

Markov Decision Processes

Chen-Yu Wei

Sequence of Actions



To win the game, the learner has to take a sequence of actions $a_1 \rightarrow a_2 \rightarrow \dots \rightarrow a_H$.

One option: view every sequence as a “meta-action”: $\bar{a} = (a_1, a_2, \dots, a_H)$

Drawback:

- The number of actions is exponential in horizon
- In stochastic environments, this does not leverage intermediate observations

Solution idea: dynamic programming

Interaction Protocol: Fixed-Horizon Case

For **episode** $t = 1, 2, \dots, T$:

For **step** $h = 1, 2, \dots, H$:

Learner observes an observation $x_{t,h}$

Learner chooses an action $a_{t,h}$

Learner receives instantaneous reward $r_{t,h}$

General case:

$$\mathbb{E}[r_{t,h}] = R(x_{t,1}, a_{t,1}, \dots, x_{t,h}, a_{t,h}), \quad x_{t,h+1} \sim P(\cdot | x_{t,1}, a_{t,1}, \dots, x_{t,h}, a_{t,h})$$

⇒ Optimal decisions may depend on the entire history $\mathcal{H}_t = (x_{t,1}, a_{t,1}, \dots, x_{t,h})$

Interaction Protocol: Fixed-Horizon Case

For **episode** $t = 1, 2, \dots, T$:

For **step** $h = 1, 2, \dots, H$:

Learner observes an observation $x_{t,h}$

Learner chooses an action $a_{t,h}$

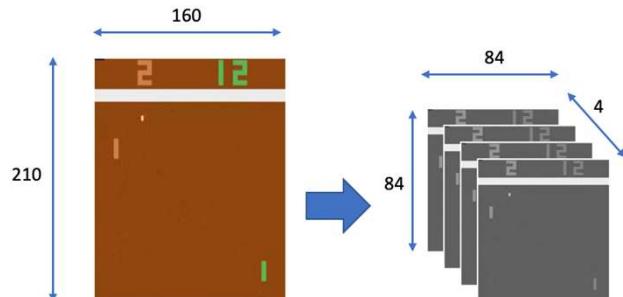
Learner receives instantaneous reward $r_{t,h}$

We assume that the history $\mathcal{H}_t = (x_{t,1}, a_{t,1}, \dots, x_{t,h})$ can be summarized as a **horizon-length-independent** representation $s_{t,h} = \Phi(x_{t,1}, a_{t,1}, \dots, x_{t,h}) \in \mathcal{S}$ so that

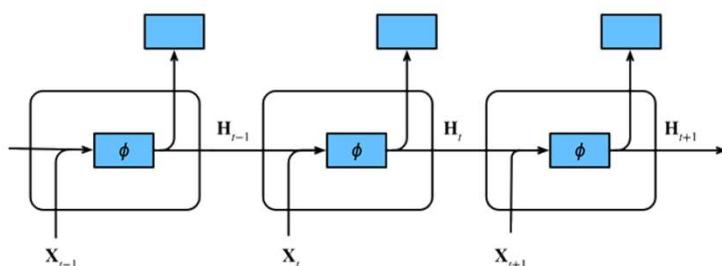
$$\mathbb{E}[r_{t,h}] = R(s_{t,h}, a_{t,h}), \quad x_{t,h+1} \sim P(\cdot | s_{t,h}, a_{t,h})$$

$s_{t,h}$ is called the “**state**” at the step h of episode t .

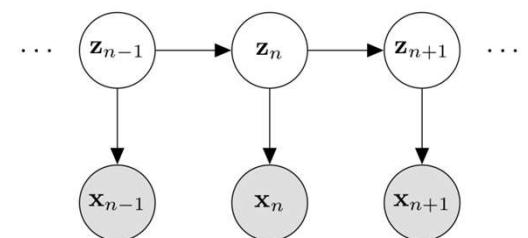
From Observations to States



Stacking recent observations



Recurrent neural network



Hidden Markov model

Interaction Protocol: Fixed-Horizon Case

For **episode** $t = 1, 2, \dots, T$:

For **step** $h = 1, 2, \dots, H$:

Environment reveals **state** $s_{t,h}$

Learner chooses an action $a_{t,h}$

Learner observes instantaneous reward $r_{t,h}$ with $\mathbb{E}[r_{t,h}] = R(s_{t,h}, a_{t,h})$

Next state is generated as $s_{t,h+1} \sim P(\cdot | s_{t,h}, a_{t,h})$

This is called the Markov decision process.

MDP as Contextual Bandits?

Viewing states as contexts, and viewing the problem as a contextual bandit problem with TH rounds (what's wrong?)

$$\text{Regret}_{\text{(Contextual bandit)}} = \sum_{t=1}^T \sum_{h=1}^H \max_a R(S_{t,h}, a) - \underbrace{\sum_{t=1}^T \sum_{h=1}^H R(S_{t,h}, a_{t,h})}_{(x_{t,1}, a_{t,1} \dots x_{t,h})}$$

$$\text{Regret}_{\text{(MDP)}} = \sum_{t=1}^T \left[\sum_{h=1}^H R(S_{t,h}^*, a_{t,h}^*) \right] - \sum_{t=1}^T \sum_{h=1}^H R(S_{t,h}, a_{t,h})$$

$$S_{t,1}^* = S_{t,1}$$

$$S_{t,h}^* \neq S_{t,h} \quad \text{for } h \geq 2$$

Formulations

- Interaction Protocol
 - Fixed-Horizon
 - Variable-Horizon (Goal-Oriented)
 - Infinite-Horizon
- Performance Metric
 - Total Reward
 - Average Reward
 - Discounted Reward
- Policy
 - History-dependent policy
 - Markov policy
 - Stationary policy

Horizon = Length of an episode

Interaction Protocols (1/3): Fixed-Horizon

Horizon length is a fixed number H

$h \leftarrow 1$

Observe initial state s_1

While $h \leq H$:

 Choose action a_h

 Observe reward r_h with $\mathbb{E}[r_h] = R(s_h, a_h)$

 Observe next state $s_{h+1} \sim P(\cdot | s_h, a_h)$

Examples: games with a fixed number of time

Interaction Protocols (2/3): Goal-Oriented

The learner interacts with the environment until reaching **terminal states** $\mathcal{T} \subset \mathcal{S}$

$h \leftarrow 1$

Observe initial state s_1

While $s_h \notin \mathcal{T}$:

 Choose action a_h

 Observe reward r_h with $\mathbb{E}[r_h] = R(s_h, a_h)$

 Observe next state $s_{h+1} \sim P(\cdot | s_h, a_h)$

$h \leftarrow h + 1$

Examples: video games, robotics tasks, personalized recommendations, etc.

