



MARKOV DECISION PROCESSES: VALUE ITERATION, POLICY ITERATION AND LINEAR PROGRAMMING METHODS

done by

Jean-Claude S. MITCHOZOUNOU

AFRICA BUSINESS SCHOOL, UM6P, MOROCCO

MASTER OF QUANTITATIVE AND FINANCIAL MODELLING (QFM)

PROJECT OF SECOND SEMESTER

Supervisor: **PROF. OMAR SAADI**

Academic Year 2023-2024

CONTENTS

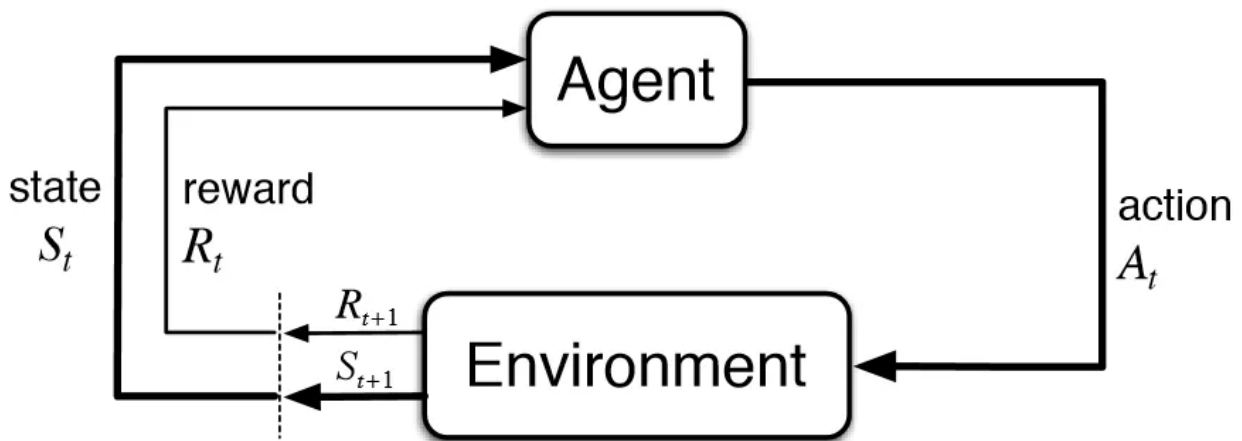
1	Reinforcement Learning and Markov Decisions processes	1
1.1	Overviews on Markov decisions processes	1
1.2	Mathematicals concepts in MDPs	2
1.2.1	Defining MDPs	2
1.2.2	Policies and Discounted Total Reward	3
1.2.3	Objective Of MDPs:	3
1.2.4	Markov Property	4
1.2.5	State Value Function of Policies	4
1.2.6	Action Value Functions of Policies	5
1.2.7	Bellman Equations	5
2	Numerical study of solving methods of MDPs	9
2.1	Value Iteration	9
2.1.1	Numerical results and interpretation	10
2.2	Policy Iteration	14
2.2.1	Numerical results and interpretation	15
2.2.2	Test of Policy iteration methods on sparse matrix	18
2.3	Linear Programming Approach	19
2.3.1	Numerical results and interpretation	21
2.4	Python code of the methods	24
3	Conclusion	26

Reinforcement Learning and Markov Decisions processes

1.1 Overviews on Markov decisions processes

Reinforcement Learning (RL) and Markov Decision Processes (MDPs) are closely related concepts, and MDPs provide a formal mathematical framework that underpins much of RL theory and practice.

Reinforcement Learning is the science of decision-making. It is about learning the optimal behavior in an environment to obtain maximum reward. MDP gives us a way to formalize sequential decision-making. In other way, a sequential decision problem for a fully observable, stochastic environment with a Markovian transition model and additive rewards is called a Markov decision process, or MDP, and consists of a set of states (with an initial state s_0); a set $ACTIONS(s)$ of actions in each state; a transition model $P(s' | s, a)$; and a reward function $R(s)$. This formalization is the basis for structuring problems that are solved with reinforcement learning.



In an MDP, we have a decision maker, called an agent, that interacts with the environment it's placed in as it is represented in the figure above. These interactions occur sequentially over time. At each time step, the agent will get some representation of the environment's state. Given this representation, the agent selects an action to take. The environment is then transitioned into a new state, and the agent is given a reward as a consequence of the previous action. Lets now discuss about the components of RL used in MDP:

1. **Agents** : Agents are a central concept in the field of RL. An RL agent is an autonomous entity or computational system that interacts with an environment to learn and make decisions in order to achieve specific goals or maximize cumulative rewards. These agents

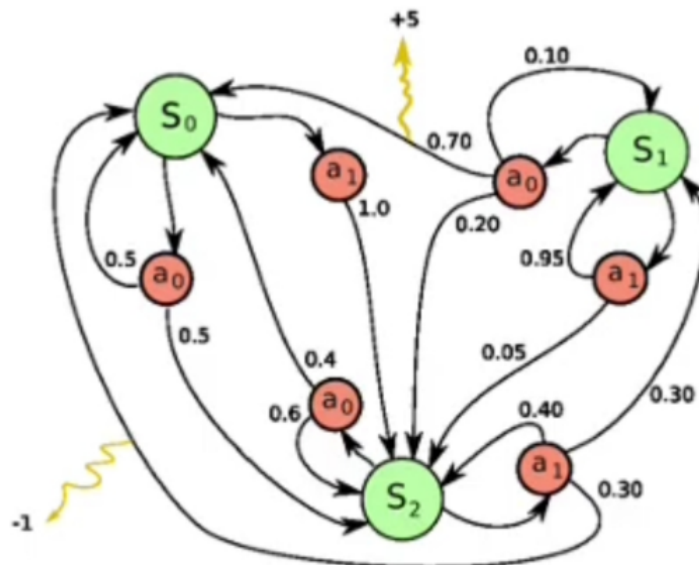
can be implemented in various domains, from robotics to game playing to recommendation systems. For instance, an agent can be a player in the game whose motive is to win the game or gain maximum rewards.

2. **Environment:** refers to what the world looks like and the rules of the world. It is a well-defined and structured system that an RL agent interacts with over time. It represents the world or domain in which the agent's actions have consequences, and it provides feedback to the agent based on those actions.
3. **States (S):** refers to the state of our agent, not the state of our environment. The environment has a set of possible states that represent different situations, configurations, or conditions it can be in. In MDP we use s to denote state and in our example we call them $s = 0, 1, 2, \dots$
4. **Action (A):** refers to the decisions that can be taken by RL agent in each state. The set of possible actions is not necessarily the same for each state.
5. **Rewards (R):** After each action, the environment provides the RL agent with a numerical signal called a reward. The reward indicates the immediate desirability or quality of the action taken. The agent's goal is often to maximize the cumulative reward over time.
6. **Trajectory:** The process of selecting an action from a given state, transitioning to a new state, and receiving a reward happens sequentially over and over again, and creates something called a trajectory that shows the sequence of states, actions, and rewards
7. **Policy:** is a decision-making strategy in which the agent chooses actions adaptively based on the history of observations;

1.2 Mathematical concepts in MDPs

1.2.1 Defining MDPs

In reinforcement learning, the interactions between the agent and the environment are often described by an infinite-horizon, discounted Markov Decision Process (MDP) $M = (\mathcal{S}, \mathcal{A}, P, r, \gamma, \mu)$, specified by:



- A state space \mathcal{S} , which may be finite or infinite. For mathematical convenience, we will assume that \mathcal{S} is finite or countably infinite.
- An action space \mathcal{A} , which also may be discrete or infinite. For mathematical convenience, we will assume that \mathcal{A} is finite.
- A transition function $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$, where $\Delta(\mathcal{S})$ is the space of probability distributions over \mathcal{S} (i.e., the probability simplex). $P(s' | s, a)$ is the probability of transitioning into state s' upon taking action a in state s . We use $P_{s,a}$ to denote the vector $P(\cdot | s, a)$.
- A reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$. $r(s, a)$ is the immediate reward associated with taking action a in state s . More generally, $r(s, a)$ could be a random variable (where the distribution depends on s, a).
- A discount factor $\gamma \in [0, 1)$, which defines a horizon for the problem.
- An initial state distribution $\mu \in \Delta(\mathcal{S})$, which specifies how the initial state s_0 is generated. In many cases, we will assume that the initial state is fixed at s_0 , i.e., μ is a distribution supported only s_0 .

1.2.2 Policies and Discounted Total Reward

In a given MDP $M = (\mathcal{S}, \mathcal{A}, P, r, \gamma, \mu)$, the agent interacts with the environment according to the following protocol:

The agent starts at some state $s_0 \sim \mu$; at each time step $t = 0, 1, 2, \dots$, the agent takes an action $a_t \in \mathcal{A}$, obtains the immediate reward $r_t = r(s_t, a_t)$, and observes the next state s_{t+1} sampled according to $s_{t+1} \sim P(\cdot | s_t, a_t)$.

The interaction record at time t is: $\tau_t = (s_0, a_0, r_0; s_1, a_1, r_1, \dots, s_t, a_t, r_t)$ is called a trajectory, which includes the observed state at time t .

A policy is a (possibly randomized) mapping from a trajectory to an action, i.e. $\pi : H \rightarrow \Delta(\mathcal{A})$ where H is the set of all possible trajectories (of all lengths) and $\Delta(\mathcal{A})$ is the space of probability distributions over \mathcal{A} . A stationary policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ specifies a decision-making strategy in which the agent chooses actions based only on the current state, i.e. $a_t \sim \pi(\cdot | s_t)$. A deterministic, stationary policy is of the form $\pi : \mathcal{S} \rightarrow \mathcal{A}$.

For a given: policy π
initial state s

SARSA sequence: $s, \pi(s), r(s, \pi(s)), s \rightarrow s', \pi(s'), r(s', \pi(s')), \dots, \infty$

State sequence: $s, s', s'', \dots, \infty$

Reward sequence: $r(s, \pi(s)), r(s', \pi(s')), r(s'', \pi(s'')), \dots, \infty$

Discounted total reward:

$$\begin{aligned} R(s, \pi) &= r(s, \pi(s)) + \gamma r(s', \pi(s')) + \gamma^2 r(s'', \pi(s'')) + \dots, \infty \\ &= \text{immediate reward} + \text{discounted future rewards} \end{aligned}$$

1.2.3 Objective Of MDPs:

The objective in an MDP is typically to find a policy π^* that maximizes the expected cumulative reward over time. This is often expressed as the expected return

$$E(R(s, \pi)) = E \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$

1.2.4 Markov Property

Andrey Andreyevich Markov (1856-1922)
Russian mathematician known for his work on
stochastic processes



The Markov Property: given the present, the future and the past are independent

- i.e., Everything you need to know about the past is included in the present state

For MDPs, the Markov property means that:

$$P(S_{t+1} = s' \mid S_t = s_t, A_t = a_t, S_{t-1} = s_{t-1}, A_{t-1} = a_{t-1}, \dots, S_0 = s_0) = P(S_{t+1} = s' \mid S_t = s_t, A_t = a_t)$$

Independent of past states and actions

1.2.5 State Value Function of Policies

For more information on State Value Function, please see ([1, 5, 13].)
Given MDP (S, A, P, R, γ) :

Value of a state s under policy π :

$V^\pi(s)$ = expected utility starting in s and acting according to π

$$V^\pi(s) = \mathbb{E} \left(\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s \right)$$

Sequence of rewards generated by following π

Optimal value of a state s :

$V^*(s)$ = expected utility starting in s and acting optimally

$$V^*(s) = V^{\pi^*}(s) = \mathbb{E} \left(\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi^*, s_0 = s \right)$$

Rewards generated by following π^* (optimal policy.)

1.2.6 Action Value Functions of Policies

For more information on action value Function, please see ([1, 5, 13]).

It is also helpful to define action-value functions.

Q-value of taking action a in the state s then following policy π :

$Q^*(s, a)$ = expected utility taking action a in the state s and then following policy π

$$Q^\pi(s, a) = \mathbb{E} \left(\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s, a_0 = a \right)$$

• The optimal Q-value of taking action a in the state s is the Q-value optaine by following the policy associated to V^*

• π^* can be greedily determined from Q^* : $\pi^*(s) = \underset{a \in A}{argmax} Q^*(s, a)$

Remark 1.1. Since $r(s, a)$ is bounded between 0 and 1, we can easily show that

$$0 \leq V^\pi(s) \leq \frac{1}{1-\gamma} \text{ and } 0 \leq Q^\pi(s, a) \leq \frac{1}{1-\gamma}.$$

1.2.7 Bellman Equations

The information and analyses presented in this section are based on previous work and research available in the literature, see([1, 13]).

Bellman's equations represent important results in MDPs. These results will be of great use throughout our paper. In particular, these equations will help us to implement the various methods for finding the optimal policy.

Lemma 1.1. Suppose that π is a stationary policy. Then V^π and Q^π satisfy the following equations: for all $s \in S$, $a \in A$,

$$V^\pi(s) = Q^\pi(s, \pi(s)).$$

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^\pi(s')].$$

Bellman consistency equations for stationnary policies

$$V^\pi(s) = Q^\pi(s, \pi(s)) \quad (i).$$

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^\pi(s')] \quad (ii).$$

Proof. (i)- By definition one has:

$$Q^\pi(s, a) = \mathbb{E} \left(\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s, a_0 = a \right)$$

. So for $a = \pi(s)$ one obtains:

$$\begin{aligned} Q^\pi(s, \pi(s)) &= \mathbb{E} \left(\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s, a_0 = \pi(s) \right) \\ &= \mathbb{E} \left(\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s \right) \text{ because } \pi \text{ is stationnary and we're starting at } s. \\ &= V^\pi(s) \end{aligned}$$

(ii)-

$$\begin{aligned}
Q^\pi(s, a) &= \mathbb{E} \left(\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s, a_0 = a \right) \\
&= \mathbb{E} \left(r(s_0, s_0) + \sum_{t=1}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s, a_0 = a \right) \\
&= r(s, a) + \mathbb{E} \left(\sum_{t=1}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s, a_0 = a \right) \\
&= r(s, a) + \gamma \mathbb{E} \left(\sum_{t=1}^{\infty} \gamma^{t-1} r(s_t, a_t) \mid \pi, s_0 = s, a_0 = a \right) \\
&= r(s, a) + \gamma \mathbb{E} \left(\sum_{t=1}^{\infty} \gamma^{t-1} r(s_{t+1}, a_{t+1}) \mid \pi, s_0 = s, a_0 = a \right) \\
&= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\mathbb{E} \left(\sum_{t=1}^{\infty} \gamma^{t-1} r(s_{t+1}, a_{t+1}) \mid \pi, s_1 = s' \right) \right] \\
&\text{by using total expected formula: } \mathbb{E}(X) = \mathbb{E}(\mathbb{E}(X|Y)) \\
&= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\mathbb{E} \left(\sum_{t=1}^{\infty} \gamma^{t-1} r(s_t, a_t) \mid \pi, s_0 = s' \right) \right] \\
&\text{by using homogeneity of Markov chain} \\
&= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V^\pi(s')]
\end{aligned}$$

Hence the results.

Remark 1.2. In matrix form, Q^π can be written as:

$$Q^\pi = r + \gamma P V^\pi$$

If we define P^π to be the transition matrix on state-action pairs induced by a stationary policy π , specifically:

$$P_{(s,a),(s',a')}^\pi := P(s'|s, a) \pi(a'|s').$$

In particular, for deterministic policies π defined by:

$$\pi(a'/s') := \begin{cases} 1 & \text{if } \pi(s') = a' \\ 0 & \text{if not} \end{cases};$$

we have:

$$P_{(s,a),(s',a')}^\pi := \begin{cases} P(s'|s, a) & \text{if } a' = \pi(s') \\ 0 & \text{if } a' \neq \pi(s') \end{cases}$$

With this notation, it is straightforward to verify:

$$Q^\pi = r + \gamma P V^\pi$$

$$Q^\pi = r + \gamma P^\pi Q^\pi.$$

Indeed, $\forall (s, a), (s', a') \in S \times A$

$$\begin{aligned}
Q^\pi(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V^\pi(s')] \\
&= r(s, a) + \gamma \sum_{s'} P(s'|s, a) V^\pi(s') \\
&= r(s, a) + \gamma (P V^\pi)(s, a) \\
&\iff Q^\pi = r + \gamma P V^\pi
\end{aligned}$$

and

$$\begin{aligned}
(P^\pi Q^\pi)(s, a) &= \sum_{s', a'} P_{(s, a), (s', a')}^\pi Q^\pi(s', a') \\
&= \sum_{s', a'} P(s'|s, a) \pi(a'|s') Q^\pi(s', a') \\
&= \sum_{s'} P(s'|s, a) \pi(\pi(s')|s') Q^\pi(s', \pi(s')) \\
&= \sum_{s'} P(s'|s, a) V^\pi(s') \\
&= (PV^\pi)(s, a)
\end{aligned}$$

$$\text{i.e. } PV^\pi = P^\pi Q^\pi \iff r + \gamma P^\pi Q^\pi = r + \gamma PV^\pi = Q^\pi$$

Corollary 1.1. *Suppose that π is a stationary policy. We have that:*

$$Q^\pi = (I - \gamma P^\pi)^{-1} r$$

where I is the identity matrix.

