

# Markov Decision Processes

*Sequential decision-making over time*

**Aditya Mahajan**  
**McGill University**

Lecture Notes for ECSE 506: Stochastic Control and Decision Theory  
February 11, 2016

# Theory of Markov decision processes

MDP functional models

Perfect state observation

MDP probabilistic models

Stochastic orders

Optimal monotone strategies in MDPs

POMDP probabilistic models

# MDP Theory: Functional models

# Functional model for stochastic dynamical systems

- Notation**
- $X_t \in \mathcal{X}$  : State of the system at time  $t$
  - $Y_t \in \mathcal{Y}$  : Observation of controller at time  $t$
  - $U_t \in \mathcal{U}$  : Control action taken by the controller at time  $t$
  - $W_t \in \mathcal{W}$  : Noise in system dynamics at time  $t$
  - $N_t \in \mathcal{N}$  : Observation noise at time  $t$

- Assumptions**
- ▶ The system runs in discrete-time until **horizon  $T$** .
  - ▶ The **primitive random variables**  $\{X_1, W_{1:T}, N_{1:T}\}$  are defined over a common probability space  $(\Omega, \mathfrak{F}, P)$ .
  - ▶ The primitive variables  $\{X_1, W_{1:T}, N_{1:T}\}$  are **mutually independent** with known probability distribution.

- System dynamics**
- ▶  $X_{t+1} = f_t(X_t, U_t, W_t)$
  - ▶ The **dynamic functions**  $\{f_t\}_{t=1}^T$  are known.

- Observations**
- ▶  $Y_t = h_t(X_t, N_t)$
  - ▶ The **observation functions**  $\{h_t\}_{t=1}^T$  are known.

# The control strategy and its performance

- Control design**
- ▶  $U_t = g_t(Y_{1:t}, U_{1:t-1})$
  - ▶ The **control strategy**  $g = \{g_t\}_{t=1}^T$  is to be determined.
  - ▶ The controller has **classical information structure** (i.e., it remembers everything that has been observed and done in the past).

- Cost**
- ▶ **Per step-cost** at time  $t \in \{1, \dots, T-1\}$ :  $c_t(X_t, U_t)$ .
  - ▶ **Terminal cost** at time  $T$ :  $c_T(X_T)$ .

- Total expected cost**
- ▶  $J(g) = \mathbb{E}^g \left[ \sum_{t=1}^{T-1} c_t(X_t, U_t) + c_T(X_T) \right]$

## Alternative formulation: Reward maximization

- ▶ In some applications, it is more natural to model per-step and terminal **reward functions**  $r_t(X_t, U_t)$  and  $r_T(X_T)$ .
- ▶ In such applications, the objective is to **maximize** the total expected reward

$$J(g) = \mathbb{E}^g \left[ \sum_{t=1}^{T-1} r_t(X_t, U_t) + r_T(X_T) \right]$$

# The problem of optimizing over time

## Objective Given

- ▶ The spaces  $(\mathcal{X}, \mathcal{Y}, \mathcal{U}, \mathcal{W}, \mathcal{N})$
- ▶ Horizon  $T$
- ▶ Probability distribution of  $\{X_1, W_{1:T}, N_{1:T}\}$
- ▶ Dynamics functions  $\{f_t\}_{t=1}^T$
- ▶ Observation functions  $\{h_t\}_{t=1}^T$
- ▶ Cost functions  $\{c_t\}_{t=1}^T$

## Choose

- ▶ **Control strategy  $g$**  to minimize the total expected cost  $J(g)$ .  
(Alternatively, to maximize the total expected reward).

## Application domains

- ▶ Systems and Control
- ▶ Communication
- ▶ Power Systems
- ▶ Artificial Intelligence
- ▶ Operations Research
- ▶ Financial Engineering
- ▶ Natural Resource Management

# Perfect and imperfect observations at the controller

**Perfect state observation** Perfect state observation refers to the scenario when  $\mathcal{Y} = \mathcal{X}$  and  $h_t(X_t, N_t) = X_t$ ; thus, at each time the controller perfectly observes the state. Such a model is also called **Markov decision process (MDP)**.

