Planning with a Known Model

(INF8250AE: Introduction to Reinforcement Learning)

Amir-massoud Farahmand

Polytechnique Montréal & Mila







Table of Contents

- 1 Some Initial Attempts
- 2 Value Iteration
- 3 Policy Iteration
 - Convergence of Policy Iteration
- 4 Linear Programming

Goal

We study some algorithms to solve the <u>Planning problem</u>: computing the (optimal) value function for a given MDP.

- Know the fundamental planning algorithms:
 - 🥦 Value Iteration (VI)
 - Policy Iteration (PI)
 - Linear Programming
- Learn the intuition behind them
- Learn the mathematics justifying them

We refer to these frequently in studying and analyzing RL/Planning algorithms.

Learning Objectives

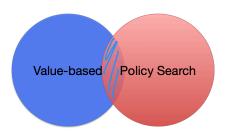
You need to

- Remember: Value Iteration, Policy Iteration, and LP-based algorithms
- Understand: Why VI, PI, LP work
- Apply: VI and PI to solve an MDP

How to Compute the Optimal Policy π^* ?

- We have defined concepts and properties such as
 - Value function for a policy π (V^{π}) and optimal value function (V^{*})
 - Relation between V^* (or Q^*) and π^* through the greedy policy
- Question: How can we find the optimal policy?
- Assumption: MDP is known, that is, we know \mathcal{R} and \mathcal{P} .
- The assumption of knowing the MDP does not hold in the RL setting.
- But designing methods for finding the optimal policy with known model provides the foundation for developing methods for the RL setting.

Different Approaches to Find π^*



- Value-based: Compute Q^* (or V^*) and then $\underline{\pi}^* \leftarrow \pi_g(Q^*)$.
- Direct policy search: Search in the space of policies without explicitly constructing the optimal value function.
- Hybrid: Explicitly construct value function to guide the search
 7 in the policy space.

Policy Evaluation vs. Control Problems

- Policy Evaluation (PE): Problem of computing the value function of a given policy $\underline{\pi}$, i.e., \underline{V}^{π} or Q^{π} .
 - Not the ultimate goal of an RL agent (finding the optimal policy is), but is often needed as an intermediate step in finding the optimal policy.
- Control: Problem of finding the optimal value function $\underline{V^*}$ or $\overline{Q^*}$ or optimal policy π^* .

Dynamic Programming (DP): Methods that benefit from the structure of the MDP, such as the recursive structure encoded in the Bellman equation, in order to compute the value function.

Policy Evaluation

Problem Statement: Given an MDP $(\mathcal{X}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$ and a policy π , we would like to compute \underline{V}^{π} or \underline{Q}^{π} .

$$V^{\pi}(x) = \mathbb{E}\left[\sum_{t=1}^{\infty} \gamma^{t-1} R_t | X_1 = x\right]$$

Policy Evaluation: A Naive Approach

$$X_1 = x$$

$$X_1$$

Idea: Expand the tree of all possible futures! Example: the expected reward at time t=2 is

$$\sum_{a,x',a'} \pi(a|x) \mathcal{P}(x'|x,a) \pi(a'|x') r(x',a').$$

Remark

- This is inefficient. The size of the tree grows very fast.
- The sample-based approximate version of this idea works.

Policy Evaluation: Linear System of Equations

Q: Can we improve the efficiency?

Key Idea: Benefit from the recursive structure of the value function

$$V^{\pi} = T^{\pi}V^{\pi}.$$

$$V(x) = r^{\pi}(x) + \gamma \sum_{x' \in \mathcal{X}} \mathcal{P}^{\pi}(x'|x)V(x'), \quad \forall x \in \mathcal{X}$$

In the discrete state-action case:

- $\underline{n} = |\mathcal{X}|$ equations
- lacksquare $|\mathcal{X}|$ unknowns $(\underline{V(x_1)},\dots,\underline{V(x_n)})$

Policy Evaluation: Linear System of Equations

We have n equations in the form of:

$$\underbrace{V(x)} - \gamma \sum_{x' \in \mathcal{X}} \underbrace{\mathcal{P}^{\pi}(x'|x)} \underbrace{V(x')} = \underbrace{r^{\pi}(x)},$$