Interaction Protocols (3/3): Infinite-Horizon

The learner continuously interacts with the environment

$h \leftarrow 1$

Observe initial state s_1

Loop forever:

Choose action a_h

Observe reward r_h with $\mathbb{E}[r_h] = R(s_h, a_h)$

Observe next state $s_{h+1} \sim P(\cdot | s_h, a_h)$

$h \leftarrow h + 1$

Examples: network management, inventory management

Formulations for Markov Decision Processes

- Interaction Protocol
 - Fixed-Horizon
 - Variable-Horizon (Goal-Oriented)
 - Infinite-Horizon
 - Performance Metric
 - Total Reward
 - Average Reward
 - Discounted Reward
 - Policy
 - History-dependent policy
 - Markov policy
 - Stationary policy
- } Episodic setting

Performance Metric

Total Reward (for episodic setting):

$$\sum_{h=1}^{\tau} r_h$$

(τ : the step where the episode ends)

Average Reward (for infinite-horizon setting):

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{h=1}^T r_h$$

If $|r_h| \leq 1$,

Discounted Total Reward (for episodic or infinite-horizon):

$$\sum_{h=1}^{\tau} \gamma^{h-1} r_h \quad \leq \frac{1}{1-\gamma}$$

τ : the step where the episode ends, or ∞ in the infinite-horizon case

$\gamma \in [0,1]$: discount factor

Interaction Protocols vs. Performance Metrics

	“natural” objective	
Fixed-Horizon	----->	Total Reward
Goal-Oriented	----->	Total Reward Could be unbonded
Infinite-horizon	----->	Average Reward Could have constant change for an infinitesimal change in policy

Discounted Total Reward?
Focusing more on the **recent** reward

There is a potential mismatch between our ultimate goal and what we optimized.

Our Focus

In most of the following lectures, we focus on the **goal-oriented / infinite-horizon** setting with **discount total reward** as the performance metric.

Policy

A mapping from observations/contexts/states to (distribution over) actions

- Contextual bandits

$$\begin{aligned} a &\sim \pi(\cdot | x) && \text{(randomized/stochastic)} \\ \text{or } a &= \pi(x) && \text{(deterministic)} \end{aligned}$$

- Multi-armed bandits

$$\begin{aligned} a &\sim \pi \\ \text{or } a &= a^* \end{aligned}$$

Policy for MDPs

History-dependent Policy

$$\begin{aligned} a_h &\sim \pi(\cdot | s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_h) \\ a_h &= \pi(s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_h) \end{aligned}$$

Markov Policy

$$\begin{aligned} a_h &\sim \pi(\cdot | s_h, h) \\ a_h &= \pi(s_h, h) \end{aligned}$$

For **fixed-horizon + total reward** setting,
there exists an optimal policy in this class

Stationary Policy

$$\begin{aligned} a_h &\sim \pi(\cdot | s_h) \\ a_h &= \pi(s_h) \end{aligned}$$

For **infinite-horizon/goal-oriented + discounted total reward**
setting, there exists an optimal policy in this class

Final Project

- I have read your proposals, and will send feedback to you soon.
- Some groups propose to perform “exploration” over recommendation system datasets. This is in general not possible. Instead, one has to use
 - Offline RL/bandits techniques
 - Synthetic data / simulators
- Try not only rely on the techniques learned in the class
 - The course aims to provide fundamental / theoretical understanding
 - Many common algorithms are introduced much later
- It is not necessary to produce “good” results. Interesting attempts and failure experiences are also valuable
 - Evaluation will be based on novelty, technicality, motivation (and writing, presentation)

Fixed-Horizon + Total Reward

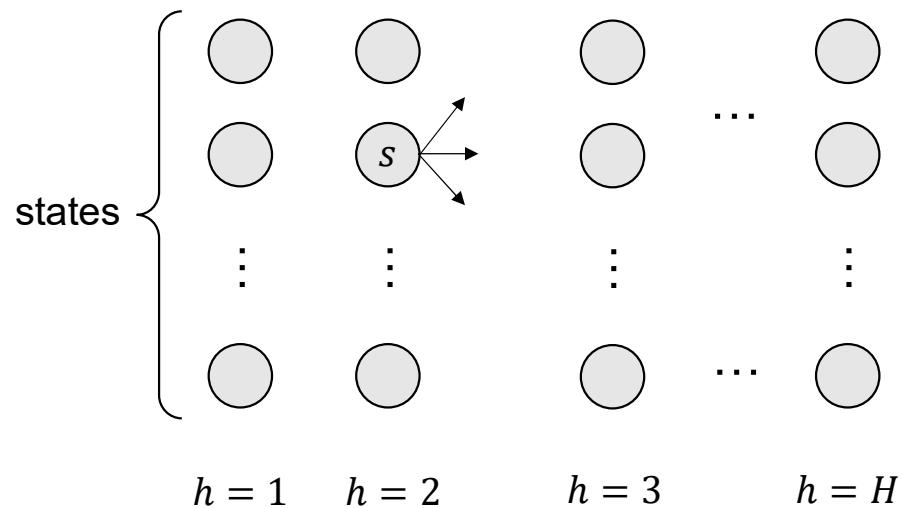
Dynamic Programming

Goal: Calculate the expected total reward of a policy

A (Markov) policy is a mapping from (state, step index) to action distribution, written as

$$\pi_h(\cdot | s) \in \Delta(\mathcal{A}) \quad \text{for } s \in \mathcal{S} \text{ and } h \in \{1, 2, \dots, H\}$$

Dynamic Programming



State transition: $P(s'|s, a)$

Reward: $R(s, a)$

$$V_h^\pi(s) = \mathbb{E} \left[\sum_{i=h}^H R(s_i, a_i) \mid s_h = s, a_i \sim \pi_i(\cdot | s_i), i \geq h \right]$$

Key quantity: $V_h^\pi(s)$ = the expected total reward of policy π starting from state s at step h .

