

# FINITE STATE MARKOV CHAINS

$$w_t = w_t(\omega) \quad \omega \in (\Omega, \mathcal{F}, \mathbb{P})$$

MDP. 1

$$n_{t+1} = f(n_t, w_t) \rightarrow \text{discrete time stochastic systems}$$

$$n_0 = n$$

$n_t$  = state

exogenous input

$w_t$  = stochastic process  $\Rightarrow n_t$  = stochastic process

Hp:  $w_t$  is an i.i.d sequence

$n_t$  = random variable

$n_t$  = random variable  $\Omega$

usually,  $n_t \in \mathbb{R}^M$ , suppose instead that  $n_t$

$$n_t \in \mathcal{X} \text{ a finite set} \quad \text{w.l.o.g. } n_t \in \{1, 2, \dots, M\}$$

$\Rightarrow n_t$  = finite state Markov chain

$n_{t+1} = f(n_t, w_t)$  induces a distribution over  $\{n_t\}$  and  
in particular the so called transition probabilities

$$P(i|j) = P(n_{t+1}=i | n_t=j) \quad + \text{invariant}$$

transition probabilities enough to characterize the evolution  
of the Markov chain

↓ thanks to

"the future is independent of  
the past given the present"

$$n_{t+1} = f(n_t, w_t) + \text{i.i.d.} \Rightarrow \text{MARKOV property: the past given the present}"$$

↳ dependence  
 $n_t$  and  $w_t$  only

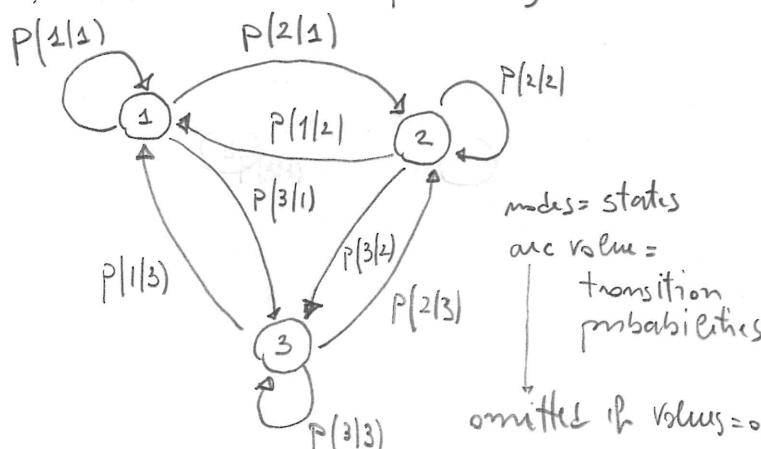
$$P(n_{t+1}=i | n_t=j \wedge n_{t-1}=k \wedge \dots) = P(n_{t+1}=i | n_t=j)$$

Clearly,  $\sum_{i=1}^M P(i|j) = 1$  and note that  $P(i|i) \neq 0$  in general

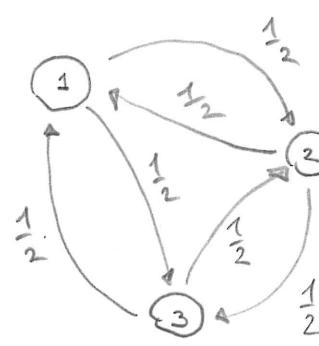
$[P(i|j)]_{M \times M}$  = transition matrix =  $P$  or conveys the whole info

$$[P(n_t=i)]_M = P^t \cdot [P(n_0=i)] \quad \text{for example}$$

example  $\mathcal{X} = \{1, 2, 3\}$



e.g.