**Imperfect state observation** Imperfect state observation refers to the general model described above (when  $Y_t \neq X_t$ ). Such a model is also called **partially observed Markov decision process (POMDP)**.

**Solution approach** First focus on problems with perfect state observation and identify the structure of optimal controllers and a recursive algorithm, called **dynamic programming** decomposition, to find an optimal strategy

Then show that an appropriate state expansion converts problems with imperfect state observations to a problem with perfect state observation. Thus, it is possible to reuse the results for models with perfect state observation in models with imperfect state observation.

# MDP Theory: Perfect state observation



# Structure of optimal strategies

## Theorem (Structural result)

A strategy  $\mathbf{g} = \{g_t\}_{t=1}^T$  is called **Markov** if it only uses  $X_t$  at time  $t$  to pick  $U_t$  i.e.,

$$U_t = g_t(X_t)$$

**Restricting attention to Markovian strategies is without any loss of optimality.**

**Implication** Let  $\mathcal{G}_{1:T}^H$  denote the family of all history dependent strategies and  $\mathcal{G}_{1:T}^M$  denote the family of all Markov strategies. The above theorem asserts that

$$\min_{\mathbf{g} \in \mathcal{G}_{1:T}^M} J(\mathbf{g}) = \min_{\mathbf{g} \in \mathcal{G}_{1:T}^H} J(\mathbf{g})$$

Note that LHS  $\leq$  RHS because  $\mathcal{G}_{1:T}^M \subset \mathcal{G}_{1:T}^H$ . The above theorem is asserting equality.

This result reduces the solution space and thereby simplifies the optimization problem.

# When is extra information irrelevant for optimal control?

## Blackwell's principle of irrelevant information

Let  $\mathcal{X}$ ,  $\mathcal{Y}$ ,  $\mathcal{U}$  be standard Borel spaces and  $X \in \mathcal{X}$  and  $Y \in \mathcal{Y}$  be random variables defined on a common probability space  $(\Omega, \mathfrak{F}, P)$ .

A decision maker observes  $(X, Y)$  and chooses  $U$  to minimize  $\mathbb{E}[c(X, U)]$  where  $c: \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$  is a measurable function.

**Then, choosing  $U$  just as a function of  $X$  is without loss of optimality.**

Formally,  $\exists g^*: \mathcal{X} \rightarrow \mathcal{U}$  such that  $\forall g: \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{U}$

$$\mathbb{E}[c(X, g^*(X))] \leq \mathbb{E}[c(X, g(X, Y))]$$

**Proof** We prove the result for the case when  $\mathcal{X}$ ,  $\mathcal{Y}$ ,  $\mathcal{U}$  are finite valued.

- ▶ Define  $g^*(x) = \arg \min_{u \in \mathcal{U}} c(x, u)$ .
- ▶ Then,  $\forall x \in \mathcal{X}$  and  $\forall u \in \mathcal{U}$ :  $c(x, g^*(x)) \leq c(x, u)$ .
- ▶ Hence,  $\forall g: \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{U}$  and  $\forall y \in \mathcal{Y}$ :  $c(x, g^*(x)) \leq c(x, g(x, y))$ .

The above point-wise inequality implies the inequality in expectation.

# How to identify irrelevant information in dynamic setups?

**Two-step Lemma** Let  $T = 2$ . For any control strategy  $g = (g_1, g_2)$  there exists a Markov control law  $g_2^*: \mathcal{X} \rightarrow \mathcal{U}$  such that  $J(g_1, g_2^*) \leq J(g_1, g_2)$ .

- Proof**
- ▶ Define  $J_1(g_1) = \mathbb{E}[c_1(X_1, U_1)]$  and  $J_2(g_1, g_2) = \mathbb{E}[c_2(X_2, U_2)]$ .
  - ▶ Then  $J(g_1, g_2) = J_1(g_1) + J_2(g_1, g_2)$
  - ▶  $J_2(g_1, g_2) = \mathbb{E}[c_2(X_2, g_2(X_2, X_1, U_1))]$ . By Blackwell's principle of irrelevant information,  $\exists g_2^*: X_2 \mapsto U_2$  such that  $J_2(g_1, g_2^*) \leq J_2(g_1, g_2)$ .

