

# Lecture #16

Stochastic (Discrete time) OCP:

$$\text{minimize}_{\gamma(\cdot) \in \Gamma} \mathbb{E} \left[ c_T(\underline{x}(T)) + \sum_{k=0}^{T-1} c_k(\underline{x}_k, \underline{u}_k) \right]$$

$$= (\underline{x}_0, \underline{u}_0, \underline{x}_1, \underline{u}_1, \dots, \underline{x}_{T-1}, \underline{u}_{T-1}, \underline{x}_T)$$

s.t.

(Discrete-time noisy case, spl. case of MDP)

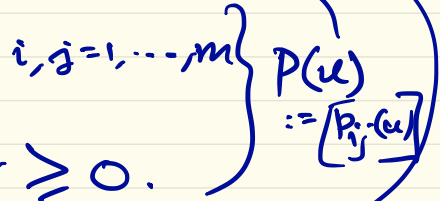
$$\underline{x}_{k+1} = \underline{f}_k(\underline{x}_k, \underline{u}_k, \underline{w}_k)$$

Controlled Markov chain  $\left. \begin{array}{l} \underline{x}_k \in \mathcal{X}, \\ \underline{u}_k \in \mathcal{U} \end{array} \right\}$

(or, the following MDP):

$$\underline{x}(k) \sim \text{Markov}([P_{ij}^k]), \quad i, j = 1, \dots, m$$

$$\mathbb{P}(\underline{x}(k+1) = s_j | \underline{x}(k) = s_i) = p_{ij}^k \geq 0.$$



(or, deterministic dynamics)

$$\underline{x}_{k+1} = \underline{f}_k(\underline{x}_k, \underline{u}_k)$$

$$\underline{w}_k \equiv 0 \quad \forall k=0, \dots, T-1$$

Given the OCP problem in any of the forms given in the prev page

Define:

$$V_k(\underline{x}) := \inf_{(\underline{\gamma}_k, \underline{\gamma}_{k+1}, \dots, \underline{\gamma}_{T-1})} \left[ \left\{ c_T(\underline{x}(T)) + \sum_{s=k}^{T-1} c_s(\underline{x}_s, \underline{u}_s) \right\} \right]_{\underline{x}_k = \underline{x}}$$

We call the above  $f_k^*$  as "value  $f_k^*$ "

Interpretation: "Cost-to-go" from state  $\underline{x}$  @ time  $k$   
(i.e.) optimal value of your action that results from applying optimal policy from  $\underline{x}_k = \underline{x}$ .

2 Main results in DP:

**Result #1** Optimal policy  $\gamma^*$  is a deterministic Markov policy, among the class of all history-dependent randomized policies " $\Gamma$ ".

**Result #2** We derive a recursion on  $V_k(\underline{x})$  [called DP eqn]

DP equations : (derived in the Chapter I emailed)

For  
Discrete-time noisy case:

$$V_T(\underline{x}) = C_T(\underline{x})$$

$$V_k(\underline{x}) = \inf_{u(\cdot) \in \mathcal{U}} \left\{ C_k(\underline{x}, \underline{u}) + \mathbb{E}_{\omega} [V_{k+1}(f_k(\underline{x}, \underline{u}, \omega))] \right\}$$

If  $C_k$  also depends on  $w_k$ , then we need to put  $\mathbb{E}[\cdot]$  outside the curly braces in RHS  $k = T-1, T-2, \dots, 0$

Markov case :  $V_T(\underline{x}) = C_T(\underline{x})$

$$V_k(\underline{x}) = \inf_{u(\cdot) \in \mathcal{U}} \left\{ C_k(\underline{x}, \underline{u}) + \sum_{j \in \mathcal{X}} p_{x,j}(u) V_{k+1}(j, \underline{x}) \right\}$$

$k = T-1, T-2, \dots, 0$

where  $[p_{ij}(u)]$  is the controlled Markov chain.

Deterministic case:

$\underline{w}_k \sim$  sample path dependency vanishes  
so the  $E_{w_k}[\cdot]$  can be dropped.

$$V_T(\underline{x}) = C_T(\underline{x})$$

$\therefore$  DP eq<sup>n</sup> becomes:

$$V_k(\underline{x}) = \inf_{u \in \mathcal{U}} \left\{ C_k(\underline{x}, \underline{u}) + V_{k+1}(f_k(\underline{x}, \underline{u})) \right\}$$

---

We call  $V_k(\underline{x})$  as "value function"

---

Example (Deterministic DP)

An investor receives annual salary/income \$  $x_k$  in the year  $k$ . He consumes \$  $u_k$ , and adds the rest  $(x_k - u_k)$  to his capital,

$0 \leq u_k \leq x_k$ . The capital is invested @ interest rate  $\theta \times 100\%$ .

