

Stochastic Dynamic Programming

V. Leclère (ENPC)

October 8, 2015

Contents

- 1 Dynamic Programming
- 2 Curses of Dimensionality
- 3 Linear-Quadratic Setting
- 4 Infinite Horizon

Contents

- 1 Dynamic Programming
- 2 Curses of Dimensionality
- 3 Linear-Quadratic Setting
- 4 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

$$\mathbf{x}_0 = \mathbf{x}_0$$

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 exogenous noise.

Examples

- Stock of water in a dam:
 - x_t is the amount of water in the dam at time t ,
 - u_t is the amount of water turbinéd at time t ,
 - w_t is the inflow of water at time t .
- Boat in the ocean:
 - x_t is the position of the boat at time t ,
 - u_t is the direction and speed chosen at time t ,
 - w_t is the wind and current at time t .
- Subway network:
 - x_t is the position and speed of each train at time t ,
 - u_t is the acceleration chosen at time t ,
 - w_t is the delay due to passengers and incident on the network at time t .

Optimization Problem

We want to solve the following optimization problem

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

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

$$\mathbf{u}_t \in U_t(\mathbf{x}_t) \quad (1c)$$

$$\sigma(\mathbf{u}_t) \subset \sigma(\mathbf{w}_0, \dots, \mathbf{w}_t) \quad (1d)$$

Dynamic Programming Principle

Assume that the noises \mathbf{w}_t are **independent** and **exogeneous**.

Then, there exists an optimal solution, called a **strategy**, of the form $\mathbf{u}_t = \pi_t(\mathbf{x}_t)$, given by

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

where (Dynamic Programming Equation)

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

Dynamic Programming Principle

Assume that the noises \mathbf{w}_t are **independent** and **exogeneous**.

Then, there exists an optimal solution, called a **strategy**, of the form $\mathbf{u}_t = \pi_t(\mathbf{x}_t)$, given by

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

where (Dynamic Programming Equation)

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

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

$$V_{t_0}(x) = \min \mathbb{E} \left[\sum_{t=t_0}^{T-1} L_t(x_t, u_t, w_{t+1}) + K(x_T) \right] \quad (2a)$$

$$s.t. \quad x_{t+1} = f_t(x_t, u_t, w_{t+1}), \quad x_{t_0} = x \quad (2b)$$

$$u_t \in U_t(x_t) \quad (2c)$$

$$\sigma(u_t) \subset \sigma(w_0, \dots, w_t) \quad (2d)$$

Information structures in the multistage setting

Open-Loop Every decision $(u_t)_{t \in \llbracket 0, T-1 \rrbracket}$ is taken before any noises $(\xi_t)_{t \in \llbracket 0, T-1 \rrbracket}$ is known. We decide a planning, and stick to it.

Decision Hazard Decision u_t is taken knowing all past noises ξ_0, \dots, ξ_t , but not knowing ξ_{t+1}, \dots, ξ_T .

Hazard Decision Decision u_t is taken knowing all past noises ξ_0, \dots, ξ_t , and the next noise ξ_{t+1} but not knowing ξ_{t+2}, \dots, ξ_T .

Anticipative Every decision $(u_t)_{t \in \llbracket 0, T-1 \rrbracket}$ is taken knowing the whole scenario $(\xi_t)_{t \in \llbracket 0, T-1 \rrbracket}$. There is one deterministic optimization problem by scenario.

With the same objective function this gives better and better value as the solution use more and more information.

Information structures: comments

Open-Loop This case can happen in practice (e.g. fixed planning). There are specific methods to solve this type of optimization problem (e.g. stochastic gradient methods).

Decision Hazard The decision \mathbf{u}_t is taken at the beginning of period $[t, t + 1[$. The decision is always implementable, and might be conservative as it does not leverage any prediction over the noise in $[t, t + 1[$.

Hazard Decision The decision \mathbf{u}_t is taken at the end of period $[t, t + 1[$. The modelization is optimistic as it assumes perfect knowledge that might not be available in practice.

Anticipative This problem is never realistic. However it is a lower bound of the real problem that can be estimated through Monte-Carlo and deterministic optimization.

Independence of noise

- The Dynamic Programming equation requires only the **time-independence of noises**.
- This can be relaxed if we consider an extended state.
- Consider a dynamic system driven by an equation

$$\mathbf{y}_{t+1} = f_t(\mathbf{X}_t, \mathbf{u}_t, \boldsymbol{\varepsilon}_{t+1})$$

where the random noise $\boldsymbol{\varepsilon}_t$ is an AR1 process :

$$\boldsymbol{\varepsilon}_t = \alpha_t \boldsymbol{\varepsilon}_{t-1} + \beta_t + \mathbf{w}_t,$$

$\{\mathbf{w}_t\}_{t \in \mathbb{Z}}$ being independent.

