory } z_i} \right) = \hat{V}^{\pi}(s)$$

i.i.d. samples

# First-visit MC method

Iterative method:

- Initialize $\hat{V}^{\pi}$

- Initialize $G(s) = \emptyset$ for all $s \in S$

- Repeat

    - Generate a sample trajectory $\tau$ using $\pi$

    - For $s \in S$ s.t. $s \in \tau$:

        state 5  2  1
        time  0  1  2   ...   $|\tau|$
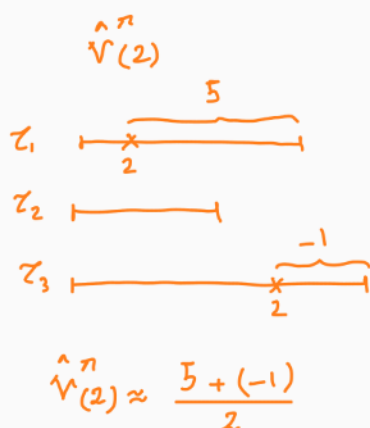
        - $t_s =$ first time $s$ is visited

        - $G = \sum_{t=t_s}^{|\tau|} \gamma^{t-t_s} R(s_t, a_t)$

        - $G(s) \leftarrow G(s) \cup \{G\}$

        - $N(s) \leftarrow N(s) + 1$    $N(s) = |G(s)|$

        - $\hat{V}^{\pi}(s) = \frac{1}{N(s)} \sum_{G \in G(s)} G$

$\hat{V}^{\pi}(2)$

$\tau_1$     5

$\tau_2$

$\tau_3$     -1

$\hat{V}^{\pi}(2) \approx \frac{5 + (-1)}{2}$

$$\hat{\mu}_{k+1} = \frac{\sum_{i=1}^{k+1} x_i}{k+1} = \frac{\sum_{i=1}^{k} x_i + x_{k+1}}{k+1} = \frac{k\hat{\mu}_k + x_{k+1}}{k+1}$$

Incremental computation of mean

$$= \frac{k}{k+1}\hat{\mu}_k + \frac{1}{k+1}x_{k+1} = \hat{\mu}_k + \frac{1}{k+1}(x_{k+1} - \hat{\mu}_k)$$

Note: By law of large numbers, $\hat{V}^{\pi}(s) \xrightarrow[N(s) \to \infty]{} V^{\pi}(s)$ .

# Every - visit MC method
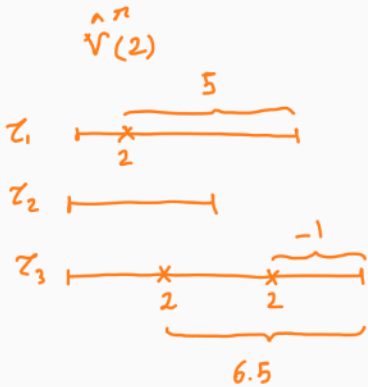
Iterative method:

- Initialize $\hat{V}^{\pi}$

- Initialize $G(s) = \emptyset$ for all $s \in S$

- Repeat

    - Generate a sample trajectory $\tau$ using $\pi$

    - For $s \in S$ s.t. $s \in \tau$:

state  5   2   1

time   0   1   2   ...   $|\tau|$

$\tau$

$\hat{V}^{\pi}(2)$



$\hat{V}^{\pi}(2) \approx \dfrac{5 + 6.5 + (-1)}{3}$

    - For every time $s$ is visited:

        - $t_s$ = time $s$ is visited

        - $G = \displaystyle\sum_{t=t_s}^{|\tau|} \gamma^{t-t_s} R(s_t, a_t)$

        - $G(s) \leftarrow G(s) \cup \{G\}$

        - $N(s) \leftarrow N(s) + 1$

    - $\hat{V}^{\pi}(s) = \dfrac{1}{N(s)} \displaystyle\sum_{G \in G(s)} G$

**Note:** It can be shown that $\hat{V}^{\pi}(s) \xrightarrow[N(s) \to \infty]{} V^{\pi}(s)$.

**MC for policy optimization:** We want to find an (approximately) optimal policy $\pi^*$ through iterative policy evaluation and policy improvement.

evaluation

$$\pi \xrightarrow{\hspace{1cm}} v^\pi$$

improvement

**Idea:** → generalized policy iteration

- Evaluate the policy by estimating the Q function using the empirical mean of the returns of the sample trajectories.

- Improve the policy using the estimated Q function in a greedy manner.