Backward calculation:

$$V_H^\pi(s) = \sum_a \pi_H(a|s) R(s, a) \quad \forall s$$

For $h = H - 1, \dots, 1$: for all s

$$V_h^\pi(s) = \sum_a \pi_h(a|s) \left(R(s, a) + \underbrace{\sum_{s'} P(s'|s, a) V_{h+1}^\pi(s')}_{\text{Expected total reward from step } h+1} \right)$$

Expected total reward from step $h + 1$

Bellman Equation

$$V_{H+1}^\pi(s) = 0$$

$$Q_h^\pi(s, a) = \mathbb{E} \left[\sum_{i=h}^H R(s_i, a_i) \mid S_h = s, a_h = a, a_i \sim \pi_i(\cdot | s_i) \forall i \geq h+1 \right]$$

$$V_h^\pi(s) = \sum_a \pi_h(a|s) \underbrace{\left(R(s, a) + \sum_{s'} P(s'|s, a) V_{h+1}^\pi(s') \right)}_{Q_h^\pi(s, a)} \quad \text{for } h = H, \dots, 1$$

$$V_h^\pi(s) = \sum_{a \in \mathcal{A}} \pi_h(a|s) Q_h^\pi(s, a)$$

$$Q_h^\pi(s, a) = R(s, a) + \sum_{s'} P(s'|s, a) V_{h+1}^\pi(s')$$

Occupancy Measures

$$d_\rho^\pi(s) = \mathbb{E} \left[\sum_{h=1}^H \mathbb{I}\{S_h=s\} \mid s_1 \sim \rho, a_i \sim \pi_i(\cdot | s_i) \forall i \geq 1 \right]$$

$d_\rho^\pi(s)$: the expected number of times state s is visited, under policy π and initial state distribution ρ

$$\boxed{d_{\rho,h}^\pi(s) = \Pr(S_h=s \mid \text{---})}, \quad d_\rho^\pi(s) = \sum_{h=1}^H d_{\rho,h}^\pi(s)$$

$$= \mathbb{E} [\mathbb{I}\{S_h=s\} \mid \text{---}]$$

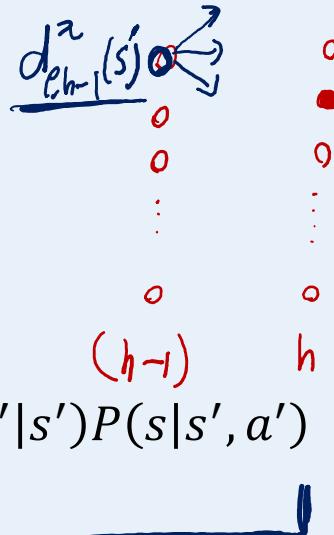
Key quantity: $d_{\rho,h}^\pi(s)$ = the probability of state s being visited **at step h** , under policy π and initial state distribution ρ

Forward calculation:

$$d_{\rho,1}^\pi(s) = \rho(s) \quad \forall s$$

For $h = 2, \dots, H$:

$$d_{\rho,h}^\pi(s) = \sum_{s'} d_{\rho,h-1}^\pi(s') \sum_{a'} \pi_{h-1}(a'|s') P(s|s', a') \quad \forall s$$



Reverse Bellman Equation

$$d_{\rho,1}^{\pi}(s) = \rho(s)$$

$$d_{\rho,h}^{\pi}(s) = \sum_{s',a'} \underbrace{d_{\rho,h-1}^{\pi}(s') \pi_{h-1}(a'|s') P(s|s',a')}_{d_{\rho,h-1}^{\pi}(s',a')} \quad \text{for } h = 2, \dots, H$$

$d_{\rho,h-1}^{\pi}(s',a') \asymp \beta_r \left(s_{h+1} = s', a_{h+1} = a' \mid s_1 \sim \rho, \pi \right)$

$$d_{\rho,h}^{\pi}(s) = \sum_{s',a'} d_{\rho,h-1}^{\pi}(s',a') P(s|s',a')$$
$$d_{\rho,h}^{\pi}(s, a) = d_{\rho,h}^{\pi}(s) \pi_h(a|s)$$

Dynamic Programming

$$V_h^*(s) = \max_{\pi} V_h^{\pi}(s)$$

Goal: Find the optimal policy

Key quantity: $V_h^*(s)$ = the optimal expected total reward starting from state s at step h .

Backward calculation:

$$V_H^*(s) = \max_a R(s, a) \quad \forall s$$

For $h = H - 1, \dots, 1$:

$$V_h^*(s) = \max_a \left(R(s, a) + \sum_{s'} P(s'|s, a) V_{h+1}^*(s') \right) \quad \forall s$$

Value Iteration

$$\pi_h^*(s) = \operatorname{argmax}_a R(s, a) + \sum_{s'} P(s'|s, a) V_{h+1}^*(s')$$

Bellman Optimality Equation

$$V_{H+1}^*(s) = 0$$
$$V_h^*(s) = \max_a \left(R(s, a) + \sum_{s'} P(s'|s, a) V_{h+1}^*(s') \right) \quad \text{for } h = H, \dots, 1$$

$Q_h^*(s, a)$

$$V_h^*(s) = \max_a Q_h^*(s, a)$$

$$Q_h^*(s, a) = R(s, a) + \sum_{s'} P(s'|s, a) V_{h+1}^*(s')$$

$$\pi_h^*(s) = \operatorname{argmax}_a Q_h^*(s, a)$$

Recap

$$V_h^\pi(s) = \sum_{a \in \mathcal{A}} \pi_h(a|s) Q_h^\pi(s, a)$$

$$Q_h^\pi(s, a) = R(s, a) + \sum_{s' \in \mathcal{S}} P(s'|s, a) V_{h+1}^\pi(s')$$

Bellman Equation
(Value Iteration for V^π)

$$d_{\rho, h}^\pi(s, a) = d_{\rho, h}^\pi(s) \pi_h(a|s)$$

$$d_{\rho, h}^\pi(s) = \sum_{s', a'} d_{\rho, h-1}^\pi(s', a') P(s|s', a')$$

Reverse Bellman Equation

$$V_h^*(s) = \max_a Q_h^*(s, a)$$

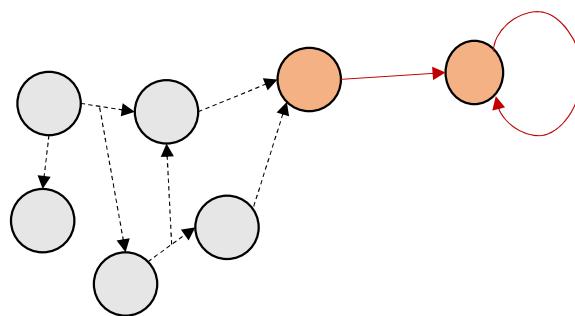
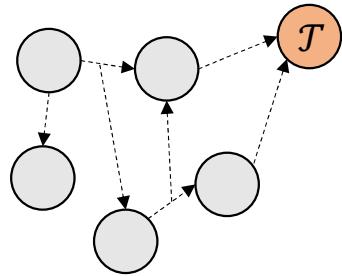
$$Q_h^*(s, a) = R(s, a) + \sum_{s' \in \mathcal{S}} P(s'|s, a) V_{h+1}^*(s')$$

