

ORF 418 : OPTIMAL LEARNING

LECTURE 6 : September 21 . 2025

1) Problem Formulation

2) Dynamic Programming



# I. PROBLEM FORMULATION

Stochastic dynamics is given by

$$x_{k+1} = Ax_k + Bu_k + w_{k+1} \quad k=0,1,\dots$$

where  $A, B$  are as before, and  $w_1, w_2, \dots$

is a sequence of **independently, identically distributed (i.i.d.)** sequence of  $\mathbb{R}^d$  valued

Gaussian random variables with mean zero and (dxd) covariance matrix  $Q$ .

**Infinite horizon** cost functional is

given by

$$J_\infty(x, u) := \mathbb{E} \left[ \sum_{k=0}^{\infty} \rho^k (x_k^T M x_k + u_k^T N u_k) \right].$$

Important restriction the control is that they must be **adapted** to the information flow. Namely,  $u_k$  can use information upto time  $k$  or equivalently  $w_1, \dots, w_k$ .

We always need  $\rho < 1$  to have a finite value function:

$$v_{\infty}(x) := \inf_{u \in \mathcal{U}} J(x, u).$$

**Example.** Consider the problem with  $d=l=1$ ,

$$x_{k+1} = 4x_k + u_k + w_{k+1}$$

where  $w_1, w_2, \dots$  have zero mean and variance  $\sigma^2$ . If we use a linear control

$$u_k = -f x_k, \text{ then}$$

$$x_{k+1} = (4-f)x_k + w_{k+1}$$

$$\Rightarrow x_n = (4-f)^n x_0 + y_n$$

and

$$y_n = \sum_{k=1}^n (4-f)^{n-k} w_k, \quad n=1, 2, \dots$$

Then,

$$\begin{aligned} \text{var}(x_n) &= \text{var}(y_n) \\ &= \sum_{k=1}^n (4-f)^{2(n-k)} \text{var}(w_k) \\ &= \sigma^2 \sum_{k=1}^n (4-f)^{2(n-k)} \geq \sigma^2. \end{aligned}$$

So no matter how we choose the gain constant  $f$ ,  $x_k$  will never converge to zero.

## II. DYNAMIC PROGRAMMING.

Proceeding exactly as in the deterministic case, we obtain the following equation for the value function  $v_\infty$ :

$$v_\infty(x) = x^T M x + \inf_u \left\{ u^T N u + \rho \mathbb{E} [v_\infty(Ax + Bu + \omega_1)] \right\}$$

We postulate that

$$v_\infty(x) = x^T V x + C$$

for some positive definite matrix, symmetric matrix  $V$  and a constant  $C$ . By substituting this form into the dynamic programming equation we obtain the following result:

THEOREM 2.4.1. Let  $V$  be the Riccati equation with matrices  $(A, B, M, N)$  and discount factor  $\rho$ .

Then,

$$v_\infty(x) = x^T V x + \frac{\rho}{1-\rho} \text{trace}(VQ)$$

and  $u^*(x) = -Fx$  (with the same gain matrix  $F$ ) is the optimal feedback control.

Proof given in the notes is a straightforward but tedious calculation. Important to note that the form of the optimal control does not change!

**Example.** Going back to the previous example, we have  $d=l=1$ ,  $A=4$ ,  $M=N=1$ , and  $\rho=\frac{1}{2}$ . Then, the dynamic programming equation is,

$$v_{\infty}(x) = x^2 + \min_u \left\{ u^2 + \frac{1}{2} \mathbb{E} [v_{\infty}(4x+u+w, 1)] \right\}.$$

We know that  $v_{\infty}(x) = Vx^2 + C$ . Hence

$$\begin{aligned} \mathbb{E} [v_{\infty}(4x+u+w, 1)] &= \mathbb{E} [V(4x+u+w)^2] \\ &= V \left[ (4x+u)^2 + 2(4x+u) \underbrace{\mathbb{E}[w, 1]}_{=0} + \underbrace{\mathbb{E}(w^2)}_{=1} \right] \\ &= V [(4x+u)^2 + 1]. \end{aligned}$$

Then,

$$Vx^2 + C = x^2 + \min_u \left\{ u^2 + \frac{1}{2} V (4x+u)^2 \right\} + 1 \cdot \frac{1}{2}$$

By calculus  $u^*(x) = -\frac{4V}{2+V}$ , exactly as in the deterministic case. Hence,

$$C = \sigma^2/2$$

and

$$\begin{aligned} \sqrt{V} &= 1 + \left( \frac{4V}{2+\sqrt{V}} \right)^2 + \frac{1}{2} \left( 4 - \frac{4V}{2+\sqrt{V}} \right)^2 \sqrt{V} \\ &= 1 + \frac{16V^2 + 32V}{(2+\sqrt{V})^2} = 1 + \frac{16V}{2+\sqrt{V}}. \end{aligned}$$

This is the same Riccati equation as in the deterministic case.