

→ Markov Chain

- $\mathcal{X}$  - set of possible states
- $P$  - transition probability matrix

→ Hidden Markov Model (HMM)

- $\mathcal{X}$  - set of possible states
- $P$  - transition probability matrix
- $\mathcal{Z}$  - set of possible observations
- $O$  - observation probability matrix

→ Markov Decision Process (MDP)

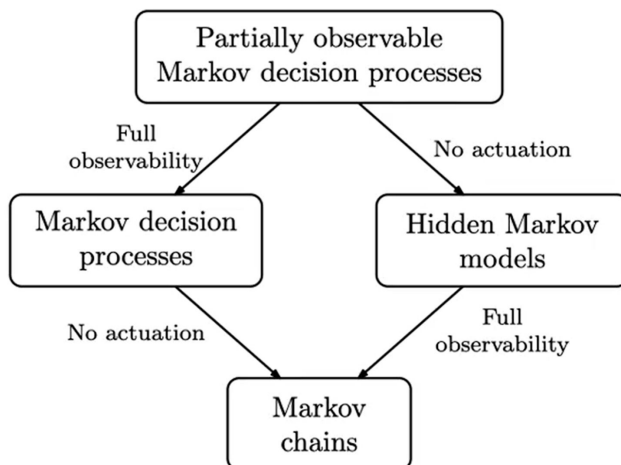
- $\mathcal{X}$  - set of possible states
- $\mathcal{A}$  - set of possible actions
- $[P_a]$  - transition probability matrices for  $a \in \mathcal{A}$
- $C$  - immediate cost function ( $\mathcal{X} \times \mathcal{A}$ )

→ Markov Decision Problem

- $\mathcal{X}$  - set of possible states
- $\mathcal{A}$  - set of possible actions
- $[P_a]$  - transition probability matrices for  $a \in \mathcal{A}$
- $C$  - immediate cost function
- $\gamma$  - discount factor

→ Partially Observable MDP (POMDP)

- $\mathcal{X}$  - set of possible states
- $\mathcal{A}$  - set of possible actions
- $\mathcal{Z}$  - set of possible observations
- $[P_a]$  - transition probability matrices for  $a \in \mathcal{A}$
- $[O_a]$  - observation probability matrices for  $a \in \mathcal{A}$
- $C$  - immediate cost function



→ Belief MDP: For a given POMDP, the belief MDP is an equivalent MDP model whose state space is the set of possible distributions over the state space of the original POMDP (also known as the belief space). Belief MDPs are thus a “bridge” between the theory of MDPs and that of POMDPs, allowing us to use MDP solution methods such as value and policy iteration to solve POMDPs.

→  $P$ , Transition probability matrix:  $\mathcal{X} \times \mathcal{X}$ , where the element  $(x_1, x_2)$  represents the probability of the agent going from state  $x_1$  to state  $x_2$ . When there are multiple actions then we have  $P_a$ , which represents the probability of the agent going from state  $x_1$  to state  $x_2$  if it selected action  $a$ .

→  $\Pi$  (ou  $\pi$ ), Policy matrix:  $\mathcal{X} \times \mathcal{A}$ , where a row represents the probability of choosing action  $a$  in state  $x$ . It serves to determine the action to choose.

→  $c$ , Cost matrix:  $\mathcal{X} \times \mathcal{A}$ , where the element  $(x, a)$  represents the cost of choosing action  $a$  when on state  $x$ . In this subject all costs are in interval  $[0, 1]$ .

→  $Q^\Pi$ , Q function:  $\mathcal{X} \times \mathcal{A}$ , where the element  $(x, a)$  represents the average cost-to-go for policy  $\Pi$ , given initial state  $x$  and action  $a$ .

→  $O$ , Observation matrix:  $\mathcal{X} \times \mathcal{Z}$ , where the element  $(x, o)$  represents the probability of witnessing Observation  $o$ , when at state  $x$ . When there are multiple actions then we have  $O_a$ , which represents the probability of the agent witnessing Observation  $o$ , when at state  $x$  if it selected action  $a$ .

→  $J^\Pi$ , Cost-to-go function:  $\mathcal{X} \times 1$ , where the element  $(x)$  represents the total cost of reaching the goal starting from state  $x$ , while following the policy  $\Pi$ .

