

## MARKOV DECISION PROCESSES (MDP)

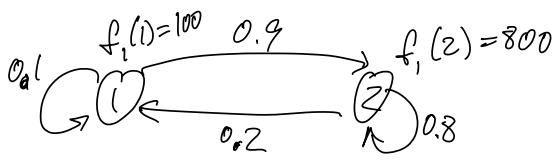
(CHAPTER 12 IN FELDMAN — RACINEZ-FLORES TEXT)

### EXAMPLE

$S = E = \{1, 2\}$  STATE SPACE

$A = \{1, 2\}$  ACTION SPACE

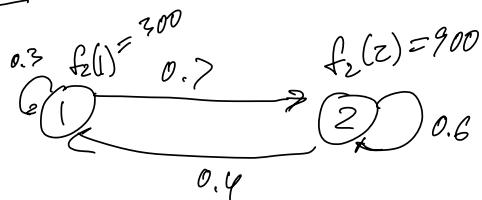
ACTION  
 $k=1:$



$$P_1 = \begin{pmatrix} 0.1 & 0.9 \\ 0.2 & 0.8 \end{pmatrix} = (P_1(i,j))$$

$$f_1 = \begin{pmatrix} 100 \\ 800 \end{pmatrix}$$

ACTION  
 $k=2:$



$$P_2 = \begin{pmatrix} 0.3 & 0.7 \\ 0.4 & 0.6 \end{pmatrix} = (P_2(i,j))$$

$$f_2 = \begin{pmatrix} 300 \\ 900 \end{pmatrix}$$

IN THIS COURSE WE ALWAYS ASSUME:  $E$  AND  $A$  FINITE

GENERAL DEFINITION OF AN MDP:

$X = (X_0, X_1, X_2, \dots)$  STATE PROCESS,  $X_n \in E = S$  = STATE SPACE;

$D = (D_0, D_1, D_2, \dots)$  ACTION PROCESS,  $D_n \in A$  = ACTION SPACE;

ACTION  $D_n$  AT TIME  $n$  MAY DEPEND ONLY ON  $X_n, D_{n-1}, X_{n-1}, \dots, D_0, X_0$ ;

$f_n(i) =$  COST INCURRED WHEN ACTION  $k$  IS CHOSEN IN STATE  $i$ ;

$P_k(i, j) =$  PROB. OF TRANSITION FROM  $i$  TO  $j$  WHEN ACTION  $k$  CHOSEN.

PROCESS  $(X, D)$  IS AN MDP IF  $\forall n, \forall i, j \in E, \forall k \in A$ :

$$\mathbb{P}\{X_{n+1}=j \mid X_0, D_0, \dots, X_n=i, D_n=k\} = \mathbb{P}\{X_{n+1}=j \mid X_n=i, D_n=k\} = P_k(i, j).$$

## FINITE HORIZON PROBLEM

$0 \leq \alpha \leq 1$  IS A DISCOUNT FACTOR.

$d \in \mathcal{D}$  A GIVEN POLICY;  $\mathcal{D} = \text{SET OF ALL POLICIES}$

POLICY  $d$ :  $a_{(m)}(i) \in A$  ACTION YOU TAKE IN STATE  $i$ ,  
WHEN AT TIME-DISTANCE  $m$   
FROM THE END OF HORIZON

EXPECTED DISCOUNTED COST UNDER A GIVEN

POLICY  $d$ , OVER TIME HORIZON  $0, 1, 2, \dots, m$ ?

$$\begin{aligned} V_{(m),d}^\alpha(i) &= \mathbb{E} \left[ \sum_{n=0}^m \alpha^n f_{\mathcal{D}_n}(X_n) \mid X_0 = i \right] = \\ &= \mathbb{E} \left[ \sum_{n=0}^m \alpha^n f_{a_{(m-n)}(X_n)}(X_n) \mid X_0 = i \right] \end{aligned}$$

EXAMPLE ( $m=2$ ):

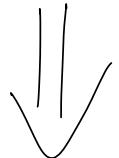
$$\begin{aligned} V_{(2),d}^\alpha(i) &= \mathbb{E} \left[ f_{a_{(2)}(X_0)}(X_0) + \alpha f_{a_{(1)}(X_1)}(X_1) + \right. \\ &\quad \left. + \alpha^2 f_{a_{(0)}(X_2)}(X_2) \mid X_0 = i \right] = \\ &= f_{a_{(2)}(i)}(i) + \alpha \sum_{j \in E} P_{a_{(2)}(i)}(i, j) \cdot \\ &\quad \cdot \mathbb{E} \left[ f_{a_{(1)}(X_1)}(X_1) + \alpha f_{a_{(0)}(X_2)}(X_2) \mid \begin{array}{l} X_0 = i, \\ a_0 = a_{(2)}(i), \\ X_1 = j \end{array} \right] \\ &= f_{a_{(2)}(i)}(i) + \alpha \sum_{j \in E} P_{a_{(2)}(i)}(i, j) \cdot \\ &\quad \cdot \mathbb{E} \left[ f_{a_{(1)}(j)}(j) + \alpha f_{a_{(0)}(X_2)}(X_2) \mid X_1 = j \right] \end{aligned}$$

