

Last class

Stochastic optimization . Newsvendor.

Info is available. State, Policies

$$\begin{aligned} \min_{\pi: \mathcal{S} \rightarrow \mathcal{A}} \mathbb{E}[c(\mathcal{S}, \pi(\mathcal{S}), \mathcal{W})] \\ \equiv \pi(\mathcal{S}) = \arg \min_{a \in \mathcal{A}} \mathbb{E}[c(\mathcal{S}, a, \mathcal{W}) | \mathcal{S} = \mathcal{S}], \quad \forall \mathcal{S} \in \mathcal{S} \end{aligned}$$

Today

- Sufficient / Irrelevant info
- Formulations / Results - Dynamic Programming + MDPs
- Example: Inventory management

IRRELEVANT INFO

$$c: S \times A \times W \rightarrow \mathbb{R}_{\geq 0}$$

(Ω, \mathcal{F}, P) Prob Space.

(S, Y, W) RVs on Prob space.

DM obs. (S, Y) and choose A .

$$A = \pi(S, Y).$$

When can the DM ignore Y .

Blackwell's Principle of Irrelevant Info (B2/G4)

$$\text{If } Y \perp\!\!\!\perp W \mid S \Rightarrow P(W|S, Y) = P(W|S)$$

Y & W are cond. indep given S

then without loss of opt.

$$\pi^0(s) = \arg \min_{a \in A} \mathbb{E}[c(s, a, W) \mid S=s] \quad \forall s \in S$$

is opt

$$\text{cf } \pi^*(s, y) = \arg \min_{a \in A} \mathbb{E}[c(s, a, W) \mid S=s, Y=y] \quad \forall s, y.$$

$$\mathbb{E}[c \mid S=s, Y=y]$$

$$= \sum_w P(W=w \mid S=s, Y=y) c(s, a, w)$$

$$= \sum_{\omega} P(W=\omega | S=s) c(s, a, \omega).$$

$$= \mathbb{E}[c(s, a, W) | S=s]$$

$$\mathbb{E}[c(s, \pi(s), W) | S=s] \leq \mathbb{E}[c(s, a, W) | S=s]$$

$\forall a \in A$

For any $\pi: S \times Y \rightarrow A$

$$\leq \mathbb{E}[c(s, \pi(s, y), W) | S=s]$$

$\forall y \in Y.$

$$= \mathbb{E}[c(s, \pi(s, y), W) \Big|_{\substack{S=s \\ Y=y}}]$$

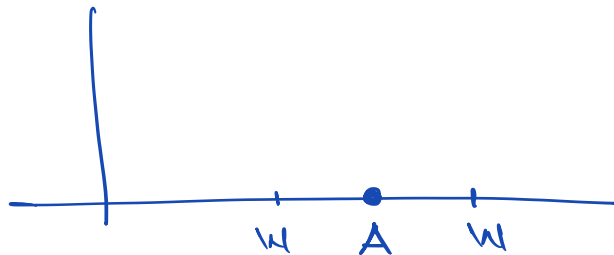
MARKOV DECISION PROCESSES

News vendor

$$S_1 = 0$$

Purchase A_1 at price per unit p .