**Estimating $v$ or $Q$?**

$$Q(s,a)$$

If $v$ known $\longrightarrow$ $\pi = \underset{a \in A}{\mathrm{argmax}} \left[ R(s,a) + \gamma \underset{s' \sim P(\cdot|s,a)}{\mathbb{E}} [v(s')] \right]$

unknown

If $Q$ known $\longrightarrow$ $\pi = \underset{a \in A}{\mathrm{argmax}} \; Q(s,a)$

$\longrightarrow$ we need to estimate $Q$ instead of $v$.

# MC method with exploring starts

Iterative method:

- Initialize $\hat{Q}$

- Initialize $\pi$

- Initialize $G(s,a) = \emptyset$ for all $s \in S$, $a \in A$

· Choose $d \in \Delta(S \times A)$ s.t. $d(s,a) > 0$ for all $s \in S$, $a \in A$

- Repeat $\underbrace{\hspace{3cm}}$ set of all possible probability distributions over $S \times A$

    - Sample $s_0 \in S$ and $a_0 \in A$ according to $d$

    - Generate a sample trajectory $\tau$ using $\pi$ starting from $(s_0, a_0)$

    - For $(s,a) \in S \times A$ s.t. $(s,a) \in \tau$:

        - $t_{s,a} = $ first time $(s,a)$ is visited

        - $G = \sum_{t=t_{s,a}}^{|\tau|} \gamma^{t-t_{s,a}} R(s_t, a_t)$

        - $G(s,a) \leftarrow G(s,a) \cup \{G\}$

        - $N(s,a) \leftarrow N(s,a) + 1$

        - $\hat{Q}(s,a) = \frac{1}{N(s,a)} \sum_{G \in G(s,a)} G$

*(top right diagram)*
evaluation
$\pi \longrightarrow \hat{Q}^{\pi}$
improvement

*(green label, left bracket)* Policy evaluation

*(orange, right)* $G(s=1, a=3)$ $= \{-5, 7, 12, \dots\}$

*(green, right)* $N(s,a) = |G(s,a)|$

- For $s \in \hat{S}$ s.t. $s \in \tau$:

  - $\pi(s) = \underset{a \in A}{\arg\max} \ \hat{Q}(s,a)$

---

**Note:** In real-world settings, it may not be possible to start from any state-action pairs.

$\longrightarrow$ Can we instead encourage continuous exploration?

**Idea:**

- Use a soft (stochastic) policy $\pi(a|s)$ such that it chooses all (possible) actions at each state with a non-zero probability, i.e.,

$$\pi(a|s) > 0 \quad \forall s \in S, \ \forall a \in A.$$

- As the estimates improve over time, reduce exploration.

# MC method with $\varepsilon$-greedy exploration

Iterative method:

- Initialize $\hat{Q}$

- Initialize $\pi(a|s)$ s.t. $\pi(a|s) > 0$ for all $s \in S, a \in A$

- Initialize $G(s,a) = \emptyset$ for all $s \in S, a \in A$

- Repeat

  - Generate a sample trajectory $\tau$ using $\pi$

  - For $(s,a) \in S \times A$ s.t. $(s,a) \in \tau$:

    - $t_{s,a} = $ first time $(s,a)$ is visited

    - $G = \sum_{t=t_{s,a}}^{|\tau|} \gamma^{t-t_{s,a}} R(s_t, a_t)$

    - $G(s,a) \longleftarrow G(s,a) \cup \{G\}$

    - $N(s,a) \longleftarrow N(s,a) + 1$

    - $\hat{Q}(s,a) = \frac{1}{N(s,a)} \sum_{G \in G(s,a)} G$

*Policy evaluation*

Policy improvement 
$\left\{\begin{array}{l}\end{array}\right.$

- For $s \in S$ s.t. $s \in \tau$ :

  - $a^* = \underset{a \in \mathcal{A}}{\text{argmax}} \ \hat{Q}(s,a)$

  - For $a \in \mathcal{A}$ :

    - $\pi(a|s) = \begin{cases} 1-\varepsilon + \dfrac{\varepsilon}{|\mathcal{A}|} & \text{if } a=a^* \\[2mm] \dfrac{\varepsilon}{|\mathcal{A}|} & \text{otherwise} \end{cases}$

$\varepsilon$-greedy

Exploration

Decaying exploration :

$\varepsilon$



$\varepsilon \propto \dfrac{1}{K}$

episode $K$