FINALLY, we obtain,

$$v_{(0), d}^\alpha(i) = f_{a_{(0)}(i)}(i) + \sum_{j \in E} P_{a_{(0)}(i)}(i, j) v_{(1), d}^\alpha(j), \quad \forall i$$

IN GENERAL:

$$\left\{ \begin{array}{l} v_{(m), d}^\alpha(i) = f_{a_{(m)}(i)}(i) + \sum_{j \in E} P_{a_{(m)}(i)}(i, j) v_{(m-1), d}^\alpha(j), \quad \forall i \\ \text{INITIAL CONDITION: } v_{(0), d}^\alpha(i) = f_{a_{(0)}(i)}(i), \quad \forall i \end{array} \right.$$



CAN RECURSIVELY COMPUTE ALL  $v_{(m), d}^\alpha(i)$ ,  $\forall i$ ,  $m = 0, 1, 2, \dots$   
FOR ANY GIVEN POLICY  $d$

PROBLEM! FIND OPTIMAL( $\min$ )  $v_{(m)}^{\alpha}(i)$  AND OPTIMAL  $a_{(m)}(i)$ !

$$v_{(m)}^{\alpha}(i) = \min_{d \in \mathcal{D}} v_{(m), d}^{\alpha}(i)$$

↓  
POLICY

SOLUTION: USE BELLMAN (DYNAMIC PROGRAMMING) EQUATIONS:

$$v_{(m)}^{\alpha}(i) = \min_{k \in A} \left\{ f_k(i) + \alpha \sum_{j \in E} P_k(i, j) v_{(m-1)}^{\alpha}(j) \right\}, \quad \forall i$$

INITIAL CONDITION:  $v_{(0)}^{\alpha}(i) = \min_{k \in A} f_k(i), \quad \forall i$

$$a_{(m)}^{\alpha}(i) = \arg \min_{k \in A} \left\{ f_k(i) + \alpha \sum_{j \in E} P_k(i, j) v_{(m-1)}^{\alpha}(j) \right\}, \quad \forall i$$

$$a_{(0)}^{\alpha}(i) = \arg \min_{k \in A} f_k(i), \quad \forall i$$

USING BELLMAN EQUATIONS, CAN RECURSIVELY, FOR  
 $m=0, 1, 2, \dots$  COMPUTE OPTIMAL VALUE FUNCTION  $v_{(m)}^{\alpha}(i)$   
 AND OPTIMAL POLICY  $a_{(m)}(i)$ .

## INFINITE HORIZON PROBLEM. $\gamma < 1$

TOTAL DISCOUNTED COST (FOR  $\gamma < 1$ ) UNDER A GIVEN POLICY  $d$ :

$$v_d^\gamma(i) = \mathbb{E} \left[ \sum_{n=0}^{\infty} \gamma^n f_{D_n}(X_n) \mid X_0 = i \right]$$

PROBLEM: FIND A POLICY  $d$  MINIMIZING  $v_d^\gamma(i)$ ,  $\forall i$ .

OPTIMAL VALUE FUNCTION:

$$v^\gamma(i) = \min_{d \in \mathcal{D}} v_d^\gamma(i)$$

SET OF ALL POLICIES

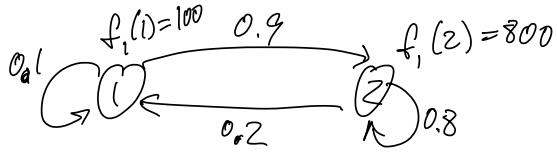
## STATIONARY POLICIES

A POLICY S.T. ACTION  $a(i)$  DEPENDS ONLY  
ON THE CURRENT STATE  $i$ ,  
BUT NOT ON TIME, IS CALLED STATIONARY.  
 $a = (a(i))$  IS CALLED ACTION FUNCTION.