# How to identify irrelevant information in dynamic setups?

**Three-step Lemma** Let  $T = 3$ . For any control strategy  $g = (g_1, g_2, g_3)$  such that  $g_3$  is Markov, there exists a Markov control law  $g_2^*: \mathcal{X} \rightarrow \mathcal{U}$  such that  $J(g_1, g_2^*, g_3) \leq J(g_1, g_2, g_3)$ .

- Proof**
- ▶ Define  $J_t(g_{1:t}) = \mathbb{E}[c_t(X_t, U_t)]$ .  
Then  $J(g_{1:3}) = J_1(g_1) + J_2(g_{1:2}) + J_3(g_{1:3})$ .
  - ▶ Define  $\tilde{c}_3(x, u; g_3) = \mathbb{E}[c_3(X_3, g_3(X_3)) \mid X_2 = x, U_2 = u]$ .
  - ▶ Then,  $J_3(g_{1:3}) = \mathbb{E}[\mathbb{E}[c_3(X_3, g_3(X_3)) \mid X_2, U_2]] = \mathbb{E}[\tilde{c}_3(X_2, U_2; g_3)]$ .
  - ▶ Define  $\tilde{c}_2(x, u; g_3) = c_2(x, u) + \tilde{c}_3(x, u; g_3)$ .
  - ▶ Then,  $J_2(g_{1:2}) + J_3(g_{1:3}) = \mathbb{E}[\tilde{c}_2(X_2, g_2(X_2, X_1, U_1); g_3)]$ .
- Use Blackwell's principle of irrelevant information, as in the two-step lemma.

# Backward induction proof of the structural result

To be written

# Dynamic programming decomposition to find optimal Markov strategy

**Definition of value functions** Define **value functions**  $\{V_t\}_{t=1}^T$ ,  $V_t: \mathcal{X} \rightarrow \mathbb{R}$  recursively as follows:

$$V_T(x) = c_T(x), \quad x \in \mathcal{X}$$

and for  $t = T - 1, T - 2, \dots, 1$ :

$$Q_t(x, y) = \mathbb{E}[c(X_t, U_t) + V_{t+1}(X_{t+1}) \mid X_t = x, U_t = u], \quad \forall x \in \mathcal{X}, u \in \mathcal{U}(x)$$

$$V_t(x) = \min_{u \in \mathcal{U}(x)} Q_t(x, u), \quad \forall x \in \mathcal{X}$$

**Verification step** A **Markov strategy**  $\{g_t^*\}_{t=1}^T$  is optimal iff

$$g_t^*(x) \in \arg \min_{u \in \mathcal{U}(x)} Q_t(x, u), \quad \forall x \in \mathcal{X} \text{ and } \forall t \in \{1, \dots, T\}$$

## Bellman's principle of optimality

An optimal policy has the property that whatever the initial state and initial decisions are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision.

# The comparison principle to prove dynamic programming

## The cost-to-go functions

For any strategy  $\mathbf{g}$ , define the cost-to-go function at time  $t$  as

$$J_t(x; \mathbf{g}) = \mathbb{E}^{\mathbf{g}} \left[ \sum_{s=t}^{T-1} c_s(X_s, \mathbf{u}_s) + c_T(X_T) \mid X_t = x \right]$$

Note that

$$J(\mathbf{g}) = \mathbb{E}[J_1(X_1; \mathbf{g})]$$

## The comparison principle

For any **Markov strategy**  $\mathbf{g}$

$$J_t(x; \mathbf{g}) \geq V_t(x)$$

with equality at  $t$  iff the **future strategy**  $\mathbf{g}_{t:T}$  satisfies the verification step.

An immediate consequence of the comparison principle is that the strategy obtained using the dynamic programming decomposition is optimal.

