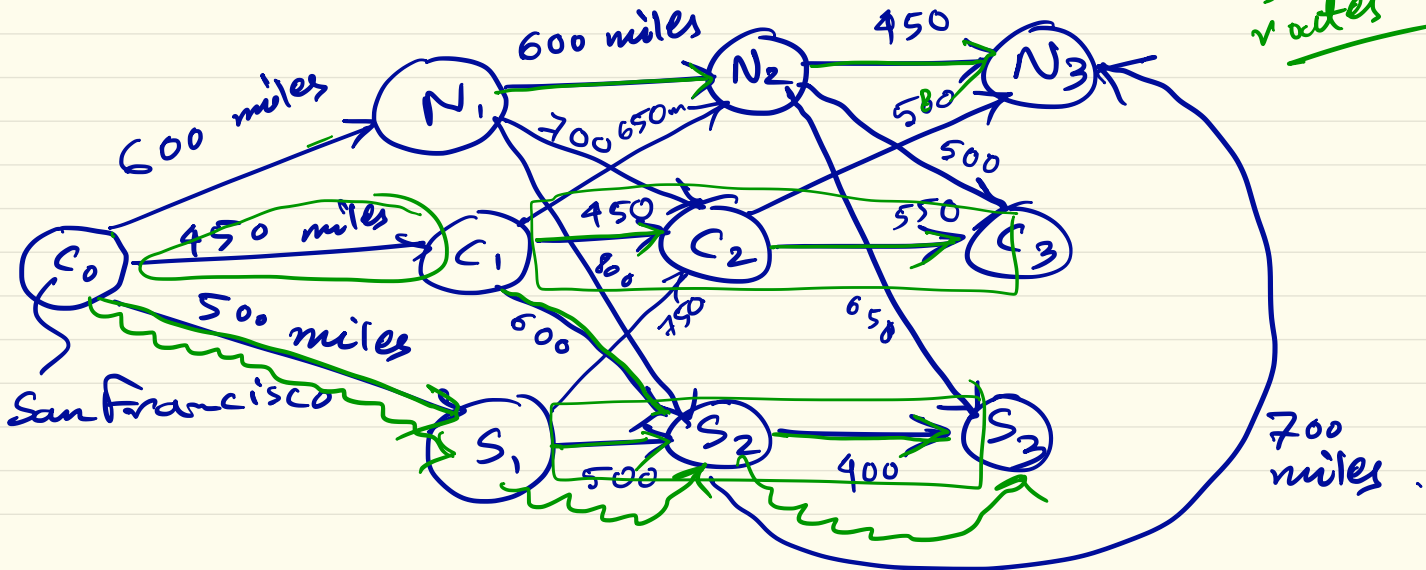# Lecture #15

Dynamic Programming (We will start with discrete time)
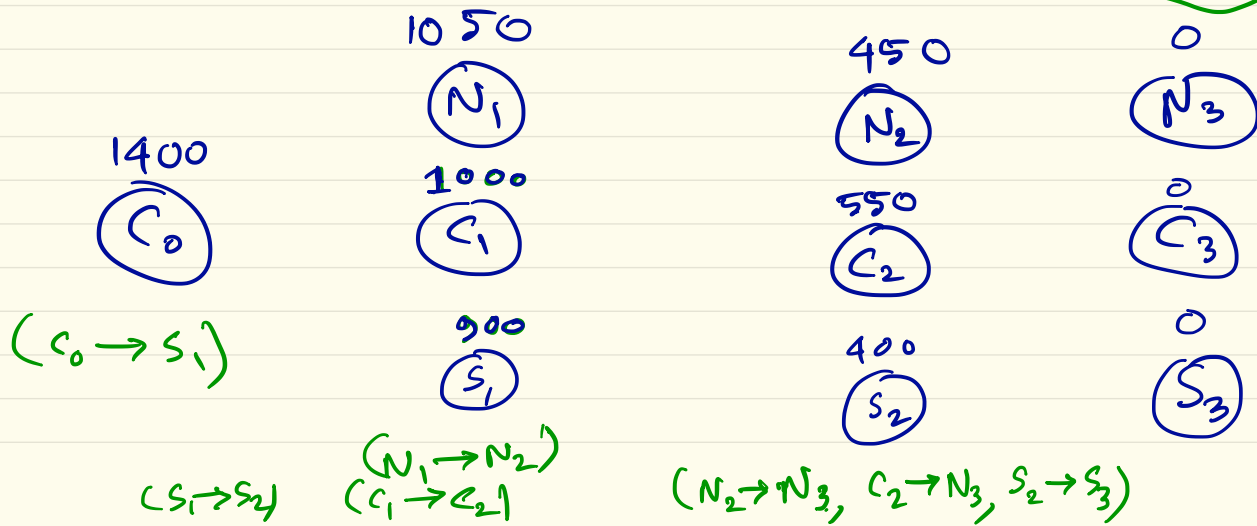
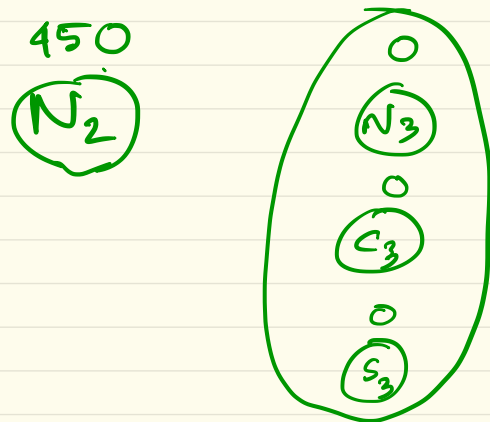## Network of roads in US

Want to find the shortest drive from SF to Chicago

all 1 way routes



600 miles

600 miles — $N_1$ — 600 miles — $N_2$ — 450 — $N_3$

700 650m

580

500

450 miles — $C_0$ — 450 — $C_1$ — 450 — $C_2$ — 550 — $C_3$

100

150

500 miles

600

650

San Francisco

500

$S_1$ — 500 — $S_2$ — 400 — $S_3$

700 miles

- from $\widehat{N_2}$, shortest path to East Cost is 950 miles.

450
$\widehat{N_2}$

$\circ$
$\widehat{N_3}$
$\circ$
$\widehat{C_3}$
$\circ$
$\widehat{S_3}$

$\circ$

- 

1050
$\widehat{N_1}$

1000
$\widehat{C_1}$

900
$\widehat{S_1}$

1400
$\widehat{C_0}$

$(C_0 \to S_1)$

450
$\widehat{N_2}$

550
$\widehat{C_2}$

400
$\widehat{S_2}$

$\widehat{N_3}$

$\circ$
$\widehat{C_3}$

$\circ$
$\widehat{S_3}$

$(S_1 \to S_2)$  $(N_1 \to N_2)$
$(C_1 \to C_2)$

$(N_2 \to N_3, \ C_2 \to N_3, \ S_2 \to S_3)$

# Observations:

(1) In order to solve the problem for one initial condition, we had to solve for all starting states/vertices/cities.

↓

This procedure is called Dynamic Programming (DP)

↓

Hence computationally (exponential complexity) difficult.

(2) DP gives you a "closed-loop-policy" / "feedback policy"

⟺ actions (controls) as $f^n$ of state

[If ever I find myself in Denver then go south].