→  $b_t$ , Belief:  $\mathcal{X} \times 1$ , which represents the stationary distribution over  $\mathcal{X}$  at the instant  $t$ . As it is a distribution all of its values must add up to 1 and are all between  $[0, 1]$

→ **Markov Property**: The state at instant  $t$  is enough to predict the state at  $t + 1$

- i.e. depends only on the last state

→ **Irreducibility**: A chain is **irreducible** if any state  $y$  can be reached from any state  $x$

- State  $y$  can be reached from a state  $x$  if  $P_t(y | x) > 0$
- We can split the state space  $X$  in **communicating classes** (sets of states that are mutually reachable) - basically componentes fortemente ligadas

→ **Aperiodicity**

- The period of a state  $x$  is the **greatest common divisor** of all time steps in which  $x$  can be visited if the chain departs from  $x$
- A state is aperiodic if its  $p_x = 1$
- A chain is **aperiodic** if all states are aperiodic, and **periodic** otherwise
- When a chain is **irreducible**, all states have the same period
- If all nodes have self loops, the chain is aperiodic

→ **Stationary** distribution

- Let  $\mu$  be a **distribution** over  $X$
- $\mu$  is a row vector ( $1 \times m$ ) such that  $\mu = [\mu(x_1) \mu(x_2) \dots \mu(x_{|X|})]$
- Component  $\mu(x)$  is the probability of  $x$  according to  $\mu$
- If the state at time  $t$  follows the distribution  $\mu$  and  $\mu$  is stationary, then the state at time  $t+1$  also follows the distribution  $\mu$ , i.e. **if  $\mu$  is stationary,  $\mu * P = \mu$**
- The stationary distribution corresponds to the stable behavior of the chain

→ An **irreducible** and **aperiodic** Markov Chain possesses a stationary distribution

- **Positive chain** - if  $x_t$  is distributed according to  $\mu$ , then so is  $x_{t+1}$

→ For an irreducible and aperiodic Markov chain with a stationary distribution will eventually converge to this distribution independently of the initial distribution

- No matter how the initial state is selected, if we let the chain run for a sufficient amount of time, or make predictions far enough in the future, then those predictions will eventually converge to the stationary distribution
- **Ergodic chain** - eventually reaches the stationary distribution

→ Filtering: given a sequence of observations, estimate the **final state**

- Forward algorithm

→ Smoothing: given a sequence of observations, estimate the **sequence of states**

- We want to estimate **the most likely sequence**, i.e.
  - $x_{0:T}^* = \operatorname{argmax} P_{\mu_0}[x_{0:T} = x_{0:T} \mid z_{0:T} = z_{0:T}]$
- Maximizing forward mapping
- Viterbi algorithm

→ Marginal smoothing: given a sequence of observations, estimate **some state** in the middle

- Forward-backward algorithm

→ Prediction: given a sequence of observations, predict **future states**

- Compute  $\mu_{T|0:T}$  using the forward algorithm
- Use the Markov property
  - $\mu_{T+1|0:T} = \mu_{T|0:T} * P$

### → Forward algorithm

- Forward mapping: maps states to real numbers
  - $a_t: X \rightarrow \mathbb{R}$
  - $a_t(x) = P\mu_0[x_t = x, z_{0:t} = z_{0:t}]$
  - Each component  $x$  represents the joint probability of being in state  $x$  at timestamp  $t$  given the observations from timestamp 0 to timestamp  $t$
- Algorithm:
  - Multiply initial distribution by  $\text{diag}(O(z_0|:))$
  - At each time step:
    - Multiply current distribution by  $P$
    - Multiply by  $\text{diag}(O(z_t|:))$
  - Normalize

Note: If there are no observations, all you have to do is multiply the by  $P_A$ , being  $A$  the action that was chosen.

