

Model-based Value methods

We will study two common Model-based, Dynamic Programming type methods

- Value Iteration
- Policy Iteration

and their associated Bellman equations.

Value iteration method

The simplest method is to

- iteratively update the Value function
 - until convergence
- derive the *final* Policy from the Value function
 - chose the action leading to highest return
 - based on the Value function

The estimate of the value of a state \state is improved iteratively.

Let

$$\text{\statevalfun}_k(\text{\state})$$

denote the *estimated* value of state \state after k iterations.

In iteration $k + 1$, the value of *each* state $\text{\state} \in \text{\States}$ is updated via the Bellman update

$$\begin{aligned}\text{\statevalfun}_{k+1}(\text{\state}) &= \max \text{\act} \\ &\sum_{\text{\state}'} \text{\transp}(\text{\state}' | \text{\state}, \text{\act})(\text{\rew}(\text{\state}, \text{\act}, \text{\state}') \\ &\quad + \gamma \text{\statevalfun}_k(\text{\state}')).\end{aligned}$$

where

$$\text{\rew}(\text{\state}, \text{\act}, \text{\state}')$$

denotes the reward returned by the environment when action \act in state \state results in successor state $\text{\state}'$

That is

- $\text{\statevalfun}_{k+1}(\text{\state})$
- learns the *previous iteration's* (k) estimate
 $\text{\statevalfun}_k(\text{\state}')$
of each successor state $\text{\state}'$
- and chooses the action leading to highest return

The backup in each iteration

- moves information about potential future (successor) states
- to the current state

The previous iteration's estimate of the successors of \state
 $\text{\statevalfun}_k(\text{\state}')$

in turn, are based on the iteration $(k - 1)$ values of the successors of $\text{\state}'$

So we need i iterations

- in order for information of states i steps ahead
- to reach state s

Hence the need to repeat until convergence.

Subtlety

Prior to convergence

- $\text{\statevalfun}_k(\text{\state})$ is a *non-random* variable
 - $\text{\transp}(\text{\state}, \text{\act})$ is known for all $\text{\state}, \text{\act}$ since method is model-based
- the estimate is *biased*
 - it is non-random, but not yet correct

We will subsequently contrast this to the iterative estimates produced by model-free methods.

Pseudo code for Value Iteration

Here is some pseudo-code

Value iteration

Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

```
| Δ ← 0
| Loop for each  $s \in \mathcal{S}$ :
|    $v \leftarrow V(s)$ 
|    $V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$ 
|   Δ ← max(Δ, |v - V(s)|)
until Δ < θ
```

Output a deterministic policy, $\pi \approx \pi_*$, such that
 $\pi(s) = \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

Attribution: <http://incompleteideas.net/book/RLbook2020.pdf#page=105>
[\(http://incompleteideas.net/book/RLbook2020.pdf#page=105\)](http://incompleteideas.net/book/RLbook2020.pdf#page=105)

Value iteration

Initialize value function $V(s)$ arbitrarily (e.g., zero for all states)
Repeat:
delta = 0
For each state s :
old_value = $V(s)$
 $V(s) = \max_{a} [R(s, a) + \gamma * \sum_{s'} [P(s' | s, a) * V(s')]]$
delta = max(delta, |old_value - V(s)|)
Until delta < threshold
Derive policy after value function converges
For each state s : $\pi(s) = \arg\max_a [R(s, a) + \gamma * \sum_{s'} [P(s' | s, a) * V(s')]]$
Return π, V

Policy iteration method

The Value Iteration method has one expensive step

- the $\max \backslash\text{act}$ is a search over all possible actions $\backslash\text{act} \in \backslash\text{Actions}$
 $\backslash\text{statevalfun}_{k+1}(\backslash\text{state}) = \max \backslash\text{act}$
$$\sum_{\backslash\text{state}'} \backslash\text{transp}(\backslash\text{state}' | \backslash\text{state}, \backslash\text{act})(\backslash\text{rew}(\backslash\text{state}, \backslash\text{act}, \backslash\text{state}') + \gamma \backslash\text{statevalfun}_k(\backslash\text{state}')).$$

Policy Iteration is a less computationally-expensive way to reach the optimal
\statevalfun.