You're in now

③ Different from Open-loop policy (time-table)
(close your eyes, drive 3 hours east, then take left turn)

④ DP gives <u>closed-loop policy.</u>

⑤ DP proceeds Backward in time
(⟸ "Backward Recursion")

Solve for   t   days   remaining
                ↓

Then (t+1) days remaining   etc.
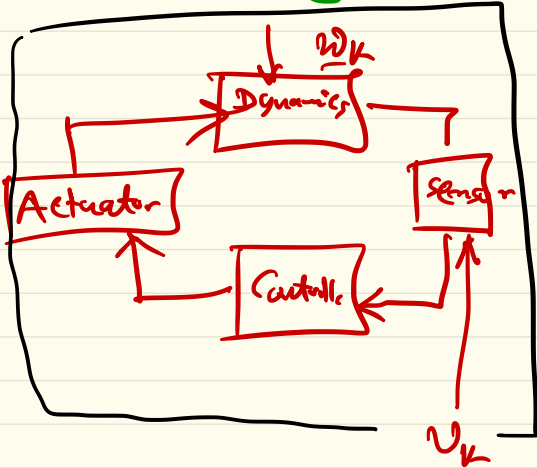
⑥ Segments / subarcs of optimal path
are   themselves optimal
(⟺ Bellman's Principle of Optimality)

(7) Recursion:

Optimal cost with $(t+1)$ days remaining $= \min\limits_{\{\text{Set of actions}\}}$ $\begin{bmatrix} \text{Immediate} \\ \text{Cost of} \\ \text{my action} \end{bmatrix} + \begin{bmatrix} \text{Optimal} \\ \text{Cost} \\ \text{from where} \\ \text{that Cost} \\ \text{takes you} \end{bmatrix}$



To fix ideas, consider discrete time.

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

$$\underline{y}_k = \underline{h}_k(\underline{x}_k, \underline{u}_k, \underline{v}_k)$$

- $\underline{x}_k \in \mathcal{X} \underset{\subset \mathbb{R}^n}{} (\text{state space})$, $\underline{u}_k \in \mathcal{U} (\text{Control space})$

- For each $k = 0, 1, 2, \ldots,$ the control values $u_k \in \mathcal{U} \subseteq \mathbb{R}^m$.

A feasible control (law) is a sequence of (policies).

Remember: $\boxed{\text{Policy/law/feedback} \neq \text{Action/control}}$

Policy/law/feedback:

$$\underline{\gamma} = \{ \underline{\gamma}_0, \underline{\gamma}_1, \underline{\gamma}_2, \dots \} \quad s.t. \quad \underline{u}_k = \underline{\gamma}_k (\underline{y}_k)$$

$\underline{\gamma}_0$ — Feedback @ time 0

$\underline{\gamma}_1$ — Feedback @ time 1

$$\in \mathcal{U}$$

$$\gamma(k, y_k)$$

- Let $\Gamma$ be the set of all possible policices:

$$\boxed{u(t) = \underline{\gamma}(\underline{x}(t), t)}$$

- We wish to find the best $\underline{\gamma}$ in $\Gamma$.