Demand W_1



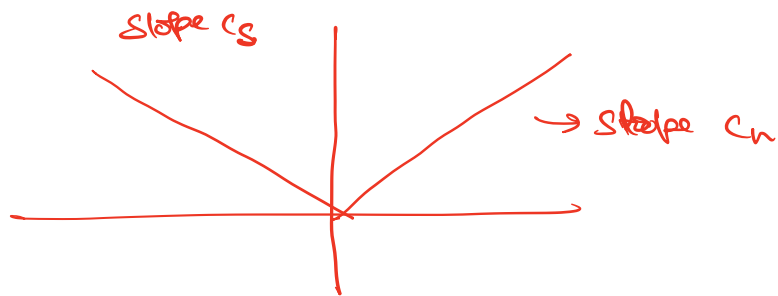
$$S_2 = S_1 + A_1 - W_1$$

Costs

$$c(s) = \begin{cases} c_h s & s \geq 0 \\ -c_s s & s < 0 \end{cases}$$

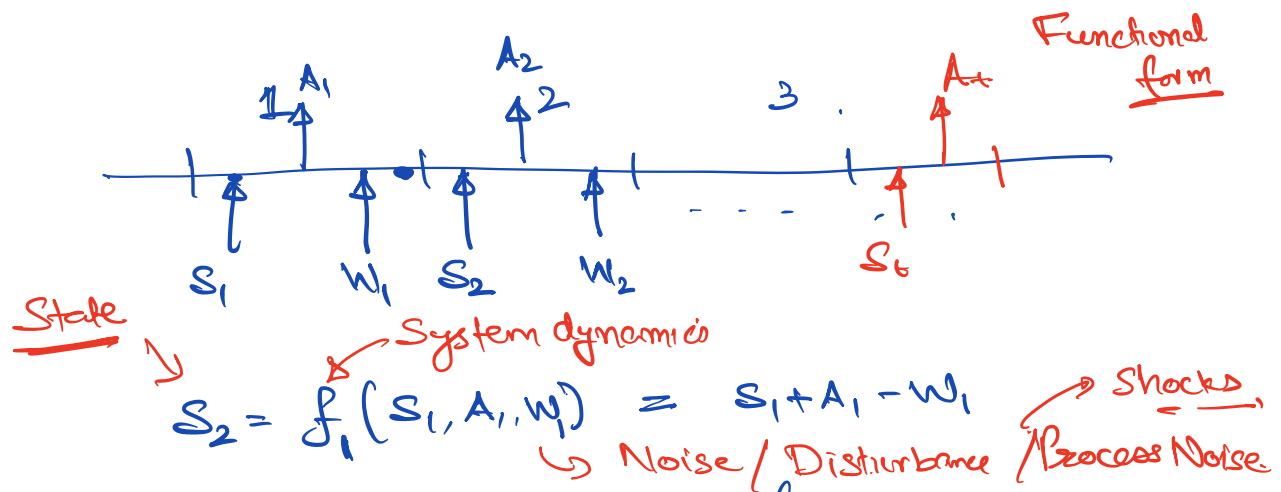
Holding

Shortage



~~$S_1 = 0$~~

$$\pi_1(s) = \arg \min_{a \in A} \mathbb{E}_{W_1} [c(s_1 + a - W_1)]$$



$$S_2 = f_1(S_1, A_1, W_1) = S_1 + A_1 - W_1$$

Primitive Random Variable (S_1, W_1, \dots, W_T)

all defined on a common prob space.

$$A_1 = \pi_1(S_1)$$

$$A_2 = \pi_2(S_1, S_2, A_1)$$

$$A_3 = \pi_3(S_1, S_2, S_3, A_1, A_2)$$

$$A_t = \pi_t(S_{1:t}, A_{1:t-1})$$

$\pi = (\pi_1, \pi_2, \dots, \pi_T)$ Policy

Performance.

$$J(\pi) = \mathbb{E} \left[\sum_{t=1}^T \underbrace{[pA_t + c(S_t, A_t)]}_{\text{Per-step cost}} \right]$$

Per-step cost $\rightarrow c(S_t, A_t)$
or $c(S_t, A_t, S_{t+1})$

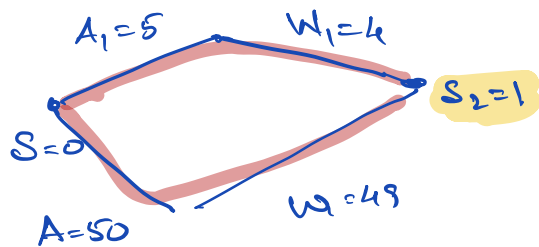
$$\min_{\pi} J(\pi)$$

Assump: All primitive random variables are indep

1 Optimality of Markov policies

WLO,

$$A_t = \pi_t(\underline{S}_t)$$



Controlled Markov Prop

Matrix form

$$P(S_{t+1} | S_{1:t}, A_{1:t}, A_t)$$

$$= P(S_{t+1} | S_t, A_t)$$

$$\mathbb{E} \quad A_1 = \pi_1(S_1) \quad |S| |A|$$

$$A_2 = \pi_2(S_2) \quad |S| |A|$$

$$A_3 = \pi_3(S_3) \quad |S| |A|$$

$$(\pi_1, \dots, \pi_T)^T = (|S| |A|)^T$$

2. Dynamic Prog. Decomp.



$$V_t: S \rightarrow \mathbb{R}$$

$$Q_t: S \times A \rightarrow \mathbb{R}$$

$$\pi_t: S \rightarrow A$$

Init $V_{T+1}(s) = 0 \quad \forall s \in S.$

Recursive for $t \in \{T, T-1, \dots, 1\}$

$$Q_t(s, a) = c_t(s, a) + E[V_{t+1}(f_t(s, a, w_t))]$$

$\underbrace{f_t(s, a, w_t)}_{S_{t+1}}$

$$V_t(s) = \min_{a \in A} Q_t(s, a)$$

$$\pi_t(s) \in \arg \min_{a \in A} Q_t(s, a)$$

Optimal Value fn (with arrow pointing to $V_t(s)$)

$$T \left(\underline{|S||A|} |w| + |S||A| \right) \quad O(T |S||A|)$$

$$Q_t(s, a) = c_t(s, a) + \sum_{\omega} P_{w_t}(\omega) V_{t+1}(f_t(s, a, \omega))$$

1) How ~~do~~ to compute efficiently

2) Can we do something beyond comp