**Marginal Smoothing:** Given a sequence of observations, estimate some state in the middle.

### → Forward-backward algorithm

- Backward mapping: maps states to real numbers
  - $b_t: X \rightarrow \mathbb{R}$
  - $b_t(x) = P\mu_0[z_{t+1:T} = z_{t+1:T} \mid x_t = x]$
- Algorithm:
  - Multiply initial distribution by  $O(z_0|:)$ ,  $b_T = 1$
  - for  $t = 0, \dots, T-1$  do **(forward update)**
    - Current distribution  $\ast P \ast \text{diag}(O(z_t|:))$
  - for  $t = T-1, \dots, 0$  do **(backward update)**
    - $b_t = P \ast \text{diag}(O(z_t|:)) \ast b_{t+1}$
  - Combine and normalize  $\rightarrow a_t \times b_t / (a_t \times b_t)$

**Joint Smoothing :** Given a sequence of observations and/or actions, estimate the whole sequence of states (most likely sequence).

### → Viterbi algorithm

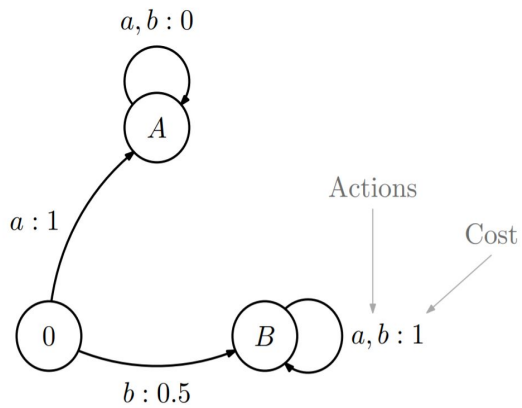
- $m_0 = \text{diag}(O_{:,z_0}) \ast \mu_0^T$  // just  $\mu_0^T$  if there is no initial observation, i.e at  $t = 0$
- for  $t = 1, \dots, T$  do
  - $m_t = \text{diag}(O_{:,z_t}) \ast \max\_row\_vector(PT \ast \text{diag}(m_{t-1}))$  // dá um row vector
  - $i_t = \text{indices dos } \max\_row\_vector() \text{ escolhidos}$  // row vector, indices de 1-N
- $x_t^* = \text{indice do } \max\_row() \text{ de } m_t$  // inteiro, um indice de 1-N
- for  $t = T-1, \dots, 0$  do
  - $x_t^* = i_{t+1}(x_{t+1}^*)$  // backtrack com indices, DM me se não perceberem
- return  $x_{0:t}^*$
- Example iteration  $m_1 = \text{diag}(O_{:,u}) \max\{P^T \text{diag}(m_0)\}$

**Decision Trees:** The objective is to minimize the cost

- $c = -u$  (where  $c$  is cost and  $u$  is utility)
- So the objective is to maximize the utility

Both policy iteration and value iteration are algorithms used to find the optimal policy for a certain MDP.

Now to explain both policy iteration and value iteration and I'm going to use the following example:



- States:
- $\mathcal{X} = \{0, A, B\}$
- Actions:
- $\mathcal{A} = \{a, b\}$

- Transition probabilities:

$$\mathbf{P}_a = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

- Cost:

$$\mathbf{P}_b = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

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

## Value Iteration

**Require:** MDP  $\mathcal{M} = (\mathcal{X}, \mathcal{A}, \{\mathbf{P}_a\}, c, \gamma)$ ; tolerance  $\varepsilon > 0$ ;

- 1: Initialize  $k = 0, J^{(0)} \equiv 0$
- 2: **repeat**
- 3:      $J^{(k+1)} \leftarrow \mathsf{T} J^{(k)}$
- 4:      $k \leftarrow k + 1$
- 5: **until**  $\|J^{(k-1)} - J^{(k)}\|_\infty < \varepsilon$ .
- 6: Compute  $\pi^*$  using  $(\dagger)$
- 7: **return**  $\pi^*$

**In step 1**, we initialize  $k = 0$  and  $J^{(0)}$  as Vector with dimensions  $X \times 1$

**In step 2**, we will now start a loop where that will stop when the distance between  $J^{(k-1)}$  and  $J^{(k)}$  is lesser than  $\varepsilon$ , which is our tolerance.

**In step 3**, we will apply the operator  $\mathsf{T}$  to  $J^{(k)}$ . The operator  $\mathsf{T}$  is actually just for each action  $a$  in de MDP's Actions do the following formula:

$$J^{(1)} = \min_a \left[ c_a + \gamma \mathbf{P}_a J^{(0)} \right]$$

In our example, we will do it for action  $a$  and for action  $b$ , since we only have 2 actions.