∴ We need critenium to compare different policies

we associate cost for each policy, and declare the best one is the one that minimizes cost.
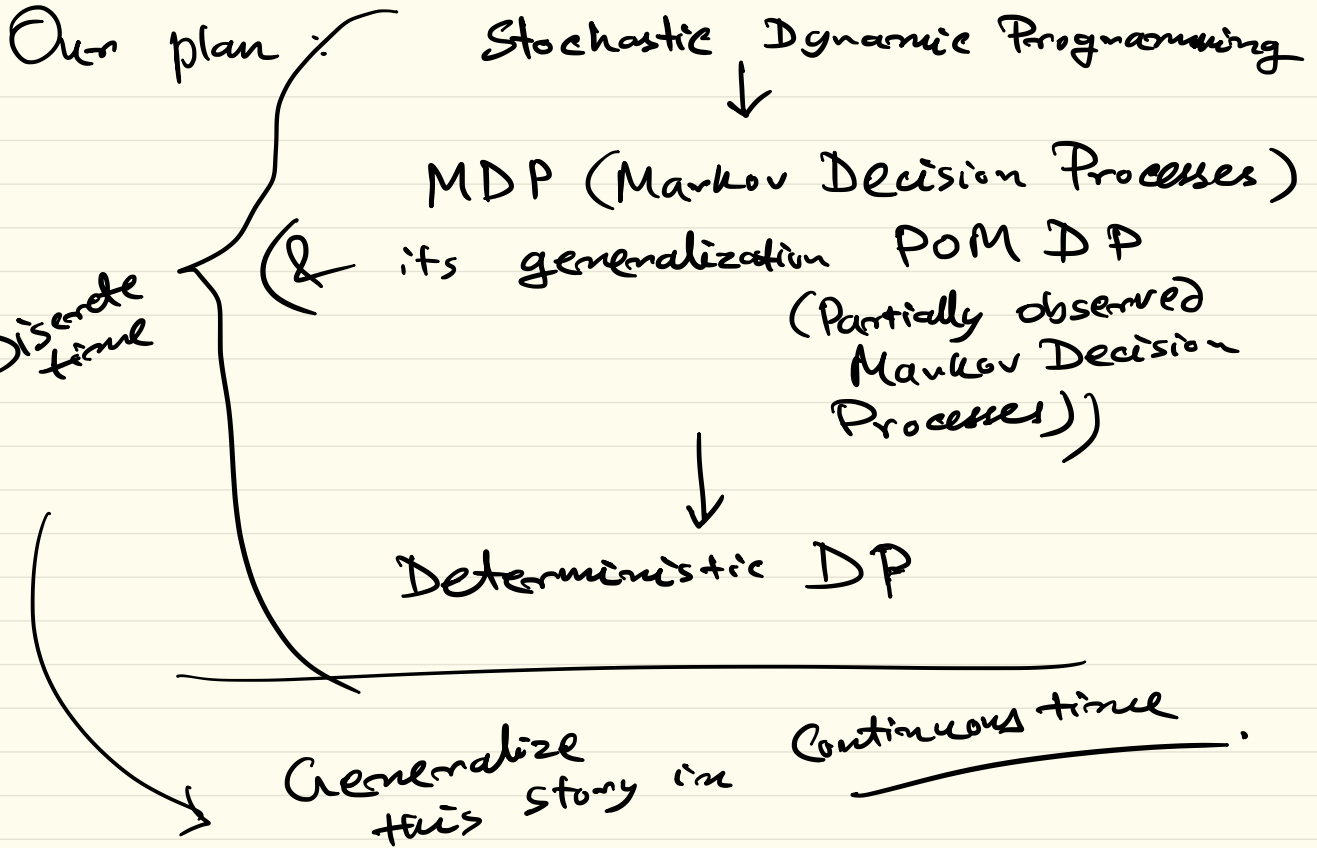
- Our cost function:

$$J(\underline{g}) := C_T \; (\underline{x}(T)) + \sum_{k=0}^{T-1} C_k(\underline{x}_k, \underline{u}_k)$$

$$\underset{\|}{\phantom{=}}$$

$$\phi(x(T), T)$$

(terminal cost)

$$\underset{\|\|\|}{\phantom{=}} L(k, \underline{x}_k, \underline{u}_k)$$

(Lagrangian)

- Finite horizon: $T < \infty$

$$\Updownarrow$$

Control law is a finite policy sequence:

$$\underline{g} = \{\underline{g}_0, \underline{g}_1, \cdots, \underline{g}_{T-1}\}$$

- The term $C_k(\underline{x}_k, \underline{u}_k)$ is called immediate/ one period cost.

Our plan: Stochastic Dynamic Programming

$\downarrow$

MDP (Markov Decision Processes)

(& its generalization POMDP
(Partially observed
Markov Decision
Processes))

Discrete time

$\downarrow$

Deterministic DP

Generalize this story in Continuous time.