Bellman Optimality Equation
(Value Iteration)

**Infinite-Horizon / Goal-Oriented +
Discounted Total Reward**

Equivalent Views

→ deterministic and zero-reward

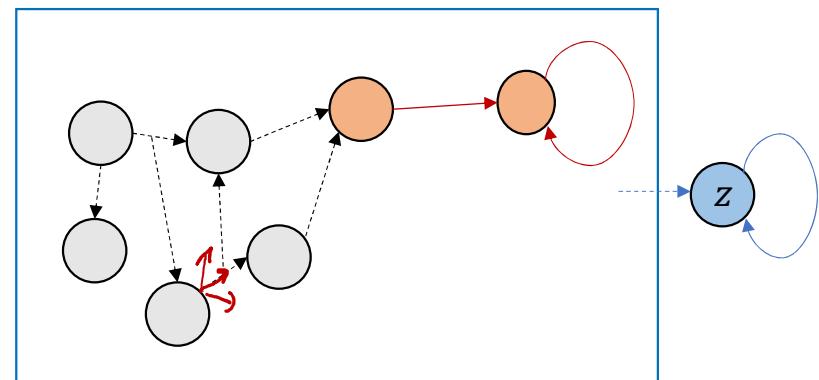


Converting goal-oriented to infinite-horizon

$$\mathbb{E}^{\text{new}} \left[\sum_{h=1}^{\infty} \gamma^{h-1} r_h \right] = \mathbb{E}^{\text{old}} \left[\sum_{h=1}^{\tau} \gamma^{h-1} r_h \right]$$

$$p(s'|s,a)^{\text{new}} = \gamma p(s'|s,a)^{\text{old}}, \quad p(z|s,a)^{\text{new}} = 1 - \gamma$$

Scale down all transitions by a factor of γ
and add probability $1 - \gamma$ transitioning to z



Converting discounted total reward to total reward

$$\mathbb{E}^{\text{new}} \left[\sum_{h=1}^{\infty} r_h \right] = \mathbb{E}^{\text{old}} \left[\sum_{h=1}^{\infty} \gamma^{h-1} r_h \right]$$

Prob of staying in original MDP at step h

Dynamic Programming

$$V_i^\pi(s) = \lim_{H \rightarrow \infty} \mathbb{E} \left[\sum_{h=1}^H \gamma^{h-1} R(s_h, a_h) \mid s_1 = s, a_h \sim \pi(\cdot | s_h) \forall h \geq 1 \right]$$

Goal: Calculate the expected discounted total reward of a stationary policy π

$V^\pi(s)$ = the expected discounted total reward starting from state s , follow π

Key quantity: $V_i^\pi(s)$ = the expected discounted total reward starting from state s supposed that i more steps can be executed

$$V_0^\pi(s) = 0 \quad \forall s$$

$$V_i^\pi(s) = \mathbb{E} \left[\sum_{h=1}^i \gamma^{h-1} R(s_h, a_h) \mid \sim \right]$$

For $i = 1, 2, 3 \dots$

$$V_i^\pi(s) = \sum_a \pi(a|s) \left(R(s, a) + \gamma \sum_{s'} P(s'|s, a) V_{i-1}^\pi(s') \right) \quad \forall s$$

$$+ (1-\gamma) \times 0$$

$$V^\pi(s) = \lim_{i \rightarrow \infty} V_i^\pi(s) \quad (\text{need to prove that the limit exists})$$

Value Iteration for V^π

$$\lim_{i \rightarrow \infty} \hat{V}_i(s) = V^\pi(s)$$

Arbitrary $\hat{V}_0(s) \quad \forall s$

For $i = 1, 2, 3 \dots$

$$\hat{V}_i(s) = \sum_a \pi(a|s) \left(R(s, a) + \gamma \sum_{s'} P(s'|s, a) \hat{V}_{i-1}(s') \right) \quad \forall s$$

To show that this algorithm converges, we prove the following statement:

For any $\epsilon > 0$, there exists a large enough N such that

$$|\hat{V}_i(s) - \hat{V}_j(s)| \leq \epsilon$$

for any $i, j \geq N$.

Proof of Convergence

$$|\hat{V}_{i+1}(s) - \hat{V}_i(s)| \leq O(\gamma^i) \quad \forall s.$$

$$|\hat{V}_i(s) - \hat{V}_j| \leq \sum_{k=i}^{j-1} |\hat{V}_k(s) - \hat{V}_{k+1}(s)| = \sum_{k=i}^{j-1} O(\gamma^k) \leq O\left(\frac{\gamma^i}{1-\gamma}\right)$$

$$\hat{V}_i(s) = \sum_a \pi(a|s) \left(R(s,a) + \sum_{s'} p(s'|s,a) \hat{V}_{i-1}(s') \right) \quad \forall s$$

$$\hat{V}_{i+1}(s) = \sum_a \pi(a|s) \left(R(s,a) + \sum_{s'} p(s'|s,a) \hat{V}_i(s') \right) \quad \forall s$$

$$\Rightarrow \hat{V}_{i+1}(s) - \hat{V}_i(s) = \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) (\hat{V}_i(s') - \hat{V}_{i-1}(s'))$$

$$\begin{aligned} \Rightarrow \max_s |\hat{V}_{i+1}(s) - \hat{V}_i(s)| &\leq \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) |\hat{V}_i(s') - \hat{V}_{i-1}(s')| \\ &\approx \gamma \underbrace{\sum_a \pi(a|s)}_{\text{constant}} \underbrace{\sum_{s'} p(s'|s,a)}_{\text{constant}} \max_{s^*} |\hat{V}_i(s^*) - \hat{V}_{i-1}(s^*)| \\ &\leq \gamma \max_{s^*} |\hat{V}_i(s^*) - \hat{V}_{i-1}(s^*)| \end{aligned}$$

Proof of Convergence

For any $\epsilon > 0$, there exists a large enough N such that

$$|\hat{V}_i(s) - \hat{V}_j(s)| \leq \epsilon$$

for any $i, j \geq N$.

$$\hat{V}(s) = \lim_{i \rightarrow \infty} \inf \{\hat{V}_j(s) : j \geq i\}$$

For any $\epsilon > 0$, there exists a large enough N such that

$$|\hat{V}_i(s) - \hat{V}(s)| \leq \epsilon$$

for any $i \geq N$.