Proof. Since P is Markovian, one has $\|Px\|_\infty \leq \|x\|_\infty \forall x$. And further $\gamma \in [0, 1)$. Consequently, $(I - \gamma P^\pi)$ is invertible.

Theorem 1.1. *Let Π be the set of all non-stationary and randomized policies. Define:*

$$\begin{aligned}
V^*(s) &:= \sup_{\pi \in \Pi} V^\pi(s) \\
Q^*(s, a) &:= \sup_{\pi \in \Pi} Q^\pi(s, a).
\end{aligned}$$

which is finite since $V^\pi(s)$ and $Q^\pi(s, a)$ are bounded between 0 and $\frac{1}{1-\gamma}$.

There exists a stationary and deterministic policy π^+ such that for all $s \in S$ and $a \in A$,

$$\begin{aligned}
V^{\pi^*}(s) &= V^*(s) \\
Q^{\pi^*}(s, a) &= Q^*(s, a).
\end{aligned}$$

We refer to such a π^+ as an optimal policy.

Proof. Indication:

Theorem 1.2. (*Bellman optimality equations*) *We say that a vector $Q \in \mathbb{R}^{|S||A|}$ satisfies the Bellman optimality equations if:*

$$Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\max_{a' \in A} Q(s', a') \right].$$

For any $Q \in \mathbb{R}^{|S||A|}$, Q is said to be optimal ($Q = Q^*$) if and only if Q satisfies the Bellman optimality equations.

To proof this theorem we will use the following lemma.

Lemma 1.2. $V^*(s) = \max_{a \in A} Q^*(s, a)$

Definition 1.2.1. Let π_Q denote the greedy policy with respect to a vector $Q \in \mathbb{R}^{|S||A|}$, i.e.,

$$\pi_Q(s) := \arg \max_{a \in A} Q(s, a).$$

With this notation, by the above theorem, the optimal policy π^* is given by:

$$\pi^* = \pi_{Q^*}.$$

Let us also use the following notation to turn a vector $Q \in \mathbb{R}^{|S||A|}$ into a vector of length $|S|$:

$$V_Q(s) := \max_{a \in A} Q(s, a).$$

Definition 1.2.2. The Bellman optimality operator $T_M : \mathbb{R}^{|S||A|} \rightarrow \mathbb{R}^{|S||A|}$ is defined as:

$$TQ := r + \gamma PV_Q.$$

This allows us to rewrite the Bellman optimality equation in the concise form:

$$Q = TQ,$$

and, so, the previous theorem states that $Q = Q^*$ if and only if Q is a fixed point of the operator T .

Remark 1.3. The equation $Q = TQ$ has solution because the operator T is γ -Lipchitz with $\gamma \in [0, 1)$. Indeed, for any two vectors $Q, Q' \in \mathbb{R}^{|S||A|}$,

$$\|TQ - TQ'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

Proof. First, let us show that for all $s \in S$, $|V_Q(s) - V_{Q'}(s)| \leq \max_{a \in A} |Q(s, a) - Q'(s, a)|$. Assume $V_Q(s) > V_{Q'}(s)$ (the other direction is symmetric), and let a be the greedy action for Q at s . Then

$$\begin{aligned} |V_Q(s) - V_{Q'}(s)| &= Q(s, a) - \max_{a' \in A} Q'(s, a') \\ &\leq Q(s, a) - Q'(s, a) \\ &\leq \max_{a \in A} |Q(s, a) - Q'(s, a)|. \end{aligned}$$

Using this,

$$\begin{aligned} \|TQ - TQ'\|_\infty &= \gamma \|PV_Q - PV_{Q'}\|_\infty \\ &= \gamma \|P(V_Q - V_{Q'})\|_\infty \\ &\leq \gamma \|V_Q - V_{Q'}\|_\infty \\ &= \gamma \max_s |V_Q(s) - V_{Q'}(s)| \\ &\leq \gamma \max_s \max_a |Q(s, a) - Q'(s, a)| \\ &= \gamma \|Q - Q'\|_\infty \end{aligned}$$

Numerical study of solving methods of MDPs

To solve an MPDs, we will use severals methods such as **Value Iteration**, **Policyiteration** and **Linear Programming algorithms**. Here is their theoretical complexity.

Algorithms	Poly?	Strongly Poly?
Value Iteration	$\frac{ S ^2 A }{1-\gamma} L(P, r, \gamma) \log \frac{1}{1-\gamma}$	\times
Policy Iteration	$\frac{ S ^3+ S ^2 A }{1-\gamma} L(P, r, \gamma) \log \frac{1}{1-\gamma}$	$(S ^3 + S ^2 A) \cdot \min \left(\frac{ A ^{ S }}{ S }, \frac{ S ^2 A }{1-\gamma} \log \frac{ S ^2}{1-\gamma} \right)$
LP-Algorithms	$ S ^3 A L(P, r, \gamma)$	$ S ^4 A ^4 \log \frac{ S }{1-\gamma}$

Table 2.1: Computational complexities of various approaches

2.1 Value Iteration

Value iteration algorithm is very well detailed in the various studies and research available in the literature as indicated in the following references (see [1, 2, 3, 7, 8, 9, 6]).

Complexity: For a fixed value of γ , Value iteration algorithm is polynomial time algorithm. Here $|S|^2|A|$ (see Tab 2.1) is the assumed runtime per iteration of value iteration

Value iteration updates the value function V using the Bellman optimality equation (Theorem 1.2). It is an algorithm used to compute the optimal policy and the value function of an MDP. It iteratively updates the value of each state until the values converge to the optimal values.

Starting at some $Q(Q = 0, \text{ generally})$, we iteratively apply $T(\text{Bellman operator})$:

$$Q \leftarrow TQ,$$

untill convergence

Pseudo-code: Value iteratoin

- 1: **input:** transition probabilities P , reward function r , Q discount factor γ and θ ,
- 2: **output:** optimal policy π^* , value function V
- 3: **initialize:** $Q(s, a) = 0, \forall (s, a) \in S \times A$
- 4: **repeat**
- 5: $\Delta \leftarrow 0$
- 6: $TQ = \text{BellmanOperator}(Q, R, P, \gamma)$
- 7: $\delta = \max(|Q - TQ|)$
- 8: $Q = TQ$
- 9: **if** $\Delta \leq \theta$ **then**
- 10: **break**

```

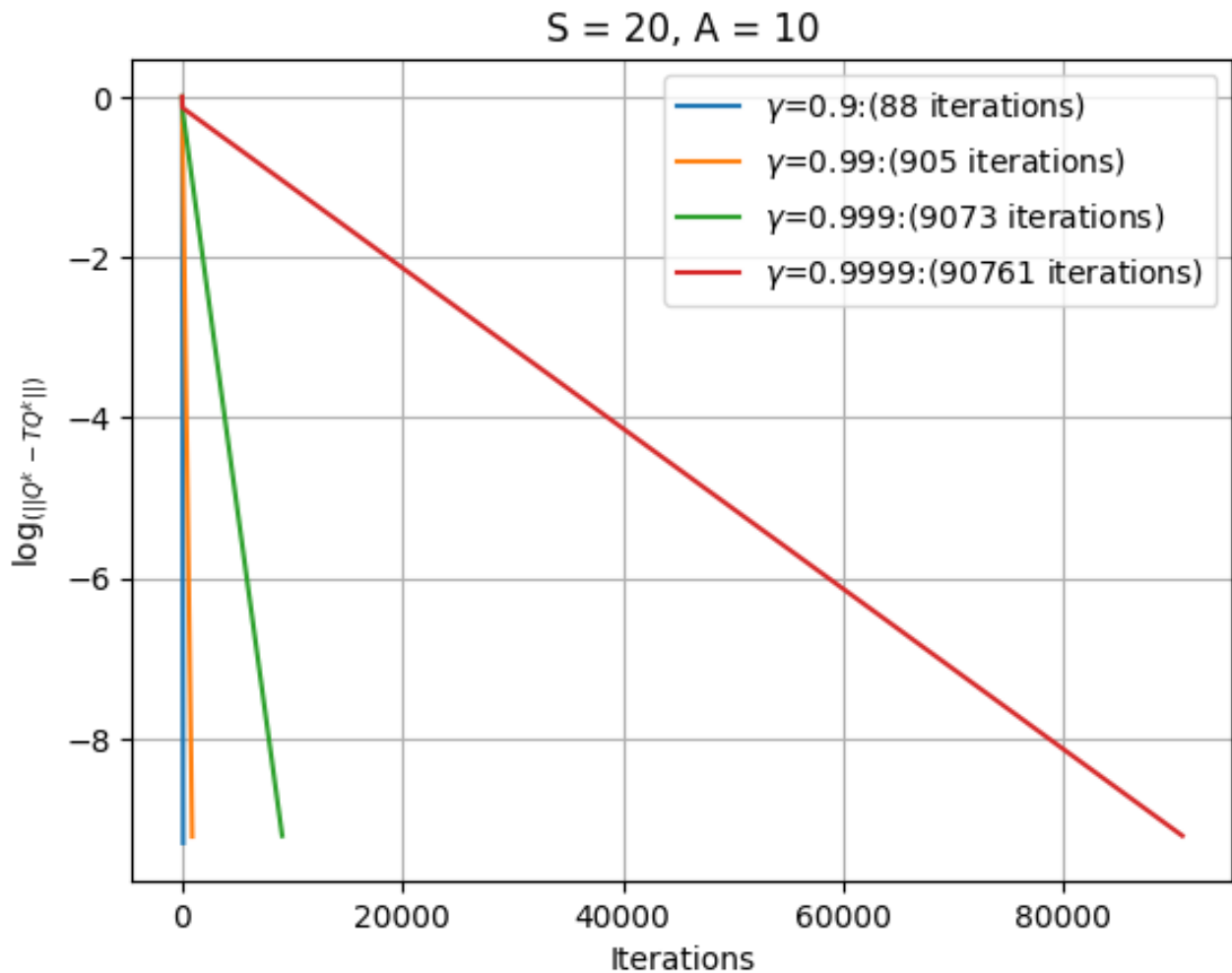
11:   end if
12: until the policy is optimal and the function converges
13: for  $s \in S$  do
14:    $V^*(s) = \max_{a \in A} Q(s, a)$ 
15:    $\pi^*(s) = \arg \max_{a \in A} Q(s, a)$ 
16: end for
17: return  $V^*, \pi^*$ 

```

At least, we have the optimal policy π^* and optimal value function V^* .

2.1.1 Numerical results and interpretation

The effects of γ on convergence.



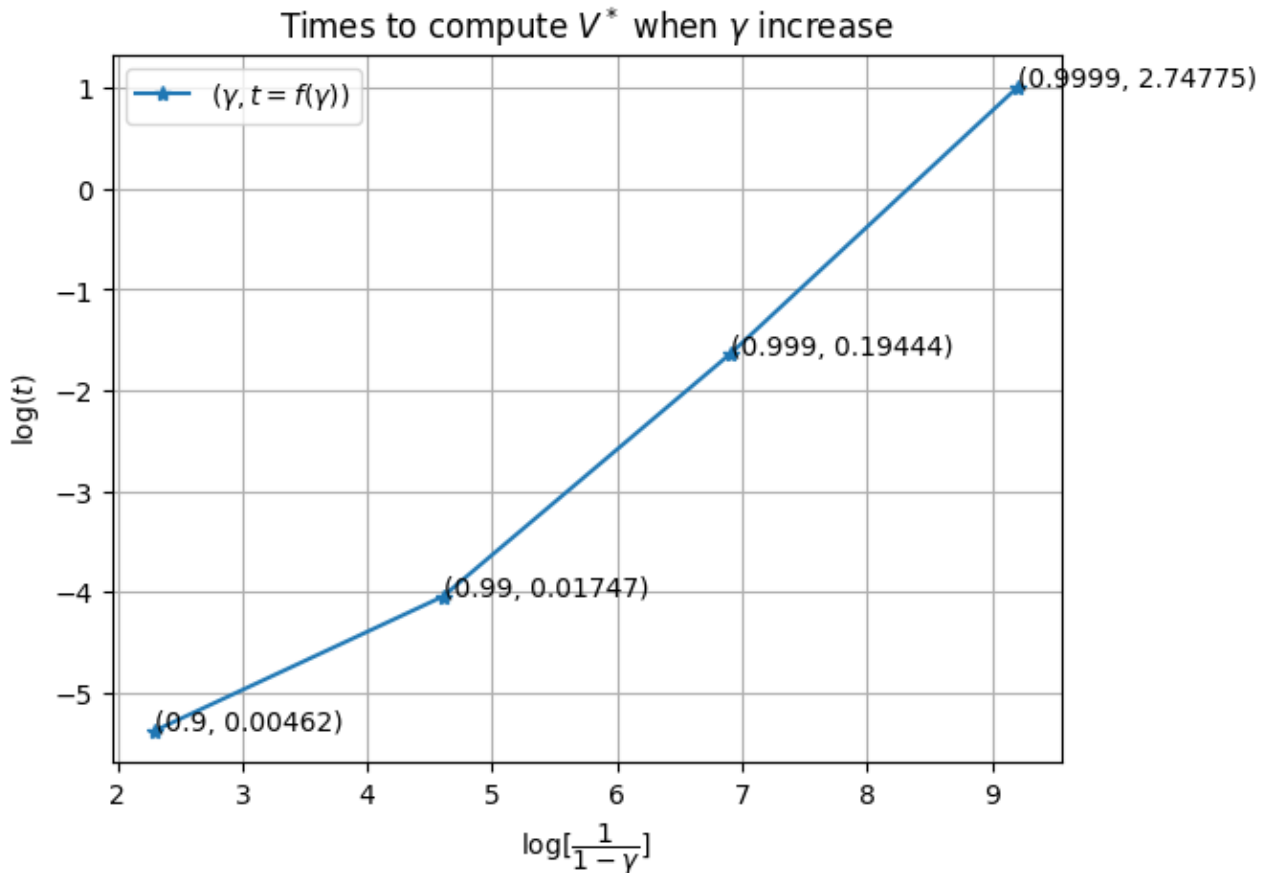
1. γ close to 0:

- The Value Iteration algorithm converges quickly because it does not look far into the future.
- When γ is close to 0, immediate rewards are much more important than future rewards.

2. γ close to 1:

- We can see that as we get closer to 1, the number of iterations is multiplied by 10. This means that the algorithm becomes 10 times slower in terms of convergence as we get closer to 1.