- For action  $a$ :
- For action  $b$ :

$$\begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} + 0.99 \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0.5 \\ 0 \\ 1 \end{bmatrix} + 0.99 \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0 \\ 1 \end{bmatrix}$$

And then we will see which one is the minimum in each line between the two, which will be the rows of our  $J^{(1)}$

$$\min \left( \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0.5 \\ 0 \\ 1 \end{bmatrix} \right) = \begin{bmatrix} 0.5 \\ 0 \\ 1 \end{bmatrix}$$

**In step 4**, we increment variable  $k$ , with  $k = k+1$ .

**In step 5**, we verify the stop condition.

**In step 6**, we will compute the policy  $\pi^{(*)}$  to do this we will look at the example:

$$\min \left( \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0.5 \\ 0 \\ 1 \end{bmatrix} \right) = \begin{bmatrix} 0.5 \\ 0 \\ 1 \end{bmatrix}$$

For each row of this matrix let's see which column has the lowest value.

In row 0, associated with state 0, column 1 has the lowest value, therefore in our policy

$$\pi^{(*)}[0][1]=1$$

In row 1, due to the fact that both have the same value, both have the lowest value, therefore in our policy,  $\pi^{(*)}[1][0]=1$  and  $\pi^{(*)}[1][1]=1$ . However all rows of a policy must add up to 1, so we must normalize this, by dividing it by the total sum of the row's values, resulting in  $\pi^{(*)}[1][0]=0.5$  and  $\pi^{(*)}[1][1]=0.5$ .

In row 2, due to the fact that both have the same value, both have the lowest value, therefore in our policy,  $\pi^{(*)}[2][0]=1$  and  $\pi^{(*)}[2][1]=1$ . However all rows of a policy must add up to 1, so we must normalize this, by dividing it by the total sum of the row's values, resulting in  $\pi^{(*)}[2][0]=0.5$  and  $\pi^{(*)}[2][1]=0.5$

$$\begin{bmatrix} 0 & 1 \\ 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

This results in  $\pi^{(*)} =$

**In step 7** we return  $\pi^{(*)}$

### Policy Iteration

---

**Require:** MDP  $\mathcal{M} = (\mathcal{X}, \mathcal{A}, \{\mathbf{P}_a\}, c, \gamma)$

1: Initialize  $k = 0$ ,  $\pi^{(0)}$  randomly

2: **repeat**

3:      $\mathbf{J} \leftarrow (\mathbf{I} - \gamma \mathbf{P}_{\pi^{(k)}})^{-1} \mathbf{c}_{\pi^{(k)}}$

4:      $\pi^{(k+1)} \leftarrow \pi_g^{\mathbf{J}}$

5:      $k \leftarrow k + 1$

6: **until**  $\pi^{(k-1)} = \pi^{(k)}$

7: **return**  $\pi^{(k)}$

Policy evaluation

Policy improvement

**In step 1**, we initialize the variable k as 0 and we initialize the policy  $\pi^{(0)}$  randomly. Remember that  $\pi^{(0)}$  is a matrix with dimension  $X \times A$ , and that each row must add up to 1.

**In step 2**, we will now start a loop that will last until  $\pi^{(k-1)} = \pi^{(k)}$ . This means that our policy can no longer be improved, the loop will stop.

**In step 3**, we will calculate the J value for this iteration.

$\mathbf{I}$  is the identity matrix, with dimension  $X \times X$ , where X is the number of states in the MDP's state space,  $\mathcal{X}$ .  $\gamma$  is the MDP's discount rate. But first, we need to calculate  $\mathbf{P}\pi^{(k)}$ . To calculate  $\mathbf{P}\pi^{(k)}$  we need to for each action A, **do the product of the diagonal matrix of  $\pi^{(k)}$ 's A column with the Probability Matrix associated with action A and then sum them all**. Second we need to calculate  $\mathbf{c}\pi^{(k)}$ , to do that we need to for each action A, **do the product of the diagonal matrix of  $\pi^{(k)}$ 's A column with the Cost Matrix's A column and then sum them all**. After we have all necessary components we can now calculate J, using the formula on the picture's step 3.