IF POLICY  $\alpha$  IS STATIONARY, GIVEN BY  
ACTION FUNCTION  $a$ , THEN OFTEN  
ACTION FUNCTION  $a$  ITSELF IS CALLED  
A (STATIONARY) POLICY

EXAMPLE:

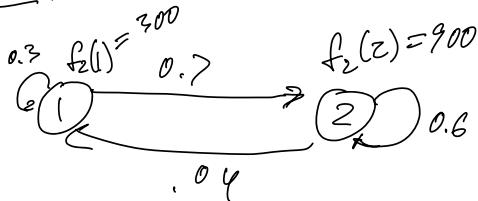
ACTION  
 $k=1$ :



$$P_1 = \begin{pmatrix} 0.1 & 0.9 \\ 0.2 & 0.8 \end{pmatrix} = (P_1(i,j))$$

$$f_1 = \begin{pmatrix} 100 \\ 800 \end{pmatrix}$$

ACTION  
 $k=2$ :



$$P_2 = \begin{pmatrix} 0.3 & 0.7 \\ 0.4 & 0.6 \end{pmatrix} = (P_2(i,j))$$

$$f_2 = \begin{pmatrix} 300 \\ 900 \end{pmatrix}$$

STAT. POLICY:  $\{a(1)=2, a(2)=1\} \iff a=(z,i)$ .

THE PROCESS BECOMES A DTMC:

$$P_a = \begin{pmatrix} 0.3 & 0.7 \\ 0.2 & 0.8 \end{pmatrix} \quad \text{TRANSITION PROBABILITIES}$$

$$f_a = \begin{pmatrix} 300 \\ 800 \end{pmatrix} \quad \text{COSTS}$$

IN GENERAL:

AN MDP UNDER A STATIONARY POLICY BECOMES A DTMC (WITH A COST ASSOCIATED WITH EACH STATE)

INFINITE HORIZON. DISCOUNT FACTOR  $\lambda < 1$

TOTAL DISCOUNTED COST UNDER

A GIVEN STATIONARY POLICY  $\alpha$

$$\begin{aligned} v_\alpha^\lambda(i) &= f_\alpha(i) + \lambda \sum_j P_\alpha(i,j) \cdot f_\alpha(j) + \\ &\quad + \lambda^2 \sum_j P_\alpha^{(2)}(i,j) f_\alpha(j) + \dots \end{aligned}$$



$$\begin{aligned} v_\alpha^\lambda &= \begin{pmatrix} \vdots \\ v_\alpha^\lambda(i) \\ \vdots \\ \vdots \end{pmatrix} = f_\alpha + \lambda P_\alpha f_\alpha + \lambda^2 P_\alpha^2 f_\alpha + \dots = \\ &= [I + \lambda P_\alpha + \lambda^2 P_\alpha^2 + \dots] f_\alpha \end{aligned}$$

$$v_\alpha^\lambda = [I - \lambda P_\alpha]^{-1} f_\alpha$$

INFINITE HORIZON. DISCOUNT FACTOR  $\alpha < 1$ .

TOTAL DISCOUNTED COST MINIMIZATION

FOR ANY  $0 \leq \alpha < 1$ , AN OPTIMAL POLICY EXISTS,  
WHICH IS A STATIONARY POLICY.

(OPT. POLICY MAY, AND TYPICALLY WILL, DEPEND ON  $\alpha$ )

THE OPTIMAL VALUE FUNCTION  $v^\alpha(i)$

SATISFIES BELLMAN EQ:

$$v^\alpha(i) = \min_{a \in A} \left\{ f_n(i) + \alpha \sum_j P_n(i, j) v^\alpha(j) \right\}.$$

MOREOVER, IT IS THE ONLY FUNCTION WHICH  
SATISFIES THIS EQUATION.

A STATIONARY POLICY  $a$  IS OPTIMAL (IF AND ONLY IF

$$a(i) = \arg \min_{a \in A} \left\{ f_n(i) + \alpha \sum_j P_n(i, j) v^\alpha(j) \right\}.$$

INFINITE HORIZON. DISCOUNT FACTOR  $\gamma < 1$ .

TOTAL DISCOUNTED COST MINIMIZATION

FINDING AN OPTIMAL STATIONARY POLICY AND  
THE OPTIMAL VALUE FUNCTION.

|| POLICY IMPROVEMENT ALGORITHM  
FOR INFINITE HORIZON DISCOUNTED COST  
(IN MDP-COMMS FILE)

## AVERAGE COST MINIMIZATION PROBLEM

(NO DISCOUNT:  $\alpha = 1$ )

$d \in \mathcal{D}$  A POLICY. AVE. COST:

$$\varphi_d = \lim_{n \rightarrow \infty} \mathbb{E} \frac{f_{d_0}(x_0) + \dots + f_{d_{n-1}}(x_{n-1})}{n} =$$
$$= \lim_{n \rightarrow \infty} \frac{f_{d_0}(x_0) + \dots + f_{d_{n-1}}(x_{n-1})}{n}$$

FINITE E  
FINITE A

## AVERAGE COST MINIMIZATION PROBLEM:

FIND POLICY  $d^*$ , FOR WHICH

$$\varphi_{d^*} = \min_d \varphi_d$$

$$\varphi^* = \varphi_{d^*} = \text{OPTIMAL (MINIMUM) COST}$$

## AVERAGE COST UNDER

---

### A GIVEN STATIONARY POLICY $\alpha = (\alpha(i))$

IN THIS COURSE:

WE WILL ALWAYS ASSUME THAT UNDER ANY STAT. POLICY  $\alpha$ , THE TRANS. PROB. ARE SUCH THAT THE RESULTING DTMC IS IRREDUCIBLE.