- Then \mathbf{y}_t is called the **physical state** of the system and DP can be used with the **information state** $\mathbf{X}_t = (\mathbf{y}_t, \boldsymbol{\varepsilon}_{t-1})$.
- Generically speaking, if the noise \mathbf{w}_t is exogeneous (not affected by decisions \mathbf{u}_t), then we can always apply Dynamic Programming with the state

$$(\mathbf{x}_t, \mathbf{w}_1, \dots, \mathbf{w}_t)$$

Contents

- 1 Dynamic Programming
- 2 Curses of Dimensionality
- 3 Linear-Quadratic Setting
- 4 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
     $V_t(x) = \infty$ ;
    for  $u \in U_t(x)$  do
       $v_u = \mathbb{E} \left[ L_t(x, u, \mathbf{w}_{t+1}) + V_{t+1} \circ f_t(x, u, \mathbf{w}_{t+1}) \right]$ ;
      if  $v_u < \underline{v}$  then
         $V_t(x) = v_u$  ;
         $\pi_t(x) = u$  ;

```

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 $|\mathbb{X}_t| = 10^3 = 1000$. 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.

Numerical considerations

- The DP equation holds in (almost) any case.
- The algorithm shown before compute a **look-up table** of control for every possible state *offline*. It is impossible to do if the state is (partly) continuous.
- Alternatively, we can focus on computing *offline* an **approximation of the value function** V_t and derive the optimal control *online* by solving a one-step problem, solved only at the current state :

$$\pi_t(x) \in \arg \min_{u \in U_t(x)} \mathbb{E} \left[L_t(x, u, \mathbf{w}_{t+1}) + V_{t+1} \circ f_t(x, u, \mathbf{w}_{t+1}) \right]$$

- The field of Approximate DP gives methods for computing those approximate value function (decomposed on a base of functions).
- The simpler one consisting in discretizing the state, and then interpolating the value function.

DP on a Markov Chain

- Sometimes it is easier to represent a problem as a controlled Markov Chain
- Dynamic Programming equation can be computed directly, without expliciting the control.
- Let's work out an example...

Contents

- 1 Dynamic Programming
- 2 Curses of Dimensionality
- 3 Linear-Quadratic Setting**
- 4 Infinite Horizon

The Linear-Quadratic setting

We assume a linear dynamic

$$\mathbf{x}_{t+1} = A_t \mathbf{x}_t + B_t \mathbf{u}_t + \mathbf{W}_{t+1}$$

associated with a quadratic cost

$$\mathbb{E} \left[\sum_{t=0}^{T-1} \left(\mathbf{x}'_t Q_t \mathbf{x}_t + \mathbf{u}'_t R_t \mathbf{u}_t \right) \right] + \mathbf{x}'_T Q_T \mathbf{x}_T.$$

A few more assumptions

- \mathbf{x}_t is of dimension n , \mathbf{u}_t of dimension m .
- Q_t is a symmetric semidefinite positive matrix, and R_t symmetric definite positive.
- \mathbf{w}_t is a centered (i.e. of mean 0) independent, exogeneous noise (i.e their law does not depend of the state or control), with finite second order moment.
- The controls are non-anticipative.

Solving the LQ case

The DP equation read

$$\begin{cases} V_T(x) &= x'Q_Tx \\ V_t(x) &= \min_u \mathbb{E} \left[x'Q_t x + u'R_t u + V_{t+1}(A_t x + B_t u + \mathbf{w}_{t+1}) \right] \end{cases}$$

Leading to

$$V_t(x_t) = x_t' K_t x_t + \sum_{\tau=t}^{T-1} \mathbb{E} [w_{\tau+1}' K_{\tau+1} w_{\tau+1}]$$

and

$$u_t^\# = \pi_t^\#(x_t) = L_t x_t .$$

Solving the LQ case

|

The DP equation read

$$\begin{cases} V_T(x) &= x'Q_Tx \\ V_t(x) &= \min_u \mathbb{E} \left[x'Q_t x + u'R_t u + V_{t+1}(A_t x + B_t u + \mathbf{w}_{t+1}) \right] \end{cases}$$

Leading to

$$V_t(x_t) = x_t' K_t x_t + \sum_{\tau=t}^{T-1} \mathbb{E} [\mathbf{w}'_{\tau+1} K_{\tau+1} \mathbf{w}_{\tau+1}]$$

and

$$\mathbf{u}_t^\# = \pi_t^\#(\mathbf{x}_t) = L_t \mathbf{x}_t .$$

Solving the LQ case



We have

$$V_t(x_t) = x_t' K_t x_t + \sum_{\tau=t}^{T-1} \mathbb{E} [\mathbf{w}'_{\tau+1} K_{\tau+1} \mathbf{w}_{\tau+1}]$$

and

$$\mathbf{u}_t^\# = \pi_t^\#(\mathbf{x}_t) = L_t \mathbf{x}_t .$$

Where

$$L_t = -(B_t' K_{t+1} B_t + R_t)^{-1} B_t' K_{t+1} A_t ,$$

and

$$\begin{cases} K_T &= Q_T \\ K_t &= A_t' \left(K_{t+1} - K_{t+1} B_t (B_t' K_{t+1} B_t + R_t)^{-1} B_t' K_{t+1} \right) A_t + Q_t \end{cases}$$

Contents

- 1 Dynamic Programming
- 2 Curses of Dimensionality
- 3 Linear-Quadratic Setting
- 4 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.

$$\min \quad \mathbb{E} \left[\sum_{t=0}^T \alpha^t L(\mathbf{x}_t, \mathbf{u}_t, \mathbf{w}_{t+1}) \right] \quad (3)$$

$$\text{s.t.} \quad \mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t, \mathbf{w}_{t+1}), \quad \mathbf{x}_0 = x_0 \quad (4)$$

$$\mathbf{u}_t \in U(\mathbf{x}_t) \quad (5)$$

$$\sigma(\mathbf{u}_t) \subset \sigma(\mathbf{w}_0, \dots, \mathbf{w}_t) \quad (6)$$

where $\alpha \in]0, 1]$. Note that the constraint and cost structure does not 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.