$P =$

$$\equiv \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix}$$

$$\pi_{t+1} = f(\pi_t, u_t, w_t)$$

$\hookrightarrow$  control action

$w_t$  = i.i.d. sequence

$\Rightarrow$  (finite state/input)

Markov Decision

Process - MDP

(controlled  
Markov chain)

transition probabilities

$$P(i|j, u) = P\{ \pi_{t+1} = i \mid \pi_t = j \wedge u_t = u \}$$

$i, j \in \mathcal{X} \quad u \in \mathcal{U}$

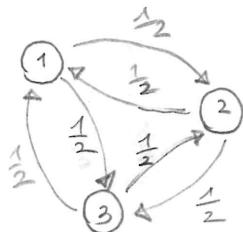
MARKOV property

$$= P\{ \pi_{t+1} = i \mid \pi_t = j \wedge u_t = u \wedge \dots \wedge \pi_{t-1} = k \wedge u_{t-1} = p \wedge \dots \}$$

$\leadsto m$  transition matrices  $P_u = [P(i|j, u)]_{m \times m}$

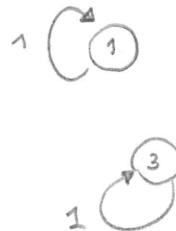
$$\text{e.g. } \mathcal{X} = \{1, 2, 3\} \quad \mathcal{U} = \{1, 2\}$$

$$u=1$$



$$P_1 = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix}$$

$$u=2$$



$$P_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$u_t$  to be chosen so as to  
impose a desired behavior  $\leadsto$  optimal control problem

every transition from  $\pi_t$  to  $\pi_{t+1}$  generates a (stage) cost

$g(\pi_t, u_t, w_t) \leadsto$  e.g.  $g(\pi_{t+1}) = g(f(\pi_t, u_t, w_t))$  but here more general to represent various situations

use the degree of freedom so as to minimize a "total" (average) cost  
 $\downarrow$   
to be precisely defined later

Example inventory control

$\pi_t$  = size of inventory at the end of day  $t$  (no of items)

$u_t$  = no. of items ordered at the end of day  $t$

$w_t$  = demand on day  $t+1$

$u_t \in \{0, 1, \dots, M\}$   $M =$  max no. of items that can be stored MDP.3

$u_t \in \{0, 1, \dots, M\}$  to fill up inventory when  $u_t=0$

$w_t \in \mathbb{Z}$

$$u_{t+1} = \min (\max(u_t + u_t, M) - w_t, 0)$$

$$g(u_t, u_t, w_t) = K \cdot \mathbf{1}_{u_t > 0} + C \cdot u_t + h \cdot u_t - p \cdot \min(\max(u_t + u_t, M), w_t)$$

fixed cost  
for placing an  
order      transportation  
costs      inventory  
maintenance  
costs      price      actually sold items      income

risk measure = expected value

expected stage cost :  $E[g(x_t, u_t, w_t)]$

we consider the evolution of the MDP up to time T

► finite horizon T ← terminal cost

$$J(x_0, \{u_t\}_{t=0}^{T-1}) = \sum_{t=0}^{T-1} E[g(x_t, u_t, w_t)] + E[g_T(x_T)]$$

t-varying

OBS : In the finite horizon problem we may admit  $f = f_t$   $g = g_t$   
results are unchanged

→ to obtain good solutions,  $u_t$  must be stochastic and adapted to  $x_t$

Hyp:  $x_t$  is observable and  $u_t$  is chosen as a function of  $x_t$   
by means of a policy  $\pi = \{M_0, M_1, \dots, M_{T-1}\}$  where  
each  $M_t: X \rightarrow U$ . ← motivated by the Markov property

In other words,  $u_t = M_t(x_t)$  (state feedback) so that

$$\begin{cases} x_{t+1} = f(x_t, M_t(x_t), w_t) \\ x_0 = \underline{x} \end{cases} \quad \text{and} \quad J(x, \{M_t(x_t)\}_{t=0}^{T-1}) =$$

$$= \sum_{t=0}^{T-1} E[g(x_t, M_t(x_t), w_t)] + E[g_T(x_T)]$$

$$= V^\pi(x) \quad \begin{matrix} \leftarrow \text{value function} \\ \text{associated to policy } \pi \end{matrix}$$

Goal:  $\min_\pi V^\pi(x)$

well posed  $\forall x \in X$  (finite no. of cases)

optimal value function  $V^*(x) = \min_\pi V^\pi(x)$

as  $x_0 = \underline{x}$  deterministic, only  $M_0(x)$  counts, while, as for  $M_1, \dots, M_{T-1}$ ,  
potentially different  $x$  could lead to different  $\pi$

⇒ there exists a unique policy  $\pi^* = \{M_0^*, M_1^*, \dots, M_{T-1}^*\}$  such that