More compactly in the matrix form:

$$(\underline{\mathbf{I} - \gamma \mathcal{P}^{\pi}})V = \underline{r^{\pi}},$$

which is the same form of a generic linear system of equations:

$$A_{n\times n}x_{n\times 1} = b_{n\times 1}$$

Policy Evaluation: Linear System of Equations

Solving $A_{n\times n}x_{n\times 1}=b_{n\times 1}$:

- \rightarrow Compute A^{-1} and then calculate $A^{-1}b$.
 - Better: Use various linear solvers.

Remark

To solve the control problem of finding \underline{V}^* , we need to solve $V=T^*V$, i.e.,

$$\underline{V(x)} = \max_{\underline{a \in \mathcal{A}}} \left\{ r(x, a) + \gamma \sum_{x' \in \mathcal{X}} \mathcal{P}(x'|x, a) V(x') \right\}.$$

This is not a linear system of equations anymore!

Value Iteration (PE)

Starting from $V_0 \in \mathcal{B}(\mathcal{X})$, we compute a sequence of $(V_k)_{k \geq 0}$ by

$$V_{\underline{k+1}} \leftarrow T^{\pi}V_{\underline{k}} [= r^{\pi} + \gamma P^{\pi}V_{\underline{k}}].$$

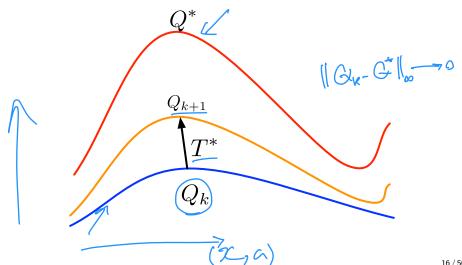
By the contraction property of the Bellman operator:

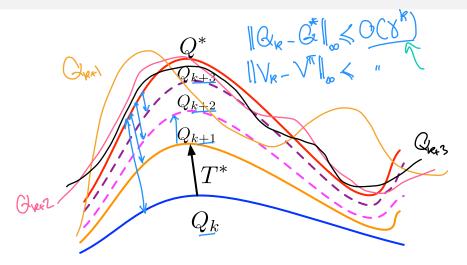
$$\lim_{k\to\infty} \|V_k - V^\pi\|_{\infty} = 0.$$

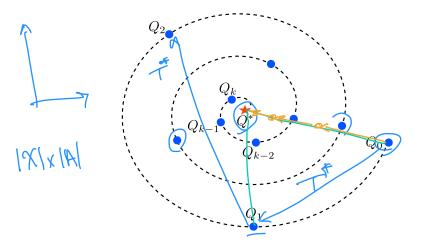
Remark

Similar procedure to compute Q^{π} , i.e., $Q_{k+1} \leftarrow T^{\pi}Q_k$.

By the contraction property of the Bellman optimality operator, it is guaranteed that $V_k \to V^*$ (or $Q_k \to Q^*$).





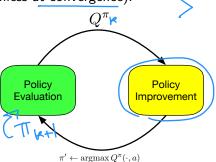


- The VI uses an existing approximation V_k of V^* (or V^{π}) to get a better approximation of V^* (or V^{π}). This idea is called bootstrapping in the in the RL literature.
 - Note: This is different from the bootstrapping method in statistics.
- VI is one of the fundamental algorithms for Planning.
- The current formulation of VI is a special case of an operator/matrix splitting-based formulation.
- Many RL algorithms are essentially the sample-based variants of VI too.

Mudineter: 1X1, 1A

A different approach is based on the iterative application of the following two steps:

- (Policy Evaluation) Given a policy $\underline{\pi_k}$, compute $\underline{V^{\pi_k}}$ (or $\underline{Q^{\pi_k}}$).
- (Policy Improvement) Find a new policy π_{k+1} that is better than π_k , i.e. $V^{\pi_{k+1}} \ge V^{\pi_k}$ (with a strict inequality in at least one state, unless at convergence).