To make step 3 clearer I'll use the example:

$$\pi^{(0)} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

In our example, our  $\pi^{(0)}$  is the uniform policy, which means  $\pi^{(0)} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$   
 Now let's calculate  $P\pi^{(k)}$ , the formula for  $P\pi$  is  $\text{diag}(\pi^{(k)}[:,a]).P_a$ , where  $\text{diag}(\pi^{(k)}[:,a])$  is the matrix where the diagonal is the a column of  $\pi^{(k)}$  and we do this for each action.  
 Let's show it on the example:

$$\text{diag}(\pi^{(0)}[:,a]) = \begin{bmatrix} 0.5 & 0 & 0 \\ 0 & 0.5 & 0 \\ 0 & 0 & 0.5 \end{bmatrix} \text{ and } \text{diag}(\pi^{(0)}[:,b]) = \begin{bmatrix} 0.5 & 0 & 0 \\ 0 & 0.5 & 0 \\ 0 & 0 & 0.5 \end{bmatrix}$$

Now let's multiply each by their respective actions's probability transition matrix:

$$\text{For action a: } \text{diag}(\pi^{(0)}[:,a]).P_a = \begin{bmatrix} 0 & 0.5 & 0 \\ 0 & 0.5 & 0 \\ 0 & 0 & 0.5 \end{bmatrix} \text{ For action b: } \text{diag}(\pi^{(0)}[:,b]).P_b = \begin{bmatrix} 0 & 0 & 0.5 \\ 0 & 0.5 & 0 \\ 0 & 0 & 0.5 \end{bmatrix}$$

$$P\pi = \begin{bmatrix} 0 & 0.5 & 0.5 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Now if we sum both we get  $P\pi$   
 Now let's calculate  $c\pi^{(k)}$ , the formula for  $c\pi$  is  $\text{diag}(\pi^{(k)}[:,a]).c[:,a]$ , where  $\text{diag}(\pi^{(k)}[:,a])$  is the matrix where the diagonal is the a column of  $\pi^{(k)}$  and  $c[:,a]$  is the a column of c.

$$\text{For action a: } \text{diag}(\pi^{(0)}[:,a]).c[:,a] = \begin{bmatrix} 0.5 \\ 0 \\ 0.5 \end{bmatrix} \text{ For action b: } \text{diag}(\pi^{(0)}[:,b]).c[:,b] = \begin{bmatrix} 0.25 \\ 0 \\ 0.5 \end{bmatrix}$$

$$c\pi = \begin{bmatrix} 0.75 \\ 0 \\ 1 \end{bmatrix}$$

Now if we sum both we get  $c\pi$   
 Now that we have both  $c\pi$  and  $P\pi$  we can finally calculate J. The formula for J is:

$$J\pi = (I - \gamma P\pi)^{-1} c\pi, \text{ whose result is } \begin{bmatrix} 50.25 \\ 0 \\ 100 \end{bmatrix}.$$

**In step 4**, We will now calculate the Q-Function for each action, using the J value we calculated in step 3. The Q-Function for the action, A is calculated by:  $c[:,A] + \gamma * P[A]J$ , where  $c[:,A]$  is the A<sup>th</sup> column of the cost matrix and  $P[A]$  is the probability transition Matrix of action A.

$$\text{For action a: } Q_a = \begin{bmatrix} 1 \\ 0 \\ 100 \end{bmatrix} \text{ For action b: } Q_b = \begin{bmatrix} 99.5 \\ 0 \\ 100 \end{bmatrix}$$



$$Q^\pi = \begin{bmatrix} 1 & 99.5 \\ 0 & 0 \\ 100 & 100 \end{bmatrix}.$$

Which will result in a  $Q^\pi = \begin{bmatrix} 1 & 99.5 \\ 0 & 0 \\ 100 & 100 \end{bmatrix}$ .

Now we will choose for each row the column that has the lowest value.

For row 0, column 0 has the lowest value, therefore  $\pi^{(1)}[0][0] = 1$ .