# Proof of comparison principle

- Proof**
- ▶ **Basis:**  $J_T(x) = V_T(x)$ . Thus, the comparison principle is true.
  - ▶ **Induction hypothesis:** Comparison principle is true for  $t + 1$ .
  - ▶ **Induction step:**

$$\begin{aligned} J_t(x; \mathbf{g}) &= \mathbb{E}^{\mathbf{g}} \left[ \sum_{s=t}^T c_s(X_s, \mathbf{u}_s) \mid X_t = x \right] \\ &= \mathbb{E}^{\mathbf{g}} \left[ c_t(x, g_t(x)) + \mathbb{E}^{\mathbf{g}} \left[ \sum_{s=t+1}^T c_s(X_s, \mathbf{u}_s) \mid X_{t+1} \right] \mid X_t = x \right] \\ &= \mathbb{E}^{\mathbf{g}} \left[ c_t(x, g_t(x)) + J_{t+1}(X_{t+1}; \mathbf{g}) \mid X_t = x \right] \end{aligned}$$

By the induction hypothesis

$$\begin{aligned} &\geq \mathbb{E}^{\mathbf{g}} \left[ c_t(x, g_t(x)) + V_{t+1}(X_{t+1}) \mid X_t = x, \mathbf{u}_t = g_t(x) \right] \\ &\geq V_t(x) \end{aligned}$$

with equality iff

- ▶ first inequality:  $g_{t+1:T}$  satisfies verification step (induction hypothesis)
- ▶ second inequality:  $g_t \in \arg \min_{\mathbf{u} \in \mathcal{U}(x)} Q_t(x, \mathbf{u})$ .



# Generalization of the basic model

**Per-step cost** Both the structural result and the dynamic programming decomposition remain valid when the per-step cost is given by

$$c_t(X_t, U_t, X_{t+1})$$

**Proof** Both the results only rely on  $\mathbb{E}[\text{per-step cost} \mid X_t, U_t]$  being independent of the control strategy.

When the per-step cost is given as above,  $\mathbb{E}[c_t(X_t, U_t, X_{t+1})]$ , is independent of the control strategy.

# MDP Theory: Probabilistic models

To be written

# Stochastic orders

# Stochastic dominance

- Notation**
- ▶  $\mathcal{X} = \{1, \dots, n\}$  and  $\mathcal{Y} = \{1, \dots, m\}$  are finite spaces.
  - ▶  $\Delta(\mathcal{X})$  is the space of probability measures (PMFs) over  $\mathcal{X}$ .

**Definition** For any  $\pi, \mu \in \Delta(\mathcal{X})$ ,  $\pi$  **stochastically dominates**  $\mu$  (denoted by  $\pi \geq_s \mu$ ) if

$$\sum_{i \geq k} \pi_i \geq \sum_{i \geq k} \mu_i, \quad \forall k.$$

Equivalently, if  $X_1 \sim \pi$  and  $X_2 \sim \mu$ , then  $\pi \geq_s \mu$  iff

$$\mathbb{P}(X_1 \geq x) \geq \mathbb{P}(X_2 \geq x), \quad \forall x \in \mathcal{X}.$$

**Example**

$$\left[ 0 \quad \frac{1}{4} \quad \frac{1}{4} \quad \frac{1}{2} \right] \geq_s \left[ \frac{1}{4} \quad 0 \quad \frac{1}{4} \quad \frac{1}{2} \right] \geq_s \left[ \frac{1}{4} \quad \frac{1}{4} \quad \frac{1}{4} \quad \frac{1}{4} \right]$$

