Machine Learning in Finance
Section I3
Lecture 14

Michael G Sotiropoulos

NYU Tandon
Deutsche Bank

10-May-2019

# Contents

Dynamic Programming (DP) is a collection of algorithms for finding optimal policies given a perfect model for the environment.

We start with the Bellman equation for the optimal state-value and action-value

$$v_\star (s) = \max_a \sum_{s',r} p\left(s', r|s, a\right) \left[r + \gamma v_\star \left(s'\right)\right] \tag{1}$$

$$q_\star (s, a) = \sum_{s',r} p\left(s', r|s, a\right) \left[r + \gamma \max_{a'} q_\star \left(s', a'\right)\right] \tag{2}$$

with states $s, s' \in \mathcal{S}$, actions $a, a' \in \mathcal{A}$, reward $r \in \mathbb{R}$, and discount $\gamma \leq 1$.
The state-value function is $v(s)$ and the action-value is $q(s, a)$.
The environment is described as an MDP via the probabilities $p(s', r|s, a)$.

Given state $s$, the future states $s'$ are called the successor states.

▶ Solving the Bellman equations means finding update rules that continually improve the approximations $v(s)$ or $q(s, a)$ and converge to $v_\star (s)$ or $q_\star (s, a)$ in the limit of infinite iterations.

The first task is policy evaluation or prediction.

▶ Given policy $\pi$, i.e. the function $\pi(a|s)$, compute the value function $v_\pi(s)$

From the definition of the value function, we need to solve

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[ r + \gamma v_\pi(s') \right] \tag{3}$$

This is a linear system of equations with unknowns $v_\pi(s)$. There are as many unknowns as states $s \in \mathcal{S}$. Brute force solution may be very inefficient.

**Iterative policy evaluation** is an algorithm that starts with an initial guess for the value function $v_0(s)$. Then, at each iteration $k$, it updates the estimate as

$$v_{k+1}(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[ r + \gamma v_k(s') \right] \tag{4}$$

The sequence $\{v_k\}$ converges to $v_\pi$ as $k \rightarrow \infty$.

▶ A two-array implementation uses two arrays for storing $v_k(s)$ and $v_{k+1}(s)$
▶ A one-array implementation or a **sweep** uses one array for $v_{k+1}(s)$ and overwrites older values with newer ones.

Iteration stops when no state update exceeds a small threshold value $\theta$.
The iterative policy algorithm with state sweep looks like this in pseudo-code

```
input: π
parameters: θ
initialize: V (s) for all s ∈ S except V (terminal) = 0
loop:
    Δ ← 0
    for each s ∈ S :
        v ← V (s)
        V (s) ← ∑_a π (a|s) ∑_{s',r} p (s', r|s, a) [r + γ V (s')]
        Δ ← max (Δ, |v − V (s) |)
until  Δ < θ
return V (s)
```

The same iterative algorithm can be applied to approximating the action-value.
We simply sweep over all possible pairs $(s, a)$ and apply the update

$$Q (s, a) \leftarrow \sum_a \pi (a|s) \sum_{s', r} p (s', r|s, a) \left[ r + \gamma Q (s', a') \right] \qquad (5)$$

The second task is to examine whether we can improve a given policy $\pi$.

The trick here is to follow policy $\pi$ for all states except for state $s$, where we switch to action $a = \pi'(s)$. This gives us a new value function

$$q_\pi(s, a) = \sum_{s', r} p(s', r | s, a) \left[ r + \gamma v_\pi(s') \right] \tag{6}$$

If it turns out that $q_\pi(s, a) \geq v_\pi(s)$ then we expect that it is always better to choose action $a$ once we are in state $s$

This is true because of the following **policy improvement theorem**

▶ If $\pi$ and $\pi'$ are deterministic policies and for all $s \in \mathcal{S}$

$$q_\pi(s, \pi'(s)) \geq v(s) \tag{7}$$

then $\pi'$ is as good or better than $\pi$, i.e for all $s \in \mathcal{S}$

$$v_{\pi'}(s) \geq v_\pi(s) \tag{8}$$

We have a procedure to check whether a policy improves if we change behavior for a single state $s$. Now we can get **greedy**.

For all states $s \in \mathcal{S}$ and all possible actions $a \in \mathcal{A}$ we select at state $s$ the action $a$ that has the highest $q_\pi(s, a)$. This generates a new policy $\pi'$ as follows

$$\pi'(s) = \arg\max_a \sum_{s',r} p(s', r | s, a) \left[ r + \gamma v_\pi(s') \right] \tag{9}$$

▶ The process of replacing a policy $\pi$ with a better policy $\pi'$ that is greedy with respect to $\pi$ is called **policy improvement**.

If $\pi'$ is as good but no better than $\pi$, then $\pi' = \pi$ and eq. (9) becomes the Bellman equation. We have found a solution!

