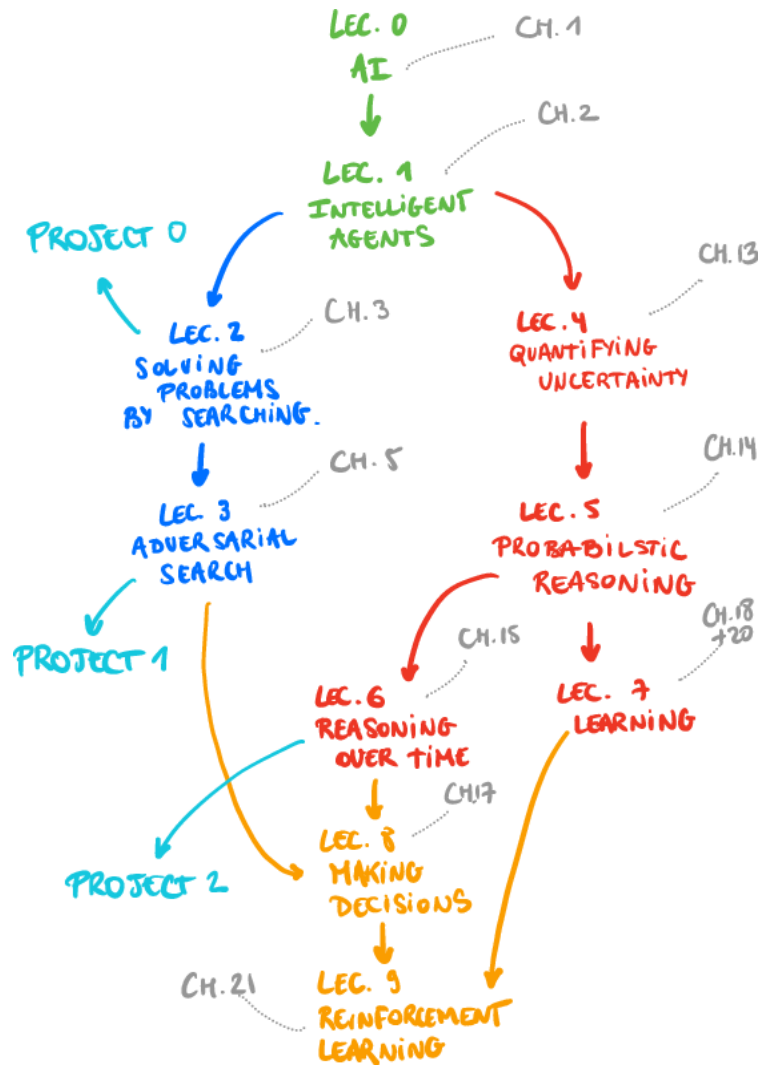


# Introduction to Artificial Intelligence

Lecture 8: Making decisions

Prof. Gilles Louppe  
[g.louppe@uliege.be](mailto:g.louppe@uliege.be)



# Today



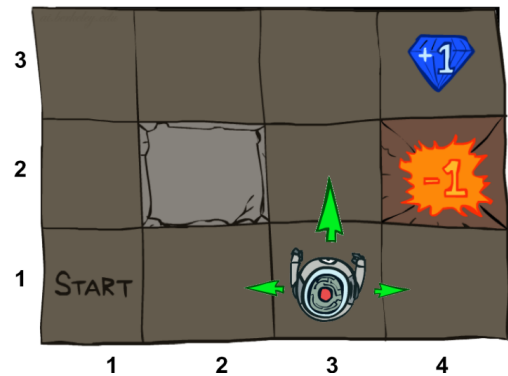
Reasoning under uncertainty and **taking decisions**:

- Markov decision processes
  - MDPs
  - Bellman equation
  - Value iteration
  - Policy iteration
- Partially observable Markov decision processes

# Grid world

Assume our agent lives in a  $3 \times 4$  grid environment.

- Noisy movements: actions do not always go as planned.
  - Each action achieves the intended effect with probability **0.8**.
  - The rest of the time, with probability **0.2**, the action moves the agent at right angles to the intended direction.
  - If there is a wall in the direction the agent would have been taken, the agent stays put.
- The agent receives rewards at each time step.
  - Small 'living' reward each step (can be negative).
  - Big rewards come at the end (good or bad).

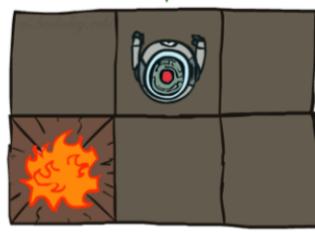
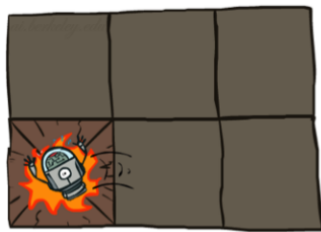
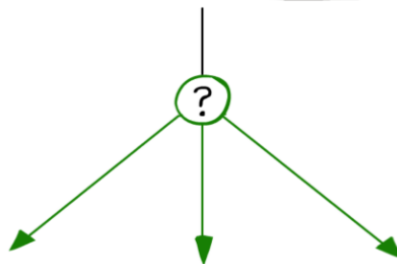


Goal: maximize sum of rewards.