# Stochastic dominance preserves monotonicity

**Lemma** Let  $\{v_i\}_{i=1}^n$  be an increasing sequence and  $\pi \geq_s \mu$ . Then,

$$\sum_{i=1}^n \pi_i v_i \geq \sum_{i=1}^n \mu_i v_i$$

Equivalently, if  $X_1 \sim \pi$ ,  $X_2 \sim \mu$ , and  $f : \mathcal{X} \rightarrow \mathbb{R}$  is an **increasing** function, then  $\pi \geq_s \mu$  implies

$$\mathbb{E}[f(X_1)] \geq \mathbb{E}[f(X_2)]$$

**Proof** Define  $v_{-1} = 0$ . Consider

$$\begin{aligned} \sum_{i=1}^{\infty} \pi_i v_i &= \sum_{i=1}^{\infty} \pi_i \sum_{j=1}^{\infty} (v_j - v_{j-1}) \\ &= \sum_{j=1}^{\infty} (v_j - v_{j-1}) \sum_{i=j}^{\infty} \pi_i \\ &\geq \sum_{j=1}^{\infty} (v_j - v_{j-1}) \sum_{i=j}^{\infty} \mu_i = \sum_{i=1}^{\infty} \mu_i v_i \end{aligned}$$

# Stochastic monotone Markov chains

**Definition** Let  $\{X_t\}_{t=1}^{\infty}$  be a time-homogeneous Markov chain with transition matrix  $P$ . The Markov chain is **stochastically monotone** if

$$P_i \geq_s P_j, \quad \forall i > j$$

where  $P_i$  denotes the row- $i$  of  $P$ .

**Implication** If  $\{X_t\}_{t=1}^{\infty}$  is stochastically monotone and  $f : \mathcal{X} \rightarrow \mathbb{R}$  is an increasing function, then

$$\mathbb{E}[f(X_{t+1}) \mid X_t = x_1] \geq \mathbb{E}[f(X_{t+1}) \mid X_t = x_2], \quad \forall x_1 > x_2.$$

# Monotone likelihood ratio (MLR) ordering

**Definition** For any  $\pi, \mu \in \Delta(\mathcal{X})$ ,  $\pi$  dominates  $\mu$  in monotone likelihood ratio (denoted by  $\pi \geq_r \mu$ ) if

$$\pi_i \mu_j \geq \mu_i \pi_j, \quad \forall i > j; \quad \text{if } \mu_i, \mu_j > 0, \text{ then } \frac{\pi_i}{\mu_i} \geq \frac{\pi_j}{\mu_j}$$

- Examples**
- ▶  $\begin{bmatrix} \frac{1}{8} & \frac{1}{8} & \frac{1}{4} & \frac{1}{2} \end{bmatrix} \geq_r \begin{bmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{bmatrix}$ .
  - ▶  $\begin{bmatrix} 0 & \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \end{bmatrix} \not\geq_r \begin{bmatrix} \frac{1}{4} & 0 & \frac{1}{4} & \frac{1}{2} \end{bmatrix}$ .

# Monotone likelihood ratio implies stochastic dominance

**Proposition** For any  $\pi, \mu \in \Delta(\mathcal{X})$ ,

$$\pi \geq_r \mu \implies \pi \geq_s \mu$$

**Proof Outline** Define

$$P_k = \sum_{i \leq k} \pi_i \quad \text{and} \quad M_k = \sum_{i \leq k} \mu_i.$$

▶ Show that

$$\frac{P_k}{M_k} \leq \frac{\pi_k}{\mu_k} \leq \frac{1 - P_k}{1 - M_k}$$

▶ Using the above relation, show that

$$\pi \geq_r \mu \implies \pi \geq_s \mu$$

**Example that the reverse implication is not true**

- ▶ Let  $\pi = [0.2 \quad 0.5 \quad 0.3]$  and  $\mu = [0.1 \quad 0.1 \quad 0.8]$ .
- ▶  $\pi \geq_s \mu$
- ▶  $\pi \not\geq_r \mu$ .



# Total positivity of order 2 (TP<sub>2</sub>) and preserving MLR

**Definition**  
(Totally positive of order 2)

- Recall for any matrix **A** and any index sets **I** and **J**
- ▶ **A**<sub>I,J</sub> denotes the submatrix corresponding to the row set **I** and the column set **J**;
  - ▶ The **(I, J) minor** of **A** is  $\det \mathbf{A}_{I,J}$ .

$$\begin{bmatrix} 4 & 3 & 2 & 1 \\ 5 & 4 & 3 & 2 \\ 6 & 5 & 4 & 3 \\ 7 & 6 & 5 & 4 \end{bmatrix} \text{ is TP}_2.$$

A  $n \times m$  matrix is **totally positive of order 2 (TP<sub>2</sub>)** if all its  $2 \times 2$  submatrices have non-negative determinant.

