# Markov chains

The uncertainty about the complex dynamics of the environment is modelled via Markov processes.

↳ Simplifying assumption: Markov property → the future dynamics of the environment is independent of the past given the present state. Formally, denoting by $S_t$ the state of the environment at time $t$, we have:

$$\mathbb{P}[S_{t+1} \mid S_t, S_{t-1} \ldots] = \mathbb{P}[S_{t+1} \mid S_t]$$

Here $S$ will be assumed to belong to a finite state space $\mathcal{S}$.

↳ Transition probabilities: $P: \mathcal{S} \times \mathcal{S} \to \mathbb{R}$ such that

* $P_{ss'} = \mathbb{P}[S_{t+1} = s' \mid S_t = s]$ : probability to transition from state $s$ to successor state $s'$.

* In general, transition can be time-dependent, but will be assumed stationary here.

* Concretely $P_{ss'}$ defines a stochastic matrix, ie a matrix with non-negative real entries which sum to one row-wise.

$$\underline{\underline{P}} = \begin{bmatrix} P_{11} & \cdots & P_{1n} \\ & \ddots & \\ P_{n1} & \cdots & P_{nn} \end{bmatrix} \text{ with } \left. \begin{array}{l} \sum_{s'} P_{ss'} = 1 \\ P_{ss'} \geq 0. \end{array} \right)$$ model for the environment

↳ Stationary dynamics: Under well-known condition (irreducibility) Markovian dynamics are ergodic, i.e., the probability distribution of $S_t$ converges toward a unique distribution $\pi_\infty$ independently from the starting distribution $\pi_0$. In other words:

$$\pi_t = \pi_0 P^t = \pi_0 (P \circ P \circ \ldots \circ P) \xrightarrow[t \to +\infty]{} \pi_\infty \qquad \text{at least weakly}$$

The invariant measure may be concentrated on one (absorbing) state. In this case, interest bears on the "transients" of the environment dynamics

## Markov reward process

Consider a Markov chain $(\mathcal{S}, P)$ with finite state space $\mathcal{S}$ and transition matrix $P$.

A Markov reward process is given by $(\mathcal{S}, P, R, \gamma)$ where

1) $R$ is a reward function, i.e., a random variable that follows the Markov property:

$$\mathbb{E}[R_{t+1} \mid S_t, S_{t-1}\cdots] = \mathbb{E}[R_{t+1} \mid S_t] = R_{S_t} \leftarrow \text{state dependent expected reward}$$

2) $\gamma$ is a discount factor $0 \le \gamma \le 1$ which measures the preference for immediate reward and is introduced to controlled the aggregated future reward / return:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+2} + \cdots = \sum_{k=0}^{+\infty} \gamma^k R_{t+k+1}$$

## Markov reward setting

Here, infinite horizon with discount rate: the performance is measured via an infinite series aggregated future reward discounted geometrically. Other settings are possible: for instance finite horizon without discounting, infinite horizon with average bounded reward, random horizon with terminal state ....

## Value Function

The value function of a Markov reward process measures the expected future return given the current state: $v: \mathcal{S} \to \mathbb{R}$

$$v(s) = \mathbb{E}[G_t \mid S_t = s] = \mathbb{E}\left[\sum_{k=0}^{+\infty} \gamma^k R_{t+k+1} \mid S_t = s\right]$$

(No concept of optimization at this stage.)

# Bellman equation for Markov reward process

Goal: given a MRP $(\mathcal{S}, P, R, \gamma)$, compute the value function $v(s)$ which measure how "beneficial" a state is.

Bellman idea: value function = immediate reward + discounted value at successor state

$$v(s) = \mathbb{E}[G_t \mid S_t = s]$$

$$= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots \mid S_t = s]$$

$$= \mathbb{E}[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \cdots) \mid S_t = s]$$

$$= \mathbb{E}[R_{t+1} \mid S_t = s] + \gamma \mathbb{E}[G_{t+1} \mid S_t = s]$$

$$= \mathbb{E}[R_{t+1} \mid S_t = s] + \gamma \mathbb{E}[\underbrace{\mathbb{E}[G_{t+1} \mid S_{t+1}]}_{v(S_{t+1})} \mid S_t = s]$$

$$\boxed{v(s) = \mathcal{R}_s + \gamma \sum_{s' \in \mathcal{S}} P_{ss'} v_{s'}} \quad \leftarrow \text{Bellman equation}$$