- The solutions found may be more optimal in the long term, but the convergence time can be significantly longer as you can see on the following figure which plots the evolution of the solution's calculation time as a function of the evolution of gamma

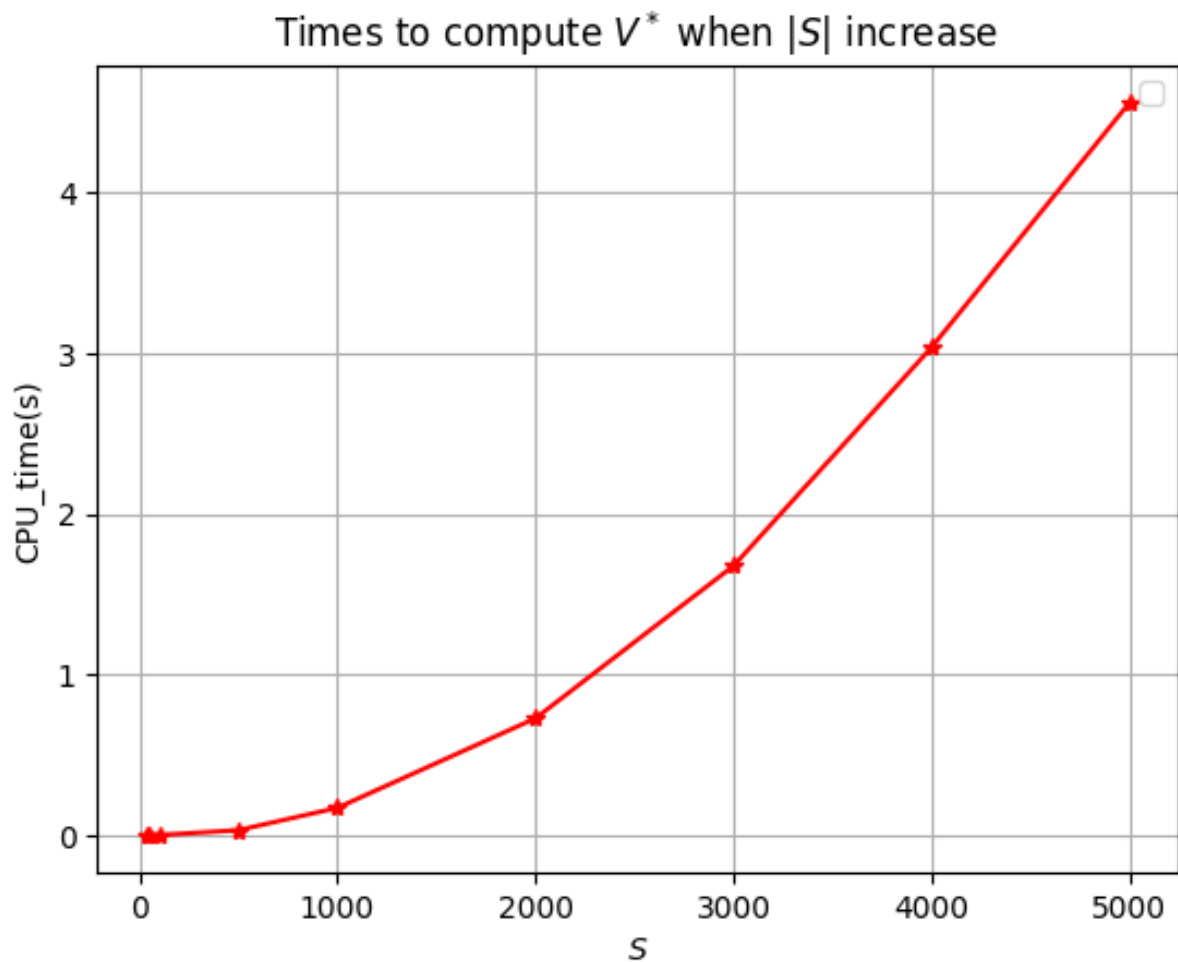


Note that, as before, the closer you get to 1, the calculation time is also multiplied by 10.

Impact of the Number of States on Convergence Speed

The number of iterations remains unchanged despite the increasing number of states in the MDPs. One might think that the variation in the number of states has no impact on the convergence of the method, but far from it. This constancy in the number of iterations in relation to the number of states hides a number of problems ranging from **computational complexity**, **convergence time** and **storage memory**.

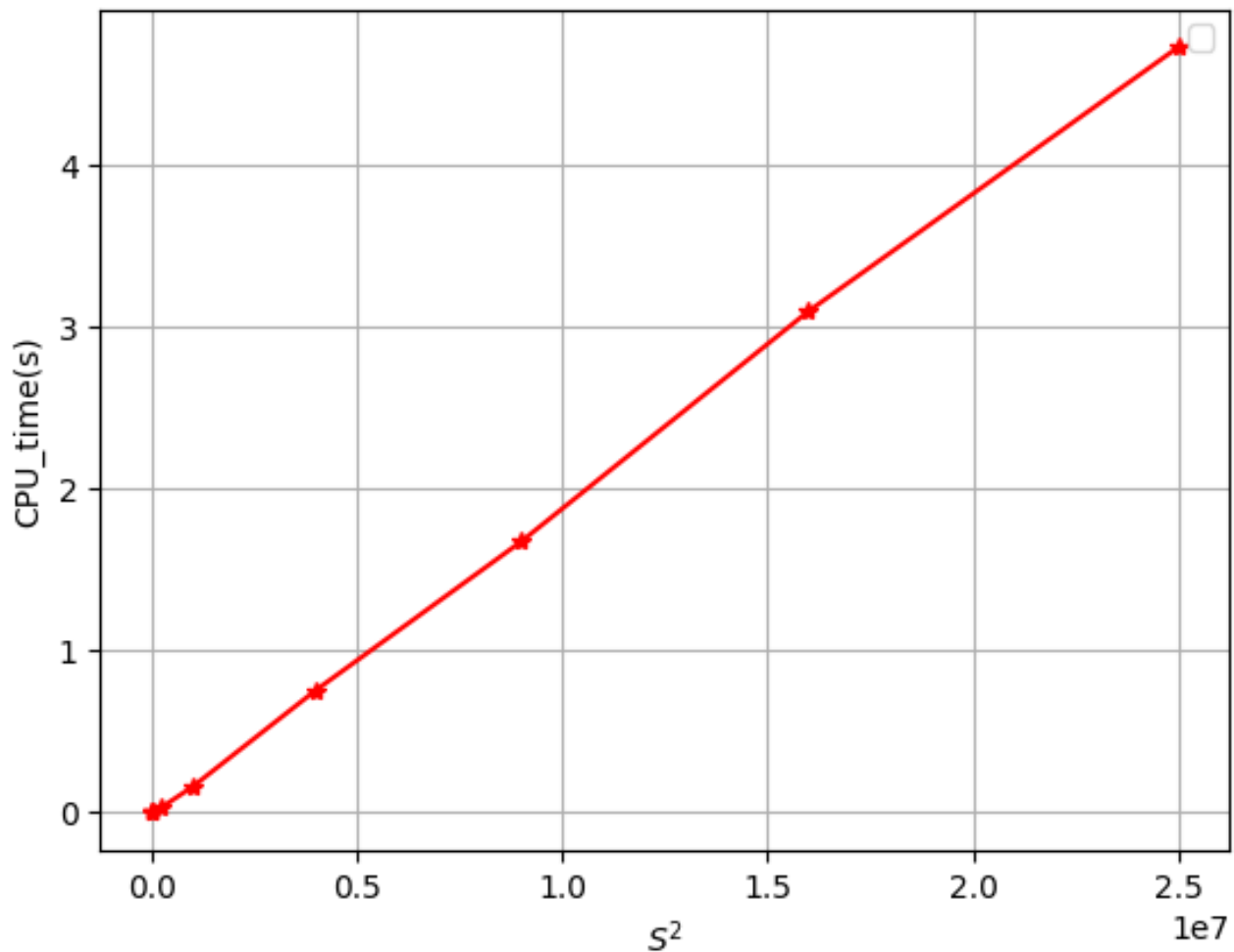
1. **Computational Complexity:** As the number of states increases, the number of computations required to update state values or policies increases. This is because each state must be updated based on all possible actions and transitions to all other possible states.
2. **Convergence Duration:**
 - **Value Propagation:** In Value Iteration algorithms, state values $V(s)$ need to propagate through the entire state space. More states mean it takes longer for values to propagate and converge to optimal values.
 - **Bellman Updates:** Each update according to the Bellman equation requires calculations that include all states and all actions, meaning the computation time for each iteration increases with the number of states as you can see in the following figure.



3. Memory and Storage:

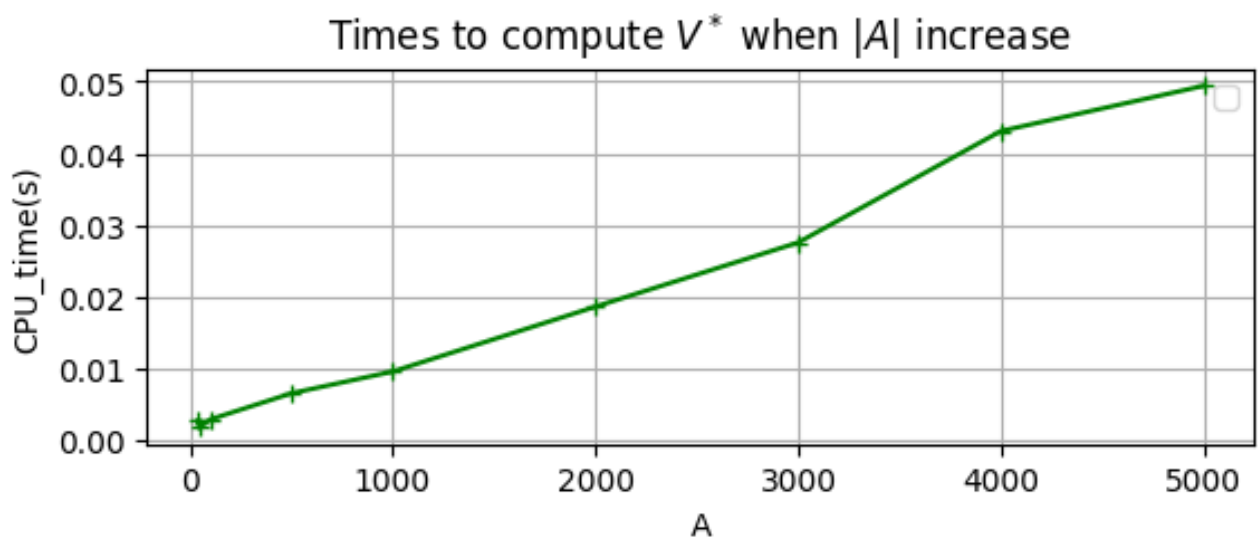
- **Value and Policy Tables:** Storing value ($V(s)$) and policy ($\pi(s)$) tables increases linearly with the number of states. More states require more memory, which can also affect computation efficiency if memory resources are limited.

By observing the curve of the evolution of the calculation time of v as the number of states increases, we can see that this evolution describes the shape of a parabolic function. More precisely, it behaves almost like the function $(x \mapsto x^2)$ shown in the comparison figure below



Impact of the Number of Actions on Convergence Speed

Increasing the number of actions in a Markov Decision Process (MDP) while keeping the number of states fixed has a significant impact on the computation time required to determine the optimal value. The relationship between the number of actions and the computation time can be described as follows:



The figure illustrates the near-linear increase in computation time as the number of actions increases. In other words, the resulting curve is almost linear, indicating that as the number

of actions increases, the computation time increases proportionally. We can therefore say that this curve behaves almost like the straight line with equation $(y = x)$.

2.2 Policy Iteration

Policy iteration algorithm is also very well detailed in the many studies and research available in the literature as indicated in the following references (see [1, 2, 3, 7, 8, 6]).

Complexity: For a fixed value of γ , Policy iteration algorithm is both polynomial and strongly polynomial time algorithm. Here $|S|^3 + |S|^2|A|$ (see Tab 2) is the assumed runtime per iteration of policy iteration

The policy iteration for discounted MDPs starts from an arbitrary policy π_0 and repeat the following steps until convergence:

1- **Policy evaluation** : keep current policy π fixed, find $V^{\pi_k}(\cdot)$

Here, we have two way to find $V^{\pi_k}(\cdot)$.

• *a*– **Exact computation of $V^{\pi_k}(\cdot)$.**

We use:

$$V^{\pi_k} = (I - \gamma P^{\pi_k})^{-1} r^{\pi_k} \text{ with } P_{s,s'}^{\pi_k} = P(s'/s, \pi_k(s))$$

Here we have to be very careful: the matrix P^{π} is a 2-dimensional matrix and r^{π} is a vector. In fact, the matrix P as we know it contains all the probabilities of taking any action in state s to transit to a state s' and the matrix r contains all the rewards associated with any action taken while in state s . These objects are obtained by extraction. This is a direct consequence of lemma 1.1 and corollary 1.1.

Remark: Using this method can be very costly as the number of states increases, since the greater the size of P^{π} , the more complicated it is to calculate its inverse, which is of complexity $\mathcal{O}(|S|^3)$. In this case, it's better to use the following iterative method to find V^{π}

• *b*– **Approximate computation of $V^{\pi_k}(\cdot)$** Iterate simplified Bellman update until values converge:

$$V_{i+1}^{\pi_k}(s) = r(s, \pi_k(s)) + \gamma \sum_{s'} P(s'/s, \pi(s)) V_i^{\pi_k}(s')$$

Here we assume that $V_0^{\pi_0}(s) = 0$

2- **Policy improvement**: find the best action for $V^{\pi_k}(\cdot)$ via one-step lookahead

$$\forall s \in \mathbb{S}, \pi_{k+1}(s) = \underset{a \in \mathbb{A}}{\operatorname{argmax}} \left[r(s, a) + \sum_{s'} P(s'/s, a) V^{\pi_k}(s') \right].$$

To implement policy iteration, we need functions for both policy evaluation and policy improvement.

Approximate Policy evaluation

- 1: **input:** reward function r , transitional model P , discounted factor γ , convergence threshold θ , policy π_k , initialized vector $v_0^{\pi_k}$
- 2: **output:** converged value v
- 3: **while** True **do**

```

4:    $\Delta \leftarrow 0$ 
5:   for  $s \in \mathcal{S}$  do
6:      $temp \leftarrow v(s)$ 
7:      $v_{i+1}^{\pi_k}(s) \leftarrow r(s, \pi_k(s)) + \gamma \sum_{s'} P(s'|s, \pi(s)) v_i^{\pi_k}(s')$ 
8:      $\Delta \leftarrow \max(\Delta, |temp - v_{i+1}^{\pi_k}|)$ 
9:   end for
10:  if  $\Delta < \theta$  then
11:    break
12:  end if
13: end while
14:  $v \leftarrow v_{i+1}^{\pi_k}$ 
15: return  $v$ 
16:

```

Policy improvement

```

1: input: transitional model  $P$ , reward function  $r$ , vector  $v$ 
2: output: updated policy  $\pi$ 
3: for  $s \in \mathcal{S}$  do
4:    $\pi(s) \leftarrow \arg \max_a [r(s, a) + \gamma \sum_{s'} P(s'|s, a) v(s')]$ 
5: end for
6: return  $\pi$ 
7:

```

Main body of the algorithm

```

1: input: reward function  $r$ , transitional model  $P$ , discounted factor  $\gamma$ , random initialized
   policy  $\pi^0$ , vector  $v$ , convergence threshold  $\theta$ 
2: output: optimal policy  $\pi^*$ 
3: initialize  $\pi(s)$  randomly
4: initialize  $v$  with zeros
5:  $change \leftarrow \text{True}$ 
6: while  $\text{True}$  do
7:    $v \leftarrow \text{policy evaluation}(r, P, \gamma, \theta, v)$ 
8:    $\pi^k \leftarrow \text{policy improvement}(r, P, v)$ 
9:   if  $\pi^k = \pi^{k+1}$  ( $change \leftarrow \text{False}$ ) then
10:    break
11:   end if
12: end while
13:  $\pi^* \leftarrow \pi$ 
14: return  $\pi^*$ 

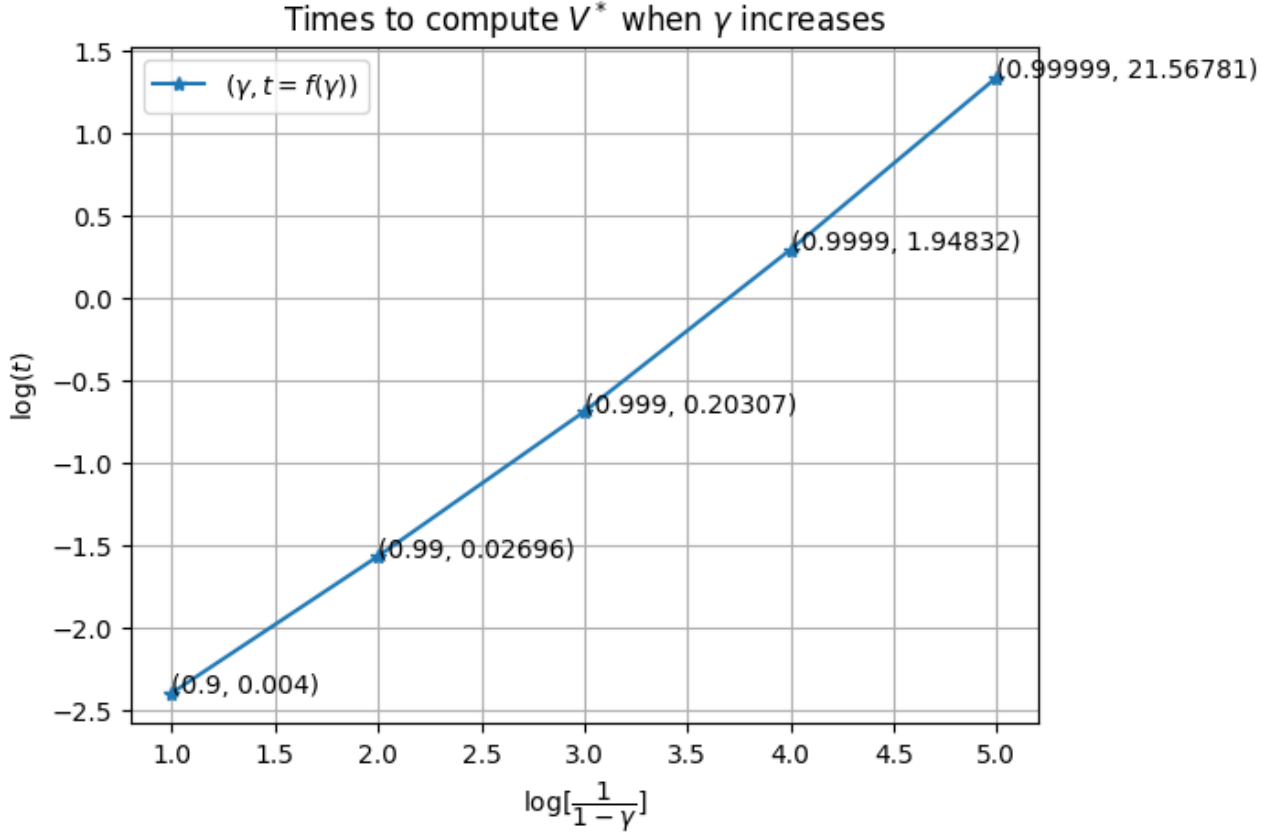
```

2.2.1 Numerical results and interpretation

The effects of γ on convergence.