## Deterministic actions



## Stochastic actions

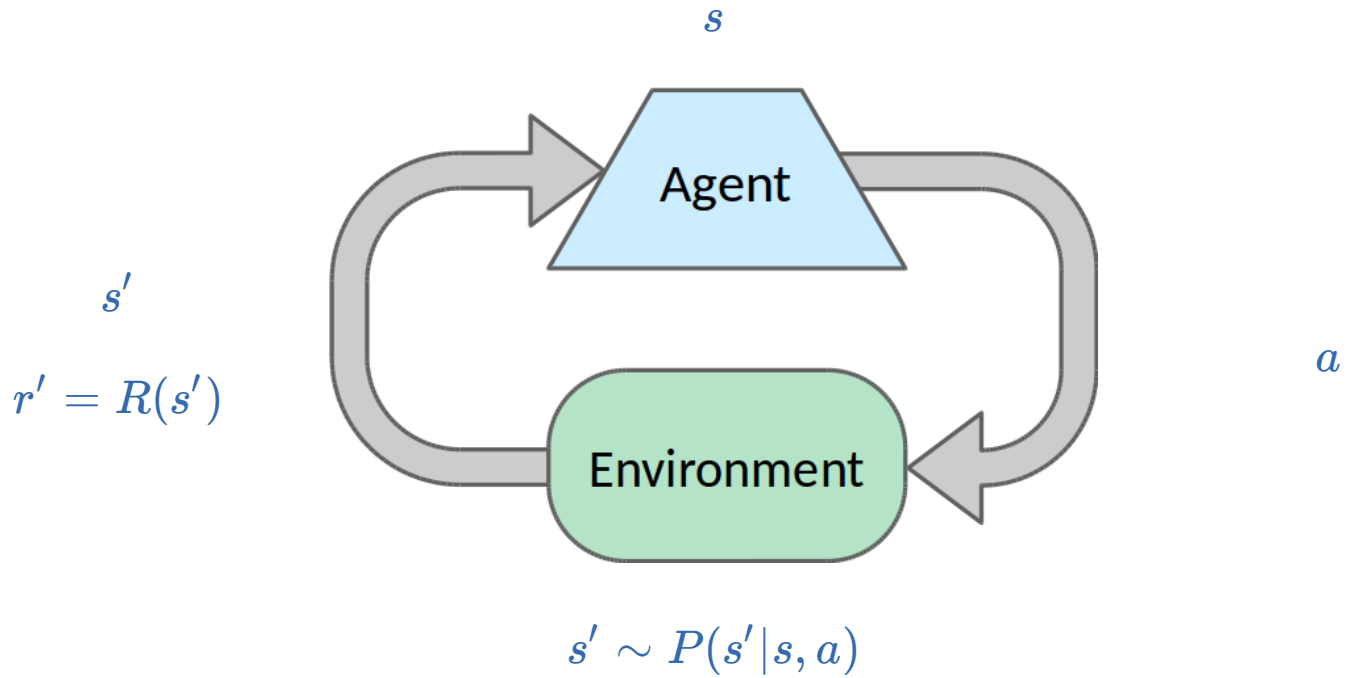


# Markov decision processes

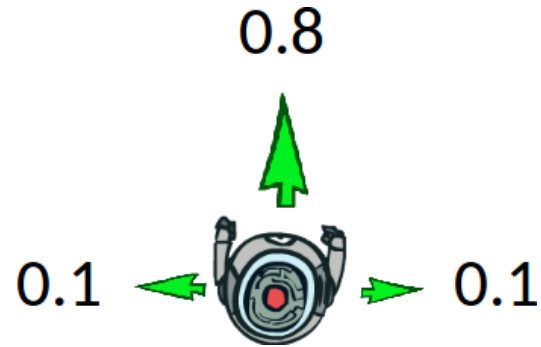
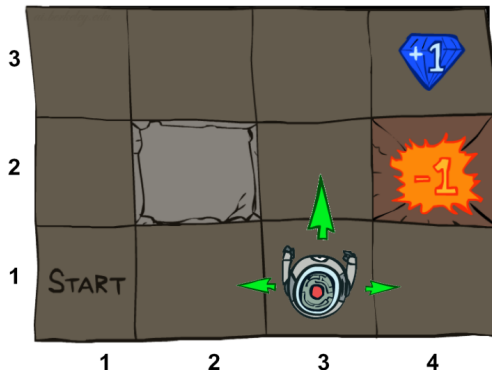
# Markov decision processes

A **Markov decision process** (MDP) is a tuple  $(\mathcal{S}, \mathcal{A}, P, R)$  such that:

- $\mathcal{S}$  is a set of states  $s$ ;
- $\mathcal{A}$  is a set of actions  $a$ ;
- $P$  is a (stationary) transition model such that  $P(s' | s, a)$  denotes the probability of reaching state  $s'$  if action  $a$  is done in state  $s$ ;
- $R$  is a reward function that maps immediate (finite) reward values  $R(s)$  obtained in states  $s$ .







## Example

- $\mathcal{S}$ : locations  $(i, j)$  on the grid.
- $\mathcal{A}$ : [Up, Down, Right, Left].
- Transition model:  $P(s' | s, a)$
- Reward:

$$R(s) = \begin{cases} -0.3 & \text{for non-terminal states} \\ \pm 1 & \text{for terminal states} \end{cases}$$

## What is Markovian about MDPs?

Given the present state, the future and the past are independent:

$$P(s_{t+1} | s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_0) = P(s_{t+1} | s_t, a_t)$$

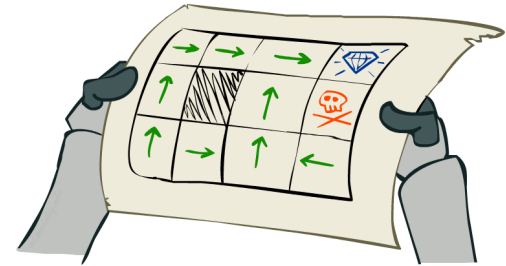
This is similar to search problems, where the successor function could only depend on the current state.



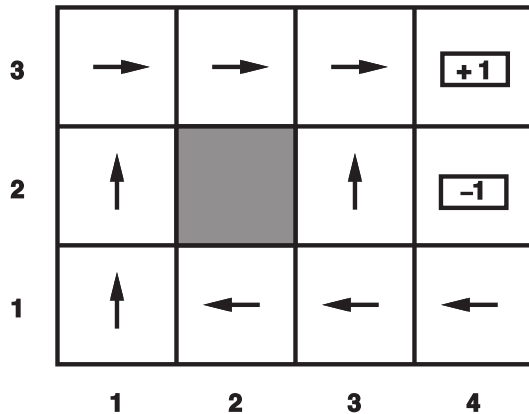
*Andrey Markov*

# Policies

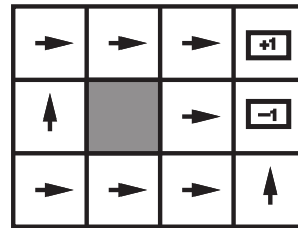
- In deterministic single-agent search problems, our goal was to find an optimal plan, or **sequence** of actions, from start to goal.
- For MDPs, we want to find an optimal **policy**  $\pi^* : \mathcal{S} \rightarrow \mathcal{A}$ .
  - A policy  $\pi$  maps states to actions.
  - An optimal policy is one that maximizes the expected utility, e.g. the expected sum of rewards.
  - An explicit policy defines a reflex agent.
- Expectiminimax did not compute entire policies, but only some action for a single state.



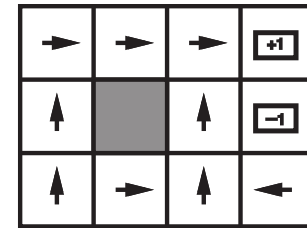
Optimal policy when  $R(s) = -0.3$  for all non-terminal states  $s$ .