**Proposition**

- If **Q** is a row stochastic matrix that is TP<sub>2</sub>, and  $\pi, \mu \in \Delta(\mathcal{X})$  then
- ▶  $Q_i \geq_r Q_j$  for  $i > j$ . Consequently,  $Q_i \geq_s Q_j$ .
  - ▶  $\pi \geq_r \mu \implies \pi Q \geq_r \mu Q$

**Proof**

- ▶ Let  $i > j$  and  $k > \ell$ . Since **Q** is TP<sub>2</sub>, the minor consisting of rows  $i, j$  and columns  $k, \ell, i > j$  is non-negative. Thus,

$$\begin{vmatrix} Q_{j\ell} & Q_{jk} \\ Q_{i\ell} & Q_{ik} \end{vmatrix} \geq 0 \implies Q_{ik}Q_{j\ell} \geq Q_{jk}Q_{i\ell} \implies Q_i \geq_r Q_j$$

- ▶ See Proposition on next page.

# TP<sub>2</sub> ordering of functions and matrices

**Definition**  
(TP<sub>2</sub> ordering)

A function  $f \geq_{tp} g$  if  $\forall x_1, x_2, y_1, y_2$

$$f(x_1 \vee x_2, y_1 \vee y_2)g(x_1 \wedge x_2, y_1 \wedge y_2) \geq f(x_1, y_1)g(x_2, y_2),$$

Note that  $a \vee b = \max(a, b)$  and  $a \wedge b = \min(a, b)$ .

- ▶ This definition extends to matrices in a natural manner.
- ▶ A matrix  $Q$  is TP<sub>2</sub> if  $Q \geq_{tp} Q$ .

**Proposition**

If  $P_1$  and  $P_2$  are row stochastic matrices such that  $P_1 \geq_{tp} P_2$ , then

$$\pi \geq_r \mu \implies \pi P_1 \geq_r \mu P_2$$

In particular,

$$\pi P_1 \geq_r \pi P_2, \quad \forall \pi$$

**Proof**

See Theorem 2.4 of Samuel Karlin and Yosef Rinott, "Classes of orderings of measures and related correlation inequalities. I. Multivariate totally positive distributions," Journal of Multivariate Analysis, vol 10, no 4 Pages 467-498, Dec 1980. [http://dx.doi.org/10.1016/0047-259X\(80\)90065-2](http://dx.doi.org/10.1016/0047-259X(80)90065-2)



# Optimal monotone strategies in MDPs

# Submodularity and supermodularity

**Definition** Let  $\mathcal{X}$  and  $\mathcal{Y}$  be partially ordered sets. A function  $f: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  is called **submodular** if for any  $x^+ \geq x^-$  and  $y^+ \geq y^-$

$$f(x^+, y^+) + f(x^-, y^-) \leq f(x^+, y^-) + f(x^-, y^+) \quad (*)$$

The function is called **supermodular** if the inequality in (\*) is reversed.

**Examples** ▶  $f(x, y) = xy$  is supermodular.

**Equivalent definition** A continuous and differentiable function is submodular iff

$$\frac{\partial^2 f(x, y)}{\partial x \partial y} \leq 0, \quad \forall x, y$$

# Submodularity and monotonicity of the arg min

**Theorem** Let  $f: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  be a submodular function and assume that for all  $x$ ,  $\arg \min_{y \in \mathcal{Y}} f(x, y)$  exists. Define,

$$g(x) = \max \{y' \in \arg \min_{y \in \mathcal{Y}} f(x, y)\}$$

Then,  $g(x)$  is (weakly) increasing in  $x$ .

**Proof** Let  $x^+ \geq x^-$ . Suppose we can show that for all  $y \leq g(x^-)$ ,

$$f(x^+, g(x^-)) \leq f(x^+, y) \quad (*)$$

Then,  $g(x^+) \geq g(x^-)$ , which concludes the proof.