# Stochastic DP

Deterministic DP is special case: $\underline{w}_k \equiv \underline{0}$, $\underline{v}_k \equiv \underline{0}$

$\underbrace{\underline{w}_k \equiv \underline{0}}_{\text{Process noise}}$, $\underline{v}_k \equiv \underline{0}$ ↑ measurement noise

- Stochastic DP:

$$\underline{w}_k \in \mathcal{W} ; \quad \underline{v}_k \in \mathcal{V} ;$$

↑ random vectors realized from (Discrete time stochastic process)

$$\mathcal{W} := (\mathbb{P}_w, \Omega_w, F_w)$$
$$\mathcal{V} := (\mathbb{P}_v, \Omega_v, F_v)$$

Then $\underline{x}_k, \underline{u}_k$ are random vectors, and hence $J(\underline{8})$ is a random variable

(i.e.) $J(\underline{8}) \equiv J(\underline{\omega}, \underline{8})$, $\underline{\omega} \in \Omega_w \times \Omega_v$

$\underbrace{\underline{\omega}}_{\text{sample path index}}$

To resolve sample path dependency, we take

$$J(\gamma) = E[J(\underline{\gamma})]$$

Let $\underline{\gamma}^* = \arg\min_{\underline{\gamma} \in \Gamma} E[J(\underline{\gamma})]$

and $\underline{J^*} = \min_{\underline{\gamma} \in \Gamma} E[J(\underline{\gamma})]$.

Deterministic
scalar $\geq 0$

We say $\underline{\gamma}^*$ is optimal policy $\vec{J}^*$ is
optimal cost.

We will now focus on : MDP
(Markov Decision Process)

Complete information/fully observed case :

$$\underline{y}_k \equiv \underline{x}_k,$$

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

Let $\underline{u}_k = \underline{g}_k(\underline{x}_0, \underline{x}_1, \ldots, \underline{x}_k)$
is allowed to depend on
previous states,

$\underline{x}_k \in \mathcal{X} \subset \mathbb{R}^n$
$\underline{u}_k \in \mathcal{U} \subset \mathbb{R}^m$
$\underline{w}_k \in \mathcal{W} \subset \mathbb{R}^p$

(i.e) $\underline{g}_k$ is history-dependent
policy.

More generally, History upto time $t =: H_t$
$:= \{\underline{x}_0, \underline{u}_0, \underline{x}_1, \underline{u}_1, \ldots, \underline{x}_{k-1}, \underline{u}_{k-1}, \underline{x}_k\}$

At each $k$, $\gamma_k(H_k) = \underline{u}_k$

$$\gamma_k : H_k \longmapsto \mathcal{U}$$

$\therefore \underline{\gamma} = (\underline{\gamma}_0, \underline{\gamma}_1, \ldots, \gamma_{T-1})$ is called history-dependent policy

$\rightarrow H_k$ is history up until time $k$.

## History Dependent Policies

Randomized

$\gamma_k : H_k \longmapsto$ Prob. over $\mathcal{U}$
(Choose $u_k$ as a sample from that Probabilits)

Non-randomized

$\gamma_k : H_k \longmapsto \underline{\mathcal{U}}$
$\quad\quad\quad\quad\quad$ returns
$\quad\quad\quad$ $u_k$ (particular action)

**Detour:** Markov process:

$$\mathbb{P}(\text{future} \mid \text{Past \& Present})$$

$$= \mathbb{P}(\text{Future} \mid \text{Present})$$

Another way to write:

$$\mathbb{P}(\text{Past \& Future} \mid \text{Present})$$

$$= \mathbb{P}(\text{Past} \mid \text{Present}) \, \mathbb{P}(\text{Future} \mid \text{Past \& Present})$$

$$\underset{\sim}{=} \mathbb{P}(\text{Past} \mid \text{Present}) \, \mathbb{P}(\text{Future} \mid \text{Present})$$

$$(\because \text{Markov}) \qquad\qquad\qquad --\cdots (A)$$

$$\left( \because \mathbb{P}(A, B) = \mathbb{P}(A) \, \mathbb{P}(B \mid A) \right)$$

One way to think this is to recall:

$$\text{(if)} \quad \mathbb{P}(A \& B) = \mathbb{P}(A) \, \mathbb{P}(B)$$

$$\text{(then)} \qquad A \& B \quad \text{are independent.}$$

So (*) means :

" Past & future are <u>conditionally independent</u>,
given the present".