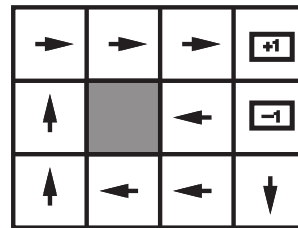
(a)



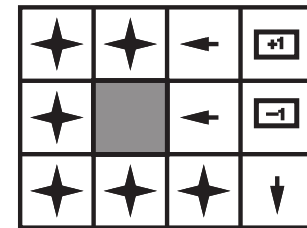
$$R(s) < -1.6284$$



$$-0.4278 < R(s) < -0.0850$$



$$-0.0221 < R(s) < 0$$



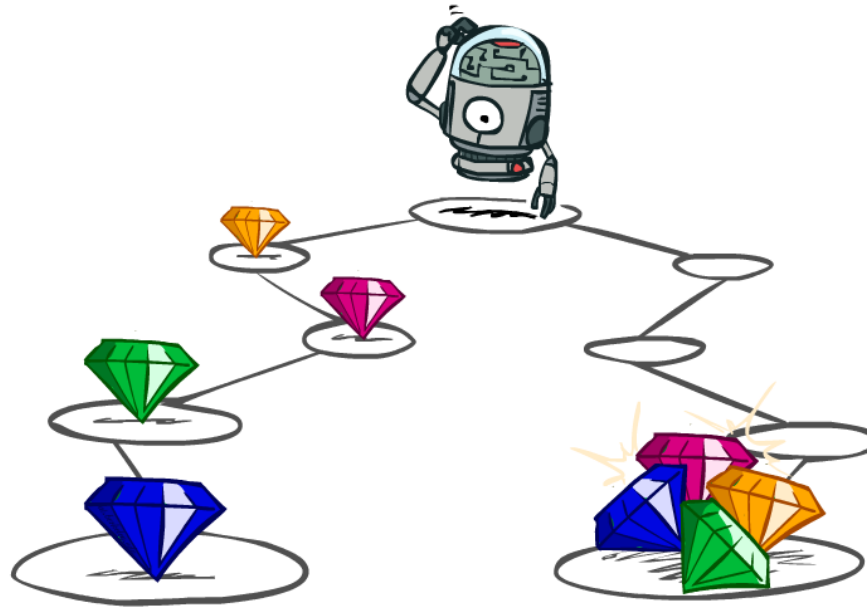
$$R(s) > 0$$

(b)

(a) Optimal policy when  $R(s) = -0.04$  for all non-terminal states  $s$ . (b) Optimal policies for four different ranges of  $R(s)$ .

Depending on  $R(s)$ , the balance between risk and reward changes from risk-taking to very conservative.

# Utilities over time



What preferences should an agent have over state or reward sequences?

- More or less?  $[2, 3, 4]$  or  $[1, 2, 2]$ ?
- Now or later?  $[1, 0, 0]$  or  $[0, 0, 1]$ ?

## Theorem

If we assume **stationary** preferences over reward sequences, i.e. such that

$$[r_0, r_1, r_2, \dots] \succ [r_0, r'_1, r'_2, \dots] \Rightarrow [r_1, r_2, \dots] \succ [r'_1, r'_2, \dots],$$

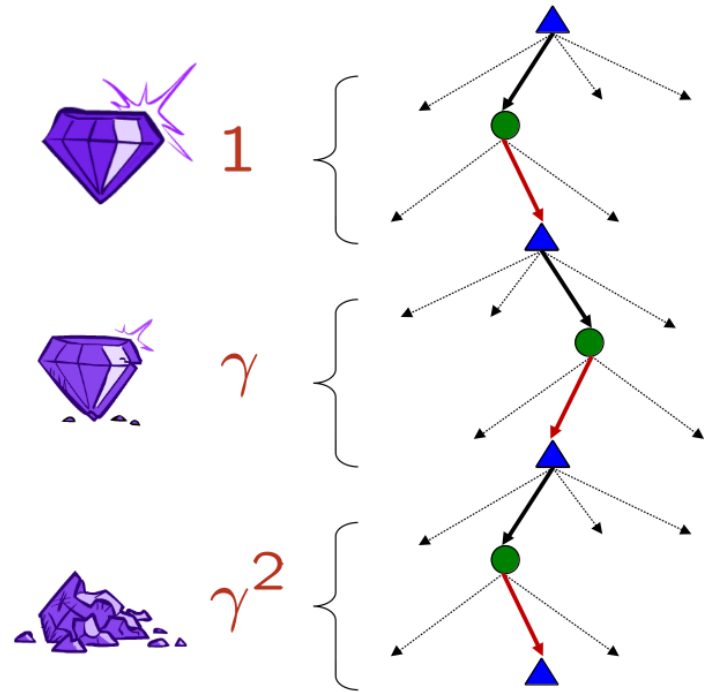
then there are only two coherent ways to assign utilities to sequences:

Additive utility:  $V([r_0, r_1, r_2, \dots]) = r_0 + r_1 + r_2 + \dots$

Discounted utility:  $V([r_0, r_1, r_2, \dots]) = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$   
( $0 < \gamma < 1$ )

# Discounting

- Each time we transition to the next state, we multiply in the discount once.
- Why discount?
  - Sooner rewards probably do have higher utility than later rewards.
  - Will help our algorithms converge.



Example: discount  $\gamma = 0.5$

- $V([1, 2, 3]) = 1 + 0.5 \times 2 + 0.25 \times 3$
- $V([1, 2, 3]) < V([3, 2, 1])$

## Infinite sequences

What if the agent lives forever? Do we get infinite rewards? Comparing reward sequences with  $+\infty$  utility is problematic.

Solutions:

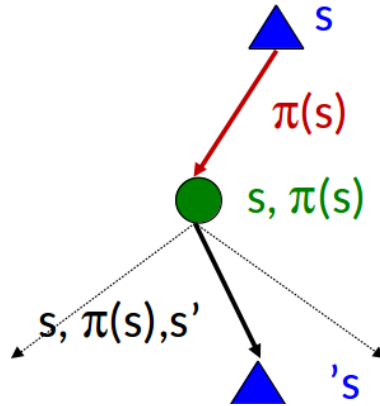
- Finite horizon: (similar to depth-limited search)
  - Terminate episodes after a fixed number of steps  $T$ .
  - Results in non-stationary policies ( $\pi$  depends on time left).
- Discounting (with  $0 < \gamma < 1$  and rewards bounded by  $\pm R_{\max}$ ):

$$V([r_0, r_1, \dots, r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_t \leq \frac{R_{\max}}{1 - \gamma}$$

Smaller  $\gamma$  results in a shorter horizon.

- Absorbing state: guarantee that for every policy, a terminal state will eventually be reached.