∴ Dynamics of income:

(i.e.) his income in the  $(k+1)^{th}$  yr:

$$x_{k+1} = f_k(x_k, u_k) = x_k + \theta(x_k - u_k)$$

Objective: to maximize total consumption over lifetime  $T$  years.

(i.e.) deterministic OCP:

maximize  $\sum_{k=0}^{T-1} u_k$

$u_k = \gamma(x_0, u_0, x_1, u_1, \dots, x_{k-1}, u_{k-1}, x_k)$

History  
History dependent policies.

Question: what is the optimal consumption policy?

Apply DP:  $C_k(x, u) = u_k$  } Time invariant dep  
 $C_T(x) = 0$  } Time homogeneous DP

Define time-to-go: (countdown clock)  
 $n = T - k$ , i.e.,  $n = 0, 1, \dots, T$

Here,  $(x, u)$  are generic values for  $(x_k, u_k)$

$\therefore$  DP eq<sup>n</sup> in forward (countdown) time:

$$W_n(x) = \sup_{u(\cdot) \in \mathcal{U}} \left\{ c_n(x, u) + W_{n-1}(x + \theta(x-u)) \right\}$$

$W(\cdot)$  is the count-down value fn.

from the  
RHS  
of dynamics

$$\text{s.t. } W_0(x) = 0$$

(nothing more to consume),  $n = 1, 2, \dots, T$

How to solve: (Let's do first few iterations of DP eq<sup>n</sup>)

$$W_1(x) = \sup_{0 \leq u \leq x} \left\{ u + W_0(x + \theta(x-u)) \right\}$$

$$= \sup_{0 \leq u \leq x} \{u + 0\} = x$$



$$W_2(x) = \sup_{0 \leq u \leq x} \left\{ u + W_1(x + \theta(x-u)) \right\}$$

$$= \sup_{0 \leq u \leq x} \left\{ u + x + \theta(x-u) \right\}$$

linear in  $u$

( $\Rightarrow$  maximum is achieved either @  $u=0$ , or @  $u=x$ )

This is like Bang-bang

$$= \max \left\{ (1+\theta)x, 2x \right\}$$

$$= \underbrace{\left( \max \{ 1+\theta, 2 \} \right)}_{p_2 \text{ (say)}} x = p_2 x$$

This motivates the guess:

$$W_{n-1}(x) = \rho_{n-1} x \quad \text{for some (to-be-determined) constant } \rho_{n-1}$$

Verification of guess:

$$W_n(x) = \max_{0 \leq u \leq x} \{u + \rho_{n-1}(x + \theta(x-u))\}$$

$$= \left( \max \{ (1+\theta)\rho_{n-1}, 1 + \rho_{n-1} \} \right) x$$
$$= \rho_n x$$

$\therefore$  Guess is correct (structurally) &  $W_n(x) = \rho_n x$   
where  $\rho_n$  itself solves a scalar recursion:

$$\rho_n = \rho_{n-1} + \max \{ \theta \rho_{n-1}, 1 \}$$

This gives: 
$$\rho_n = \begin{cases} n, & \text{if } n \leq n^* \\ (1+\theta)^{n-n^*} n^*, & \text{if } n \geq n^* \end{cases}$$



where  $n^*$  is the least integer

$$\text{s.t. } (1+\theta)n^* \geq 1+n^*$$



$$n^* \geq 1/\theta$$



$$n^* = \lceil 1/\theta \rceil.$$

$\therefore$  Optimal policy is to invest entire income in yrs.

0, 1, 2, ...,  $T-n^*-1$  (to build up enough capital)

& then consume whole of the income in yrs.

$T-n^*, T-n^*+1, \dots, T-1.$



# Example of Stochastic DP:

(Inventory Control)  
Dynamics of stock:

Suppose you are the store manager for Walmart

$$x_{k+1} = x_k + u_k - w_k$$

$k=0, 1, 2, \dots, T-1$

$x_k$  = Stock available @ the beginning of the  $k^{\text{th}}$  period

$u_k$  = Stock ordered (and immediately delivered) at the beginning of the  $k^{\text{th}}$  period

Assume that demands  $w_0, w_1, \dots, w_{T-1}$  are indep. r.v.s

Assume  $u_k \geq 0$

$w_k$  = Demand during  $k^{\text{th}}$  period with known probability distribution

Negative  $x$  (stock) means backlogged demand

cost has 2 components:

- (i)  $r(x_k)$ : penalty for either positive stock (i.e., holding cost for excess inventory) OR negative stock (i.e., shortage cost for unfulfilled demand)
- (ii)  $c u_k$ : purchasing cost (here,  $c$  is the cost for per unit ordered)

want to solve the OCP:

$$\min_{x \in \Gamma} \mathbb{E} \left\{ r(x_0) + \sum_{k=0}^{T-1} (r(x_k) + c u_k) \right\}, \text{ where } r(\cdot) \text{ is a piecewise linear cost}$$

piecewise linear cost:  $v(x_k) := p \max(0, -x_k) + h \max(0, x_k)$

with slope  $h$  when stock  $> 0$

" " " " " "  $< 0$ .