In matrix form:

$$\begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix} = \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \cdots & P_{1n} \\ \vdots & & \vdots \\ P_{n1} & \cdots & P_{nn} \end{bmatrix} \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix}$$

$$\underline{v} = \underline{R} + \gamma \underline{\underline{P}} \, \underline{v} \quad \longrightarrow \text{linear matrix equation}$$

$$\underline{v} = (\underline{\underline{I}} - \gamma \underline{\underline{P}})^{-1} \underline{R} \quad \longrightarrow \text{matrix inversion in } O(N^3)$$
$$\text{where } N = |\mathcal{S}|.$$

Impractical for large N : alternative methods
  ↳ dynamical programming
  ↳ Monte Carlo evaludation

# Markov decision process

The agent collecting reward can perform actions that will biased future reward collection. This means the transition probabilities of the environment can be affected by the agent's decision.

Formally, a Markov decision process is given by $(S, A, P, R, \gamma)$
 * $S$ is a finite set of states.
 * $A$ is a finite set of actions.
 * $P$ is a state transition matrix: $P_{ss'}^a = P[S_{t+1}=s' \mid S_t=s, A_t=a]$
 * $R$ is a reward function: $R_s^a = E[R_{t+1} \mid S_t=s, A_t=a]$
 * $\gamma$ is a discount factor: $0 \leq \gamma \leq 1$.

# Action policy

A policy $\pi$ is a distribution of actions on $A$ given states in $S$.
$$\pi(a|s) = P[A_t=a \mid S_t=s]$$
↳ policy may be stochastic or deterministic (here deterministic is enough)
↳ policy may be stationary or time dependent (here stationary is enough)

Restriction to stationary deterministic policies can be justified by the memoryless setting of finite state reward Markov process.

# Remark about Markovianity

Given a MDP $M = (S, A, P, R, \gamma)$ and fixed policy $\pi$, the joint sequence of states and rewards $(S_t, R_t)$ defines a Markov chain. Justification via averaging over policy actions:

$$P_{s,s'}^\pi = \sum_{a \in A} \pi(a|s) P_{s,s'}^a \qquad \leftarrow \text{transition of environment process}$$

$$R_{s,s'}^\pi = \sum_{a \in A} \pi(a|s) R_{s,s'}^a \qquad \leftarrow \text{transition of reward process}$$

## Action-value function

The action-value function of a MDP measures the expected return of taking the action $a$ in state $s$ under policy $\pi$:

$$q_\pi(s,a) = \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] \quad \leftarrow \text{policy dependent!}$$

Related state-value function:

$$v_\pi(s) = \mathbb{E}_\pi[G_t \mid S_t = s] = \mathbb{E}_\pi\Big[\mathbb{E}_\pi[G_t \mid S_t = s, A_t]\Big]$$

$$= \mathbb{E}_\pi[q_\pi(s, A_t)] = \sum_{a \in \mathcal{A}} \pi(a,s)\, q_\pi(s,a)$$

Reciprocally:

$$q_\pi(s,a) = \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a]$$

$$= \mathbb{E}_\pi[R_{t+1} \mid S_t = s, A_t = a] + \gamma\, \mathbb{E}_\pi\Big[\mathbb{E}_\pi[G_{t+1} \mid S_{t+1}] \mid S_t = s, A_t = a\Big]$$

$$= \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a\, v_\pi(s')$$

Bellman equation:

$$\boxed{v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a\mid s)\left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a\, v_\pi(s')\right)}$$

Matrix equation:

$$\underline{v}_\pi = \underline{\mathcal{R}}^\pi + \gamma\, \underline{\underline{P}}^\pi\, \underline{v}_\pi \left.\vphantom{\Big|}\right\} \text{linear equation}$$

$$\underline{v}_\pi = \left(\underline{\underline{I}} - \gamma\, \underline{\underline{P}}^\pi\right)^{-1} \underline{\mathcal{R}}^\pi$$

## Optimal functions

The optimal value function $v_*(s)$ is obtained by maximization over all policies $\pi$:

$$v_*(s) = \max_\pi v_\pi(s) \qquad \pi \in \text{convex, simplex space of distributions}$$

The optimal action-value function is similarly defined as

$$q_*(s,a) = \max_\pi q_\pi(s,a)$$

A MDP is solved when the optimal value $v^*$ is known. The optimal value function specifies the best possible performance of the MDP.

## Notion of optimality

Partial ordering on policies $\pi \geqslant \pi' \iff v_\pi(s) \geqslant v_{\pi'}(s)$ for all $s$.