Q: How to perform Policy Evaluation and Policy Improvement?

- Policy Evaluation: This is clear. We can either solve a linear system of equations or even perform VI (PE) to compute the value of a policy π_k .
- Policy Improvement: Choose the greedy policy, i.e.,

$$\underline{\pi_{k+1}(x)} \leftarrow \underline{\pi_g(x; Q^{\pi_k})} = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q^{\pi_k}(\underline{x}, \underline{a}), \quad \forall x \in \mathcal{X}.$$

Remark

The Policy Iteration (PI) algorithm refers to the specific case that we pick the new policy π_{k+1} as $\pi_g(Q^{\pi_k})$.

Why Greedy Policy for Policy Improvement? (Intuition)

Assume that at state
$$x$$
, we act according to $\pi_g(x;Q^{\pi_k})$, and afterwards, we follow π_k . The value of this new policy is
$$Q^{\pi_k}(x,\pi_g(x;Q^{\pi_k})) = Q^{\pi_k}(x,\arg\max_{a\in\mathcal{A}}Q^{\pi_k}(x,a)) = \max_{a\in\mathcal{A}}Q^{\pi_k}(x,a).$$
 Comparing $\max_{a\in\mathcal{A}}Q^{\pi_k}(x,a)$ with $V^{\pi_k}(x) = Q^{\pi_k}(x,\pi_k(x))$, we see
$$Q^{\pi_k}(x,\pi_g(x;Q^{\pi_k})) = \max_{a\in\mathcal{A}}Q^{\pi_k}(x,a).$$

So this new policy is equal to or better than π_k at state x.

Recall that:

- $lackbox{V}^{\pi_k}$ is the unique fixed point of T^{π_k} .
- The greedy policy satisfies $T^{\pi_{k+1}}Q^{\pi_k} = T^*Q^{\pi_k}$.

We can summarize each iteration of the Policy Iteration algorithm:

greedy policy

- (Policy Evaluation) Given $\underline{\pi_k}$, compute Q^{π_k} , i.e., find a Q that satisfies $Q = T^{\pi_k}Q$.
- (Policy Improvement) Obtain $\underline{\pi_{k+1}}$ as a policy that satisfies $T^{\pi_{k+1}}Q^{\pi_k}=T^*Q^{\pi_k}$

- 1: Initialize π_0 arbitrarily
- 2: $k \leftarrow 0$
- 3: repeat
- 4: $Q^{\pi_k} \leftarrow \text{ solution of } Q = T^{\pi_k}Q$

▶ Policy Evaluation:

compute
$$Q^{\pi_k}$$

5:
$$\pi_{k+1} \leftarrow \text{ policy s.t.} T^{\pi_{k+1}} Q^{\pi_k} = T^* Q^{\pi_k}$$

Improvement

6:
$$k \leftarrow k + 1$$

7: **until**
$$\pi_k = \pi_{k-1}$$



⊳ Policy

Approximate Policy Iteration

popper. Value Hention _ VRA = TVR

We also have approximate policy iteration algorithms too, where policy evaluation or improvement steps are performed approximately:

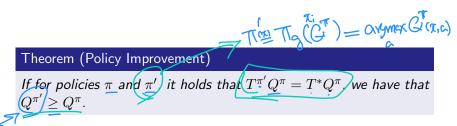
- $Q \otimes T^{\pi_k}Q$
- $T^{\pi_{k+1}}Q^{\pi_k} \otimes T^*Q^{\pi_k}$

We discuss this later when we get to function approximation.



- The Policy Iteration algorithm converges to the optimal policy.
- For finite MDPs, the convergence happens in a finite number of iterations.

Policy Improvement Theorem



In other words, the greedy policy is a proper policy improvement step.