Proof of Uniqueness

No matter what the initial values of $\hat{V}_0(s)$ are, the limit $\lim_{i \rightarrow \infty} \hat{V}_i(s)$ is the same.

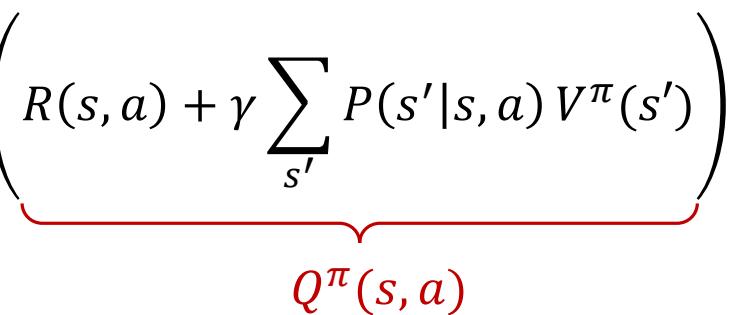
(This value is $V^\pi(s)$)

Assume $V^{(1)}(s)$ and $V^{(2)}(s)$ are different convergence point.

$$\begin{aligned} V^{(1)}(s) &= \sum_a \pi(a|s) \left(R(s,a) + \gamma \sum_{s'} p(s'|s,a) V^{(1)}(s') \right) \\ V^{(2)}(s) &= \underbrace{\qquad\qquad\qquad}_{V^{(1)}(s) - V^{(2)}(s)} \underbrace{\qquad\qquad\qquad}_{V^{(2)}(s')} \\ V^{(1)}(s) - V^{(2)}(s) &= \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) (V^{(1)}(s') - V^{(2)}(s')) \\ \max_s |V^{(1)}(s) - V^{(2)}(s)| &\leq \gamma \max_{s'} |V^{(1)}(s') - V^{(2)}(s')| \rightarrow |V^{(1)}(s) - V^{(2)}(s)| = 0 \end{aligned}$$

Bellman Equation

$$V^\pi(s) = \sum_a \pi(a|s) \left(R(s, a) + \gamma \sum_{s'} P(s'|s, a) V^\pi(s') \right)$$



$$Q^\pi(s, a)$$

$$V^\pi(s) = \sum_a \pi(a|s) Q^\pi(s, a)$$

$$Q^\pi(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) V^\pi(s')$$

Approximate Bellman Equations

$$\begin{aligned}
 & \left| \gamma \sum_a \pi(a|s) \mathbb{E} \left[\hat{V}(s') - V^\pi(s') \right] \right| \\
 & \leq \left| \gamma \sum_a \pi(a|s) \mathbb{E} \left[\max_{s' \sim p(\cdot|s,a)} |\hat{V}(s') - V^\pi(s')| \right] \right| \\
 & = \gamma \sum_a \pi(a|s) \max_{s'} |\hat{V}(s') - V^\pi(s')| \\
 & = \gamma \max_{s'} |\hat{V}(s') - V^\pi(s')|
 \end{aligned}$$

If $\left| \hat{V}(s) - \sum_a \pi(a|s) \left(R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [\hat{V}(s')] \right) \right| \leq \epsilon$

then $|\hat{V}(s) - V^\pi(s)| \leq \frac{\epsilon}{1-\gamma} \quad \forall s$

$$\begin{aligned}
 V^\pi(s) &= \sum_a \pi(a|s) \left(R(s,a) + \gamma \mathbb{E}_{s' \sim p(\cdot|s,a)} [V^\pi(s')] \right) \\
 \Rightarrow \max_s |\hat{V}(s) - V^\pi(s)| &= \left| \hat{V}(s) - \sum_a \pi(a|s) \left(R(s,a) + \gamma \mathbb{E}[V^\pi(s')] \right) \right| \\
 &\leq \left| \underbrace{\sum_a \pi(a|s) (R(s,a) + \gamma \mathbb{E}[\hat{V}(s')])}_{\hat{V}(s)} - \sum_a \pi(a|s) (R(s,a) + \gamma \mathbb{E}[V^\pi(s)]) \right| + \varepsilon \\
 &\leq \gamma \max_{s'} |\hat{V}(s') - V^\pi(s')| + \varepsilon \Rightarrow (1-\gamma) \max_s |\hat{V}(s) - V^\pi(s)| \leq \varepsilon
 \end{aligned}$$

Occupancy Measures

$$d_{\rho}^{\pi}(s) = \mathbb{E} \left[\sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{I}\{S_h=s\} \mid s_1 \sim \rho, a_h \sim \pi(\cdot | S_h) \text{ for } h \geq 1 \right]$$

$d_{\rho}^{\pi}(s)$: the expected discounted number of times state s is visited, under policy π and initial state distribution ρ

$$d_{\rho,h}^{\pi}(s) = \mathbb{E} \left[\gamma^{h-1} \mathbb{I}\{S_h=s\} \mid \text{---} \right]$$

Key quantity: $d_{\rho,h}^{\pi}(s)$ = the discounted probability of state s being visited at step h , under policy π and initial state distribution ρ

Forward calculation:

$$d_{\rho,1}^{\pi}(s) = \rho(s) \quad \forall s$$

For $h = 2, 3, \dots$

$$d_{\rho,h}^{\pi}(s) = \sum_{s'} d_{\rho,h-1}^{\pi}(s') \sum_{a'} \pi(a'|s') P(s|s', a')$$