Policy iteration algorithm converges in few iterations generally. For this reason, the effect of varying γ values is not noticeable if we decide to reason in terms of the number of iterations. In fact, the first step of *Policy Iteration* (**Policy Evaluation**) is actually an execution of the *Value Iteration's algorithm*, and this takes place at each iteration of the outer loop. These operations bring us very close to the optimal solution very quickly, which reduces the number of iterations of the *Policy Iteration's algorithm* to converge. It's therefore more convenient to look at the

execution time of the algorithm to see more clearly the impact of variations in γ values. The following figure give us a description about it.

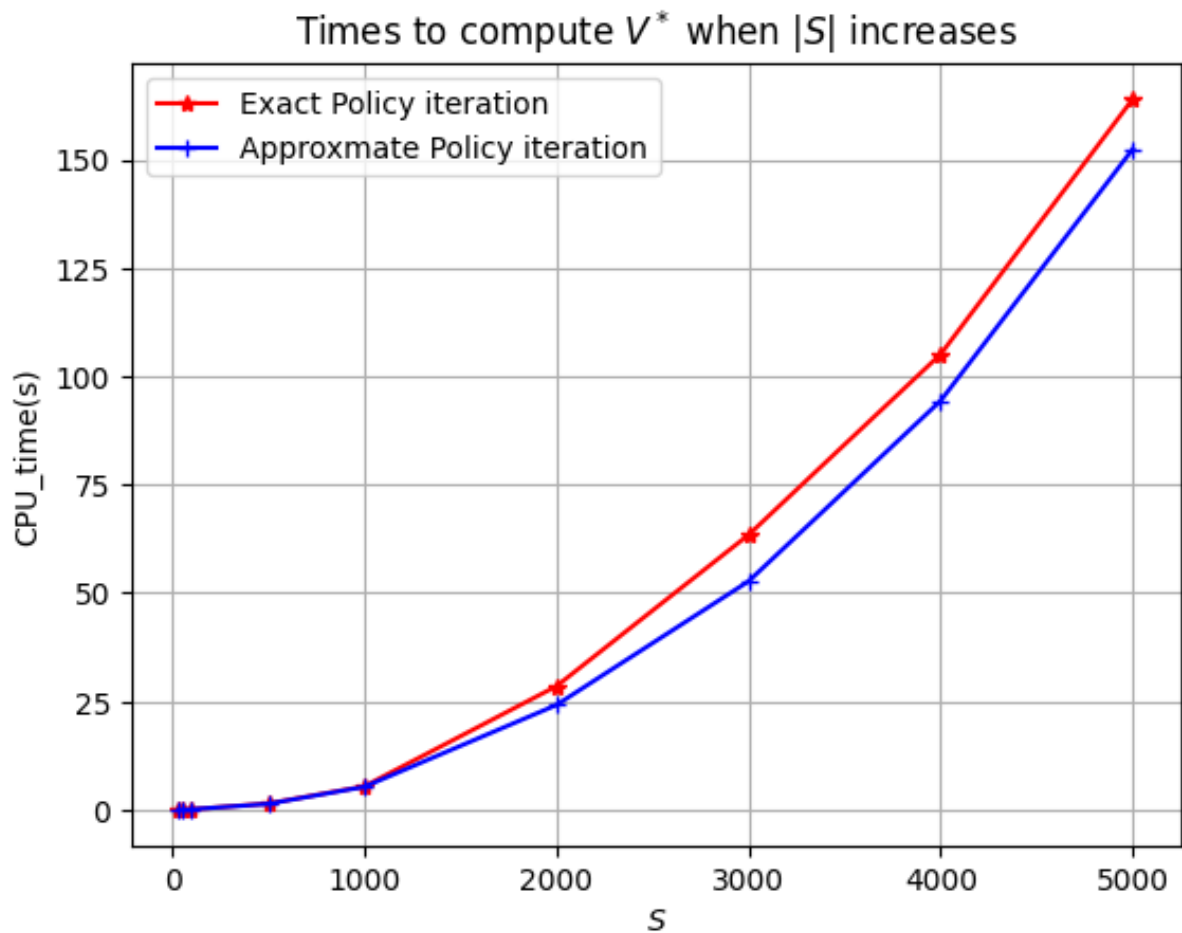


As you can see from the figure, as gamma increases, so does calculation time. The smaller the gamma, the faster the algorithm converges. But as gamma gets closer and closer to 1, the algorithm takes longer to converge; and in particular, convergence time increases almost tenfold as it gets closer and closer to 1.

Impact of the Number of States on Convergence Speed

Policy Iteration's complexity being $\mathcal{O}(|S|^3 + |S|^2|A|)$, then increasing the number of states could have several consequences on the algorithm's performance. As mentioned in case of *Value Iteration*, these problems concern, among other things, computational complexity, convergence time and memory storage.

1. **Computational Complexity:** When the step of *Policy Evaluation* is executed, if the number of states increases, the number of computations required to update state values or to invert the matrix $(I - \gamma P^{\pi_k})$ for each policy π_k increases.
2. **Convergence Duration:**
 - **Value Propagation:** In Policy Iteration algorithms, state values $V(s)$ need to propagate through the entire state space at each step of *Policy Evaluation*. More states mean it takes longer for values to propagate.
 - **Bellman Updates during *Policy Evaluation*:** Each update according to the Bellman equation requires calculations that include all states and their corresponding actions $\pi_k(s)$, meaning the computation time for each iteration increases with the number of states.

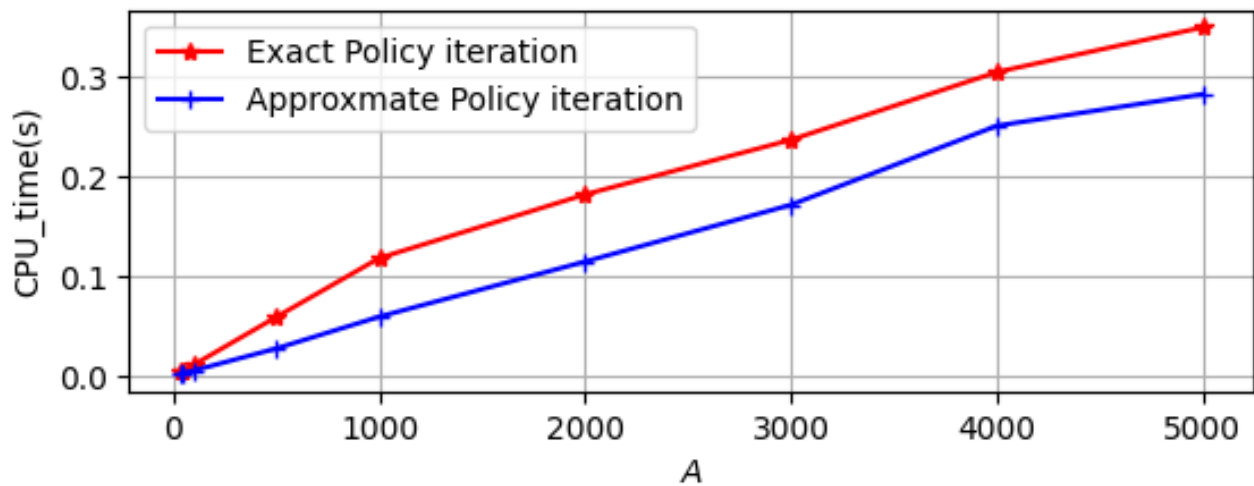


3. Memory and Storage:

- **Value and Policy Tables:** Storing value ($V(s)$) and policy ($\pi(s)$) tables increases linearly with the number of states. More states require more memory, which can also affect computation efficiency if memory resources are limited.

Impact of the Number of Actions on Convergence Speed

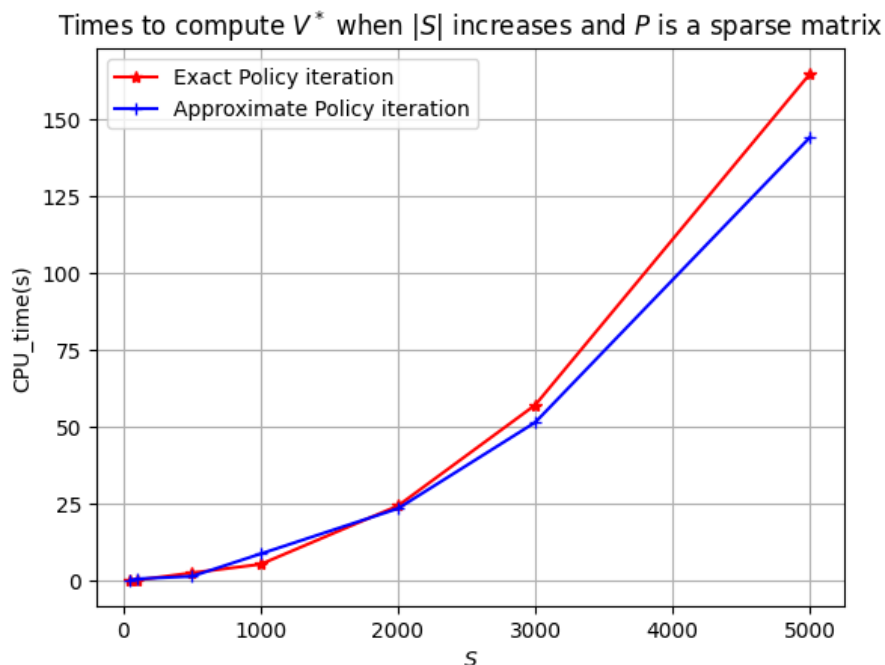
Increasing the number of available actions at each state also complicates calculations. In effect, to update the policy at the Policy Improvement stage, we need to propagate through all the action space for each state. The following figure shows the evolution of calculation time as the number of actions increases.

Times to compute V^* when $|A|$ increases

Remark 2.1. As can be seen from the two last previous figure above, the curve of the exact method is always above that of the approximate method. This implies that the exact method takes longer to converge than the approximate method as the number of actions or states increases. This is the cost of the matrix inversion that takes place at each iteration in the exact method during the PE step.

2.2.2 Test of Policy iteration methods on sparse matrix

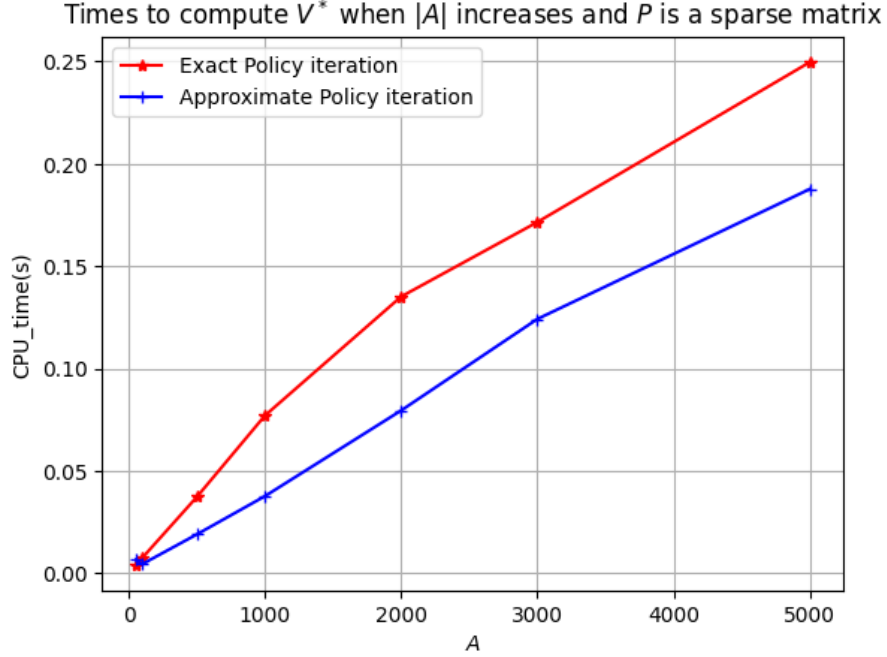
Impact of the Number of States on Convergence Speed when P is a sparse matrix



As you can see, the Approximate Policy Iteration method performs better than the Exact Policy Iteration method. Approximate Policy Iteration converges in less time on sparse matrices than on dense matrices (almost 10 seconds less than on dense matrices). On the other hand, the exact method shows no improvement despite the transition matrix being sparse. This is because the inversion of a matrix does not depend on its density. A sparse matrix generally has a dense

inverse. However, with sparse matrices, the gap between the exact and approximate methods becomes much wider as the number of states increases. We can easily see that from 5000 states upwards.

Impact of the Number of Action on Convergence Speed when P is a sparse matrix



2.3 Linear Programming Approach

For more details on Linear Programming approach to solve MDPs, please consult the following references(see [1, 7, 8, 6, 10]).

Complexity: For a fixed value of γ , Linear programming algorithm is both polynomial and strongly polynomial time algorithm. Here $|S^3||A|$ (see Tab 2) is the assumed runtime per iteration of Linear programming

As we've broached, the algorithms of PI and VI can be used to solve MDPs. But instead of using those approach, we can remark that MDPs can be reformulated as a linear programming problem. Let's describe this approach very quickly.