Policy Improvement Theorem (Proof)

We first show that
$$T^{\pi'}Q^{\pi} \geq Q^{\pi}$$
 Notice that $T^{\pi'}Q^{\pi} = T^*Q^{\pi}$, by the assumption. We have $T^*Q^{\pi} \geq T^*Q^{\pi} = Q^{\pi}$ because for any $(x,a) \in \mathcal{X} \times \mathcal{A}$, it holds that
$$r(x,a) + \gamma \int \mathcal{P}(\mathrm{d}x'|x,a) \max_{a' \in \mathcal{A}} Q^{\pi}(x',a') \geq r(x,a) + \gamma \int \mathcal{P}(\mathrm{d}x'|x,a) Q^{\pi}(x',\pi(x'))$$
 Therefore, $T^{\pi'}Q^{\pi} = T^*Q^{\pi} \geq T^{\pi}Q^{\pi} = Q^{\pi}$.

Policy Improvement Theorem (Proof)

The second step is to use $T^{\pi'}Q^{\pi} \geq Q^{\pi}$ to show that $Q^{\pi'} \geq Q^{\pi}$. Apply $T^{\pi'}$ to both sides of $T^{\pi'}Q^{\pi} \geq Q^{\pi}$, and use the monotonicity property of the Bellman operator to conclude

$$T^{\pi'}(T^{\pi'}Q^{\pi}) \ge T^{\pi'}Q^{\pi} = T^*Q^{\pi} \ge Q^{\pi}.$$

So we also have $(T^{\pi'})^{\underline{2}}Q^{\pi} \geq Q^{\pi}$.

By repeating this argument, we get that for any $m \ge 1$,

$$(T^{\pi'})^{\underline{m}}Q^{\pi} \ge \underline{T}^*\underline{Q}^{\pi} \ge \underline{Q}^{\pi}. \tag{1}$$

Policy Improvement Theorem (Proof)

$$(T^{\pi'})^m \underline{Q}^{\pi} \ge T^* Q^{\pi} \ge \underline{Q}^{\pi}.$$

Take the limit of $m \to \infty$.

Because of the contraction property of the Bellman operator $T^{\pi'}$:

$$\lim_{m \to \infty} (T^{\overline{\pi'}})^m Q^{\overline{\pi}} = \underline{Q}^{\overline{\pi'}}. \tag{2}$$

By combining (1) and (2), we get that

$$Q^{\pi'} = \lim_{m \to \infty} (T^{\pi'})^m Q^{\pi} \nearrow T^* Q^{\pi} \nearrow Q^{\pi}, \tag{3}$$

- The Policy Improvement theorem shows that if we are given π_k , the new policy π_{k+1} is at least as good as the previous one.
- We can show that the PI algorithm converges to an optimal policy. We shall prove this.
- \blacksquare If $|\mathcal{X}\times\mathcal{A}|<\infty,$ this happens in a finite number of iterations.

Theorem (Convergence of the Policy Iteration Algorithm)

Let $(\pi_k)_{k\geq 0}$ be the sequence generated by the PI algorithm.

- For all \underline{k} , we have that $\underline{V}^{\pi_{k+1}} \geq \underline{V}^{\pi_k}$, with equality if and only if $V^{\pi_k} = V^*$.
- If the set of policies is finite, the PI algorithm converges in a finite number of iterations.

Remark

We follow the line of proof of Proposition 2.4.1 of Bertsekas 2018.

Convergence of Policy Iteration (Proof)

The basic idea behind the proof is that either we can strictly improve the policy, or if we cannot, we are already at the optimal policy.

Convergence of Policy Iteration (Proof)

Proof of $V^{\pi_{k+1}} > V^{\pi_k}$:

By the Policy Improvement Theorem (Theorem 1), we have that

$$V^{\pi_{k+1}} \ge V^{\pi_k}$$

$$V^{\pi_{k+1}} \ge V^{\pi_k}$$
Proof of $V^{\pi_{k+1}} = V^{\pi_k} \Rightarrow V^{\pi_k} = V^*$:

Suppose that instead of a strict inequality, we have an equality of

$$V^{\pi_{k+1}} = V^{\pi_k}.$$

Apply $T^{\pi_{k+1}}$ to both side to get

$$T^{\pi_{k+1}}V^{\pi_{k+1}} = T^{\pi_{k+1}}V^{\pi_k}.$$

As $T^{\pi_{k+1}}V^{\pi_k}=T^*V^{\pi_k}$ by the definition of the PI algorithm, we get that

$$T^{\pi_{k+1}}V^{\pi_{k+1}} = T^{\pi_{k+1}}V^{\pi_k} = T^*V^{\pi_k} = T^*V^{\pi_{k+1}},$$

where in the last step we used $V^{\pi_{k+1}} = V^{\pi_k}$ again.