## Policy evaluation

The expected utility obtained by executing  $\pi$  starting in  $s$  is given by

$$V^\pi(s) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \right] \Big|_{s_0=s}$$

where the expectation is with respect to the probability distribution over state sequences determined by  $s$  and  $\pi$ .

## Optimal policies

Among all policies the agent could execute, the **optimal policy** is the policy  $\pi_s^*$  that maximizes the expected utility:

$$\pi_s^* = \arg \max_{\pi} V^{\pi}(s)$$

Because of discounted utilities, the optimal policy is **independent** of the starting state  $s$  (see later). Therefore we simply write  $\pi^*$ .

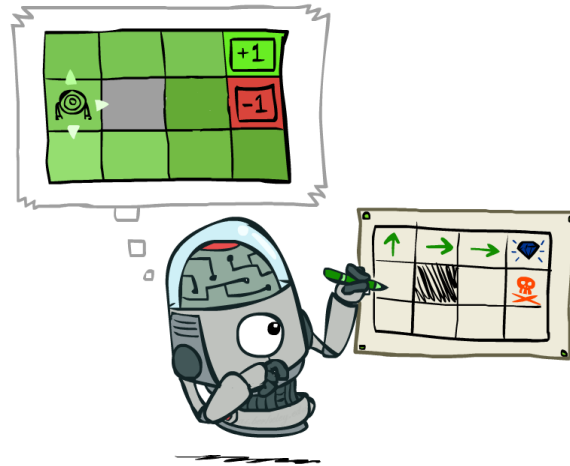
# Values of states

The utility, or value,  $V(\mathbf{s})$  of a state is now simply defined as  $V^{\pi^*}(\mathbf{s})$ .

- That is, the expected (discounted) reward if the agent executes an optimal policy starting from  $\mathbf{s}$ .
- Notice that  $R(\mathbf{s})$  and  $V(\mathbf{s})$  are quite different quantities:
  - $R(\mathbf{s})$  is the short term reward for having reached  $\mathbf{s}$ .
  - $V(\mathbf{s})$  is the long term total reward from  $\mathbf{s}$  onward.

<b>3</b>	<b>0.812</b>	<b>0.868</b>	<b>0.918</b>	<b>+1</b>
<b>2</b>	<b>0.762</b>		<b>0.660</b>	<b>-1</b>
<b>1</b>	<b>0.705</b>	<b>0.655</b>	<b>0.611</b>	<b>0.388</b>
	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>

Utilities of the states in Grid World, calculated with  $\gamma = 1$  and  $R(s) = -0.04$  for non-terminal states.



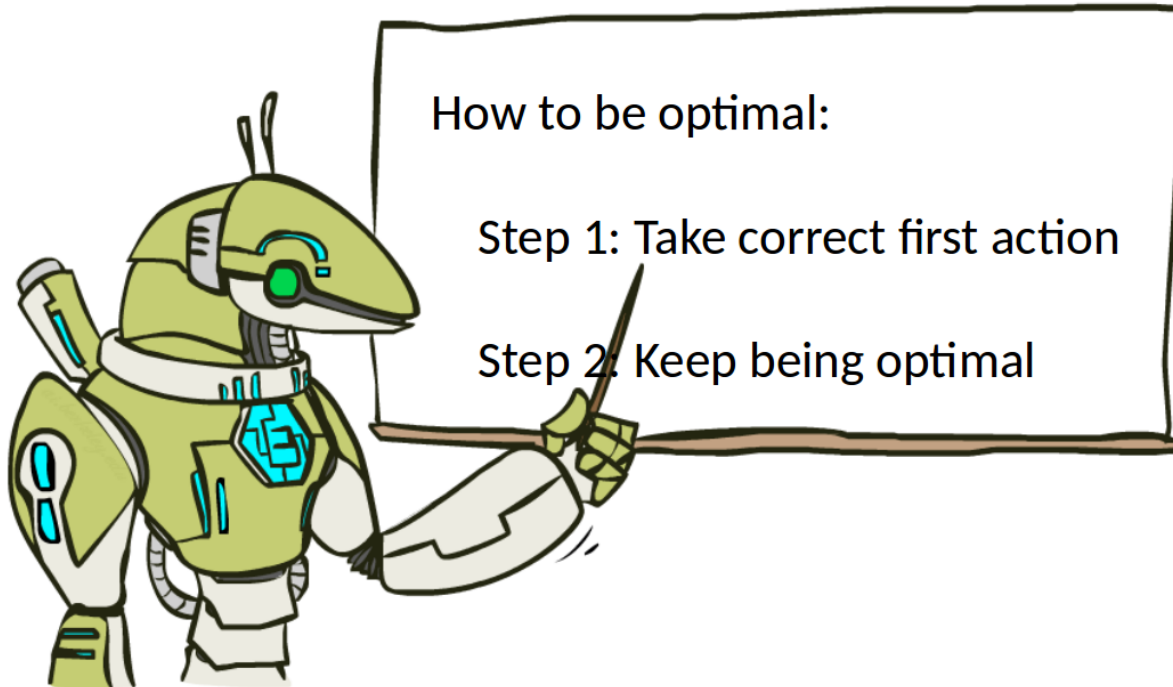
## Policy extraction

Using the principle of maximum expected utility, the optimal action maximizes the expected utility of the subsequent state. That is,

$$\pi^*(s) = \arg \max_a \sum_{s'} P(s'|s, a) V(s').$$

Therefore, we can extract the optimal policy provided we can estimate the utilities of states.

$$\pi^*(s) = \arg \max_a \sum_{s'} P(s'|s, a) V(s')$$



# The Bellman equation

The utility of a state is the immediate reward for that state, plus the expected discounted utility of the next state, assuming that the agent chooses the optimal action:

$$V(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V(s').$$

- These equations are called the **Bellman equations**. They form a system of  $n = |\mathcal{S}|$  non-linear equations with as many unknowns.
- The utilities of states, defined as the expected utility of subsequent state sequences, are solutions of the set of Bellman equations.

## Example

$$V(1, 1) = -0.04 + \gamma \max[0.8V(1, 2) + 0.1V(2, 1) + 0.1V(1, 1), \\ 0.9V(1, 1) + 0.1V(1, 2), \\ 0.9V(1, 1) + 0.1V(2, 1), \\ 0.8V(2, 1) + 0.1V(1, 2) + 0.1V(1, 1)]$$



# Value iteration

Because of the **max** operator, the Bellman equations are non-linear and solving the system is problematic.

The **value iteration** algorithm provides a fixed-point iteration procedure for computing the state utilities  $V(\mathbf{s})$ :

- Let  $V_i(\mathbf{s})$  be the estimated utility value for  $\mathbf{s}$  at the  $i$ -th iteration step.
- The **Bellman update** consists in updating simultaneously all the estimates to make them **locally consistent** with the Bellman equation:

$$V_{i+1}(\mathbf{s}) := R(\mathbf{s}) + \gamma \max_a \sum_{s'} P(s' | \mathbf{s}, a) V_i(s')$$

- Repeat until convergence.