TRANS. PROB. MATRIX  $\underline{P}_\alpha$ ; COSTS (COLUMN VECTOR)  $\underline{f}_\alpha$

$$\text{Ave. cost } \varphi_\alpha = \sum_{i \in E} \pi_i f_\alpha(i), \text{ WHERE}$$

$\pi = (\pi_i)$  IS THE (UNIQUE) STATIONARY DISTR.

$$\begin{cases} \pi = \pi \underline{P}_\alpha \\ \sum_i \pi_i = 1 \end{cases}$$

UNDER OUR ASSUMPTIONS, THERE ALWAYS EXISTS  
AN OPTIMAL POLICY  $\alpha$ , WHICH IS STATIONARY.

$\exists$  NUMBER  $\varphi^*$  AND A FUNCTION (VECTOR)  $h = (h(i))$   
SUCH THAT :

$$(i) \quad \varphi^* + h(i) = \min_{\alpha \in A} \left\{ f_\alpha(i) + \sum_{j \in E} P_\alpha(i, j) h(j) \right\}, \forall i;$$

(ii)  $\varphi^*$  IS THE OPT. COST ;

(iii) STAT. POLICY  $\alpha$  IS OPTIMAL IF AND ONLY IF

$$\alpha(i) = \arg \min_{\alpha \in A} \left\{ f_\alpha(i) + \sum_{j \in E} P_\alpha(i, j) h(j) \right\}, \forall i;$$

(iv) SUCH FUNCTION  $h = (h(i))$  IS UNIQUE  
UP TO AN ADDITIVE CONSTANT.

FOR EXAMPLE, IF YOU FIX  $i_0 \in E$ , AND SET  $h(i_0) = 0$ ,  
THEN  $h$  IS JUST UNIQUE.

INTUITIVELY :

$$1) \quad \varphi^* = \lim_{\alpha \uparrow 1} ((1-\alpha) v^\alpha(i)) , \quad \forall i.$$

$\uparrow$

$$(1-\alpha) v^\alpha(i) = (1-\alpha) \alpha^0 \mathbb{E} f_{D_0}(X_0) + (1-\alpha)\alpha^1 \mathbb{E} f_{D_1}(X_1) + (1-\alpha)\alpha^2 \mathbb{E} f_{D_2}(X_2) + \dots$$

$$2) \quad h(i) - h(j) = \lim_{\alpha \uparrow 1} \{ v^\alpha(i) - v^\alpha(j) \}$$

$$v^\alpha(i) = \min_k \left\{ f_k(i) + \alpha \sum_j P_k(i, j) v^\alpha(j) \right\} , \quad \forall i$$

$$(1-\alpha) v^\alpha(i) = \min_k \left\{ f_k(i) + \alpha \sum_j P_k(i, j) v^\alpha(j) \right\} - \alpha v^\alpha(i) =$$

$$= \min_k \left\{ f_k(i) + \alpha \sum_j P_k(i, j) (v^\alpha(j) - v^\alpha(i)) \right\}$$

$\underbrace{\quad}_{\alpha \uparrow 1}$

$$\varphi^* = \min_k \left\{ f_k(i) + \sum_j P_k(i, j) [h(j) - h(i)] \right\}$$

$$\varphi^* = \min_k \left\{ f_k(i) + \sum_j P_k(i, j) h(j) - h(i) \right\}$$

$$\varphi^* + h(i) = \min_{k \in A} \left\{ f_k(i) + \sum_{j \in E} P_k(i, j) h(j) \right\}$$

HOW TO CHECK THAT A GIVEN STATIONARY

POLICY  $\pi$  IS OPTIMAL

→ SEE EXAMPLE ON PAGE 338 OF MDP-COMMS

## HOW TO FIND AN OPTIMAL STATIONARY POLICY

|| POLICY IMPROVEMENT ALGORITHM (FOR AVERAGE COST) :

|| PROPOSITION 12.11 ON PAGE 340 OF MDP-COMMS

— SEE EXAMPLE OF THE POLICY IMPROVEMENT

ALGORITHM APPLICATION,

PAGE 341 OF MDP-COMMS