$\vdots \quad \bullet s$
 $\cdot \quad h$

$$d_{\rho,h}^{\pi}(s) = \gamma \sum_{s'} d_{\rho,h-1}^{\pi}(s') \sum_{a'} \pi(a'|s') P(s|s', a') \quad \forall s$$

$$d_{\rho,h}^{\pi}(s) = \gamma \sum_{s'} d_{\rho,h-1}^{\pi}(s') \sum_{a'} \pi(a'|s') p(s|s',a')$$

$$\Rightarrow \sum_{h=2}^{\infty} d_{\rho,h}^{\pi}(s) = \gamma \sum_{s'} \underbrace{\sum_{h=2}^{\infty} d_{\rho,h-1}^{\pi}(s')}_{\rho(s)} \sum_{a'} \pi(a'|s') p(s|s',a')$$

$$\Rightarrow d_{\rho}^{\pi}(s) - \underbrace{d_{\rho,1}^{\pi}(s)}_{\rho(s)} = \gamma \sum_{s'} d_{\rho}^{\pi}(s') \sum_{a'} \pi(a'|s') p(s|s',a')$$

Reverse Bellman Equation

$$d_\rho^\pi(s, a) = \mathbb{E} \left[\sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{I}\{s_h=s, a_h=a\} \mid s_1 \sim \rho \text{ follows } \right]$$

$$d_\rho^\pi(s) = \rho(s) + \gamma \sum_{s', a'} d_\rho^\pi(s') \underbrace{\pi(a'|s') P(s|s', a')}_{d_\rho^\pi(s', a')}$$

$$\sum_s \underline{d_\rho^\pi(s)} = \frac{1}{1-\gamma}$$

$$d_\rho^\pi(s) = \rho(s) + \gamma \sum_{s', a'} d_\rho^\pi(s', a') P(s|s', a')$$

$$d_\rho^\pi(s, a) = d_\rho^\pi(s) \pi(a|s)$$

$$\boxed{\sum_s d_\rho^\pi(s) = 1}$$

Another (more common) version makes $d_\rho^\pi(s)$ a distribution over s
 → Just change the $\rho(s)$ in the first equation by $(1 - \gamma)\rho(s)$

Dynamic Programming

$$V^*(s) = \max_{\pi} V^\pi(s)$$

Goal: find optimal policy

Key quantity: $V_i^*(s)$ = the optimal discounted total reward starting from state s **supposed that i more steps can be executed**

$$V_0^*(s) = 0 \quad \forall s$$

For $i = 1, 2, 3 \dots$

$$V_i^*(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s'|s, a) V_{i-1}^*(s') \right) \quad \forall s$$

Value Iteration

$$V^*(s) = \lim_{i \rightarrow \infty} V_i^*(s) \quad \pi^*(s) = \operatorname{argmax}_a R(s, a) + \gamma \sum_{s'} P(s'|s, a) V^*(s')$$

Bellman Optimality Equation

$$V^*(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s'|s, a) V^*(s') \right)$$

$Q^*(s, a)$

$$V^*(s) = \max_a Q^*(s, a)$$

$$Q^*(s, a) = R(s, a) + \sum_{s'} P(s'|s, a) V^*(s')$$

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

Approximate Bellman Optimality Equations

Suppose that $\left| \hat{V}(s) - \max_a \left(R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [\hat{V}(s')] \right) \right| \leq \epsilon \quad \forall s$

Then

$$(1) \quad |\hat{V}(s) - V^*(s)| \leq \frac{\epsilon}{1 - \gamma} \quad \forall s$$

$$(2) \quad V^*(s) - V^{\hat{\pi}}(s) \leq \frac{2\epsilon}{1 - \gamma} \quad \forall s$$

where $\hat{\pi}(s) = \operatorname{argmax}_a \left(R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [\hat{V}(s')] \right)$

Summary

Guarantees for approximate solutions

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) Q^\pi(s, a)$$

$$Q^\pi(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V^\pi(s')$$

$$\left| \hat{V}(s) - \sum_{a \in \mathcal{A}} \pi(a|s) \left(R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \hat{V}(s') \right) \right| \leq \epsilon \quad \forall s$$

$$\Rightarrow |\hat{V}(s) - V^\pi(s)| \leq \frac{\epsilon}{1-\gamma} \quad \forall s$$

$$d_\rho^\pi(s, a) = d_\rho^\pi(s) \pi(a|s)$$

$$d_\rho^\pi(s) = (1 - \gamma)\rho(s) + \gamma \sum_{s', a'} d_\rho^\pi(s', a') P(s|s', a')$$

<https://www.youtube.com/watch?v=XVuRQWXtxLA>

$$V^*(s) = \max_a Q^*(s, a)$$

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V^*(s')$$

$$\left| \hat{V}(s) - \max_a \left(R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \hat{V}(s') \right) \right| \leq \epsilon \quad \forall s$$

$$\Rightarrow |\hat{V}(s) - V^*(s)| \leq \frac{\epsilon}{1-\gamma} \text{ and } V^*(s) - V^{\hat{\pi}}(s) \leq \frac{2\epsilon}{1-\gamma} \quad \forall s$$

Policy Iteration

Policy Iteration

value iteration for V^π (Q^π)

for V^*

Policy Iteration

For $k = 1, 2, \dots$

$$\forall s, \quad \pi^{(k+1)}(s) \leftarrow \operatorname{argmax}_a Q^{\pi^{(k)}}(s, a)$$

$$Q^\pi(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} [V^\pi(s')], \quad V^\pi(s) = \sum_a \pi(a|s) Q^\pi(s, a)$$

Theorem (monotonic improvement). Policy Iteration ensures

$$\forall s, \quad V^{\pi^{(k+1)}}(s) \geq V^{\pi^{(k)}}(s)$$

Below, we will establish a more general lemma (not only show monotonic improvement, but also quantify *how much* the improvement is).