**function** VALUE-ITERATION( $mdp, \epsilon$ ) **returns** a utility function

**inputs:**  $mdp$ , an MDP with states  $S$ , actions  $A(s)$ , transition model  $P(s' | s, a)$ ,  
rewards  $R(s)$ , discount  $\gamma$

$\epsilon$ , the maximum error allowed in the utility of any state

**local variables:**  $U, U'$ , vectors of utilities for states in  $S$ , initially zero

$\delta$ , the maximum change in the utility of any state in an iteration

**repeat**

$U \leftarrow U'; \delta \leftarrow 0$

**for each** state  $s$  **in**  $S$  **do**

$U'[s] \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s' | s, a) U[s']$

**if**  $|U'[s] - U[s]| > \delta$  **then**  $\delta \leftarrow |U'[s] - U[s]|$

**until**  $\delta < \epsilon(1 - \gamma)/\gamma$

**return**  $U$

(Step-by-step code example)

## Convergence

Let  $V_i$  and  $V_{i+1}$  be successive approximations to the true utility  $V$ .

**Theorem.** For any two approximations  $V_i$  and  $V'_i$ ,

$$\|V_{i+1} - V'_{i+1}\|_\infty \leq \gamma \|V_i - V'_i\|_\infty.$$

- That is, the Bellman update is a contraction by a factor  $\gamma$  on the space of utility vector.
- Therefore, any two approximations must get closer to each other, and in particular any approximation must get closer to the true  $V$ .

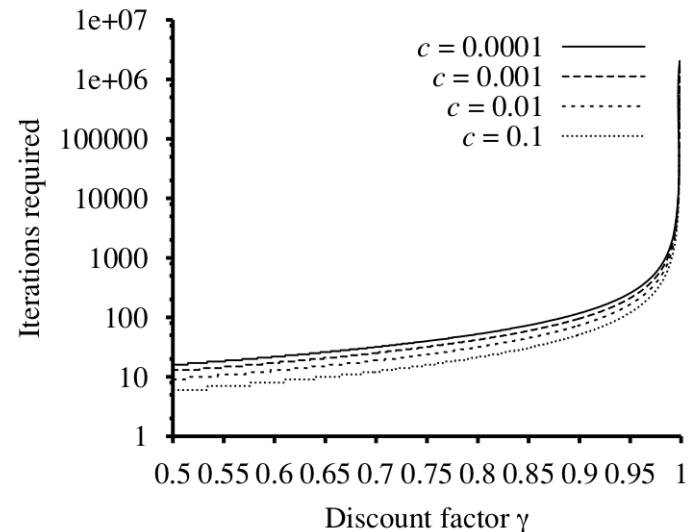
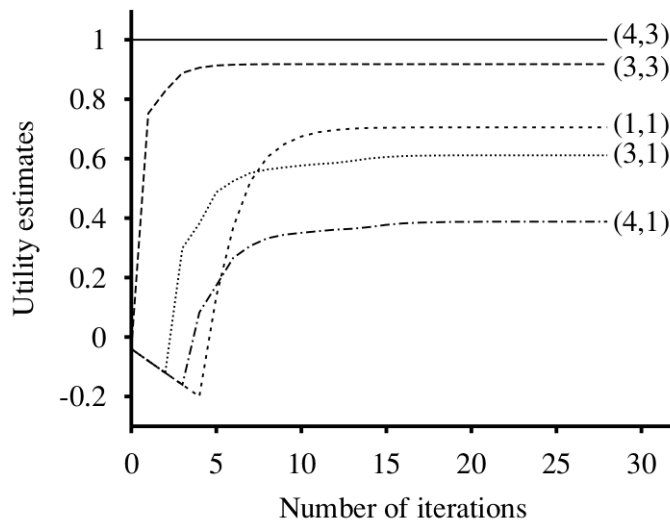
$\Rightarrow$  Value iteration always converges to a unique solution of the Bellman equations whenever  $\gamma < 1$ .

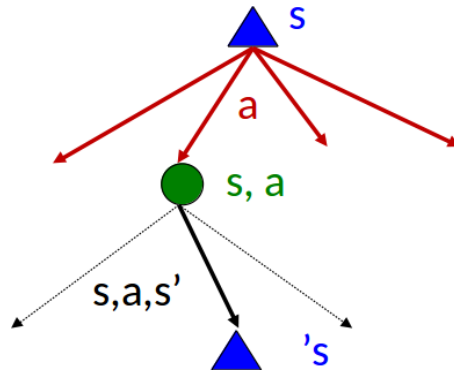
# Performance

Since  $\|V_{i+1} - V\|_\infty \leq \gamma \|V_i - V\|_\infty$ , the error is reduced by a factor of at least  $\gamma$  at each iteration.

Therefore, value iteration converges exponentially fast:

- The maximum initial error is  $\|V_0 - V\|_\infty \leq 2R_{\max}/(1 - \gamma)$ .
- To reach an error of at most  $\epsilon$  after  $N$  iterations, we require  $\gamma^N 2R_{\max}/(1 - \gamma) \leq \epsilon$ .





## Problems with value iteration

Value iteration repeats the Bellman updates:

$$V_{i+1}(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V_i(s')$$

- Problem 1: it is slow –  $O(|\mathcal{S}|^2 |\mathcal{A}|)$  per iteration.
- Problem 2: the **max** at each state rarely changes.
- Problem 3: the policy  $\pi_i$  extracted from the estimate  $V_i$  might be optimal even if  $V_i$  is inaccurate!

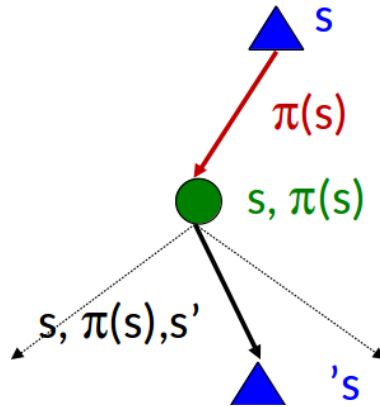
# Policy iteration

The **policy iteration** algorithm instead directly computes the policy (instead of state values). It alternates the following two steps:

- Policy evaluation: given  $\pi_i$ , calculate  $V_i = V^{\pi_i}$ , i.e. the utility of each state if  $\pi_i$  is executed.
- Policy improvement: calculate a new policy  $\pi_{i+1}$  using one-step look-ahead based on  $V_i$ :

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} P(s'|s, a) V_i(s')$$

This algorithm is still optimal, and might converge (much) faster under some conditions.