Rather than updating the Policy once (after Value function convergence)

- we introduce a method that periodically updates the Policy.

Policy iteration is an algorithm that improves π_p to π_{p+1} by alternating two steps during round p

The algorithm alternates between

- Policy evaluation
 - update the estimate of \statevalfun $_{\pi_p}$ to \statevalfun $_{\pi_{p+1}}$
- Policy improvement
 - update π_p to π_{p+1} using the newly updated \statevalfun $_{\pi_{p+1}}$

It uses the Bellman update

$$\begin{aligned} \text{\statevalfun}_{\pi,k+1}(\text{\state}) = \\ \sum_{\text{\state}'} \text{\transp}(\text{\state}' | \text{\state}, \pi(\text{\state})) (\text{\rew}(\text{\state}, \pi(\text{\state}), \text{\state}') \\ + \gamma \text{\statevalfun}_{\pi,k}(\text{\state}')). \end{aligned}$$

where

- we add an explicit subscript π
- to the value function

$$\text{\statevalfun}_{\pi,k+1}(\text{\state})$$

- to indicate the dependence of the update on the *current policy* π

Rather than

- evaluating *all* actions $\text{\act} \in \text{\Actions}$
- it chooses a single action $\text{\state}, \pi(\text{\state})$ according to the current policy

Here is some pseudo-code

Policy iteration

Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$; $V(\text{terminal}) \doteq 0$

2. Policy Evaluation

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_{s',r} p(s', r | s, \pi(s)) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement

policy-stable \leftarrow true

For each $s \in \mathcal{S}$:

$$\text{old-action} \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$$

If $\text{old-action} \neq \pi(s)$, then *policy-stable* \leftarrow false

If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Attribution: <http://incompleteideas.net/book/RLbook2020.pdf#page=102>
[\(http://incompleteideas.net/book/RLbook2020.pdf#page=102\)](http://incompleteideas.net/book/RLbook2020.pdf#page=102).

Pseudo code for Policy Iteration

Policy iteration

Initialize policy π arbitrarily (e.g., random policy) Initialize value function $V(s)$ arbitrarily (e.g., zero for all states)
Repeat: # Policy Evaluation (compute V for current policy π)
Repeat: For each state s : $V(s) = R(s, \pi(s)) + \gamma * \text{sum over } s' [P(s' | s, \pi(s)) * V(s')]$ Until $V(s)$ converges (changes smaller than threshold)
Policy Improvement (update policy based on current V)
 $\text{policy_stable} = \text{True}$
For each state s : $\text{old_action} = \pi(s)$ $\pi(s) = \text{argmax over } a [R(s, a) + \gamma * \text{sum over } s' [P(s' | s, a) * V(s')]]$
if $\text{old_action} \neq \pi(s)$: $\text{policy_stable} = \text{False}$
Until policy_stable is True
Return π, V

Both Value Iteration and Policy Iteration will converge to the same policy.

The advantage of alternating between Policy Evaluation and Policy Improvement

- faster convergence
 - the Value function under the current policy is fully evaluated
 - before the Policy is updated
- more stable convergence
 - Policy doesn't change until Value function (under current policy) is fully known

Finding the best action in a Value-based method

The Value-based methods don't directly give you a policy

- the Value function gives you the best successor state $\text{\textbackslash state}'$ from current state $\text{\textbackslash state}$
- but it doesn't directly tell you the action $\text{\textbackslash act}$ that leads to $\text{\textbackslash state}'$

In order to find $\text{\textbackslash act}$ you either

- need a model
 - search over all possible actions, using the model to determine the value of action's successor state
- use actual experience to estimate the effect of performing action $\text{\textbackslash act}$ in state $\text{\textbackslash state}$

Here is some pseudo-code that uses a model to determine the best action:

```
# For each possible action a in current state s:  
for each action a in actions:  
    # Take action a from state s in the environment  
    observe next state s', reward r  
    # Estimate value for taking action a from s  
    A[a] = r + gamma * V[s']  
# Find the maximum estimated action value  
A_max = max over a of A[a]  
  
# Update the value function for state s  
V[s] = V[s] + alpha * (A_max - V[s])
```

Q-learning: From Value function to State-Action function

We can more readily identify the optimal *action* to take from state $\backslash\text{state}$

- the one leading to $\backslash\text{state}'$ with maximal $\backslash\text{statevalfun}(\backslash\text{state}')$

with slightly more detailed bookkeeping.

