

Further Considerations on Dynamic Programming

V. Leclère (ENPC)

October 16, 2014

Contents

- 1 Dynamic Programming
- 2 Curses of Dimensionality
- 3 Infinite Horizon

Contents

- 1 Dynamic Programming
- 2 Curses of Dimensionality
- 3 Infinite Horizon

Stochastic Controlled Dynamic System

A stochastic controlled dynamic system is defined by its *dynamic*

$$\mathbf{X}_{t+1} = f_t(\mathbf{X}_t, \mathbf{U}_t, \mathbf{W}_{t+1})$$

and initial state

$$\mathbf{X}_0 = x_0$$

The variables

- \mathbf{X}_t is the *state* of the system,
- \mathbf{U}_t is the *control* applied to the system at time t ,
- \mathbf{W}_t is an exogeneous noise.

Optimization Problem

We want to solve the following optimization problem

$$\min \quad \mathbb{E} \left[\sum_{t=0}^{T-1} L_t(\mathbf{x}_t, \mathbf{u}_t, \mathbf{w}_{t+1}) + K(\mathbf{x}_T) \right] \quad (1)$$

$$\text{s.t.} \quad \mathbf{x}_{t+1} = f_t(\mathbf{x}_t, \mathbf{u}_t, \mathbf{w}_{t+1}), \quad \mathbf{x}_0 = x_0 \quad (2)$$

$$\mathbf{u}_t \in U_t(\mathbf{x}_t) \quad (3)$$

$$\mathbf{u}_t \preceq \sigma(\mathbf{w}_0, \dots, \mathbf{w}_t) \quad (4)$$

Dynamic Programming Principle

Assume that the noises \mathbf{W}_t are **independent**.

Then, there exists an optimal solution of the form $\mathbf{U}_t = \pi_t(\mathbf{X}_t)$, given by

$$\pi_t(x) = \arg \min_{u \in U_t(x)} \mathbb{E} \left[\underbrace{L_t(x, u, \mathbf{W}_{t+1})}_{\text{current cost}} + \underbrace{V_{t+1} \circ f_t(x, u, \mathbf{W}_{t+1})}_{\text{future costs}} \right],$$

where

$$\begin{cases} V_T(x) &= K(x) \\ V_t(x) &= \min_{u \in U_t(x)} \mathbb{E} \left[L_t(x, u, \mathbf{W}_{t+1}) + V_{t+1} \circ f_t(x, u, \mathbf{W}_{t+1}) \right] \end{cases}$$

Interpretation of Bellman Value

The Bellman's value function $V_{t_0}(x)$ can be interpreted as the value of the problem starting at time t_0 from the state x . More precisely we have

$$V_{t_0}(x) = \min \mathbb{E} \left[\sum_{t=t_0}^{T-1} L_t(\mathbf{X}_t, \mathbf{U}_t, \mathbf{W}_{t+1}) + K(\mathbf{X}_T) \right] \quad (5)$$

$$s.t. \quad \mathbf{X}_{t+1} = f_t(\mathbf{X}_t, \mathbf{U}_t, \mathbf{W}_{t+1}), \quad \mathbf{X}_{t_0} = x \quad (6)$$

$$\mathbf{U}_t \in U_t(\mathbf{X}_t) \quad (7)$$

$$\mathbf{U}_t \preceq \sigma(\mathbf{W}_0, \dots, \mathbf{W}_t) \quad (8)$$

Contents

- 1 Dynamic Programming
- 2 Curses of Dimensionality
- 3 Infinite Horizon

Dynamic Programming Algorithm

Data: Problem parameters

Result: optimal control and value;

$V_T \equiv K$;

for $t : T \rightarrow 0$ **do**

for $x \in \mathbb{X}_t$ **do**

$\underline{v} = -\infty$;

for $u \in U_t(x)$ **do**

$\underline{v} = \min \left\{ \underline{v}, \mathbb{E} \left[L_t(x, u, \mathbf{W}_{t+1}) + V_{t+1} \circ f_t(x, u, \mathbf{W}_{t+1}) \right] \right\}$;

$V_t(x) = \underline{v}$;

Algorithm 1: Dynamic Programming Algorithm (discrete case)

Number of flops: $O(T \times |\mathbb{X}_t| \times |U_t| \times |W_t|)$.

3 curses of dimensionality

- 1 **State.** If we consider 3 independent states each taking 10 values, then $|\mathbb{X}_t| = 1000$. In practice DP is not applicable for states of dimension more than 5.
- 2 **Decision.** The decision are often vector decisions, that is a number of independent decision, hence leading to huge $|U_t(x)|$.
- 3 **Expectation.** In practice random information came from large data set. Without a proper statistical treatment computing an expectation is costly. Monte-Carlo approach are costly too, and unprecise.