## Policy evaluation

At the  $i$ -th iteration we have a simplified version of the Bellman equations that relate the utility of  $s$  to the utilities of its neighbors:

$$V_i(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) V_i(s')$$

These equations are now **linear** because the **max** operator has been removed.

- for  $n$  states, we have  $n$  equations with  $n$  unknowns;
- this can be solved exactly in  $O(n^3)$  by standard linear algebra methods.



In some cases  $O(n^3)$  is too prohibitive. Fortunately, it is not necessary to perform exact policy evaluation. An approximate solution is sufficient.

One way is to run  $k$  iterations of simplified Bellman updates:

$$V_{i+1}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) V_i(s')$$

This hybrid algorithm is called **modified policy iteration**.

**function** POLICY-ITERATION( $mdp$ ) **returns** a policy

**inputs:**  $mdp$ , an MDP with states  $S$ , actions  $A(s)$ , transition model  $P(s' | s, a)$

**local variables:**  $U$ , a vector of utilities for states in  $S$ , initially zero

$\pi$ , a policy vector indexed by state, initially random

**repeat**

$U \leftarrow$  POLICY-EVALUATION( $\pi, U, mdp$ )

$unchanged? \leftarrow$  true

**for each** state  $s$  **in**  $S$  **do**

**if**  $\max_{a \in A(s)} \sum_{s'} P(s' | s, a) U[s'] > \sum_{s'} P(s' | s, \pi[s]) U[s']$  **then do**

$\pi[s] \leftarrow \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s' | s, a) U[s']$

$unchanged? \leftarrow$  false

**until**  $unchanged?$

**return**  $\pi$

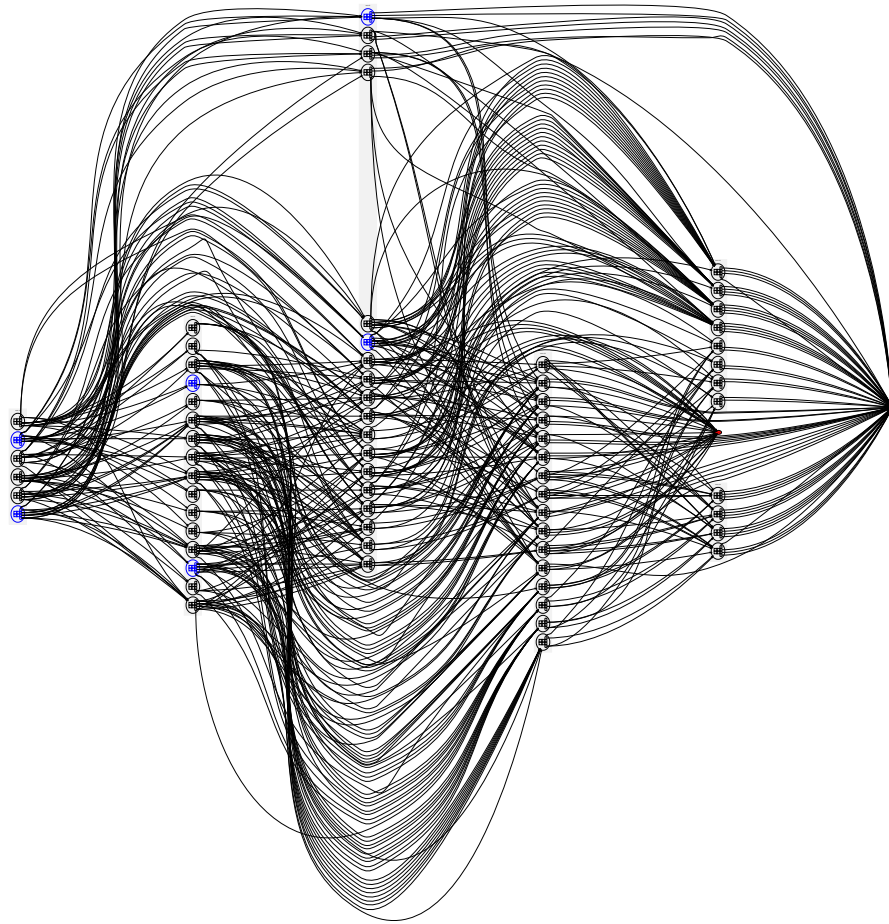
(Step-by-step code example)

# Recap example: 2048

The game 2048 is a Markov decision process!

- $\mathcal{S}$ : all possible configurations of the board (huge!)
- $\mathcal{A}$ : swiping left, right, up or down.
- $P(s'|s, a)$ : encodes the game's dynamic
  - collapse matching tiles
  - place a random tile on the board
- $R(s) = 1$  if  $s$  is a winning state, and 0 otherwise.

256	2		
16	8		2
64	2	4	2



The transition model for a  $2 \times 2$  board and a winning state at 8.

Optimal play for a **3** × **3** grid and a winning state at **1024**.

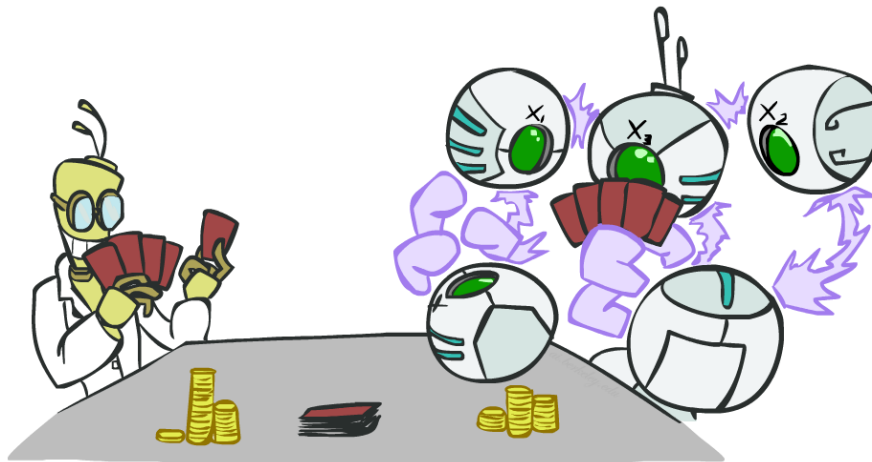
See [jd1m.info](http://jd1m.info): The Mathematics of 2048.

# Partially observable Markov decision processes

# POMDPs

What if the environment is only **partially observable**?