Single-Step Policy Modification under Fixed Horizon

Assume $\pi'_h(\cdot|s) = \pi_h(\cdot|s)$ for all $h \neq h^*$

 $1 \quad h^* \quad H$	$\mathbb{E}_{s \sim \rho} [V_1^{\pi'}(s)] - \mathbb{E}_{s \sim \rho} [V_1^{\pi}(s)] = ?$
\vdots	$= \mathbb{E} \left[\sum_{h=1}^H R(s_h, a_h) \mid s_1 \sim \rho, \pi' \right] - \mathbb{E} \left[\sum_{h=1}^H R(s_h, a_h) \mid s_1 \sim \rho, \pi \right]$
\vdots	$= \mathbb{E} \left[\sum_{h=h^*}^H R(s_h, a_h) \mid s_1 \sim \rho, \pi' \right] - \mathbb{E} \left[\sum_{h=h^*}^H R(s_h, a_h) \mid s_1 \sim \rho, \pi \right]$
\vdots	$= \mathbb{E} \left[\sum_{h=h^*}^H R(s_h, a_h) \mid s_{h^*} \sim d_{\rho, h^*}^{\pi_{\text{in}}}, \pi' \right] - \mathbb{E} \left[\sum_{h=h^*}^H R(s_h, a_h) \mid s_h \sim d_{\rho, h^*}^{\pi_{\text{in}}}, \pi \right]$
\vdots	$= \mathbb{E}_{s_{h^*} \sim d_{\rho, h^*}^{\pi_{\text{in}}} \mathbb{E}_{a_{h^*} \sim \pi'_{h^*}(\cdot s_{h^*})} \left[Q_{h^*}^{\pi_{\text{out}}}(s_{h^*}, a_{h^*}) \right]} - \mathbb{E}_{s_{h^*} \sim d_{\rho, h^*}^{\pi_{\text{in}}} \mathbb{E}_{a_{h^*} \sim \pi_{h^*}(\cdot s_{h^*})} \left[Q_{h^*}^{\pi_{\text{out}}}(s_{h^*}, a_{h^*}) \right]}$

$$\pi'_h(\cdot|s) = \pi_h(\cdot|s) = \pi_{\text{in}}(\cdot|s) \text{ for } h < h^*$$

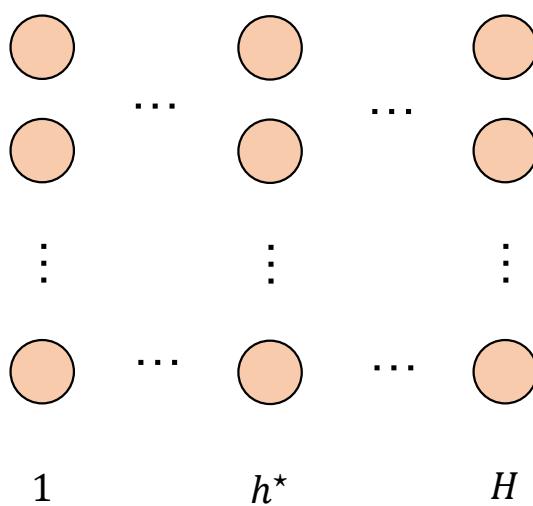
$$\pi'_h(\cdot|s) = \pi_h(\cdot|s) = \pi_{\text{out}}(\cdot|s) \text{ for } h > h^*$$

$$= \sum_{s,a} d_{\rho, h^*}^{\pi_{\text{in}}}(s) \pi'_{h^*}(a|s) Q_{h^*}^{\pi_{\text{out}}}(s, a) - \sum_{s,a} d_{\rho, h^*}^{\pi_{\text{in}}}(s) \pi_{h^*}(a|s) Q_{h^*}^{\pi_{\text{out}}}(s, a)$$

$$= \sum_{s,a} d_{\rho, h^*}^{\pi_{\text{in}}}(s) (\pi'_{h^*}(a|s) - \pi_{h^*}(a|s)) Q_{h^*}^{\pi_{\text{out}}}(s, a)$$

All-Step Policy Modification under Fixed Horizon

$$\sup_{\lambda} \mathbb{E} \left[V_1^{\lambda'}(s) - V_1^{\lambda}(s) \right] \text{ for arbitrary } \lambda', \lambda$$



Let $\pi^{(h)}$ be a Markov policy such that it is

$\left\{ \begin{array}{l} \text{same as } \pi' \text{ in steps 1 to } h-1 \\ \text{same as } \pi \text{ in steps } h \text{ to } H \end{array} \right.$

$\pi' = \pi^{(H+1)}$ and $\pi = \pi^{(1)}$

$$\pi^{(1)}, \pi^{(2)}, \pi^{(3)}.$$

$$\pi^{(H)}; \quad \bar{\pi}^{(H+I)};$$

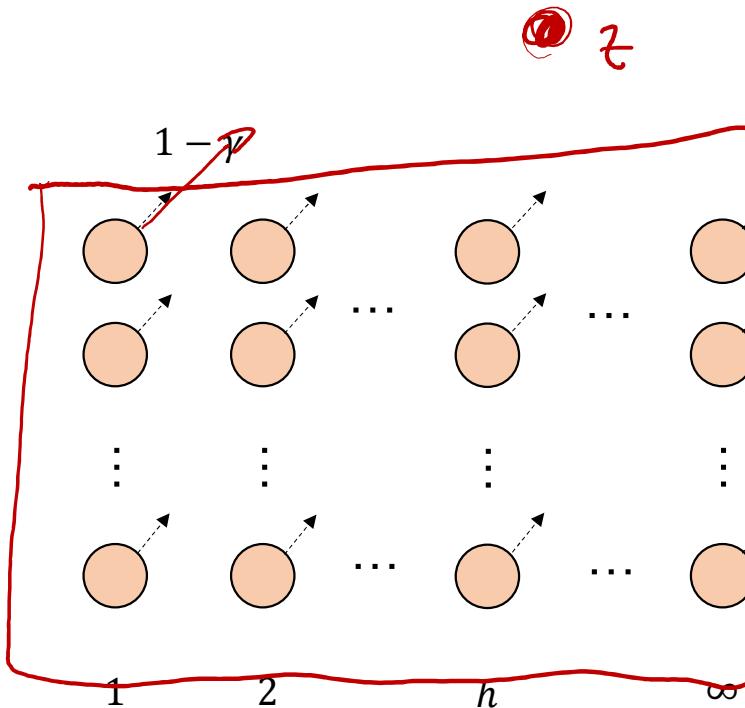
$$1 = E \left[V^{(t+1)}(s) - V^{(t)}(s) \right]$$

$$= \mathbb{E} \left[\sum_{h=1}^H \left(V^{\pi^{(h+1)}}(s) - V^{\pi^{(h)}}(s) \right) \right]$$

The diagram illustrates two sets of points, labeled x and x' , arranged in horizontal rows. The set x is on the left, and the set x' is on the right. Each set has several points, with specific points highlighted by red dots. Below the sets, indices h and $h+1$ are shown above groups of points. The index h is positioned above the second point from the left in both sets. The index $h+1$ is positioned above the third point from the left in both sets. Red arrows point from the labels h and $h+1$ to their respective points in both sets.

$$= \mathbb{E} \left[\sum_{h=1}^H \sum_{s,a} d_{ph}^{\pi'}(s) \left(\bar{\pi}_h'(a|s) - \pi_h(a|s) \right) Q_h^\pi(s,a) \right]$$