From Theorem 1.2 and Definition 1.2.2, we search for Q^* such that $Q^* = TQ^*$ and

$$V^*(s) = \max_{a \in A} Q^*(s, a) \textcircled{1}.$$

$\forall (s, a) \in S \times A$, one has:

$$\begin{aligned} Q^*(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\max_{a' \in A} Q(s', a') \right] \\ &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V^*(s')] \textcircled{2} \end{aligned}$$

$\textcircled{2}$ in $\textcircled{1}$ gives:

$$\begin{aligned} V^*(s) &= \max_{a \in A} \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V^*(s')] \right] \\ &= \max_{a \in A} \left[r(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V^*(s') \right] \end{aligned}$$

Then we can define a new operator (also named *Bellman Operator* τ defined by:

$$\tau : \mathbb{R}^{|S|} \rightarrow \mathbb{R}^{|S|}$$

$$V \mapsto \tau(V)$$

with $\tau(V)(s) = \max_{a \in A} [r(s, a) + \gamma \sum_{s' \in S} P(s'|s, a)V(s')]$

Now, the problem to solve is to find V such that $V = \tau(V)$. This problem can be reformulated as the following optimization problem with variables $V \in \mathbb{R}^{|S|}$:

$$\begin{aligned} \min \quad & \sum_s \mu(s)V(s) \\ \text{subject to} \quad & V(s) \geq r(s, a) + \gamma \sum_{s'} P(s' | s, a)V(s') \\ & \forall a \in A, \forall s \in S \end{aligned}$$

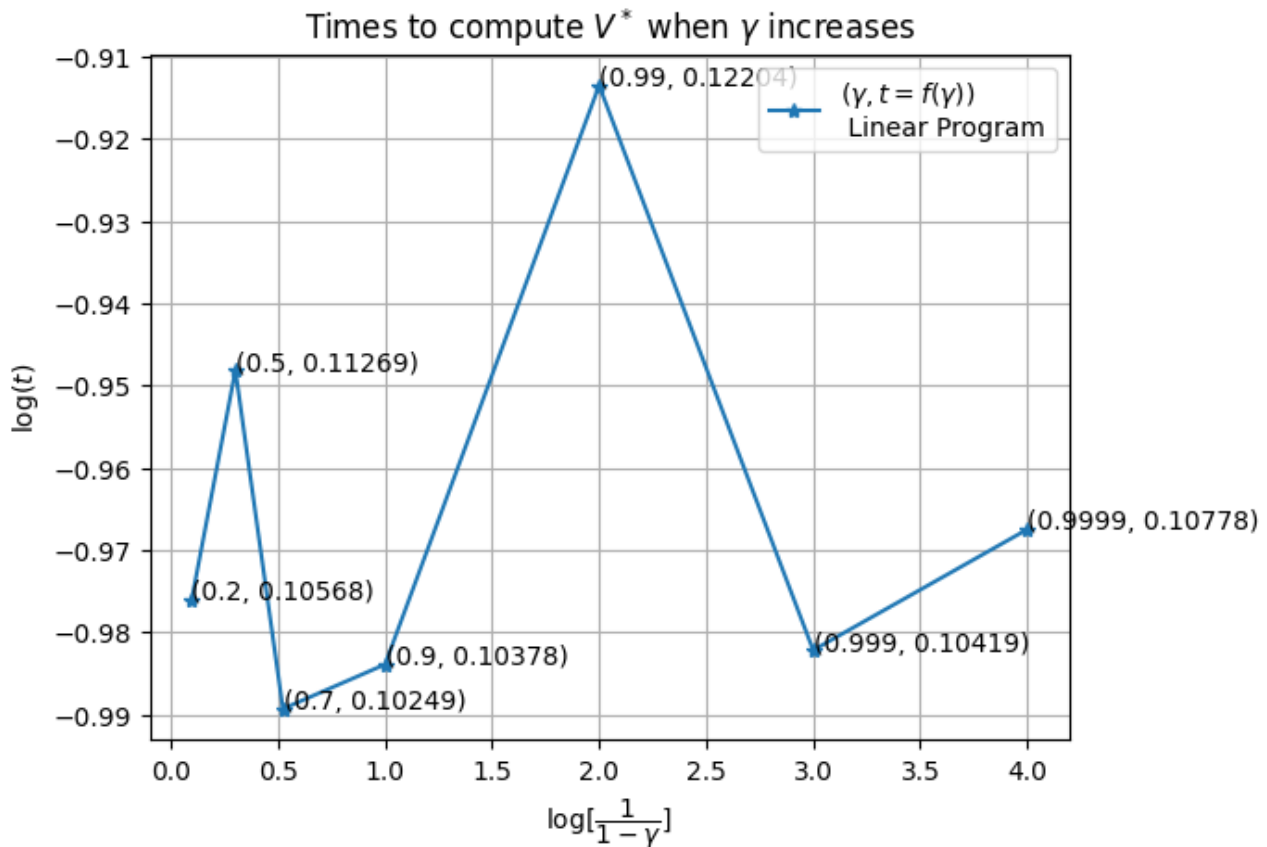
This is a linear programming problem.

To solve this linear programming (LP) problem, we are going to use **PuLP** library in python as follows:

- **Import PuLP:** We import the PuLP library.
- **Define the Problem:** We create an LP problem instance with `pulp.LpProblem`.
- **Variables:** Define the decision variables using `pulp.LpVariable`.
- **Objective Function:** Add the objective function to the LP problem using `+=`.
- **Constraints:** Add constraints similarly using `+=`.
- **Solve the Problem:** Call the `solve()` method to solve the LP.
- **Results:** Print the status, values of the decision variables, and the objective function.

2.3.1 Numerical results and interpretation

Impact of γ on Convergence Speed



The impact of the discount factor γ on the computation time of the optimal solution V^* is not linear. Variations in γ influence the complexity of the problem in a non-trivial manner. Indeed ;

- Between 0.9 and 0.99, We observe a significant increase in calculation time up to 0.12204 indicating that very high values of γ significantly increase computation time.
- But, between 0.99 to 0.999, there here is a significant reduction in calculation time to 0.10419, which may seem counter-intuitive, but perhaps indicates reduced complexity in some specific algorithm and data configurations.

The behavior observed in the curve could indeed be influenced by some of the open problems and well-known challenges in the field of linear programming, particularly those related to the simplex algorithms and other solution methods. Here are some specific points to consider:

- **Complexity of the Simplex**[11, 12, 16, 15, 17]

The simplex algorithm, while efficient in practice for many problems, has exponential complexity in the worst case. This means that, for certain specific configurations of γ , the computation time can increase unpredictably.

- **Conditioning of the Matrix**[11, 12]

High values of γ can affect the conditioning of the constraint matrix of the linear program. Poor conditioning can lead to convergence difficulties and increases in computation time. This is a common problem in the solution of linear systems and directly affects the performance of linear programming algorithms.

- **Cycling and Degenerate Pivoting**[11, 12]

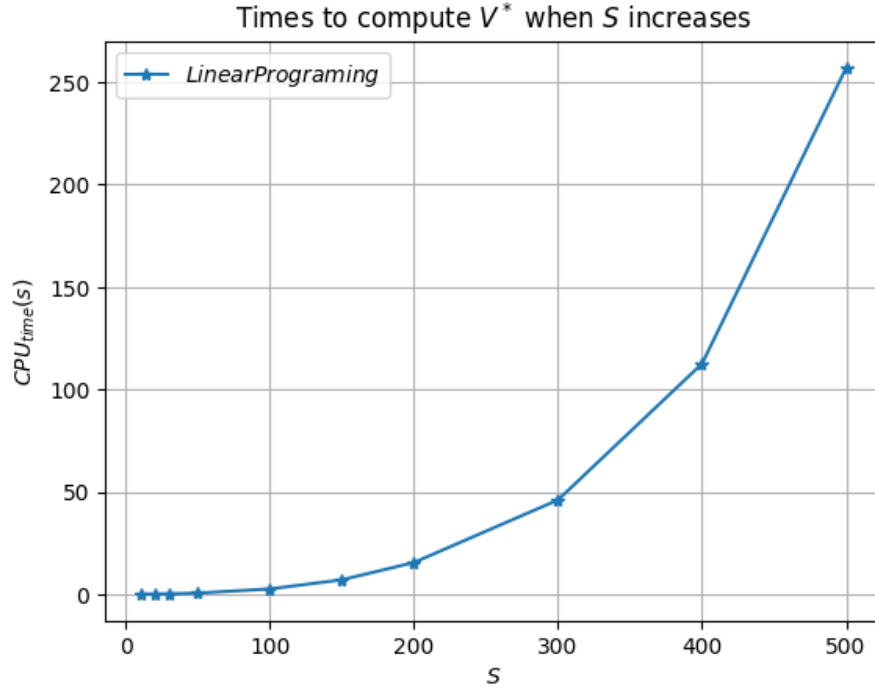
The phenomenon of cycling, where the simplex algorithm can revisit the same basic solutions without making progress, is a known problem that can significantly lengthen computation time. Specific values of γ can exacerbate these conditions, especially if they lead to degenerate solutions (where multiple optimal solutions exist).

- **Oscillations in Computation Time**

Non-linear variations and oscillations in computation time with changes in γ can also be due to complex interactions between the specific characteristics of the solution algorithm and the particular structure of the MDP being considered. For example, certain states and transitions may become more or less important at different levels of γ , affecting the difficulty of the problem in a non-linear way.

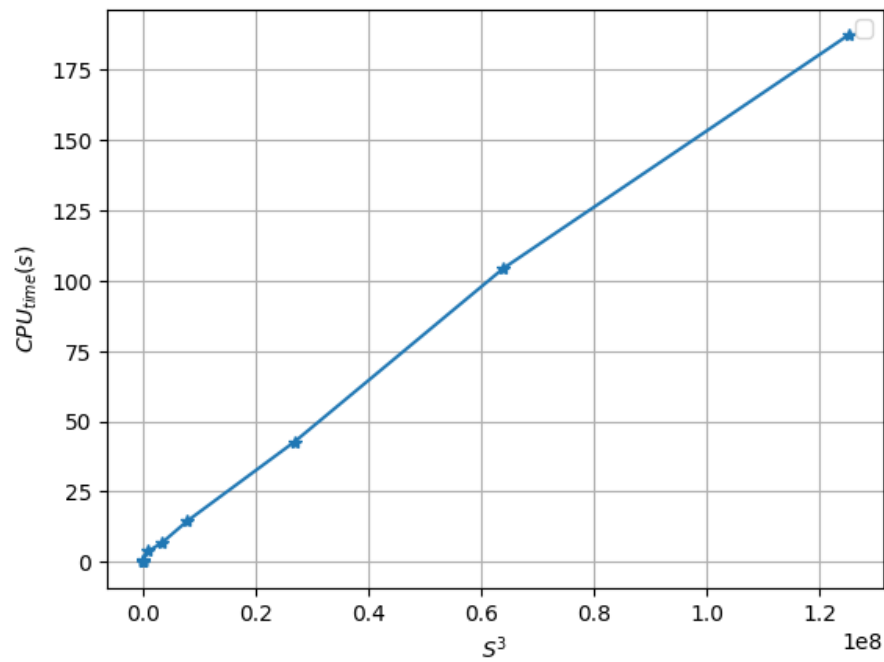
This suggests that, in practice, it may be useful to test several values of γ to strike a balance between future reward accuracy and computational efficiency. Algorithms for solving MDPs by linear programming may react differently depending on the specific parameters of the problem

Impact of the Number of States on Convergence Speed

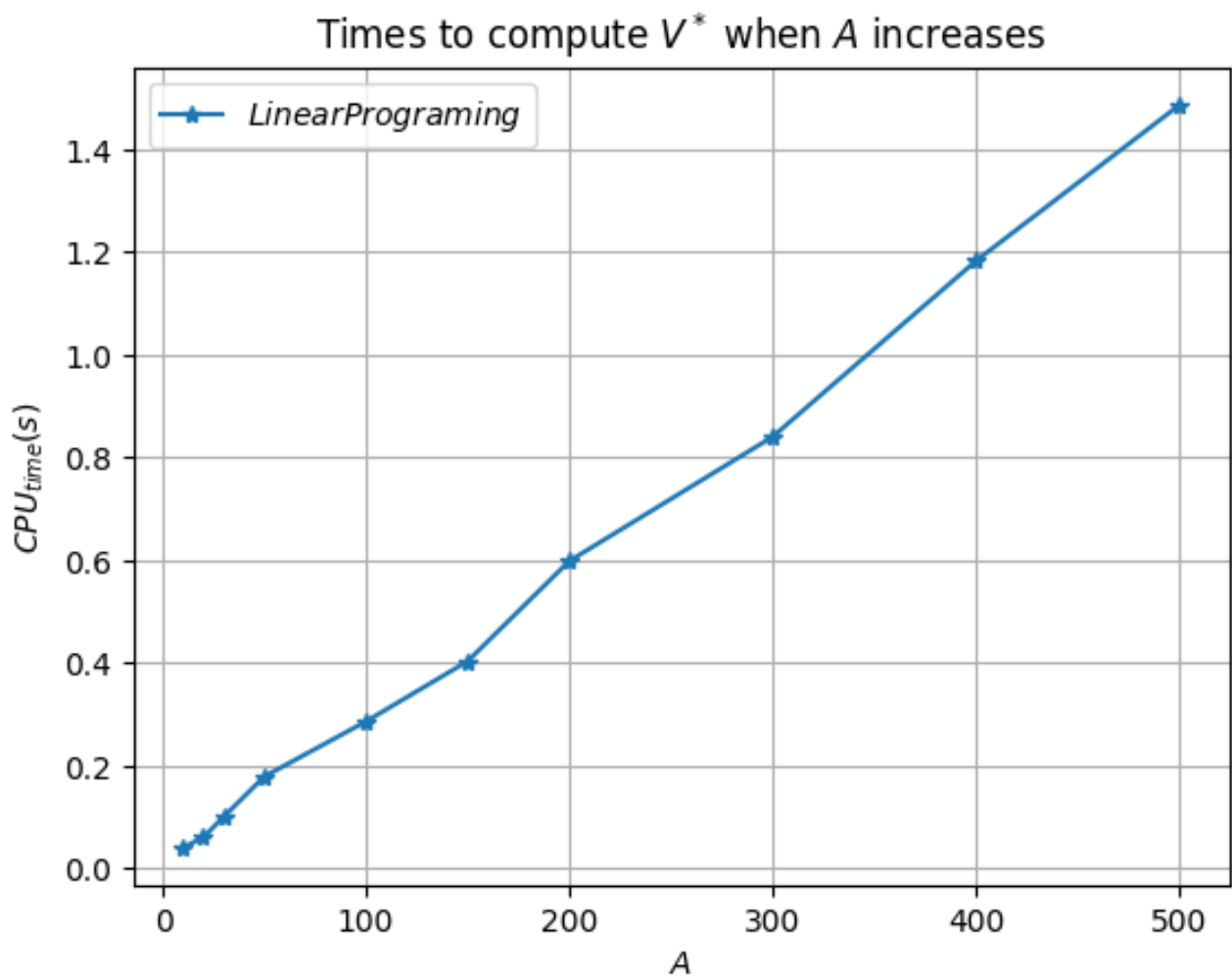


LP -algorithm becomes much more complex to solve as the number of states increases, so it takes longer to converge. Indeed, as the state space increases, so does the number of degrees of freedom (unknowns: $|S|$) in the cost function and the number of constraints ($|S| \cdot |A|$) in the problem, making the solution much more complex and more complicated to find.

The calculation time behaves like a cubic function ($x \mapsto x^3$: see figure below). This corresponds exactly to the theoretical results (refer to Table 2)



Impact of the Number of Actions on Convergence Speed



Value Iteration vs Linear Programming vs Policy Iteration

Comparison table						
States	Value Iteration		Approximate Policy Iteration		Linear Programming	
S_i	Value $V(s_i)$	Policy π^*	Value $V(s_i)$	Policy π^*	Value $V(s_i)$	Policy π^*
s_0	9.1055857	a_3	9.100605	a_3	9.1064532	a_3
s_1	9.2387334	a_3	9.229618	a_3	9.2396009	a_3
s_2	9.2370624	a_2	9.202891	a_2	9.2379299	a_2
s_3	9.2802963	a_3	9.290088	a_3	9.2811638	a_3
s_4	9.2639249	a_2	9.256704	a_2	9.2647924	a_2
s_5	9.1822204	a_4	9.169880	a_4	9.1830879	a_4
s_6	9.3022003	a_1	9.290648	a_2	9.3030678	a_1
s_7	9.2214848	a_0	9.211004	a_0	9.2223523	a_0
s_8	9.2665669	a_0	9.258590	a_0	9.2674345	a_0
s_9	9.0516525	a_4	9.043167	a_4	9.05252	a_4

Exact Policy Iteration vs Approximate Policy Iteration

Comparison table				
States	Exact Policy Iteration		Approximate Policy Iteration	
S_i	Value V	Policy π^*	Value V	Policy π^*
s_0	9.101454	a_3	9.100605	a_3
s_1	9.230468	a_3	9.229618	a_3
s_2	9.203740	a_2	9.202891	a_2
s_3	9.290938	a_3	9.290088	a_3
s_4	9.257554	a_2	9.256704	a_2
s_5	9.170729	a_4	9.169880	a_4
s_6	9.291498	a_2	9.290648	a_2
s_7	9.211854	a_0	9.211004	a_0
s_8	9.259439	a_0	9.258590	a_0
s_9	9.044017	a_4	9.043167	a_4

The values $V(s_i)$ from these methods are very close to each other, indicating that these methods are effectively converging to similar optimal values. Slight differences in values are due to differences in numerical precision and the algorithms' convergence criteria.