- The agent does not know in which state  $s$  it is in.
  - Therefore, it cannot evaluate the reward  $R(s)$  associated to the unknown state.
  - Also, it makes no sense to talk about a policy  $\pi(s)$ .
- Instead, the agent collects percepts  $e$  through a sensor model  $P(e|s)$ , from which it can reason about the unknown state  $s$ .

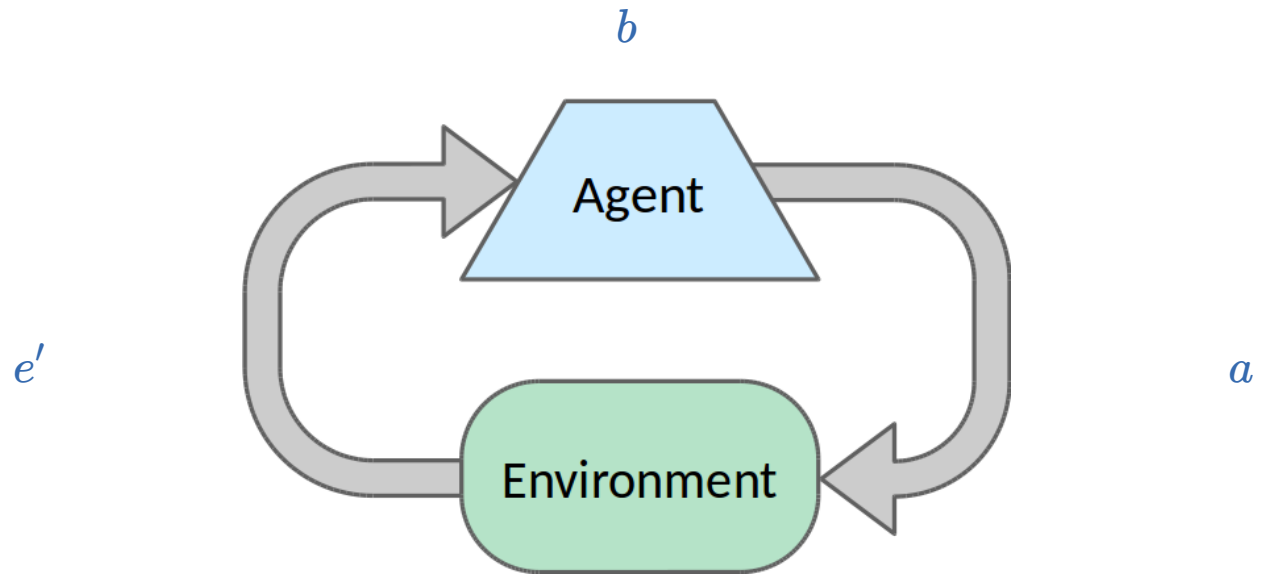




We will assume that the agent maintains a belief state  $b$ .

- $b$  represents a probability distribution  $\mathbf{P}(S)$  of the current agent's beliefs over its state;
- $b(s)$  denotes the probability  $P(S = s)$  under the current belief state;
- the belief state  $b$  is updated as evidence  $e$  are collected.

This is filtering!



$$s' \sim P(s'|s, a)$$

$$e' \sim P(e'|s')$$

# Belief MDP

**Theorem (Astrom, 1965).** The optimal action depends only on the agent's current belief state.

- The optimal policy can be described by a mapping  $\pi^*(b)$  from beliefs to actions.
- It does not depend on the actual state the agent is in.

In other words, POMDPs can be reduced to an MDP in belief-state space, provided we can define a transition model  $P(b'|b, a)$  and a reward function  $\rho$  over belief states.

If  $b$  was the previous belief state and the agent does action  $a$  and perceives  $e$ , then the new belief state over  $S'$  is given by

$$b' = \alpha \mathbf{P}(e|S') \sum_s \mathbf{P}(S'|s, a) b(s) = \alpha \text{forward}(b, a, e).$$

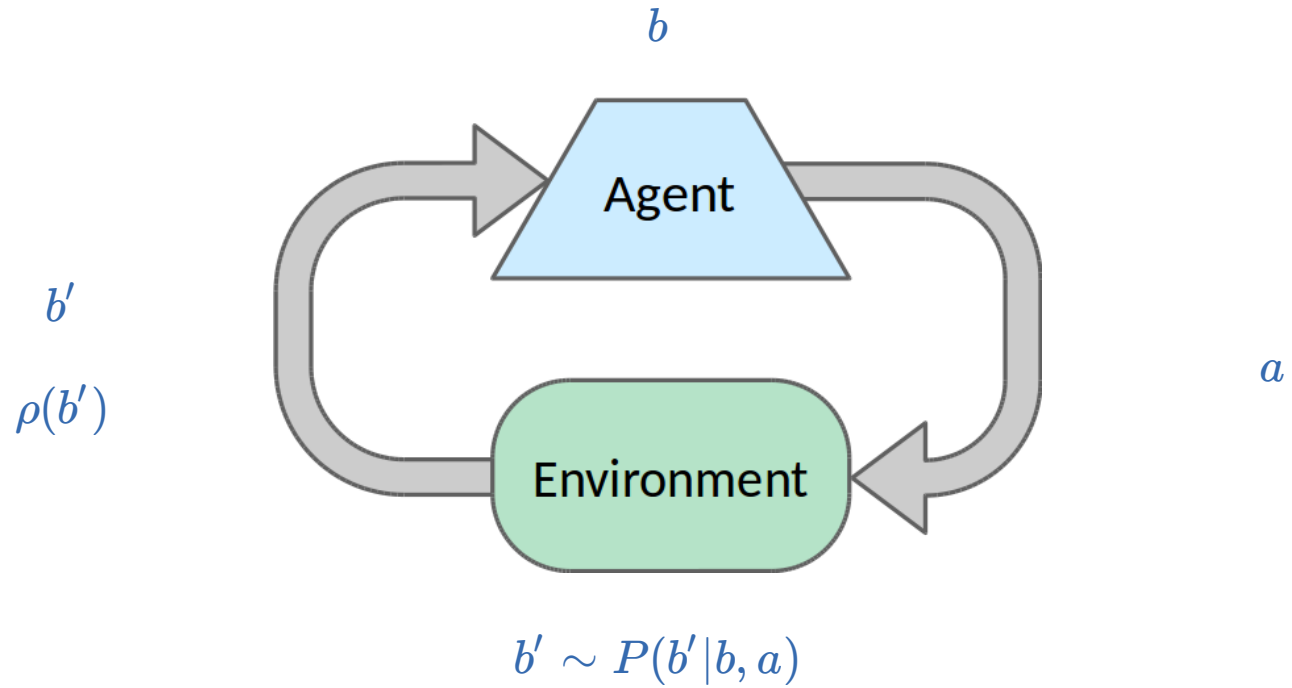
Therefore,

$$\begin{aligned} P(b'|b, a) &= \sum_e P(b', e|b, a) \\ &= \sum_e P(b'|b, a, e) P(e|b, a) \\ &= \sum_e P(b'|b, a, e) \sum_{s'} P(e|b, a, s') P(s'|b, a) \\ &= \sum_e P(b'|b, a, e) \sum_{s'} P(e|s') \sum_s P(s'|s, a) b(s) \end{aligned}$$

where  $P(b'|b, a, e) = 1$  if  $b' = \text{forward}(b, a, e)$  and  $0$  otherwise.