Define a *State Value function* $\text{actvalfun}_\pi(\text{\state}, \text{\act})$

$$\text{actvalfun}_\pi : \text{\States} \times \text{\Actions} \rightarrow \text{\Reals}$$

to map a state/action pair $(\text{\state}, \text{\act})$ to the expected value of the successor state.

Re-write the Value function backup equation

$$\begin{aligned} \text{\statevalfun}_{k+1}(\text{\state}) &= \max \text{\act} \\ \sum_{\text{\state}'} \text{\transp}(\text{\state}' | \text{\state}, \text{\act})(\text{\rew}(\text{\state}, \text{\act}, \text{\state}')) \\ &\quad + \gamma \text{\statevalfun}_k(\text{\state}')). \end{aligned}$$

as

$$\text{\statevalfun}_{k+1}(\text{\state}) = \max \text{\act} \text{\actvalfun}_k(\text{\state}, \text{\act})$$

where

$$\begin{aligned} \text{\actvalfun}(\text{\state}, \text{\act}) &= \\ \sum_{\text{\state}'} \text{\transp}(\text{\state}' | \text{\state}, \text{\act})(\text{\rew}(\text{\state}, \text{\act}, \text{\state}')) \\ &\quad + \gamma \text{\statevalfun}_k(\text{\state}')) \end{aligned}$$

Key Differences Between Value Iteration and Policy Iteration

Here is a comparison of the Value and Policy Iteration methods.

Feature	Value Iteration	Policy Iteration
Approach	Updates value function until convergence	Alternates between value evaluation and improvement
Convergence	When value function $V(s)$ stabilizes	When policy $\pi(s)$ stops improving
Computational Cost	Higher per iteration, simpler logic	Lower per iteration, more complex
Speed	Requires more iterations	Fewer iterations; often faster
Policy Output	Extracted after value convergence	Updated during each iteration

Feature	Value Iteration	Policy Iteration
Approach	Updates value function iteratively using the Bellman Optimality Equation in one step	Alternates between full policy evaluation and policy improvement steps
Policy Handling	Policy is implicitly updated after value convergence	Explicitly maintains and updates policy each iteration
Initialization	Starts with an initial value function	Starts with an initial policy
Iteration Steps	Single step combining evaluation and improvement	Two-step process: separate evaluation and improvement
Convergence Criterion	Value function converges (changes below a threshold)	Policy stabilizes (no change between iterations)
Computation per Iteration	Potentially more expensive per iteration (max over all actions for each state)	More computationally intensive due to full policy evaluation, but fewer total iterations needed
Number of Iterations	Typically more iterations	Usually fewer iterations
Complexity	Simpler to implement	More complex implementation
Suitability	Suitable for smaller state spaces or when full policy evaluation is expensive	Can handle larger state spaces more efficiently when full evaluation is feasible
Policy Updates	Policy derived after convergence of value function	Policy updated after each evaluation phase

Model-based Value methods: Classical vs Deep Learning

In "Classical" Model-base Value methods

- the state and action spaces are finite
- leading to methods based on tables

Dynamic Programming style solutions are efficient methods for exploring the finite possibilities.

But non-finite action and state spaces may be accommodated

- by using Neural Networks (NN)
- to *learn*
 - an approximation of the model
 - the value function

Aspect	Classical Model-based RL (Tabular)	Model-based Value Method with NN
State/action space	finite	continuous/infinite (via approximation)
Model representation	table or analytic function	neural network (NN), e.g., $f\phi f \setminus \phi$
Value function representation	table	$NN: V\theta(s) \setminus \theta(s), Q\theta(s,a) \setminus \theta(s,a) Q\theta(s,a)$
Planning/Improvement	Bellman backup over all states	Rollouts/approximate backups, often via sampling
Example use cases	gridworlds, small MDPs	robotics, control, games with pixel input

In [2]: `print("Done")`

Done