To prove  $(*)$ , note that since  $f$  is submodular,

$$f(x^+, y) - f(x^+, g(x^-)) \geq f(x^-, y) - f(x^-, g(x^-)) \geq 0$$

where the last inequality follows because  $g(x^-)$  is an arg min of  $f(x^-, y)$ .

# Monotonicity of the value function

**Theorem** Suppose that the arg min at each step of the dynamic program is attained and that

C1.  $c_t(x, u)$  is (weakly) increasing in  $x$  for all  $u \in \mathcal{U}$ .

C2.  $P(u)$  is stochastic monotone for all  $u \in \mathcal{U}$ .

Then,  $V_t(x)$  is weakly increasing in  $x$

**Proof** We proceed by backward induction.

▶ **Basis:**  $V_{T+1}(x)$  is weakly increasing in  $x$ .

▶ **Hypothesis:** Assume that  $V_{t+1}(x)$  is weakly increasing in  $x$ .

▶ **Induction step:** Consider  $x' \leq x$  and let  $u^*$  be the optimal action at  $x$ .

$$\begin{aligned} \text{Thus, } V_t(x) &= Q_t(x, u^*) \\ &= c_t(x, u^*) + \mathbb{E}[V_{t+1}(X_{t+1} \mid X_t = x, U_t = u^*)] \\ &\stackrel{(a)}{\geq} c_t(x', u^*) + \mathbb{E}[V_{t+1}(X_{t+1} \mid X_t = x', U_t = u^*)] \\ &= Q_t(x', u^*) \geq \min_{u \in \mathcal{U}} Q_t(x', u) = V_t(x') \end{aligned}$$

where (a) follows from the assumptions in the theorem.

# Monotonicity of optimal strategy

**Theorem** Suppose that the arg min at each step of the dynamic program is attained and that

C1.  $c_t(x, u)$  is (weakly) increasing in  $x$  for all  $u \in \mathcal{U}$ .

C2.  $P(u)$  is stochastic monotone for all  $u \in \mathcal{U}$ .

C3. For any increasing function  $v: \mathcal{X} \rightarrow \mathbb{R}$ , the function

$$c_t(x, u) + \mathbb{E}[v(X_{t+1}) \mid X_t = x, U_t = u] \quad (*)$$

is submodular in  $(x, u)$  for all  $t$ .

Then,  $g_t^*(x)$  is weakly increasing in  $x$ .

- Proof**
- ▶ Conditions C1 and C2 imply that  $V_{t+1}(x)$  is increasing in  $x$ .
  - ▶ Condition C3 implies that  $Q_t(x, u)$  is submodular in  $(x, u)$ .
  - ▶ Therefore, the arg min  $g_t^*(x)$  is increasing in  $x$ .

# Sufficient condition for C3

**Theorem** Suppose conditions C1 and C2 of the previous theorem are satisfied. In addition,

C3a.  $c_t(x, u)$  is submodular in  $(x, u)$ .

C3b.  $q(y | x, u) := \sum_{x' \geq y} P(x' | x, u)$  is submodular in  $(x, u)$  for all  $y \in \mathcal{X}$ .

Then,  $g_t^*(x)$  is weakly increasing in  $x$ .

**Proof** We will show that C3a and C3b imply C3 of the previous theorem.

Since  $q(y|x, u)$  is submodular,

$$q(y | x^-, u^-) + q(y | x^+, u^+) \leq q(y | x^-, u^-) + q(y | x^+, u^+)$$

$$\implies \sum_{x' \geq y} [P(y | x^-, u^-) + P(y | x^+, u^+)] \leq \sum_{x' \geq y} [P(x' | x^-, u^-) + P(x' | x^+, u^+)]$$