Finally, we use sequence of alternating evaluation and improvement updates to obtain a monotonically improving policy and value function

$$\pi_0 \xrightarrow{\text{E}} v_{\pi_0} \xrightarrow{\text{I}} \pi_1 \xrightarrow{\text{E}} v_{\pi_1} \xrightarrow{\text{I}} \pi_2 \xrightarrow{\text{E}} \cdots \xrightarrow{\text{I}} \pi_* \xrightarrow{\text{E}} v_*.$$

Figure: The policy evaluation-improvement sequence. From [SB]

The policy iteration/evaluation algorithm looks like this in pseudo-code

```
parameters: θ
1. initialize: V (s), π (s) for all s ∈ S
2. policy evaluation
     loop:
          Δ ← 0
          for each s ∈ S :
           v ← V (s)
           V (s) ← ∑_a π (a|s) ∑_{s',r} p (s', r|s, a) [r + γV (s')]
           Δ ← max (Δ, |v − V (s)|)
     until Δ < θ
3. policy improvement
     stable ← true
     for each s ∈ S :
          a_old ← π (s)
          π (s) ← arg max_a ∑_{s',r} p (s', r|s, a) [r + γV (s')]
          if a_old ≠ π (s) then stable ← false
     if stable :
          return V (s), π (s)
     else :
          go to 2
```

In policy iteration, evaluation and improvement alternate in sequence. One affects the other through its output, but they do not "step on each other".

**Generalized Policy Iteration** (GPI) <u>allows evaluation and improvement to communicate</u> (or overwrite each other's results) while they both run. The final solution is returned once it becomes stable w.r.t. both evaluation and improvement.

- ▶ In GPI, evaluation and improvement both cooperate and compete
- ▶ Greedy search makes the value function inaccurate w.r.t. the new policy, but encourages exploration
- ▶ Making the value function consistent w.r.t. the current policy increases accuracy but decreases exploration

Most RL algorithms are some form of GPI.

Dynamic programming is the cleanest way of formulating RL problems. However it has two major drawbacks:

1. It requires full information about the environment
2. It is computationally inefficient (curse of dimensionality)

Below we examine ways to improve on DP algorithms.

Monte Carlo (MC) methods try to overcome the limitations of DP, by using simulated experience. Their advantages are

- ▶ No need to know the complete specification of the environment, i.e. the functional form of $p\left(s', r|s, a\right)$
- ▶ It is much easier to generate sample transitions from target distributions
- ▶ Learning occurs via sample averaging

In an MC method, we generate an episode and then for each state-action pair $(s, a)$ that we saw during the episode, we compute the average reward.

Similarly to DP, here also we need to solve for policy evaluation (prediction), policy improvement (control) and combine the two in a GPI algorithm.

- ▶ Since we do not assume knowledge of the environment, we have to use the action-value function $q\left(s, a\right)$.

$$\pi_0 \xrightarrow{\text{ E }} q_{\pi_0} \xrightarrow{\text{ I }} \pi_1 \xrightarrow{\text{ E }} q_{\pi_1} \xrightarrow{\text{ I }} \pi_2 \xrightarrow{\text{ E }} \cdots \xrightarrow{\text{ I }} \pi_* \xrightarrow{\text{ E }} q_*$$

Figure: The policy evaluation-improvement sequence. From [SB]

With an initialization that explores states and actions (exploring starts) the MC algorithm looks like this in pseudo-code

```
parameters: θ
initialize: π(s) ∈ 𝒜, Q(s, a) ∈ ℝ for all s ∈ 𝒮, a ∈ 𝒜
            Returns(s, a) ←empty_list
loop for every episode:
    choose S₀, A₀ randomly so that all pairs (S₀, A₀) are candidates
    generate episode  S₀, A₀, R₁, S₁, A₁, ... Rᴛ following π
    G ← 0
    loop for step t = T − 1, T − 2, ... 0 :
        G ← γG + Rₜ₊₁
        if (Sₜ, Aₜ) not in sequence S₀, A₀, R₁, S₁, A₁, ... Sₜ₋₁, Rₜ₋₁:
            append G to Returns(Sₜ, Aₜ)
            Q(Sₜ, Aₜ) ← average(Returns(Sₜ, Aₜ))
            π(Sₜ) ← arg maxₐ Q(Sₜ, Aₜ)
```

▶ A method is called **on-policy** if it uses the same policy for prediction and control

▶ A method is called **off-policy** if it uses one policy for prediction and a different policy for control.
The policy used for prediction is also called *behavior*. It is used for exploring state-action pairs.
The policy used for control is also called *target*. This is the policy that is being improved.

The assumption of exploring starts is too restrictive.
A better method is on-policy MC with $\epsilon$-greedy policy.
At each new episode, we do not randomly re-initialize, but follow a greedy policy to explore states.
This is similar to the multi-arm bandit problem of the previous lecture.