$V^{\pi^*}(x) = V^*(x) \quad \forall x \in X \quad \hookrightarrow$  simultaneously optimal for  
all initial conditions

⇒ consequence of the Bellman optimality principle : given any  $x \in X$ , let

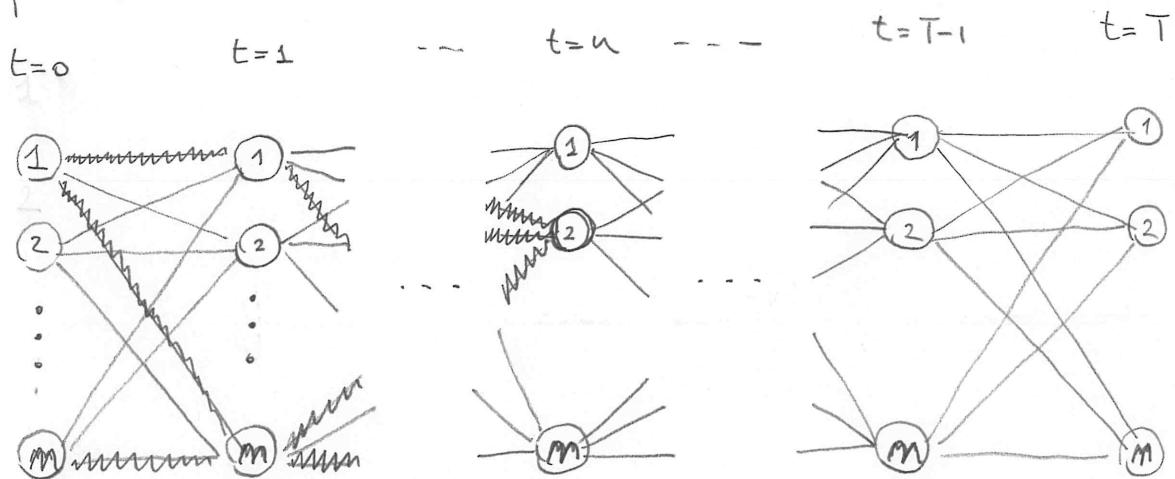
$\pi^* = \{M_0^*, M_1^*, \dots, M_{T-1}^*\}$  optimal for problem  $\min_\pi V^\pi(x)$ . and assume

that using  $\pi^*$ ,  $x_n = x'$  arises. Consider the sub-problem where

$x_n = x'$  and we wish to minimize the "cost-to-go" from  $t=n$  to MDP.5

$$t=T : V_K^{Tn}(x') = \sum_{t=n}^{T-1} \mathbb{E}[g(x_t, u_t, w_t)] + \mathbb{E}[g_T(x_T)] \quad \text{where } x_n = x'$$

and  $\Pi_n = \{u_0, \dots, u_{T-1}\}$ . Then  $\Pi_n^* = \{u_n^*, \dots, u_{T-1}^*\}$  is optimal for this sub-problem



once at  $t=n$  I'm arrived at  $x_n = 2$ , the future evolution doesn't depend on the past (Markov), if I'm not optimal we can switch to an optimal policy and simultaneously improve all the situations

$\Rightarrow$  Dynamic Programming algorithm: break down multistage optimization into single stage optimizations; generate backwards  $U_{T-1}^*$ , solving the cost-to-go from  $T-1$  to  $T$  subproblem for all  $x \in X$ ,  $x_{T-1} = x$ , then  $U_{T-2}^*$ , solving the cost-to-go from  $T-2$  to  $T$  subproblem where  $U_{T-1}^*$  is already implemented for all  $x \in X$  and  $x_{T-2} = x$ , and so forth and so on till  $U_0^*$

$$\text{init: } V_T^*(x) = g_T(x) \quad \forall x \in X$$

for  $t = T-1, \dots, 0$  compute  $\forall x \in X$

$$V_t^*(x) = \min_{u \in U} \mathbb{E}[g(x, u, w_t) + V_{t+1}^*(f(x, u, w_t))] \quad (*)$$

$$U_t^*(x) = \arg \min_{u \in U} \mathbb{E}[g(x, u, w_t) + V_{t+1}^*(f(x, u, w_t))]$$