↖ can be taken as alternative def⁼
of Markov process.

Discrete Time : Markov Chain (have states $\{s_1, \ldots, s_m\}$)

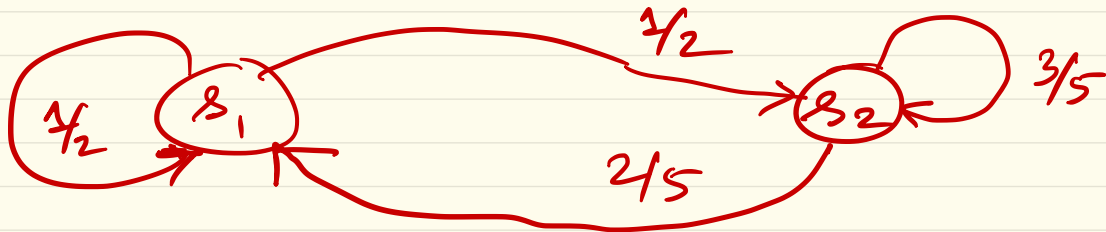$$\mathbb{P}\left(\underline{x}(t+1) = s_j \mid \underline{x}(t) = s_i\right)$$

$$= p_{ij} \in [0, 1]$$

This defines an $m \times m$ matrix

$$P = [p_{ij}] \text{ where } 0 \leq p_{ij} \leq 1, \quad \& \quad \sum_{j=1}^{m} p_{ij} = 1$$

↖ called (row) stochastic matrix

**Example:** 2 state Markov Chain:



$$\therefore \quad P = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{2}{5} & \frac{3}{5} \end{bmatrix}$$

**Example:** 3 state Markov Chain

$$\{S, C, R\}$$

$$P = \begin{array}{c} \\ S \\ C \\ R \end{array} \begin{array}{ccc} S & C & R \\ \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \\ 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix} \end{array} \longrightarrow \text{sum to 1.}$$

$$\text{etc.}$$

Coming back to feedback policy:

__Def:__ A feasible policy $\underline{\delta}_k$ is called "Markovian" or "Markov Policy" if $\underline{\delta}_k$ only depends on $\underline{y}_k \equiv \underline{x}_k$ (MDP)

("action now" depends on "state now")

Set of Markov policies: $\Gamma_M \subset \underline{\Gamma}$

all history dependent randomized policies

Intuition suggests:

$$\underline{\delta}^* \in \Gamma_M.$$

# Dynamic Programming Sol<sup>n</sup>:

Let
$$V_K(\underline{x}) = \text{Optimal remaining expected}$$
$$\text{cost from state } \underline{x} \text{ at time } k.$$
$$\text{(generic)}$$

$$= \inf_{(\underline{g}_K, \underline{g}_{K+1}, \dots, \underline{g}_{T-1})} \mathbb{E}\left[\left\{c_T(\underline{x}(T)) + \sum_{s=K}^{T-1} c_s(\underline{x}_s, \underline{u}_s)\right\} \,\middle|\, \underline{x}_K = \underline{x}\right]$$

Under a Markovian policy, can show that:

$$V_K^{\underline{g}}(\underline{x}_K^{\underline{g}}) = \mathbb{E}\left[\text{copy} \,\middle|\, \underline{x}_K^{\underline{g}}\right]$$

(we're using: If $\underline{g} \in \Gamma_M$, then $\{\underline{x}_K^{\underline{g}}\}$ is a Markov process)

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

Nothing random here

and $V_k(\underline{x}) :=$

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

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

This minimization is over actions (NOT over policies)
(Even if policies are randomized, there is nothing random about this minimization)

minimize
$\gamma(\cdot)$

$$\mathbb{E}\left[ C_T(\underline{x}(T)) + \sum_{k=0}^{T-1} C_k(\underline{x}_k, u_k) \right]$$