We can also define a reward function for belief states as the expected reward for the actual state the agent might be in:

$$\rho(b) = \sum_s b(s)R(s)$$



Although we have reduced POMDPs to MDPs, the Belief MDP we obtain has a **continuous** (and usually high-dimensional) state space.

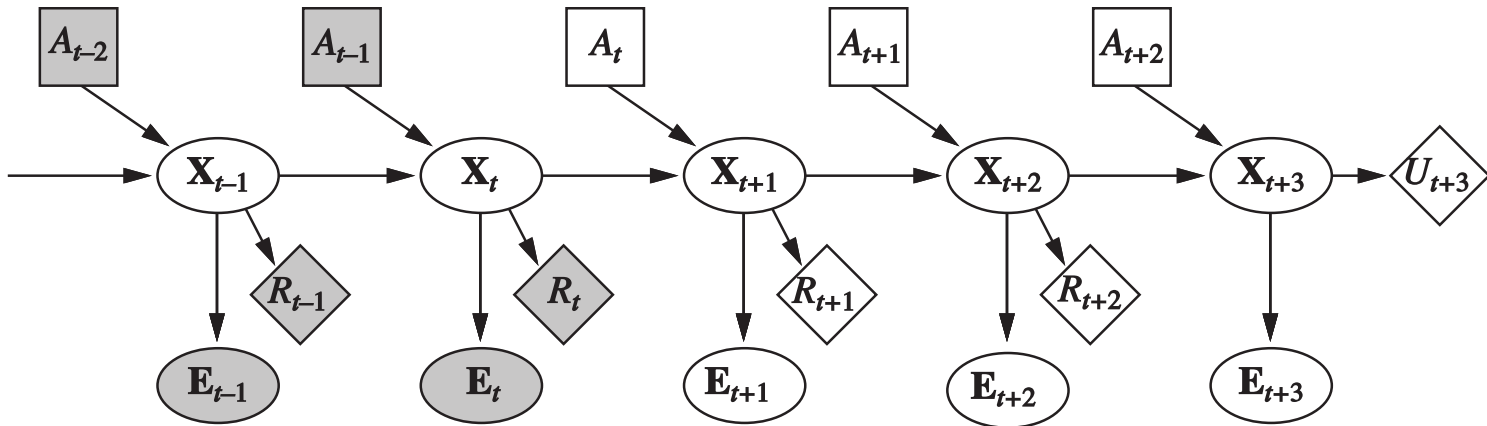
- None of the algorithms described earlier directly apply.
- In fact, solving POMDPs remains a difficult problem for which there is no known efficient exact algorithm.
- Yet, Nature is a POMDP.

# Online agents

While it is difficult to directly derive  $\pi^*$ , a decision-theoretic agent can be constructed for POMDPs:

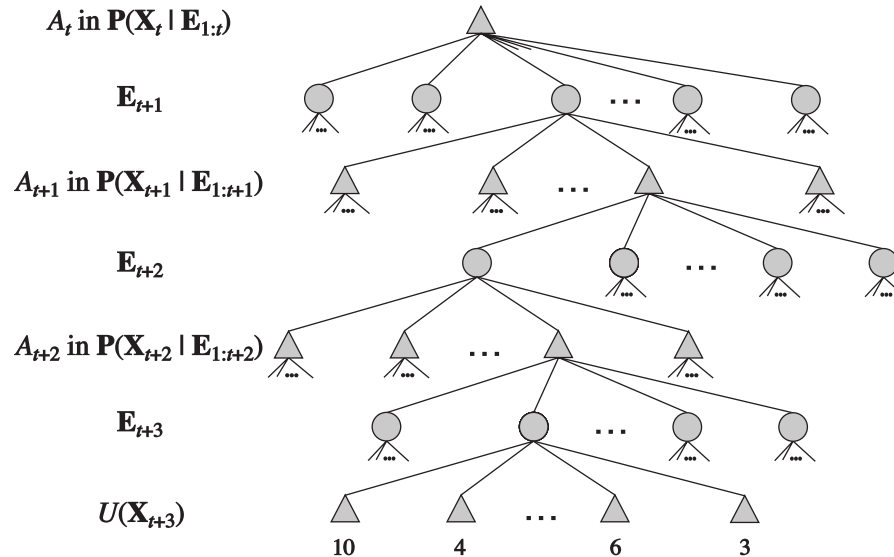
- The transition and sensor models are represented by a [dynamic Bayesian network](#);
- The dynamic Bayesian network is extended with decision ( $A$ ) and utility ( $R$  and  $U$ ) nodes to form a dynamic decision network;
- A [filtering algorithm](#) is used to incorporate each new percept and action and to update the belief state representation;
- Decisions are made by projecting forward possible action sequences and choosing (approximately) the best one, in a manner similar to a truncated [Expectiminimax](#).





At time  $t$ , the agent must decide what to do.

- Shaded nodes represent variables with known values.
- The network is unrolled for a finite horizon.
- It includes nodes for the reward of  $\mathbf{X}_{t+1}$  and  $\mathbf{X}_{t+2}$ , but the (estimated) utility of  $\mathbf{X}_{t+3}$ .



Part of the look-ahead solution of the previous decision network:

- Each triangular node is a belief state in which the agent makes a decision.
  - The belief state at each node can be computed by applying a filtering algorithm to the sequence of percepts and actions leading to it.
- The round nodes correspond to choices by the environment.

A decision can be extracted from the search tree by backing up the (estimated) utility values from the leaves, taking the average at the chance nodes and taking the maximum at the decision nodes.

# Summary

- Sequential decision problems in uncertain environments, called MDPs, are defined by transition model and a reward function.
- The utility of a state sequence is the sum of all the rewards over the sequence, possibly discounted over time.
  - The solution of an MDP is a policy that associates a decision with every state that the agent might reach.
  - An optimal policy maximizes the utility of the state sequence encountered when it is executed.
- Value iteration and policy iteration can both be used for solving MDPs.
- POMDPs are much more difficult than MDPs. However, a decision-theoretic agent can be constructed for those environments.