The optimal policies π^* (the actions associated with each state that maximize the expected reward) are also provided for each state. For most states, the policies derived from both methods are identical. This indicates a strong agreement between the two methods on the best actions to take from each state. There are a few differences in policies for some states (e.g., s_6) where **Policy Iteration** takes action a_2 and **Linear programming & Value Iteration** take action a_1 .

One might think that the optimal policy provided by these methods are not equivalent, but this is not true. In fact, these are two optimum policies

2.4 Python code of the methods

You can access the Google Colab notebook of all methods by clicking on the following links:

- [Value Iteration's algorithm](#)
- [Policy Iteration's algorithm](#)
- [Linear Programming' algorithm](#)

For those who will print the document

Value Iteration:

https://colab.research.google.com/drive/1M-Bxsr-Xwh0_Q1M-0Ps5ub7y5GjUXJLX?usp=sharing

Policy Iteration:

https://colab.research.google.com/drive/1r8W5EpijgxuojrJ9V6cN7pIsAC_zpB9X?usp=sharing

Linear Programming' algorithm:

<https://colab.research.google.com/drive/1dBiOB2JzkZ6uffLmBeXbU8nZ6UXaUtO7?usp=sharing>

Conclusion

Markov Decision Processes (MDPs) offer a powerful framework for modeling decision-making in stochastic environments. Among the various methods for solving MDPs, value iteration, policy iteration and linear programming stand out as the most widely used. In our work, we have studied and verified the theoretical complexity of these different methods through numerical simulations by varying the parameters of the algorithms. This practical approach allowed us to observe how changes in parameters impact the performance and convergence of each method. Summarizing each of the three MDP resolution methods studied, it can be seen that:

Value Iteration is an iterative algorithm that focuses on updating the value function until convergence. It is relatively simple to implement and guarantees convergence to the optimal value function. However, it can require a large number of iterations to achieve an accurate solution, especially in complex MDPs.

Policy Iteration combines policy evaluation and policy improvement steps to iteratively refine the policy. Starting with an arbitrary policy, it alternates between evaluating the current policy and improving it. This method typically converges faster than Value Iteration in terms of the number of iterations, but each iteration is computationally more expensive due to the need for solving systems of linear equations during policy evaluation.

Linear Programming offers a different approach by formulating the MDP as a linear optimization problem. This method directly solves for the optimal value function or policy by leveraging linear programming solvers.

Each method has its own strengths and weaknesses, making them suitable for different types of MDPs and computational environments. Value Iteration is advantageous for its simplicity and ease of implementation, Policy Iteration for its efficiency in practice, and Linear Programming for its precision.

In conclusion, the choice of method depends on the specific requirements of the problem at hand, including the size of the state and action spaces, the desired accuracy, and the available computational resources. Understanding the trade-offs and applications of these methods allows practitioners to effectively leverage MDPs for decision-making in uncertain and dynamic environments.

BIBLIOGRAPHY

- [1] Alekh Agarwal, Nan Jiang, Sham M. Kakade, Wen Sun, *Reinforcement Learning: Theory and Algorithms*, January 31, 2022
- [2] Amir massoud Farahmand, Remi Munos, Csaba Szepesvari, *Error Propagation for Approximate Policy and Value Iteration*
- [3] Alekh Agarwal, Mikael Henaff, Sham Kakade, and Wen Sun. *Pc-pg: Policy cover directed exploration for provable policy gradient learning*. *NeurIPS*, 2020a.
- [4] Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. *Flambe: Structural complexity and representation learning of low rank mdps*. *NeurIPS*, 2020b.
- [5] Sutton and Barto, *Reinforcement Learning: An Introduction*
- [6] Michael L. Littman, Thomas L. Dean, and Leslie Pack Kaelbling. *Learning and Planning in Markov Decision Processes*
- [7] Michael L. Littman, Thomas L. Dean, and Leslie Pack Kaelbling. *On the Complexity of Solving Markov Decision Problems*
- [8] Christos H. Papadimitriou and John N. Tsitsiklis. *The complexity of Markov decision processes*. *Math. Oper. Res.*, 12(3):441–450, 1987. URL: <https://doi.org/10.1287/moor.12.3.441>, doi:10.1287/moor.12.3.441.
- [9] Brafman and Tennenholtz, *On the Complexity of Value Iteration*
- [10] Dimitri P. Bertsekas and John Tsitsiklis, *Neuro-Dynamic Programming*
- [11] Donald Goldfarb, *On the Complexity of the Simplex Method*
- [12] Y. Disser and M. Skutella. *In defense of the simplex algorithm’s worst-case behavior*. *CoRR*, abs/1311.5935, 2013.
- [13] Martin L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*
- [14] O. Friedmann, *An exponential lower bound for the parity game strategy improvement algorithm as we know it*. In *Proc. of LICS*, pages 145–156, 2009.
- [15] O. Friedmann, *An exponential lower bound for the latest deterministic strategy iteration algorithms*. *Logical Methods in Computer Science*, 7(3), 2011.

- [16] O. Friedmann. *A subexponential lower bound for Zadeh's pivoting rule for solving linear programs and games.*In *Proc. of IPCO*, pages 192–206, 2011.
- [17] O. Friedmann, T. D. Hansen, and U. Zwick. *A subexponential lower bound for the random facet algorithm for parity games.*In *Proc. of SODA*, pages 202–216, 2011.
- [18] O. Friedmann, T. D. Hansen, and U. Zwick. *Subexponential lower bounds for randomized pivoting rules for the simplex algorithm.*In *Proc. of STOC*, pages 283–292, 2011.

Codes of different algorithm

- Value Iteration

```

1  # -*- coding: utf-8 -*- #
2
3
4  import numpy as np
5  import matplotlib.pyplot as plt
6
7  ##-----Define the MDP
   ↪ components-----##
8
9  num_states = 20
10 num_actions = 10
11 states = [f's{i}' for i in range(num_states)]
12 actions = [f'a{i}' for i in range(num_actions)]
13
14 ###----- Transition probability matrix P[s, a,
   ↪ s']-----###
15
16 # Random generation of matrix
17 def Markov_decision_matrix(num_states, num_actions):
18     Markov = []
19     for s in range(num_states):
20         matrix = np.random.rand(num_actions, num_states)
21         # Normalization of each row
22         matrix = matrix/matrix.sum(axis=1, keepdims=True)
23         Markov.append(matrix)
24     return np.array(Markov)
25 np.random.seed(42)
26 P = Markov_decision_matrix(num_states, num_actions)
27
28
29
30
31 ##-----Reward function R[s,
   ↪ a]-----##
32 R = np.round(np.random.rand(num_states, num_actions), 2) # random
   ↪ generation of R
33
34
35 # Action_value Function Q(initial value)
36 Q = np.zeros((num_states, num_actions))
37
38 # Discount factor

```

```

39 gamma = 0.7
40
41 Q.shape, R.shape, P.shape
42
43 V = np.zeros(num_states)
44 #np.dot(P[a,:,:],V)
45 np.dot(P,V).shape # The product make really sense
46
47 ##----- Definition of Bellman operator-----##
48
49 def Bellman_Operator(Q,R,P,gamma):
50     V_Q = np.max(Q, axis=1)
51     return R + gamma*np.dot(P,V_Q)
52
53 #Bellman_Operator(Q,R,P,gamma)# Test pour Q = Zeros_23
54
55 ##-----Value iteration
56 ↪ algorithm-----##
57
58 def Matrix_value_iteration(Q, P, R, gamma, theta=0.0001):
59     iter = 0
60     err = []
61     while True:
62         TQ=Bellman_Operator(Q,R,P,gamma)
63         delta = np.max(np.abs(Q - TQ))
64         err.append(delta)
65         Q = TQ
66         iter+=1
67         if delta < theta:
68             break
69     V_star = np.max(Q, axis=1)
70     return V_star, iter, err, R.shape
71
72 # Commented out IPython magic to ensure Python compatibility.
73 # %timeit Matrix_value_iteration(Q,P, R, gamma)
74
75 ##-----Compute the optimal value
76 ↪ function-----##
77
78 optimal_values, k, err, _ = Matrix_value_iteration(Q,P, R, gamma)
79 print("Optimal Value Function:")
80 for i, value in enumerate(optimal_values):
81     print(f"V({states[i]}) = {value:.2f}")
82
83 print("Nombre d'itérations:",k)
84 print("Erreur",err)
85
86 """> *After convergence, the optimal policy is derived by selecting the
87 ↪ action that maximizes the expected value for each state*.
88
89 $$$\forall s \in \mathbb{S}, \pi^*(s) =
90 ↪ \underset{a \in \mathbb{A}}{\operatorname{argmax}} Q^*(s,a) =
91 ↪ \underset{a \in \mathbb{A}}{\operatorname{argmax}} \left[ r(s,a) + \gamma \sum_{s'} P(s'/s,a) V^*(s') \right]
92 """
93
94

```

```

88  ##-----Optimal policy-----##
89
90  def optimal_policy(P, R, V, gamma):
91      num_states, num_actions, _ = P.shape
92      policy = np.zeros(num_states, dtype=int)
93      for s in range(num_states):
94          policy[s] = np.argmax([
95              R[s,a] + sum(gamma * P[s, a, s_next] * (
96                  ↪ V[s_next])#np.argmax(Q[:,s])
97                  for s_next in range(num_states))
98                  for a in range(num_actions)])
99      return policy
100
101  policy_star = optimal_policy(P, R, optimal_values, gamma)
102  print("\nOptimal Policy:")
103  for i, action in enumerate(policy_star):
104      print(f"({states[i]}) = {actions[action]}")
105
106  ##-----Plotting function-----##
107
108  def loss_Plot(k,err,gamma,sh,type_err):
109      K = np.arange(k)
110      if type_err == "log":
111          plt.plot(K, np.log(err), label = f"$\gamma$={gamma}:({k}
112              ↪ iterations)")
113          plt.xlabel("Iterations")
114          plt.ylabel("$\log_{10}(|Q^k - TQ^k|)$")
115          plt.title(f"S = {sh[0]}, A = {sh[1]}")
116          plt.grid(visible=True)
117          plt.legend()
118      else:
119          plt.plot(K, err, label = f"$\gamma$={gamma}:({k} iterations)")
120          plt.xlabel("Iterations")
121          plt.ylabel("Error")
122          plt.title(f"S = {sh[0]}, A = {sh[1]}")
123          plt.grid(visible=True)
124          plt.legend()
125      return 0
126
127  import time
128  ###-----Analyze of convergence based on gamma's
129  ↪ value-----###
130  Gamma = [1-1e-1, 1-1e-2, 1-1e-3, 1-1e-4]
131  Times = []
132  for gam in Gamma:
133      start = time.time()
134      _, k, err,sh = Matrix_value_iteration(Q,P, R, gam)
135      end = time.time()
136      Times.append((end-start))
137      loss_Plot(k,err,gam,sh,"log")
138
139  ###-----Plotting of Gamma/times curve-----###
140  Gamma_log = [np.log10(1/(1-g)) for g in Gamma] # Log(times)
141  Times_log = [np.log10(t) for t in Times] # Log(1/(1-gamma))

```



```

139 plt.figure()
140 for i in range(len(Times)):
141     plt.text(Gamma_log[i],Times_log[i],f'{Gamma[i],round(Times[i],5)}')
142
143 plt.plot(Gamma_log,Times_log, "-*", label="$\\gamma, t=f(\\gamma)$")
144 plt.xlabel("$\\log[\\frac{1}{1-\\gamma}]$")
145 plt.ylabel("$\\log(t)$")
146 plt.title("Times to compute $V^*$ when $\\gamma$ increase ")
147 plt.legend()
148 plt.grid(1)
149
150 ###-----Analyze of convergence by varying the number of
151 ↪ states-----###
152
153 gam = 0.7 # Here gamma is constant
154 actions = 10 # Number of action is constant
155 Times_state = []
156 States = [35, 50, 100, 500, 1000, 2000, 3000, 4000, 5000] # List of states
157
158 for state in States:
159     np.random.seed(42)
160     P_s = Markov_decision_matrix(state, actions)
161     R_s = np.round(np.random.rand(state, actions),2)
162     Q_s = np.zeros((state, num_actions))
163     start = time.time()
164     _, k, err, sh = Matrix_value_iteration(Q_s, P_s, R_s, gam)
165     end = time.time()
166     Times_state.append((end-start))
167     #loss_Plot(k, err, gam, sh, "")
168
169 plt.figure()
170 plt.plot(States, Times_state, "-*r", label="")
171 plt.xlabel("$S$")
172 plt.ylabel("CPU_time(s)")
173 plt.title("Times to compute $V^*$ when $|S|$ increase ")
174 plt.grid(1)
175 plt.legend()
176
177 ##-----Proof that CPU_time of VI when number of state behave like
178 ↪ quaratic function-----#
179
180 plt.figure()
181 ss = [s**2 for s in States]
182
183 plt.plot(ss, Times_state, "-*r", label="")
184 plt.xlabel("$S^2$")
185 plt.ylabel("CPU_time(s)")
186 plt.grid(1)
187 plt.legend()
188
189 ###-----Analyze of convergence by varying the number of
190 ↪ actions-----###

```

```

190 gam = 0.7 # Here gamma is constant
191 state = 10 # Number of state is constant
192 Times_action = []
193 Actions = [35, 50, 100, 500, 1000, 2000, 3000, 4000, 5000] # List of states
194
195 for i, action in enumerate(Actions):
196     np.random.seed(42)
197     P_a = Markov_decision_matrix(state, action)
198     R_a = np.round(np.random.rand(state, action), 2)
199     Q_a = np.zeros((state, action))
200     start = time.time()
201     _, k, err, sh = Matrix_value_iteration(Q_a, P_a, R_a, gam)
202     end = time.time()
203     Times_action.append(end-start)
204     #loss_Plot(k, err, gam, sh, "")
205
206 plt.figure()
207 #for i in range(len(Times_state)):
208     ↪ #plt.text(Times_action[i], States[i], f'{States[i], round(Times_action[i], 5)}')
209 plt.subplot(2, 1, 1)
210 plt.plot(Actions, Times_action, "-+g", label="")
211 plt.xlabel("A")
212 plt.ylabel("CPU_time(s)")
213 plt.title("Times to compute  $V^*$  when  $|A|$  increase")
214 plt.legend()
215 plt.grid(1)
216
217 #-----
218 plt.figure()
219 plt.subplot(2, 1, 1)
220 plt.plot(Actions, Times_action, "-*g", label="")
221 plt.ylabel("CPU_time(s)")
222 plt.legend()
223 plt.grid(1)
224 plt.subplot(2, 1, 2)
225 plt.plot(Actions, Actions, "-b", label="$y=x$")
226 plt.xlabel("A")
227 plt.ylabel("A")
228 plt.legend()
229 plt.grid(1)

```

- Policy Iteration

```

1      # -*- coding: utf-8 -*-
2
3
4  import numpy as np
5  import matplotlib.pyplot as plt
6  import time
7  from scipy import sparse
8
9