## Theorem:

For all MDP, there is an optimal policy $\pi_*$ such that $\pi_* \geqslant \pi$ for any $\pi$. All optimal policies achieve the optimal value function $v_*(s) = v_{\pi_*}(s)$. All optimal policies achieve the optimal action-value function $q_*(s,a) = q_{\pi_*}(s,a)$.

## Idea of the proof

Bellman's optimality principle:

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

More formally, the optimal value function is recursively determined by the Bellman optimality equations. Intuitively, we have

1) $v_*(s) = \max_a q_*(s,a) \leftarrow$ value given most beneficial choice

2) $q_*(s,a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s') \leftarrow$ value of action $a$ given optimal value

1) + 2) $\implies$

$$\boxed{v_*(s) = \max_a \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s') \right)}$$

Bellman optimality equation

Similarly: $q_*(s,a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a \max_{a'} q_*(s',a')$.

## Proof of Bellman optimal equation

For any policy $\pi = (\pi_0, \pi')$ $\leftarrow$ not necessarily stationary

$\quad\quad\quad\quad\quad\quad\quad\quad \hookrightarrow$ can be assumed deterministic $\pi_0 = \delta_a$.

$v_*(s) = \max_\pi \mathbb{E}_\pi [G_t | S_t = s] = \max_\pi \mathbb{E}_\pi \left[ \sum_{k \geqslant 0} \gamma^k R_{t+k+1} | S_t = s \right]$

$$v_*(s) = \max_\pi \mathbb{E}_\pi \left[ R_{t+1} + \gamma \left( \sum_{k \geq 0} \gamma^k R_{t+k+2} \right) \Big| S_t = s \right]$$

$$= \max_\pi \mathbb{E}_\pi \left[ R_{t+1} + \gamma G_{t+1} \Big| S_t = s \right]$$

$$= \max_{(\pi_0, \pi')} \left\{ \mathbb{E}_{\pi_0} [R_{t+1} | S_t = s] + \gamma \mathbb{E}_{\pi_0} \left[ \mathbb{E}_{\pi'} [G_{t+1} | S_{t+1}] \Big| S_t = s \right] \right\}$$

$$= \max_{\pi_0} \left\{ R_s^{\pi_0} + \gamma \mathbb{E}_{\pi_0} \left[ \underbrace{\max_{\pi'} \mathbb{E}_{\pi'} [G_{t+1} | S_{t+1}]}_{v_*(S_{t+1})} \Big| S_t = s \right] \right\}$$

$$= \max_a \left\{ R_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a v_*(s') \right\}$$

## Proof of theorem:

Interpreting Bellman optimality equation as a fixed point equation suggests considering the Bellman operator:

$$T : \mathbb{R}^N \longrightarrow \mathbb{R}^N, \quad N = |\mathcal{S}|$$

$$TV(s) = \max_\pi \left\{ \pi(a|s) \left( R_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V(s') \right) \right\} \quad \Big| \quad \text{$\pi$ can be assumed to be deterministic } \pi = \delta_a$$

The operator $T$ is monotone, i.e., if $V \leq V'$ then $TV \leq TV'$.
More importantly $T$ is a contraction for the $L_\infty$ norm on $\mathbb{R}^N$.
Indeed:
$$|TV(s) - TV'(s)| = \Big| \max_{a \in A} \left[ R_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V(s') \right.$$
$$\left. - \max_{a' \in A} \left[ R_s^{a'} + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^{a'} V'(s') \right] \right.$$

$$(*) \qquad \leq \max_{a \in A} \Big| \cancel{R_s^a} + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V(s')$$
$$- \cancel{R_s^a} - \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V'(s') \Big|$$

$$\leq \gamma \underbrace{\|V - V'\|_\infty}_{} \underbrace{\max_{a \in A} \Big| \sum_{s' \in \mathcal{S}} P_{ss'}^a \Big|}_{\leq 1}$$
$$\underset{< 1}{\uparrow}$$

$(*)$ from $\big| \max_a f(a) - \max_{a'} g(a') \big| \leq \max_a |f(a) - g(a)|$

Justification:

$$\max_a F(a) - \max_{a'} g(a) = F(a^*) - \max_{a'} g(a')$$
$$\leq F(a^*) - g(a^*)$$
$$\leq \max_a |F(a) - g(a)|$$

The contraction property of $T$ allows one to use the Banach fixed point theorem on $(\mathbb{R}^N, \|\ \|_\infty)$ to show that there is a unique solution $v^*$, the optimal value function.

Moreover: $\lim_{k \to +\infty} T^k v = v^*$.