The on-policy $\epsilon$-greedy MC algorithm looks like this in pseudo-code

```
parameter: small ε > 0
initialize: π ← some ε-greedy policy
             Q (s, a) ∈ ℝ for all s ∈ 𝒮, a ∈ 𝒜
             Returns (s, a) ← empty_list
loop for every episode:
    generate episode S₀, A₀, R₁, S₁, A₁, . . . Rₜ following π
    G ← 0
    loop for every step t = T − 1, T − 2, . . . 0 :
        G ← γG + Rₜ₊₁
        if (Sₜ, Aₜ) not in sequence S₀, A₀, R₁, S₁, A₁, . . . Sₜ₋₁, Rₜ₋₁:
            append G to Returns (Sₜ, Aₜ)
            Q (Sₜ, Aₜ) ← average (Returns (Sₜ, Aₜ))
            A* ← arg maxₐ Q (Sₜ, Aₜ)
            for all a ∈ 𝒜(Sₜ):
```

$$\pi\left(a|S_t\right) \leftarrow \begin{cases} 1 - \epsilon + \epsilon/\left|\mathcal{A}\left(S_t\right)\right| & \text{if} \quad a = A^\star \\ \epsilon/\left|\mathcal{A}\left(S_t\right)\right| & \text{if} \quad a \neq A^\star \end{cases}$$

In the last step the algorithm updates the policy using an $\epsilon$-greedy method.
The notation $\left|\mathcal{A}\left(S_t\right)\right|$ means the number of actions taken from state $S_t$.

Temporal Learning (TD) is a pivotal method in Reinforcement Learning.

▶ TD combines the concepts of DP and MC

▶ Like DP, TD updates the current action-value estimates based on other concurrent action-value estimates, i.e. it **bootstraps**.
Therefore, it does not wait until the end of an episode to do the updates.

▶ Like MC, TD learns from experience without requiring a model for the environment.

Let's compare the prediction methods in MC and TD within the simple context of updating the value function

In MC, we run to the end time $T$ to obtain the payoff $G_t$, and then we update

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)] \tag{10}$$

In TD, just after time $t$ we update

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_t) - V(S_t)] \tag{11}$$

This is called TD(0) or one-step TD. Its generalization is $n$-step TD or TD($\lambda$).

**Sarsa** is on-policy TD control.
It continually estimates the action value $q_\pi$ for the behavior policy $\pi$ and at the same time it changes $\pi$ towards a target that is greedy with respect to $q_\pi$.
The algorithm in pseudo-code looks like this

```
parameters: step size α ∈ (0,1], small ε > 0
initialize: Q(s,a) ∈ ℝ for all s ∈ 𝒮,a ∈
𝒜 with Q(terminal,.) = 0
loop for every episode:
    initialize S
    choose A from S using policy derived from Q (ε-greedy)
    loop for each step:
        take A from S, observe R, S′
        choose A′ from S′ using policy derived from Q (ε-greedy)
        Q(St,At) ← Q(St,At) + α[R + γQ(S′,A′) − Q(St,At)]
        S ← S′; A ← A′
    until terminal S
```

Sarsa converges almost surely to the optimal policy and action-value function in the limit when all state-action pairs are visited an infinite number of times.

**Q-learning** is off-policy TD control.
This algorithm was a major breakthrough in RL (Watkins 1989).

▶ The main idea in Q-learning is to approximate directly the optimal action-value $q_\star$ independently of the policy being followed during learning. This is achieved via the update

$$Q\left(S_t, A_t\right) \leftarrow Q\left(S_t, A_t\right) + \alpha\left[R_{t+1} + \gamma \max_a Q\left(S_{t+1}, a\right) - Q\left(S_t, A_t\right)\right] \quad (12)$$

▶ The policy that we follow during learning (behavior) affects the results to the extent that it determines which state-action pairs are visited and therefore updated.

▶ If all the state-action pairs are visited and updated, Q-learning converges to the optimal target $q_\star$ independent of the behavior policy.

The Q-learning algorithm is structurally quite similar to Sarsa.
In pseudo-code looks like this

```
parameters: step size α ∈ (0, 1], small ε > 0
initialize: Q (s, a) ∈ ℝ for all s ∈ S, a ∈ A with Q (terminal, .) = 0
loop for every episode:
    initialize S
    loop for each step:
        choose A from S using policy derived from Q (ε-greedy)
        take A from S, observe R, S'
        Q (S_t, A_t) ← Q (S_t, A_t) + α [R + γ max_a Q (S', a) − Q (S_t, A_t)]
        S ← S'
    until terminal S
```

In TD(0) methods within one time step we update (bootstrap) and move to the
next state. In TD($n$) methods, the bootstrapping happens every $n$ time steps.

Sarsa and Q-learning can be generalized to TD($n$). In fact, this generates a
whole spectrum of methods, with TD(0) on one end and MC on the other.

## Further Reading

1. For a deeper discussion of the topics in this lecture read chapters 5 and 6 in [SB]
2. Become familiar with OpenAI Gym [GYM]. It provides ready-made components for simulating environments and coding RL algorithms.

**References**

[SB] Sutton R.S. and Barto A.G. "Reinforcement Learning An Introduction" 2nd edition. MIT, 2018

[GYM] Open AI Gym. A toolkit for Developing and Comparing RL Algorithms. gym.openai.com