end

then  $\Pi^* = \{U_0^*, \dots, U_{T-1}^*\}$  is optimal for  $\min_{\Pi} V^{\Pi}(x) \quad \forall x \in X$  and

$$V_0^*(n) = V^*(n).$$

$= T-1, \dots, 0$  MDP 6

FORMAL PROOF: It is enough to show that  $\forall k$  and  $\forall n \in X$

$$V_n^*(n) = \min_{\pi_k} V_n^{\pi_k}(n) = \min_{\pi_k} \sum_{t=n}^{T-1} \mathbb{E}[g(n_t, M_t(n_t), w_t)] + \mathbb{E}[g_T(n_T)]$$

$n_h = n$

as defined in  $(*)$

and that  $\pi_n^* = \{M_n^*, \dots, M_{T-1}^*\}$  as defined in  $(*)$  is optimal for  $\min_{\pi_n} V_n^{\pi_n}(n)$ . This can be done by induction

$$\begin{aligned} \min_{\pi_{T-1}} V_{T-1}^{\pi_{T-1}}(n) &= \min_{M_{T-1}} \mathbb{E}[g(n, M_{T-1}(n), w_{T-1})] + \mathbb{E}[g_T(f(n, M_{T-1}(n), w_{T-1}))] \\ &= \min_{M_{T-1}} \mathbb{E}[g(n, M_{T-1}(n), w_{T-1}) + V_T^*(f(n, M_{T-1}(n), w_{T-1}))] \\ &= \min_{u \in U} \mathbb{E}[g(n, u, w_{T-1}) + V_T^*(f(n, M_{T-1}(n), w_{T-1}))] \end{aligned}$$

for  $n \in X$ , we have  $m$  distinct opt. problems

Clearly,  $M_{T-1}^*(n) = \arg \min_{u \in U} \dots$  defines an optimal policy  $\pi_{T-1}^*$

Suppose, property is true for  $t = n+1$

$$\begin{aligned} \min_{\pi_n} V_n^{\pi_n}(n) &= \min_{M_n, \pi_{n+1}} \sum_{t=n}^{T-1} \mathbb{E}[g(n_t, M_t(n_t), w_t)] + \mathbb{E}[g_T(n_T)] \\ &= \min_{M_n, \pi_{n+1}} \mathbb{E}[g(n, M_n(n), w_n)] + \sum_{t=n+1}^{T-1} \mathbb{E}[g(n_t, M_t(n_t), w_t)] + \mathbb{E}[g_T(n_T)] \\ &= \min_{M_n, \pi_{n+1}} \mathbb{E}[g(n, M_n(n), w_n)] + \mathbb{E}[V_{n+1}^{\pi_{n+1}}(f(n, M_n(n), w_n))] \\ &= \min_{M_n, \pi_{n+1}} \mathbb{E}[g(n, M_n(n), w_n)] + V_{n+1}^{\pi_{n+1}}(f(n, M_n(n), w_n)) \\ &\geq \min_{M_n} \mathbb{E}[g(n, M_n(n), w_n)] + \min_{\pi_{n+1}} V_{n+1}^{\pi_{n+1}}(f(n, M_n(n), w_n)) \\ &= \min_{M_n} \mathbb{E}[g(n, M_n(n), w_n) + V_{n+1}^*(f(n, M_n(n), w_n))] \rightarrow m \text{ distinct problems} \\ &= \min_{u \in U} \mathbb{E}[g(n, u, w_n) + V_{n+1}^*(f(n, M_n(n), w_n))] \rightarrow \arg \min \text{ defines a feedback } M_n^* \end{aligned}$$

$$= \min_{u \in U} \mathbb{E} \left[ g(u, u, w_u) + V_{u+1}^{\pi_{k+1}^*} (f(u, u, w_u)) \right] =$$

$$= \mathbb{E} \left[ g(u, \mu_n^*(u), w_u) + V_{u+1}^{\pi_{k+1}^*} (f(u, \mu_n^*(u), w_u)) \right] = V_n^{\pi_n^*}(u)$$

altogether gives a policy  $\pi_n^*$

$$\min_{\pi_n} V_n^{\pi_n}(u) \geq V_n^*(u) = V_n^{\pi_n^*}(u)$$

must be an equality