Contents

- 1 Dynamic Programming
- 2 Curses of Dimensionality
- 3 Infinite Horizon**

Introducing the Bellman operators

We define the Bellman operator associated to our optimisation problem

$$T_t(J) : x \mapsto \min_{u \in U_t(x)} \mathbb{E} \left[L_t(x, u, \mathbf{W}_{t+1}) + J \circ f_t(x, u, \mathbf{W}_{t+1}) \right].$$

The Dynamic Programming equation can then be written

$$\begin{cases} V_T = K \\ V_t = T_t(V_{t+1}) \end{cases}$$

We also construct the *policy-dependent Bellman operator*

$$T_t^\pi(J) : x \mapsto \mathbb{E} \left[L_t(x, \pi(x), \mathbf{W}_{t+1}) + J \circ f_t(x, \pi(x), \mathbf{W}_{t+1}) \right].$$

Discounted fixed cost case

We now consider the following specific case problem, where $(\mathbf{W}_t)_{t \in \mathbb{N}}$ is i.i.d.

$$\min \quad \mathbb{E} \left[\sum_{t=0}^T \alpha^t L(\mathbf{X}_t, \mathbf{U}_t, \mathbf{W}_{t+1}) \right] \quad (9)$$

$$s.t. \quad \mathbf{X}_{t+1} = f(\mathbf{X}_t, \mathbf{U}_t, \mathbf{W}_{t+1}), \quad \mathbf{X}_0 = x_0 \quad (10)$$

$$\mathbf{U}_t \in U(\mathbf{X}_t) \quad (11)$$

$$\mathbf{U}_t \preceq \sigma(\mathbf{W}_0, \dots, \mathbf{W}_t) \quad (12)$$

where $\alpha \in]0, 1]$. Note that the constraint and cost structure doesnot depend on t .

The Bellman operator is given by

$$T(J) : x \mapsto \min_{u \in U(x)} \mathbb{E} \left[L(x, u, \mathbf{W}_{t+1}) + \alpha J \circ f(x, u, \mathbf{W}_{t+1}) \right]$$

Infinite horizon problems

There is different ways of considering the above problem in an “infinite horizon” setting.

- 1 **Discounted case.** This is the case where $\alpha < 1$. It is especially easy to treat if the cost L is bounded.
- 2 **Stochastic shortest path.** In this case $\alpha = 1$ but there is a “cemetery state” such that once reached the system remains there with null cost. Moreover, we assume that the system always reach the cemetery state in a finite time.
- 3 **Average cost per stage problems.** This approach is mainly taken if the infinite time cost isn't finite (for example $\alpha = 1$ and $L > 0$). We consider

$$\lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\sum_{t=0}^{T-1} L(\mathbf{x}_t, \mathbf{u}_t, \mathbf{w}_{t+1}) \right].$$

An overview of typical infinite horizon results

Here are the main results that can be shown in infinite horizon problems (under the right set of assumptions)

- 1 the sequence of value function $V_{n+1} = T(V_n)$, converges toward the value function of the infinite horizon problem:
 $\lim_{n \rightarrow \infty} V_n = V^\#$.
- 2 The optimal value of the infinite horizon problem is a fixed point of the Bellman operator: $V^\# = T(V^\#)$.
- 3 If π is such that $V^\# = T^\pi V^\#$ then the stationary policy π is optimal.

Value iteration algorithm

```

Data: Initial value  $V^{(0)}$ 
Result: optimal policy and value;
repeat
  |   for  $x \in \mathbb{X}$  do
  |   |
  |   |            $V^{(k+1)}(x) = T(V^{(k)})(x)$ 
  |   |
  |   └──
until  $\|V^{(k+1)} - V^{(k)}\|_\infty < \varepsilon$ ;

```

Algorithm 2: Value iteration algorithm

- Each step takes $O(|\mathbb{X}| \times |\mathbb{U}| \times |\Omega|)$ flops.
- The error $|V_n(x) - V^\#(x)|$ is bounded by $C\alpha^n$.

Policy iteration algorithm

Data: Initial policy $\pi^{(0)}$

Result: optimal policy and value;

repeat

policy evaluation step: solve $V = T^{\pi^{(k)}}(V)$ which gives $V^{(k)}$;
 policy improvement step :

for $x \in \mathbb{X}$ **do**

$$\pi^{(k+1)}(x) = \arg \min_{u \in U(x)} \mathbb{E} \left[L(x, u, \mathbf{W}_{t+1}) + \alpha V^{(k)} \circ f(x, u, \mathbf{W}_{t+1}) \right]$$

until $V^{(k)} = T(V^{(k)})$;

Algorithm 3: Policy iteration algorithm

The policy iteration algorithm terminate in a finite number of step.