```

```

10  ##-----Define the MDP
    ↪ components-----##
11
12  num_states = 10
13  num_actions = 5
14  states = [f's{i}' for i in range(num_states)]
15  actions = [f'a{i}' for i in range(num_actions)]
16
17
18  ###----- Transition probability matrix P[s, a,
    ↪ s']-----###
19
20  # Random generation of matrix
21  def Markov_decision_matrix(num_states, num_actions):
22      Markov = []
23      for s in range(num_states):
24          matrix = np.random.rand(num_actions, num_states)
25          # Normalization of each row
26          matrix = matrix/matrix.sum(axis=1, keepdims=True)
27          Markov.append(matrix)
28      return np.array(Markov)
29  np.random.seed(42)
30  P = Markov_decision_matrix(num_states, num_actions)
31
32
33
34
35  ##-----Reward function R[s,
    ↪ a]-----##
36  R = np.round(np.random.rand(num_states, num_actions), 2) # random
    ↪ generation of R
37
38
39  # Action_value Function Q(initial value)
40  Q = np.zeros((num_states, num_actions))
41
42  # Discount factor
43  gamma = 0.9
44
45  Q.shape, R.shape, P.shape
46
47  """
48  # Policy Iteration
49  Repeat steps until convergence:
50
51  1-  $\text{Policy evaluation}$ : keep current policy  $\pi$  fixed, find
52   $V^{\pi_k}(\cdot)$ 
53  + Approximative evaluation of  $V^{\pi_k}$ :
54
55  Iterate simplified Bellman update until values converge:
56  
$$V_{i+1}^{\pi_k}(s) = r(s, \pi_k(s)) + \gamma \sum_{s'} P(s'/s, \pi_k(s)) V_i^{\pi_k}(s')$$

    ↪  $\text{Le faire juste apr\`es comme Approximate Policy Iteration, bien}$ 
    ↪  $\text{choisir le theta de l'ordre de } (1-\gamma)^2 \epsilon/4$ 
57  Here we assume that  $V_0^{\pi_0}(s) = 0$ 

```

```

58
59
60 + Exact evaluation of  $V^{\pi_k}$ :
61
62  $\boxed{V^{\pi_k} = (I - \gamma P^{\pi_k})^{-1} r^{\pi_k} \mid \text{with} \mid$ 
63  $\hookrightarrow P^{\pi_k}_{s,s'} = P(s'/s, \pi_k(s))$ 
64
65
66 2-  $\text{Policy improvement}$ : find the best action for  $V^{\pi_k}(\cdot)$ 
67  $\hookrightarrow$  via one-step lookahead
68
69  $\text{forall } s \in \mathbb{S}, \pi_{k+1}(s) =$ 
70  $\hookrightarrow \underset{a \in \mathbb{A}}{\operatorname{argmax}} \left[ R(s, a) + \underset{s'}{\sum} P(s'/s, a) V^{\pi_k}(s') \right]$ 
71
72 """
73
74 np.random.seed(42)
75
76 ##-----Random initialization of pi-----##
77 pi = np.random.randint(num_actions, size = num_states)
78
79 pi
80
81 I = np.eye(num_states)
82 #I
83
84 ##-----Function to evaluate the utility V on the the arbitrary
85  $\hookrightarrow$  policy-----##
86
87 def Exact_Policy_evaluation(pi, P, R, gamma):
88     num_states, num_actions, _ = P.shape
89     P_pi = np.array([P[s, pi[s], :] for s in range(num_states)])
90     R_pi = np.array([R[s, pi[s]] for s in range(num_states)])
91
92     I = np.eye(num_states)
93     V = np.dot(np.linalg.inv(I - gamma*P_pi), R_pi)
94
95     return V
96
97 def Approximate_Policy_evaluation(pi, V, P, R, gamma, theta):
98     count = 0
99     num_states, num_actions, _ = P.shape
100     P_pi = np.array([P[s, pi[s], :] for s in range(num_states)])
101     R_pi = np.array([R[s, pi[s]] for s in range(num_states)])
102
103     while True:
104         delta = 0
105         TV = R_pi + gamma*sparse.coo_array.dot(P_pi, V)
106         delta = np.max(np.abs(V - TV))
107         V = TV
108         count += 1
109         #print(dt)
110
111         if delta <= theta:
112             break
113     #print("Number of iteration loop intern approximate", count)

```

```

108     return V
109
110     ##-----Function to compute the greedy
111     ↪ polyci-----##
112
113 def Policy_improvement(V, P, R, gamma):
114     num_states, num_actions, _ = P.shape
115     policy = np.zeros(num_states, dtype=int)
116     for s in range(num_states):
117         policy[s] = np.argmax([R[s,a]+sum(gamma *P[s, a, s_next] * (
118             ↪ V[s_next])
119             for s_next in range(num_states))#R(s,a)+...
120             for a in range(num_actions)
121         ])
122     return policy
123
124     ##-----The main
125     ↪ algorithm-----##
126
127 def Policy_iteration(pi,V,P,R, gamma, method="Exact", theta = 0.0001):
128     iter = 0
129
130     if method=="Exact":
131         while True:
132             V_Q = Exact_Policy_evaluation(pi, P, R, gamma)
133             greedy_pi = Policy_improvement(V_Q, P, R, gamma)
134             iter+=1
135             if greedy_pi.all() == pi.all():
136                 break
137             pi = greedy_pi
138
139     elif method=="Approx":
140         V_Q = np.copy(V)
141
142         while True:
143             V_old = np.copy(V_Q)
144             V_Q = Approximate_Policy_evaluation(pi,V_old, P, R, gamma, theta)
145             greedy_pi = Policy_improvement(V_Q, P, R, gamma)
146             iter+=1
147
148             if greedy_pi.all() == pi.all():
149                 break
150
151             pi = greedy_pi
152     return pi, V_Q, iter
153
154 V = np.zeros(num_states)
155
156 optimal_policy, optimal_values , k= Policy_iteration(pi,V,P,
157     ↪ R,gamma,"Approx", theta=0.0001)
158 print("Optimal Value Function:")
159 for i, value in enumerate(optimal_values):
160     print(f"V({states[i]}) = {value:7f}")

```

```

158 print("Nombre d'itérations:",k)
159
160 print("\nOptimal Policy:")
161 for i, action in enumerate(optimal_policy):
162     print(f"({states[i]}) = {actions[action]}")
163
164 ###-----Analyze of convergence based on gamma's
165 ↪ value-----###
166
167 Gamma = [1-1e-1, 1-1e-2, 1-1e-3, 1-1e-4, 1-1e-5]
168 Policy_times = []
169 for gam in Gamma:
170     start = time.time()
171     _, _, k = Policy_iteration(pi,V,P, R,gam,"Approx", theta=0.0001)
172     end = time.time()
173     Policy_times.append(end-start)
174
175 Gamma_log = [np.log10(1/(1-g)) for g in Gamma] # Log(1/(1-gamma))
176 Times_log = [np.log10(t) for t in Policy_times] # Log(times)
177 plt.figure()
178 for i in range(len(Policy_times)):
179     ↪ plt.text(Gamma_log[i],Times_log[i],f'{Gamma[i],round(Policy_times[i],5)}')
180
181 plt.plot(Gamma_log,Times_log, "-*", label="$\\gamma, t=f(\\gamma)$")
182 plt.xlabel("$\\log[\\frac{1}{1-\\gamma}]$")
183 plt.ylabel("$\\log(t)$")
184 plt.title("Times to compute $V^*$ when $\\gamma$ increases")
185 plt.legend()
186 plt.grid(1)
187
188 print(Policy_times)
189
190 ###-----Analyze of convergence by varying the number of
191 ↪ states-----###
192
193 gam = 0.7 # Here gamma is constant
194 actions = 10 # Number of action is constant
195 Times_state = []
196 States = [35, 50, 100, 500, 1000, 2000, 3000, 4000, 5000] # List of states
197 Methods = ["Exact", "Approx"]
198 for m in Methods:
199     T = []
200     for state in States:
201         np.random.seed(42)
202         P_s = Markov_decision_matrix(state, actions)
203         R_s = np.round(np.random.rand(state, actions),2)
204         pi_s = np.random.randint(actions, size = state)
205         V_s = np.zeros(state)
206         start = time.time()
207         _, _, k = Policy_iteration(pi_s,V_s,P_s, R_s,gam,m, theta=0.0001)
208         end = time.time()
209         T.append((end-start))
210     Times_state.append(T)

```

```

209
210 plt.figure()
211 plt.plot(States, Times_state[0], "-*r", label=f"{Methods[0]} Policy
↪ iteration")
212 plt.plot(States, Times_state[1], "-+b", label=f"{Methods[1]}imate Policy
↪ iteration")
213 plt.xlabel("$S$")
214 plt.ylabel("CPU_time(s)")
215 plt.title("Times to compute  $V^*$  when  $|S|$  increases")
216 plt.grid(1)
217 plt.legend()
218
219
220
221 ###-----Analyze of convergence by varying the number of
↪ actions-----###
222
223 gam = 0.7 # Here gamma is constant
224 state = 10 # Number of state is constant
225 Times_action = []
226 Actions = [35, 50, 100, 500, 1000, 2000, 3000, 4000, 5000] # List of states
227
228 Methods = ["Exact", "Approx"]
229 for m in Methods:
230     T = []
231     for actions in Actions:
232         np.random.seed(42)
233         P_a = Markov_decision_matrix(state, actions)
234         R_a = np.round(np.random.rand(state, actions), 2)
235         pi_a = np.random.randint(actions, size = state)
236         V_a = np.zeros(state)
237         start = time.time()
238         _, _, k = Policy_iteration(pi_a, V_a, P_a, R_a, gam, m, theta=0.0001)
239         end = time.time()
240         T.append((end-start))
241     Times_action.append(T)
242
243 plt.figure()
244 plt.subplot(2,1,1)
245 plt.plot(States, Times_action[0], "-*r", label=f"{Methods[0]} Policy
↪ iteration")
246 plt.plot(States, Times_action[1], "-+b", label=f"{Methods[1]}mate Policy
↪ iteration")
247 plt.xlabel("$A$")
248 plt.ylabel("CPU_time(s)")
249 plt.suptitle("Times to compute  $V^*$  when  $|A|$  increases")
250 plt.grid(1)
251 plt.legend()
252
253 """> # Test of Policy iteration methods on Sparse matrix"""
254
255 #-----Function to create sparse transition
↪ matrix-----#
256

```

```

257 def create_random_sparse_3d_transition_matrix(nb_states, nb_actions,
↪ nb_next_states, density=0.3):
258     """
259     Create a random sparse 3D transition matrix for an MDP.
260
261     Parameters:
262     - nb_states: Number of states.
263     - nb_actions: Number of actions.
264     - nb_next_states: Number of next states.
265     - density: The proportion of non-zero elements (default is 0.1).
266
267     Returns:
268     - A 3D numpy array representing the transition matrix.
269     """
270     data = []
271     row_indices = []
272     col_indices = []
273
274     rng = np.random.default_rng()
275
276     for s in range(nb_states):
277         for a in range(nb_actions):
278             # Determine how many transitions to generate for this
↪ state-action pair
279             num_transitions = int(density * nb_next_states)
280             if num_transitions < 1:
281                 num_transitions = 1
282
283             # Randomly choose next states and probabilities
284             next_states = rng.choice(nb_next_states, size=num_transitions,
↪ replace=False)
285             probabilities = rng.random(size=num_transitions)
286             probabilities /= probabilities.sum() # Normalize to sum to 1
287
288             for s_prime, p in zip(next_states, probabilities):
289                 row_index = s * nb_actions + a
290                 row_indices.append(row_index)
291                 col_indices.append(s_prime)
292                 data.append(p)
293
294     # Convert to numpy arrays
295     data = np.array(data)
296     row_indices = np.array(row_indices)
297     col_indices = np.array(col_indices)
298
299     # Create sparse matrix in COO format
300     transition_matrix_2d = sparse.coo_matrix((data, (row_indices,
↪ col_indices)),
301                                             shape=(nb_states *
↪ nb_actions,
↪ nb_next_states))
302
303     # Convert to a dense 3D numpy array
304     transition_matrix_3d = np.zeros((nb_states, nb_actions,
↪ nb_next_states))

```



```

304
305     for row, col, val in zip(row_indices, col_indices, data):
306         s = row // nb_actions
307         a = row % nb_actions
308         transition_matrix_3d[s, a, col] = val
309
310     return transition_matrix_3d
311
312 # Example usage
313 nb_states = 5
314 nb_actions = 2
315 nb_next_states = 5
316 density = 0.3 # 30% of transitions will be non-zero
317
318 transition_matrix_3d =
319     ↪ create_random_sparse_3d_transition_matrix(nb_states, nb_actions,
320     ↪ nb_next_states, density)
321
322 # Display the 3D transition matrix
323 print("3D Transition Matrix:\n", transition_matrix_3d[:, :, 0].shape)
324
325 """## When $|S|$ increases"""
326
327 gam = 0.7 # Here gamma is constant
328 actions = 10 # Number of action is constant
329 Density = 0.3
330 Sp_Times_state = []
331 States = [50,100,500,1000,2000, 3000,4000,5000] # List of states
332 Methods = ["Exact", "Approx"]
333 for m in Methods:
334     T = []
335     for state in States:
336         np.random.seed(42)
337         P_s = create_random_sparse_3d_transition_matrix(state, actions, state,
338         ↪ Density) # Sparse transition matrix
339         R_s = np.round(np.random.rand(state, actions),2)
340         pi_s = np.random.randint(actions, size = state)
341         V_s = np.zeros(state)
342         start = time.time()
343         _ , _ , k = Policy_iteration(pi_s,V_s,P_s, R_s,gam,m, theta=0.0001)
344         end = time.time()
345         T.append((end-start))
346     Sp_Times_state.append(T)
347
348 plt.figure()
349 plt.plot(States, Sp_Times_state[0], "-*r", label=f"{Methods[0]} Policy
350     ↪ iteration")
351 plt.plot(States, Sp_Times_state[1], "-+b", label=f"{Methods[1]}imate
352     ↪ Policy iteration")
353 plt.xlabel("$S$")
354 plt.ylabel("CPU_time(s)")
355 plt.title("Times to compute $V^*$ when $|S|$ increases and $P$ is a sparse
356     ↪ matrix")
357 plt.grid(1)

```

```

352 plt.legend()
353
354 """## When  $|A|$  increases"""
355
356 gam = 0.7 # Here gamma is constant
357 state = 10 # Number of action is constant
358 Density = 0.6
359 ASp_Times_state = []
360 Action = [50,100,500,1000,2000,3000,4000,5000] # List of states
361 Methods = ["Exact", "Approx"]
362 for m in Methods:
363     T = []
364     for actions in Action:
365         np.random.seed(42)
366         P_a = create_random_sparse_3d_transition_matrix(state, actions, state,
367             ↪ Density) # Sparse transition matrix
368         R_a = np.round(np.random.rand(state, actions),2)
369         pi_a = np.random.randint(actions, size = state)
370         V_a = np.zeros(state)
371         start = time.time()
372         _, _, k = Policy_iteration(pi_a,V_a,P_a, R_a,gam,m, theta=0.0001)
373         end = time.time()
374         T.append((end-start))
375     ASp_Times_state.append(T)
376
377 plt.figure()
378 plt.plot(Action, ASp_Times_state[0], "-*r", label=f"{Methods[0]} Policy
379     ↪ iteration")
380 plt.plot(Action, ASp_Times_state[1], "--+b", label=f"{Methods[1]}imate
381     ↪ Policy iteration")
382 plt.xlabel("$A$")
383 plt.ylabel("CPU_time(s)")
384 plt.title("Times to compute  $V^*$  when  $|A|$  increases and  $P$  is a sparse
385     ↪ matrix")
386 plt.grid(1)
387 plt.legend()

```

- Linear Program

```

1         # -*- coding: utf-8 -*-
2
3     # Please install by tapping: !pip install pulp
4
5     import numpy as np
6     import matplotlib.pyplot as plt
7     import pulp
8     import time
9     # Test
10    P1 = np.array([[0.07200801, 0.18278161, 0.14073106, 0.11509637, 0.0299957
11        ↪ ,
12        0.02999106, 0.01116699, 0.16652855, 0.11556865, 0.13613201],
13        [0.00520773, 0.24538041, 0.21060217, 0.05372031, 0.04600045,
14        0.04640006, 0.07697116, 0.13275971, 0.10927907, 0.07367894],
15        [0.15281528, 0.03483974, 0.07296552, 0.09150188, 0.11390722,