For row 1, both columns have the same value, therefore  $\pi^{(1)}[1][0] = 1$  and  $\pi^{(1)}[1][1] = 1$ , however because rows need to add up to zero, we need to normalize it by dividing each value on the row by the sum of the row's value, so in actually  $\pi^{(1)}[1][0] = 0.5$  and  $\pi^{(1)}[1][1] = 0.5$ .

For row 2, both columns have the same value, therefore  $\pi^{(1)}[2][0] = 1$  and  $\pi^{(1)}[2][1] = 1$ , however because rows need to add up to zero, we need to normalize it by dividing each value on the row by the sum of the row's value, so in actually  $\pi^{(1)}[2][0] = 0.5$  and  $\pi^{(1)}[2][1] = 0.5$ .

$$\pi^{(1)} = \begin{bmatrix} 1 & 0 \\ 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

So in the end we get  $\pi^{(1)} = \begin{bmatrix} 1 & 0 \\ 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$ , concluding step 4.

**In step 5**, we increment the value of  $k$ , using  $k = k + 1$

**In step 6**, we verify if we reached our stopping point. If we didn't we will go back to step 3.

**And in step 7**, we return  $\pi^{(k)}$

→ How to verify if a cost-to-go function is optimal?

To verify if a cost-to-go is optimal, do one iteration of value iteration using it. If it remains unchanged then it's optimal, otherwise it isn't.

→ How to verify if a policy is optimal?

To verify if a policy is optimal, do one iteration of policy iteration using it. If it remains unchanged then it's optimal, otherwise it isn't.

# Heuristics

In time step  $t$ , we have a belief  $b$ , and also the optimal MDP policy

$$b = [b(x_1) \quad b(x_2) \quad \dots \quad b(x_N)]$$

## → MLS (Most Likely State)

- We choose the most likely state from the belief
- After that we will execute the corresponding action:  $\pi_{\text{mdp}}(x_2)$ .
- Basically  $\pi_{\text{MLS}}(b) = \pi_{\text{MDP}}(\underset{x \in \mathcal{X}}{\operatorname{argmax}} b(x))$

## → AV

- States vote proportionally to their probability
- Each action  $a$  will have a number of votes.
- Each state  $x$  will calculate the corresponding action  $\pi_{\text{mdp}}(x)$  and add  $b(x)$  to action's number of votes. The action that has the highest number of votes will be executed.

## → Q-MDP

- No partial observability
- Start by creating a row vector  $Q_{\text{MDP}}$  with dimensions  $A \times 1$ , with zeros in every element
- for each  $b(x)$  in  $b$ 
  - multiply it by  $Q_{\text{MDP}}(x, :)$
  - $Q_{\text{MDP}} += b(x) * Q_{\text{MDP}}(x, :)$
- Select  $\operatorname{argmin}(Q_{\text{MDP}})$

Another way to calculate it directly is to transpose  $b$ , obtaining  $b^T$  and multiply it by  $Q_{\text{mdp}}$ . In that case we have  $q_{\text{mdp}} = b^T Q_{\text{mdp}}$

## → FIB

- Factor out the belief but include partial observability

# Sauce

→ Compute **optimal policy** for MDP

- $\pi^*(x) = \operatorname{argmin}_{a \in \mathcal{A}} Q^*(x, a)$
- For each row in  $Q^*$ , select the smallest value
- i.e. for  $Q^* = \begin{bmatrix} 1.45 & 1.9 \\ 1.31 & 0.0 \\ 1.0 & 1.9 \end{bmatrix}$ , the optimal policy is  $\begin{bmatrix} 1.0 & 0.0 \\ 0.0 & 1.0 \\ 1.0 & 0.0 \end{bmatrix}$

→ Suppose that, at time step  $t = 0$ , agent chooses by Q-MDP. Compute the **action selected**

- Q-MDP prescribes  $a_0 = \operatorname{argmin}_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} b_0(x) Q^*(x, a)$ , and i.e.  $b_0 Q^* = \begin{bmatrix} 1.25 & 1.27 \end{bmatrix}$ , we choose  $a_1$

→ Suppose that agent observes  $z_1 = u$ . Compute the **resulting belief** at  $t = 1$