Assume  $p > c$ , otherwise  $u_k^* = 0 \forall k = 0, \dots, T-1$

$\therefore$  DP eq<sup>n</sup>:

$$V_T(\underline{x}) = 0.$$

Recall: OCP is equivalent to  

$$\min_{p \in \Gamma} \mathbb{E} \left[ \sum_{k=0}^{T-1} \left( v(x_k + u_k - w_k) + c u_k \right) \right]$$

Since  $\mathbb{E}[v(x_0)]$  cannot be influenced  
 so, we can subtract that.

$$V_k(\underline{x}) = \inf_{u_k \geq 0} \mathbb{E}_{w_k} \left\{ \underbrace{v(x + u - w_k) + c u}_{c_k(x, u, w_k)} + V_{k+1} \left( \underbrace{x + u - w_k}_{f_k(x, u, w_k)} \right) \right\}$$

$$= \inf_{u \geq 0} \mathbb{E}_w \left\{ v(x + u - w) + c u + V_{k+1}(x + u - w) \right\}$$

where  $k = T-1, T-2, \dots, 0$

(assuming  $w_k$  as iid)

Can show by induction that  $V_k(\underline{x})$  is a non-neg. convex f<sup>n</sup> that  $\rightarrow \infty$  as  $|x_k| \rightarrow \pm \infty$

Write  $\mathbb{E}_w[r(x+u-w)]$  as  $H(x+u)$ , where  $H(y) := \mathbb{E}[r(y-w)]$

Then the DP recursion in the prev. page can be re-written as:

$$V_T(x) = 0$$

$$V_k(x) = \inf_{u \geq 0} \left\{ cu + H(x+u) + \mathbb{E}[V_{k+1}(x+u-w)] \right\}$$

From the def<sup>n</sup>,  $H(y) := \mathbb{E}[r(y-w)]$ , notice that since  $r(\cdot)$  is a convex f<sup>n</sup>, so is  $H(\cdot)$ . Furthermore, slope of  $H(\cdot)$  approaches " $h$ " as  $x \rightarrow \infty$  and " $-p$ " as  $x \rightarrow -\infty$

Let  $y := x+u$ . Then the constraint in "inf"  $u \geq 0 \Leftrightarrow y \geq x$ , and  $cu \equiv cy - cx$ . Also,  $G(y) := cy + H(y) + \mathbb{E}[V_{k+1}(y-w)]$

$$\text{Then, } V_k(x) = \inf_{y \geq x} (G(y) - cx) = \left( \inf_{y \geq x} G(y) \right) - cx.$$

Suppose, we can show that " $G$ " is convex and  $\lim_{|x| \rightarrow \infty} G(x) \rightarrow \infty$ , and let  $x_k := \arg \min_{y \geq x} (G_k(x_k))$

Then,  $V_k(x_k) = G_k(y_k^*) - cx_k$ , where

$$y_k^* = \begin{cases} s_k & \text{if } s_k \geq x_k, \\ x_k & \text{otherwise.} \end{cases}$$

(i.e.) Optimal policy is a time-varying threshold policy (i.e., restock to  $s_k$  iff  $x_k \leq s_k$ ).

To finish this example, we now prove that  $V_k(x_k)$  is convex with slope  $> 0$  as  $x_k \rightarrow \infty$ , and with slope  $< 0$  as  $x_k \rightarrow -\infty$  for  $k = 0, 1, \dots, T-1$ . Backward induction proof

Notice that  $G_{T-1}(y) = cy + H_{T-1}(y)$  is convex since  $H_{T-1}(y)$  is convex, and its slope approaches  $(h+c) > 0$  as  $y \rightarrow \infty$ , and approaches  $-p+c < 0$  as  $y \rightarrow -\infty$ .

$$\text{So, } V_{T-1}(x_{T-1}) = c(s_{T-1} - x_{T-1})^+ + H_{T-1}(\max\{s_{T-1}, x_{T-1}\})$$

where  $z^+ := \max(z, 0) \forall z \in \mathbb{R}$ .