Discounted Total Reward Setting



Define Markov policy

$$\pi^{(h)} = \begin{cases} \pi' & \text{in step } 1 \rightarrow h-1 \\ \pi & \text{in step } h \rightarrow \infty \end{cases}$$

$$\pi^{(1)} = \pi, \quad \pi^{(\infty)} = \pi'$$

$$\begin{aligned} \mathbb{E}_{\pi'} \left[V^{\pi'}(s) - V^{\pi}(s) \right] &= \mathbb{E}_{\pi'} \left[\sum_{h=1}^{\infty} \gamma^{(h+1)} (V^{\pi'}(s) - V^{\pi}(s)) \right] \\ &= \sum_{h=1}^{\infty} \sum_{s,a} d_{\rho_h}^{\pi'}(s) \underbrace{(\pi'(a|s) - \pi(a|s))}_{\rho_h} Q^{\pi}(s,a) \\ &= \sum_{s,a} d_{\rho}^{\pi'}(s) (\pi'(a|s) - \pi(a|s)) Q^{\pi}(s,a) \end{aligned}$$

Performance / Value Difference Lemma

For any two stationary policies π' and π in the discounted total reward setting,

$$\begin{aligned}
 \underbrace{\mathbb{E}_{s \sim \rho} [V^{\pi'}(s)] - \mathbb{E}_{s \sim \rho} [V^\pi(s)]}_{\text{e.g. } \pi'(s) = \arg \max Q^\pi(s, a)} &= \sum_{s,a} d_\rho^{\pi'}(s) (\pi'(a|s) - \pi(a|s)) Q^\pi(s, a) \\
 &\quad \stackrel{Q^\pi(s, a) = \sum_a d_\rho^{\pi'}(s, a) \pi'(a|s)}{=} \sum_{s,a} d_\rho^{\pi'}(s) \pi'(a|s) Q^\pi(s, a) - \sum_s \underbrace{d_\rho^{\pi'}(s)}_{\sum_a d_\rho^{\pi'}(s, a)} V^\pi(s) \\
 &= \sum_{s,a} \underbrace{d_\rho^{\pi'}(s, a)}_{d_\rho^{\pi'}(s) \pi'(a|s)} (Q^\pi(s, a) - V^\pi(s)) \underbrace{\sum_{s,a} d_\rho^{\pi'}(s) \pi'(a|s)}_{V^\pi(s)}
 \end{aligned}$$

$$\begin{aligned}
 -\mathbb{E}_{s \sim \cdot} [V^\pi(s)] + \mathbb{E}_s [V^{\pi'}(s)] &= \sum_{s,a} d_\rho^{\pi'}(s) \left(\cancel{\pi(s,a)} - \cancel{\pi'(a|s)} \right) Q^{\pi'}(s, a) \\
 &\quad \stackrel{\cancel{\pi(s,a)} - \cancel{\pi'(a|s)} = A^{\pi'}(s, a)}{=} \sum_{s,a} d_\rho^{\pi'}(s, a) A^{\pi'}(s, a) \\
 &= Q^{\pi'}(s, a) - \sum_{a'} \cancel{\pi(a'|s)} Q^{\pi'}(s, a')
 \end{aligned}$$

Modified Policy Iteration $\underline{V} \rightarrow \underline{\mathcal{T}^\pi V}$

$$\begin{aligned} V^\pi &= \mathcal{T}^\pi V^\pi \quad (\text{Bellman eq}) \\ V^{(k+1)} &\leftarrow \mathcal{T}^\pi V^{(k)} \xrightarrow{\text{convergence}} V^{(k)} \rightarrow V^\pi \end{aligned}$$

Bellman Operator \mathcal{T}^π

$$\underline{(\mathcal{T}^\pi V)(s)} = \sum_a \pi(a|s) \left(R(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s') \right)$$

\checkmark

Greedy Policy Operator \mathcal{G}

$$(\mathcal{G}V)(s) = \operatorname{argmax}_a \left(R(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s') \right)$$

$$V_{k+1}(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} p(s'|s, a) V_k(s') \right)$$

Value Iteration:

$$\pi_{k+1} = \mathcal{G}V_k$$

$$V_{k+1} = \mathcal{T}^{\pi_{k+1}} V_k$$

Policy Iteration:

$$\pi_{k+1} = \mathcal{G}V_k = \underline{V^\pi} = \underline{V^{\pi_{k+1}}}$$

$$V_{k+1} = \underline{(\mathcal{T}^{\pi_{k+1}})^\infty V_k}$$

MPI:

$$\pi_{k+1} = \mathcal{G}V_k$$

$$V_{k+1} = (\mathcal{T}^{\pi_{k+1}})^m V_k$$

Policy update

Value update

Difference:
Relative speed
between policy and
value updates

Summary for the Basics of MDPs

- MDPs model decision-making problems where the return depends on sequences of actions.
- “State” summarizes all the information needed to make decisions (in the fixed-horizon setting, the step index is also important).
- Interaction Protocols: fixed-horizon, goal-oriented, infinite-horizon
- Performance Metrics: total reward, average reward, discounted total reward
- Policies: history-dependent, Markov, stationary
- While the **number of action sequence is exponential** in the horizon length, the optimal policy can be computed in **poly(#state, #actions, horizon length)** time using dynamic programming techniques (Value Iteration).
- The dynamic programming here is slightly more complicated since it involves infinite horizon and recursive states.
- Bellman equation, Reverse Bellman equation, Bellman optimality equation
- Approximate Bellman optimality → Approximate optimal policy
- Policy Iteration and Performance Difference Lemma
- Unifying Value Iteration and Policy Iteration

Summary for the Basics of MDPs

- MDPs model decision-making problems where the return depends on sequences of actions.
- “State” summarizes all the information needed to make decisions (in the fixed-horizon setting, the step index is also important).
- Interaction Protocols: fixed-horizon, goal-oriented, infinite-horizon
- Performance Metrics: total reward, average reward, discounted total reward
- Policies: history-dependent, Markov, stationary
- While the **number of action sequence is exponential** in the horizon length, the optimal policy can be computed in **poly(#state, #actions, horizon length)** time using dynamic programming techniques (**Value Iteration**).
- The dynamic programming here is slightly more complicated since it involves infinite horizon and recursive states.
- Bellman equation, Reverse Bellman equation, Bellman optimality equation
- Approximate Bellman optimality → Approximate optimal policy
- **Policy Iteration** and Performance Difference Lemma
- Unifying Value Iteration and Policy Iteration