- Belief update rule
- $b_1 = \rho b_0 P_a \operatorname{diag}(O_{a,u})$ , where  $\rho$  is a normalization constant

→ Compute **optimal action** at  $t = 0$  ( $\Gamma$  is gifted to us by the elves,  $\alpha$ -vectors)

- $J^*(b_0) = \min_{\alpha \in \Gamma} b_0 \cdot \alpha$
- Cada  $b_0 * \alpha$  dá um número real; escolhemos o mais pequeno e a ação correspondente

→ Get the **distribution over states** at  $t = 2$ , assuming that the agent makes no observations:

- $\mu_2 = b_0 * P_{a0} * P_{a1}$

→ Compute the **most likely state** at  $t = 2$

- If they give us the first state,  $\mu_0$  is a column vector with 1 on that state
- $\mu_2 = \mu_0 * P_{a0} * P_{a1}$
- Forward algorithm
- If they ask for  $t = 0$ , just need to do it backwards (don't forget  $\mu_{0|0:3} = \alpha_0 \otimes \beta_0 / (\alpha_0^\top \beta_0)$ )

→ Compute the **most likely sequence of states** given actions and observations

- Viterbi

→ Show that  $b_0$  is a **stationary distribution** for the policy

- Obtain the transition probability matrix  $P$ 
  - For each row of the policy, see which action it chooses and get the corresponding row from said action
- Then we just have to prove that  $b_0 P = b_0$

→ Calculate chain **periodicity**

- Draw it from the transition probability matrix

→ Is the chain **ergodic**?

- “The chain is periodic with period N. Therefore, for any initial distribution  $b \neq b_0$ , the state distribution will oscillate with a period of N and will not converge to  $b_0$ . It follows that the chain is not ergodic.”

→ Is the chain **irreducible**?

- No, since there are multiple communicating classes

→ Indicate one advantage of Q-MDP over MLS and AV

- Both the MLS and AV heuristics use directly the optimal MDP policy and, therefore, prescribe only actions that are optimal for the underlying MDP. Q-MDP, on the other hand, may prescribe actions that the underlying MDP policy never selects. This was observed, for example, in the Tiger problem, where the Q-MDP did select the action “Listen”, while neither MLS nor AV were able to select such action.

→ Compute the **cost-to-go** associated with policy  $\pi$

- $J^\pi = (I - \gamma P_\pi)^{-1} c_\pi$

→ Compute the **greedy policy** with respect to the cost-to-go function...

- We compute the Q-values for each action  $a \in A$
- $Q_{:,A}^\pi = c_{:,A} + \gamma P_A J^\pi$  for action A, for example
- The resulting greedy policy is, therefore,  $\pi_g^{J^\pi}$  (same dimension as each of the Q-values)
  - To get it, just put probability 1 on the smallest value of the actions for a given row
  - And 0 in all others

→ Do you believe that  $\pi$  is the policy  $\pi$  is optimal?

- We have  $J\pi$ , so we perform one step of value iteration (or policy iteration as alternative)
  - $Q_{:,a}^\pi = c_{:,a} + \gamma P_a J^\pi$  for each action
  - $J_{\text{new}} = \min_a Q_{:,a}^\pi$ , we create a row vector with the smallest values for each of the rows in the Q things computed for each action
  - Since  $J_{\text{new}} = J\pi$ , we can conclude that  $J\pi = J^*$  and thus  $\pi = \pi^*$
  - And vice versa if  $J_{\text{new}} \neq J\pi$
- “The policy  $\pi$  is, most likely, not optimal, since the greedy policy obtained from  $J\pi$  is different from  $\pi$ , which can be observed by noting that the transition probabilities resulting from  $\pi_g^{J^\pi}$  are different from  $P_\pi$ . Therefore, since the greedy policy is what we would get after one step of policy iteration,  $\pi$  is, most likely, not optimal.”

→ Perform one step of **policy iteration** (using policy  $\pi$ )

- $Q_a = c_a + \gamma \mathbf{P}_a \mathbf{J}^\pi$  for each action
- $\pi_{\text{new}}$ , we create a row vector with the smallest values for each of the rows in the Q things computed for each action
- If  $\pi_{\text{new}} = \pi$ , then  $\pi$  is optimal