This  $f^*$  is convex with slope  $> 0$  as  $x_{T-1} \rightarrow \infty$  and slope approaching  $-c < 0$  as  $x_{T-1} \rightarrow -\infty$ . Now show this recursively backward.

## Back to DP Theors

General result:

$\gamma^* \in \Gamma$ , (optimal policy) is a non-randomized (deterministic)

Markovian policy

(i.e)  $\gamma^* \in \Gamma_M \subset \Gamma$   class of all history dependent randomized policies.

Infinite Horizon Problems:

$\min_{\gamma \in \Gamma} \mathbb{E} \left[ \sum_{k=0}^{+\infty} c_k(x_k, u_k, w_k) \mid x(0) = \bar{x} \right]$

called Total cost

→ may NOT be well defined  
(series sum may NOT converge or may diverge even when "c" is bdd.)

NOTE :

$$\sum_{k=0}^{+\infty} a_k \text{ converges} \implies \lim_{k \rightarrow \infty} a_k = 0$$

↳ (contrapositive)

$$\lim_{k \rightarrow \infty} a_k \neq 0 \implies \sum_{k=0}^{+\infty} a_k \text{ diverges.}$$

This motivates introducing discount factor

$$0 < \beta < 1$$

a number

Total discounted cost

$$\min_{\gamma \in \Gamma} \mathbb{E} \left[ \sum_{k=0}^{+\infty} \beta^k c_k(x_k, u_k^{i\beta}) \mid \underline{x}(0) = \underline{x} \right]$$

--- (\*)

$\beta \approx 0$  : Tomorrow/future is NOT imp.

If  $\beta \approx 1$  : All days are important (approx.)

Total discounted cost is ALWAYS well defined if  $f^k = c_k(\cdot)$  is bdd.

"MYOPIA"

Today \$1 = \$1

Tomorrow \$1 = 90¢  
today

geometric fall of USD

Day after tomorrow : 81¢

$$\beta = 0.9$$

Finite horizon Discounted cost :

$$\text{Let } W(x, n) = \min_{\delta \in \Gamma} \mathbb{E} \left[ \sum_{k=0}^{T-1} \beta^k c_k(x_k, u_k) \mid x(0) = x \right]$$



DP eq<sup>n</sup> : (when  $x$  is Markov CP)

$$W(x, u) = \min_{u \in U} \left\{ c(x, u) + \sum_{j \in X} p_{x_j}(u) \cdot \beta W(j, u) \right\}$$

$$W(x, 0) = 0.$$

Infinite horizon Discounted Cost:

$$W(x, \infty) = \min_{u \in U} \left\{ \underline{c(x, u)} + \beta \sum_{j \in X} p_{x_j}(u) W(j, \infty) \right\}$$

--- (\*\*)

same function.

Algebraic nonlinear eq<sup>n</sup>

nonlinear since  $c(x, u) \neq 0$ .

One eq<sup>n</sup> for each  $x$  (when  $x$  is enumerable)

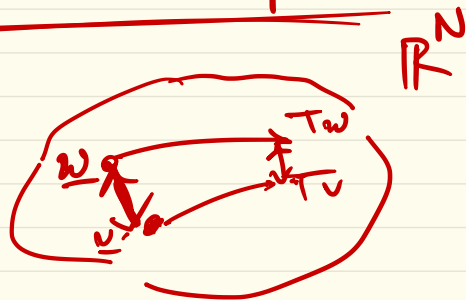
$\therefore x$  eq<sup>n</sup>s in  $x$  unknowns.

Question :

Solution  $x$

- Exists?
- Unique?
- Algorithm to solve it?

Contraction Map.



$$F \subseteq \mathbb{R}^N$$

closed subset  $0 < \beta < 1$

$$d(Tx, Ty) \leq \beta d(x, y)$$

The map  $T$  is called contractive.

$$T: \begin{matrix} F \\ x, y \end{matrix} \mapsto \begin{matrix} F \\ Tx, Ty \end{matrix}$$

# Contraction Mapping Thm.

$F \subseteq \mathbb{R}^N$   
closed set.

Suppose  $\exists$  a scalar  $\beta$  where  $0 < \beta < 1$ , s.t.

$$\|T\underline{w} - T\underline{v}\| \leq \beta \|\underline{w} - \underline{v}\|$$

$\forall \underline{v}, \underline{w} \in F$

Then

①  $\exists \underline{z} \in F$  s.t.  $T\underline{z} = \underline{z}$ , called fixed pt. of map.  $T(\cdot)$

② The fixed pt.  $\underline{z}$  is unique

③ start with any  $\underline{w} \in F$ , then

$$\lim_{n \rightarrow \infty} T^n \underline{w} = \underline{z} \quad (\text{fixed pt. iteration})$$