Convergence of Policy Iteration (Proof)

By these equalities, we have

$$T^{\pi_{k+1}}V^{\pi_{k+1}} = T^*V^{\pi_{k+1}}.$$

 $T^{\pi_{k+1}}V^{\pi_{k+1}}=T^*V^{\pi_{k+1}}.$ As $V^{\pi_{k+1}}$ is the value function of π_{k+1} it satisfies the Bellman equation $T^{\pi_{k+1}}V^{\pi_{k+1}}=V^{\pi_{k+1}}$ Therefore, we also have $V^{\pi_{k+1}} = T^* V^{\pi_{k+1}}.$

This means that $V^{\pi_{k+1}}$ is a fixed point of T^* . But the fixed point of T^* is unique and is equal to V^* . So we must have that

$$V^{\pi_{k+1}} = V^*$$

Convergence of Policy Iteration

Convergence of Policy Iteration (Proof)

Proof of $V^{\pi_k}=V^*\Rightarrow V^{\pi_{k+1}}=V^{\pi_k}$: If $V^{\pi_k}=V^*$, then π_k is an optimal policy. The greedy policy of $V^{\pi_k}=V^*$ is still an optimal policy, hence $V^{\pi_{k+1}}=V^*=V^{\pi_k}$.

Convergence of Policy Iteration (Proof)

Proof of $\lim_{k\to\infty} \|V^{\pi_k} - V^*\|_{\infty} = 0.$

To prove the convergence, recall from (3) that

$$Q^{\pi_{k+1}} \ge T^* Q^{\pi_k} \ge Q^{\pi_k}.$$

By induction,

$$Q^{\pi_{k+1}} \ge T^*Q^{\pi_k} \ge T^*(T^*Q^{\pi_{k-1}}) \ge \cdots \ge (T^*)^k Q^{\pi_0}.$$

By the definition of the optimal policy, we have $Q^{\pi} \leq Q^*$ for any π , including all π_k generated during the iterations of the PI algorithm. So $Q^{\pi_{k+1}}$ is sandwiched between Q^* and $(T^*)^k Q^{\pi_0}$, i.e.,

$$Q^* \ge Q^{\pi_{k+1}} \ge (T^*)^k Q^{\pi_0}.$$

This entails that $\|Q^{\pi_{k+1}} - Q^*\|_{\infty} \le \|(T^*)^k Q^{\pi_0} - Q^*\|_{\infty}$

Convergence of Policy Iteration (Proof)

We show that the RHS of $\|Q^{\pi_{k+1}} - Q^*\|_{\infty} \le \|\underline{(T^*)^k Q^{\pi_0} - Q^*}\|_{\infty}$ converges to zero.

By the contraction property of the Bellman optimality operator, we have that

$$\lim_{k\to\infty} \left\| (T^*)^k Q^{\pi_0} - Q^* \right\|_{\infty} = 0.$$

Therefore,

$$\lim_{k \to \infty} \|Q^{\pi_k} - Q^*\|_{\infty} = 0.$$

This implies the convergence of V^{π_k} too.

Convergence of Policy Iteration

Convergence of Policy Iteration (Proof)

Proof of finite convergence:

If the number of policies is finite, the number of times (4) can be a strict inequality is going to be finite too.