Consider two measures  $\pi$  and  $\mu$  where

$$\pi(x) = 0.5P(x | x^-, u^-) + 0.5P(x | x^+, u^+)$$

$$\mu(x) = 0.5P(x | x^-, u^+) + 0.5P(x | x^+, u^-)$$

Then, the above equation implies that  $\pi \leq_s \mu$ . Therefore, for any



increasing function  $v: \mathcal{X} \rightarrow \mathbb{R}$ ,

$$\sum_{x' \in \mathcal{X}} \pi(x') v(x) \leq \sum_{x' \in \mathcal{X}} \mu(x') v(x)$$

or, equivalently,

$$H(x^-, u^-) + H(x^+, u^+) \leq H(x^-, u^+) + H(x^+, u^-)$$

where  $H(x, u) = \mathbb{E}[v(X_{t+1}) \mid X_t = x, U_t = u]$ .

Therefore,  $c_t(x, u) + H(x, u)$  is submodular.

# Constraint on actions

- Note** The results on monotonicity of the value function and the optimal strategy remain valid if  $\mathcal{U}$  depends on  $x$  provided:
- ▶  $\mathcal{U}(x) \subseteq \mathcal{U}(x')$  for all  $x' \geq x$ .
  - ▶ For any  $x \in \mathcal{X}$ ,  $u, u' \in \mathcal{U}$  such that  $u' \leq u$ , if  $u \in \mathcal{U}(x)$  then  $u' \in \mathcal{U}(x)$ .

# Optimal threshold strategies in optimal stopping problems

# Some definitions

**Benefit function**  $B_t(x) = \mathbb{E}[V_{t+1}(X_{t+1} | X_t = x) + c_t(x) - s_t(x)]$ .

Note that the value function may be written as

$$V_t = \min\{B_t(x) + s_t(x), s_t(x)\} = s_t(x) + \min\{B_t(x), 0\}.$$

Therefore, it is optimal to stop when  $B_t(x) \geq 0$ .

**One-step benefit function**  $M_t(x) = \mathbb{E}[s_{t+1}(X_{t+1} | X_t = x) + c_t(x) - s_t(x)]$ .

The benefit function and the one-step benefit function are closely related:

$$B_T(x) = M_T(x)$$

and

$$\begin{aligned} B_t(x) &= \mathbb{E}[s_{t+1}(X_{t+1} + \min\{B_{t+1}(X_{t+1}), 0\} | X_t = x) + c_t(x) - s_t(x)] \\ &= M_t(x) + \mathbb{E}[\min\{B_{t+1}(X_{t+1}), 0\} | X_t = x]. \end{aligned}$$

# Optimality of threshold strategies

**Theorem** Suppose the following conditions hold:

- S1.  $M_t(x)$  is (weakly) increasing in  $x$  for all  $t$ .
- S2.  $\{X_t\}_{t \geq 1}$  is stochastic monotone.

Then,  $B_t(x)$  is (weakly) increasing in  $x$  for all  $t$  and there exists a sequence  $\{\lambda_t\}_{t \geq 1}$  such that it is optimal to stop at time  $t$  if  $X_t \geq \lambda_t$ .

**Proof** Proceed by backward induction:

- ▶ **Basis:**  $B_T(x) = M_T(x)$  is increasing in  $x$ .
- ▶ **Hypothesis:**  $B_{t+1}(x)$  is increasing in  $x$ .
- ▶ **Induction:** By induction hypothesis,  $\min\{B_{t+1}(x), 0\}$  is increasing in  $x$ .  
Due to (S2),  $\mathbb{E}[\min\{B_{t+1}(X_{t+1}), 0\} | X_t = x]$  is increasing in  $x$ .  
Therefore,  $B_t(x) = M_t(x) + \mathbb{E}[\min\{B_{t+1}(X_{t+1}), 0\} | X_t = x]$  is increasing in  $x$ .

Recall that it is optimal to stop if  $B_t(x) \geq 0$ . Hence, the optimal decision rule is of a threshold type.

# POMDP Theory: Probabilistic models

To be written

**MDP Theory: Functional models**

**MDP Theory: Perfect state observation**

**MDP Theory: Probabilistic models**

To be written

**Stochastic orders**

**Optimal monotone strategies in MDPs**