```

```

15      0.19610414, 0.04987017, 0.12843427, 0.1479604 , 0.01160137],
16      [0.11929704, 0.03348399, 0.01277348, 0.18632243, 0.18961076,
17      0.15873628, 0.05981372, 0.01917882, 0.13435547, 0.086428 ],
18      [0.03016589, 0.12239978, 0.00850029, 0.22476941, 0.06396626,
19      0.16376487, 0.07704997, 0.12855247, 0.135138 , 0.04569305]],
20
21      [[0.16851056, 0.13471549, 0.16328177, 0.15551799, 0.10391301,
22      0.16021865, 0.0153797 , 0.03406116, 0.00786035, 0.05654132],
23      [0.08316254, 0.05805864, 0.17731912, 0.07633199, 0.06010957,
24      0.11611685, 0.03015256, 0.17164043, 0.01595108, 0.21115723],
25      [0.16981899, 0.04369819, 0.00121433, 0.17932246, 0.15544009,
26      0.16031091, 0.16960471, 0.01628265, 0.07882771, 0.02547996],
27      [0.16460084, 0.11886802, 0.06310494, 0.01212109, 0.05930686,
28      0.0620151 , 0.13914183, 0.12158739, 0.16919868, 0.09005523],
29      [0.02655733, 0.15838451, 0.16894139, 0.12463829, 0.17120245,
30      0.1096532 , 0.11607906, 0.09494058, 0.00564462, 0.02395858]],
31      [[0.16790231, 0.04428548, 0.02678531, 0.09048039, 0.18220764,
32      0.04474641, 0.12425118, 0.14079323, 0.04392975, 0.1346183 ],
33      [0.0868652 , 0.14934175, 0.14963081, 0.12654245, 0.02132517,
34      0.19728671, 0.07576374, 0.04405305, 0.00963052, 0.13956061],
35      [0.1538449 , 0.00376636, 0.11627368, 0.05142717, 0.14649021,
36      0.0395909 , 0.1568814 , 0.08781049, 0.21268995, 0.03122492],
37      [0.06412921, 0.02133593, 0.17386608, 0.16496227, 0.04849963,
38      0.1240939 , 0.1536587 , 0.10439197, 0.09958787, 0.04547442],
39      [0.01431886, 0.1379885 , 0.138481 , 0.09736869, 0.05214154,
40      0.05370715, 0.11164933, 0.13797227, 0.13643065, 0.119942 ]],
41
42      [[0.19263895, 0.02524585, 0.04849603, 0.26960749, 0.18195654,
43      0.00275954, 0.03044612, 0.19908098, 0.00151871, 0.04824979],
44      [0.0973444 , 0.12274097, 0.11565676, 0.03978498, 0.12633932,
45      0.04208756, 0.05772533, 0.13242624, 0.11524372, 0.15065072],
46      [0.12930379, 0.11174424, 0.01841889, 0.07230249, 0.05214568,
47      0.04797471, 0.19131917, 0.07729323, 0.17539954, 0.12409826],
48      [0.18013797, 0.11391889, 0.1307509 , 0.11162541, 0.04425035,
49      0.16373831, 0.06363493, 0.00551103, 0.14629141, 0.0401408 ],
50      [0.12822126, 0.13005775, 0.12473177, 0.0504671 , 0.00210734,
51      0.1265661 , 0.05837824, 0.13179283, 0.13137907, 0.11629853]],
52
53      [[0.19263895, 0.02524585, 0.04849603, 0.26960749, 0.18195654,
54      0.00275954, 0.03044612, 0.19908098, 0.00151871, 0.04824979],
55      [0.0973444 , 0.12274097, 0.11565676, 0.03978498, 0.12633932,
56      0.04208756, 0.05772533, 0.13242624, 0.11524372, 0.15065072],
57      [0.12930379, 0.11174424, 0.01841889, 0.07230249, 0.05214568,
58      0.04797471, 0.19131917, 0.07729323, 0.17539954, 0.12409826],
59      [0.18013797, 0.11391889, 0.1307509 , 0.11162541, 0.04425035,
60      0.16373831, 0.06363493, 0.00551103, 0.14629141, 0.0401408 ],
61      [0.12822126, 0.13005775, 0.12473177, 0.0504671 , 0.00210734,
62      0.1265661 , 0.05837824, 0.13179283, 0.13137907, 0.11629853]],
63
64      [[0.06042058, 0.07902161, 0.17465227, 0.06503203, 0.03477972,
65      0.11425498, 0.19209789, 0.14282452, 0.1169759 , 0.01994051],
66      [0.10363778, 0.16683866, 0.02360622, 0.08734618, 0.14785029,
67      0.12483043, 0.11745742, 0.11837892, 0.06057956, 0.04947455],
68      [0.1099875 , 0.11008973, 0.11783012, 0.12410411, 0.06948848,

```

```

69         0.06815317, 0.1084837 , 0.08832634, 0.09539324, 0.10814362],
70     [0.21203704, 0.08052479, 0.0894798 , 0.02239049, 0.13777087,
71     0.00856297, 0.11092521, 0.12928098, 0.06826629, 0.14076156],
72     [0.00848645, 0.01039183, 0.22888188, 0.1002201 , 0.03535355,
73     0.14530992, 0.21424441, 0.06005044, 0.17331418, 0.02374723]],
74
75     [[0.00962717, 0.09897969, 0.10070844, 0.11873918, 0.13525485,
76     0.18177979, 0.0961754 , 0.06015969, 0.14812571, 0.05045009],
77     [0.10902578, 0.01948593, 0.00629627, 0.23908958, 0.20762943,
78     0.17285665, 0.1015702 , 0.0430405 , 0.03885372, 0.06215194],
79     [0.09584904, 0.12470868, 0.11521524, 0.04885304, 0.16663961,
80     0.12877508, 0.09674385, 0.10675528, 0.07322707, 0.04323311],
81     [0.1028186 , 0.2188951 , 0.00415739, 0.03352624, 0.01328733,
82     0.01176404, 0.24708992, 0.20324345, 0.13695964, 0.02825827],
83     [0.11520358, 0.11095176, 0.04058753, 0.10166731, 0.09338424,
84     0.14431621, 0.14882567, 0.01061639, 0.08778544, 0.14666187]],
85
86     [[0.10018177, 0.17053963, 0.13115552, 0.03244262, 0.01405127,
87     0.12791506, 0.00527879, 0.11663647, 0.18721357, 0.11458531],
88     [0.07064585, 0.11707668, 0.08340076, 0.09930075, 0.17134399,
89     0.07026961, 0.17493403, 0.16477132, 0.03563345, 0.01262357],
90     [0.03099827, 0.00560485, 0.02904968, 0.21008584, 0.02189689,
91     0.09811361, 0.25987493, 0.00715821, 0.2505221 , 0.08669562],
92     [0.02010703, 0.11855739, 0.10702145, 0.14931138, 0.12508031,
93     0.136721 , 0.04799124, 0.03019326, 0.12772524, 0.13729169],
94     [0.14971506, 0.06236725, 0.05623061, 0.11735498, 0.05151252,
95     0.14068416, 0.12974926, 0.06484254, 0.11349432, 0.11404931]],
96
97     [[0.0208345 , 0.18234617, 0.10207804, 0.16697232, 0.06466083,
98     0.18092594, 0.07863188, 0.00218957, 0.18291774, 0.01844299],
99     [0.05393382, 0.16047065, 0.16056274, 0.09685679, 0.10672076,
100     0.0757449 , 0.0495249 , 0.05551324, 0.11359204, 0.12708017],
101     [0.1731693 , 0.17274033, 0.01995265, 0.10816155, 0.01259181,
102     0.12021735, 0.09659115, 0.19419806, 0.07676771, 0.02561009],
103     [0.03495161, 0.18613685, 0.1511117 , 0.02471752, 0.02055832,
104     0.17133863, 0.01778554, 0.20088813, 0.17262754, 0.01988417],
105     [0.01468904, 0.17082958, 0.06480231, 0.06417404, 0.14073043,
106     0.16400931, 0.17071903, 0.13044204, 0.06514666, 0.01445755]],
107
108     [[0.19799489, 0.1422655 , 0.10807968, 0.23091328, 0.02832995,
109     0.12550684, 0.00289259, 0.11940138, 0.01434447, 0.03027142],
110     [0.02035325, 0.11243055, 0.12920041, 0.10102808, 0.1666295 ,
111     0.06492027, 0.04947976, 0.15042441, 0.03872243, 0.16681134],
112     [0.00218722, 0.17453113, 0.00776669, 0.16036252, 0.0949606 ,
113     0.17868549, 0.0132798 , 0.0996669 , 0.17442743, 0.09413223],
114     [0.10386453, 0.11481374, 0.07500921, 0.10356079, 0.09642463,
115     0.14871076, 0.00749965, 0.04636507, 0.15683866, 0.14691298],
116     [0.09149112, 0.12451615, 0.05569525, 0.03777276, 0.0931058 ,
117     0.07094944, 0.11719205, 0.0156083 , 0.1956483 , 0.19802082]]])
118
119 P1.shape
120
121 #Reward function
122 R1 = np.array([

```

```

123     [0.7 , 0.54, 0.31, 0.81, 0.68],
124     [0.16, 0.91, 0.82, 0.95, 0.73],
125     [0.61, 0.42, 0.93, 0.87, 0.05],
126     [0.03, 0.38, 0.81, 0.99, 0.15],
127     [0.59, 0.38, 0.97, 0.84, 0.84],
128     [0.47, 0.41, 0.27, 0.06, 0.86],
129     [0.81, 1. , 1. , 0.56, 0.77],
130     [0.94, 0.85, 0.25, 0.45, 0.13],
131     [0.95, 0.61, 0.23, 0.67, 0.62],
132     [0.36, 0.11, 0.67, 0.52, 0.77]])
133
134 R1.shape
135
136 ##-----Define the MDP
137 ↪ components-----##
138
139 num_states = 20
140 num_actions = 10
141
142 ###----- Transition probability matrix P[s, a,
143 ↪ s']-----###
144
145 # Random generation of matrix
146 def Markov_decision_matrix(num_states, num_actions):
147     Markov = []
148     for s in range(num_states):
149         matrix = np.random.rand(num_actions, num_states)
150         # Normalization of each row
151         matrix = matrix/matrix.sum(axis=1, keepdims=True)
152         Markov.append(matrix)
153     return np.array(Markov)
154
155 #-----Transition probabilities-----#
156
157 np.random.seed(42)
158 P = Markov_decision_matrix(num_states, num_actions)
159
160
161 ##-----Reward function R[s,
162 ↪ a]-----##
163
164 R = np.round(np.random.rand(num_states, num_actions),2) # random
165 ↪ generation of R
166
167 # Discount factor
168 gamma = 0.9
169
170 # Coefficient of v in objective function
171 mu = np.random.randint(1,num_actions, size = num_states)
172
173 ##-----Function to solve MDPs using linear
174 ↪ programm-----##

```

```

172
173 def Linear_Programming(P,R,gamma, nb_states,nb_actions,disp = "F", mu =
    ↪ None):
174     states = [f's{i}' for i in range(nb_states)]
175     actions = [f'a{i}' for i in range(nb_actions)]
176
177     # Initializing a linear programming problem
178     lp = pulp.LpProblem("MDP", pulp.LpMinimize)
179
180     # Definition of the value variables
181     V = pulp.LpVariable.dicts("V", states, lowBound=None, cat='Continuous')
182
183     # Objective function: minimize the sum of state values
184     if mu is not None:
185         lp += pulp.lpSum([mu[s] * V[states[s]] for s in range(nb_states)])
186     else:
187         lp += pulp.lpSum([V[s] for s in states]), "Objective"
188
189     # Adding constraints for each state and action
190     for s in range(nb_states):
191         for a in range(nb_actions):
192             lp += V[states[s]] >= R[s,a] + gamma *
                ↪ pulp.lpSum([P[s,a,s_prime] * V[states[s_prime]] for s_prime
                ↪ in range(nb_states)]) #, f"Constraint_{s}_{a}"
193
194     # Solving the linear program
195     cplex_path = "/path/to/CPLEX_StudioXXX/cplex/bin/x86-64_linux/cplex"
196
197     solver1 = pulp.CPLEX_CMD(path=cplex_path)
198
199     solver2 = pulp.PULP_CBC_CMD()
200     #lp.solve(solver1)
201     lp.solve(solver2)
202
203
204
205     if disp=="T":
206
207         # Print the optimal values
208         for s in states:
209             print(f"V({s})= {pulp.value(V[s])}")
210
211         # Deriving the optimal policy
212         optimal_policy = {}
213         for s in range(nb_states):
214             best_action = None
215             best_value = float('-inf')
216             for a in range(nb_actions):
217                 value = R[s, a] + gamma * sum(P[s,a,s_prime] *
                    ↪ pulp.value(V[states[s_prime]]) for s_prime in
                    ↪ range(nb_states))
218                 if value > best_value:
219                     best_value = value
220                     best_action = a

```

```

221         optimal_policy[states[s]] = best_action
222
223     print("\n Optimal policy:")
224     for s in range(nb_states):
225         print(f"{s}: a{optimal_policy[states[s]]}")
226
227     return lp,V
228
229 pulp.listSolvers(onlyAvailable=False)
230
231 #Prog_L,V = Linear_Programming(P,R,gamma, num_states,num_actions,"F")
232
233 Pl, _ = Linear_Programming(P1,R1,gamma, 10,5,"T",mu)
234
235 Pl
236
237 ##----- Impact of gamma on the resolution of Linear Program
238 ↪ -----##
239 Gamma = [0.2,0.5,0.7,1-1e-1, 1-1e-2, 1-1e-3,1-1e-4]
240 LP_times = []
241
242 for gam in Gamma:
243     start = time.time()
244     _ ,_ = Linear_Programming(P,R,gam, num_states,num_actions)
245     end = time.time()
246     LP_times.append(end-start)
247
248 Gamma_log = [np.log10(1/(1-g)) for g in Gamma] # Log(1/(1-gamma))
249 Times_log = [np.log10(t) for t in LP_times] # Log(times)
250 plt.figure()
251 for i in range(len(LP_times)):
252     plt.text(Gamma_log[i],Times_log[i],f'{Gamma[i],round(LP_times[i],5)}')
253
254 plt.plot(Gamma_log,Times_log, "-*", label="$\\gamma, t=f(\\gamma)$ \n
255 ↪ Linear Program")
256 plt.xlabel("$\\log[\\frac{1}{1-\\gamma}]$")
257 plt.ylabel("$\\log(t)$")
258 plt.title("Times to compute $V^*$ when $\\gamma$ increases")
259 plt.legend()
260 plt.grid(1)
261
262 ###-----Analyze of convergence by varying the number of
263 ↪ state-----###
264
265 LPState_times = []
266 S = [10,20,30,50,100,150,200,300,400,500] #states
267 a = 10 # fixed action
268 gamma = 0.7
269
270 for s in S:
271     Pr = Markov_decision_matrix(s, a)
272     Rr = np.round(np.random.rand(s, a),2)
273     start = time.time()

```

```

272     _ , _ = Linear_Programming(Pr,Rr,gamma, s,a)
273     end = time.time()
274     LPState_times.append(end-start)
275
276 plt.figure()
277
278 plt.plot(S,LPState_times, "-*", label="$Linear Programing$")
279 plt.xlabel("$S$")
280 plt.ylabel("$CPU_{time}(s)$")
281 plt.title("Times to compute $V^*$ when $$ increases")
282 plt.legend()
283 plt.grid(1)
284
285 plt.figure()
286 Sss = [s**3 for s in S]
287 plt.plot(Sss,LPState_times, "-*", label="")
288 plt.xlabel("$S^3$")
289 plt.ylabel("$CPU_{time}(s)$")
290 plt.grid(1)
291 plt.legend()
292
293 ###-----Analyze of convergence by varying the number of
294 ↪ action-----###
295
296 LPAction_times = []
297 A = [10,20,30,50,100,150,200,300,400,500] #states
298 s = 10 # fixed action
299 gamma = 0.7
300
301 for a in A:
302     Pr = Markov_decision_matrix(s, a)
303     Rr = np.round(np.random.rand(s, a),2)
304     start = time.time()
305     _ , _ = Linear_Programming(Pr,Rr,gamma, s,a)
306     end = time.time()
307     LPAction_times.append(end-start)
308
309 plt.figure()
310
311 plt.plot(A,LPAction_times, "-*", label="$Linear Programing$")
312 plt.xlabel("$A$")
313 plt.ylabel("$CPU_{time}(s)$")
314 plt.title("Times to compute $V^*$ when $A$ increases")
315 plt.legend()
316 plt.grid(1)
317

```