Convergence of Policy Iteration: Some Remarks



- The PI algorithm converges to the optimal policy in a finite number of iterations whenever the number of policies is finite.
- If the state space \mathcal{X} and the action space \mathcal{A} are finite, the number of policies are finite and is $|\mathcal{A}|^{|\mathcal{X}|}$
- Even though finite, this can be very large.
 - Example: A 10×10 grid world problem with 4 actions at each state has $4^{100} \approx 1.6 \times 10^{60}$ possible policies.
- In practice, PI converges much faster.
- This suggest that the previous analysis might be crude.

$$\pi(\alpha_{l}|x_{l}) = P_{l} \rightarrow \pi(\alpha_{l}|x_{l}) = 1 - P_{l}$$

$$P_{l} \in [\alpha_{l}|x_{l}] \rightarrow P_{l} = 0, |x_{l}|$$

Fast Convergence of Policy Iteration

It can be shown that the PI algorithm converges in

$$O\left(\frac{|\mathcal{X}||\mathcal{A}|}{1-\gamma}\log\left(\frac{1}{1-\gamma}\right)\right)$$

fce = 4000

iterations.

The proof is in the Foundations of Reinforcement Learning [Farahmand, 2025].

This is a significant quantitative improvement over the previous result.

Remark

This is a relatively recent result, which in various forms have been proven by Ye [2011]; Hansen et al. [2013]; Scherrer [2016].

Linear Programming

Linear Programming for Finding V^*



We can find V^* by solving a Linear Program (LP) too. Consider the set of all V that satisfy $V \geq T^*V$, i.e.,

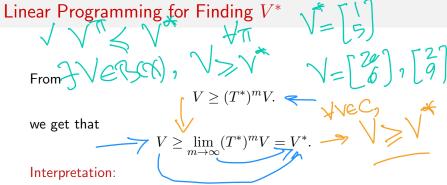
$$C = \{V : V \ge T^*V\}.$$

Interesting property: For any $V \in C$, we have

$$V \ge T^*V \Rightarrow T^*V \ge T^*(T^*V) = (T^*)^2V.$$

Repeating this argument, we get that for any $m \geq 1$,

$$V \ge (T^*)^m V$$



- Any $V \in C$ is lower bounded by V^* .
- (OR) V^* is the function in C that is smaller or equal to any other function in C (pointwise sense).



Choose a strictly positive vector $\mu > 0$ with the dimension of \mathcal{X} .

$$\frac{\min_{v \in \mathcal{O}} \mu(x)}{\sqrt{x}} \sqrt{x}$$

Can be written as

$$\min_{V} \mu^{\top} V,$$
s.t. $V(x) \ge (T^*V)(x)$ $\forall x \in \mathcal{X}$.

Linear objective; nonlinear constraints.

Linear Programming for Finding V^*

Each nonlinear constraint:

$$V(x) \ge \max_{a \in \mathcal{A}} \left\{ r(x, a) + \gamma \sum_{y} \mathcal{P}(y|x, a) V(y) \right\}$$

is equivalent to

$$V(x) \ge r(x, a) + \gamma \sum_{y} \mathcal{P}(y|x, a)V(y), \quad \forall a \in \mathcal{A}.$$

Linear Programming for Finding
$$V^*$$

$$\mathcal{L} = \mathcal{L} \times \mathcal{L}$$

This is a linear program with $|\mathcal{X} \times \mathcal{A}|$ constraints.

Summary



- Three methods for computing the optimal value function
 - Value Iteration
 - Policy Iteration
 - Linear Programming
- Established convergence of VI and PI
- These methods have variants for the RL setting.

References

- Dimitri P. Bertsekas. *Abstract dynamic programming*. Athena Scientific Belmont, 2nd edition, 2018.
- Amir-massoud Farahmand. Foundations of Reinforcement Learning (draft). 2025.
- Thomas Dueholm Hansen, Peter Bro Miltersen, and Uri Zwick. Strategy iteration is strongly polynomial for 2-player turn-based stochastic games with a constant discount factor. *Journal of the ACM (JACM)*, 60(1):1–16, 2013.
- Bruno Scherrer. Improved and generalized upper bounds on the complexity of policy iteration. *Mathematics of Operations Research*, 41(3):758–774, 2016.
- Yinyu Ye. The simplex and policy-iteration methods are strongly polynomial for the Markov decision problem with a fixed discount rate. *Mathematics of Operations Research*, 36(4): 593–603, 2011.