

CS 482/682 – AI: Making Complex Decisions

Fall 2021 - Chapter 17

Making Decisions

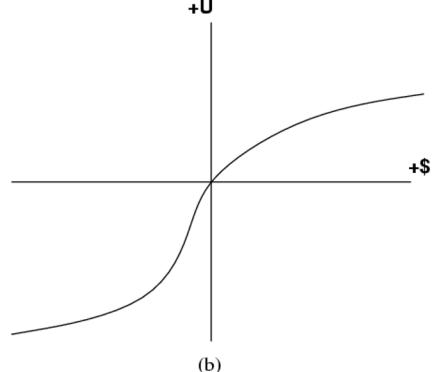
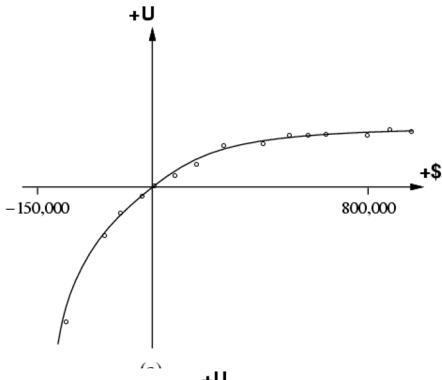
- Probability theory describes what an agent should believe on the basis of evidence
- Utility theory describes what an agent wants
- Decision theory puts the two together to describe what an agent should do

Expected Utility

- Consider a non-deterministic action A that has possible outcome states $\text{Result}_i(A)$
- Compute the probability of ending up in each of these outcome states given that you perform action A : $P(\text{Result}_i(A) | \text{Do}(A))$
- Weight the utility of each possible state by its probability and find the expected utility:

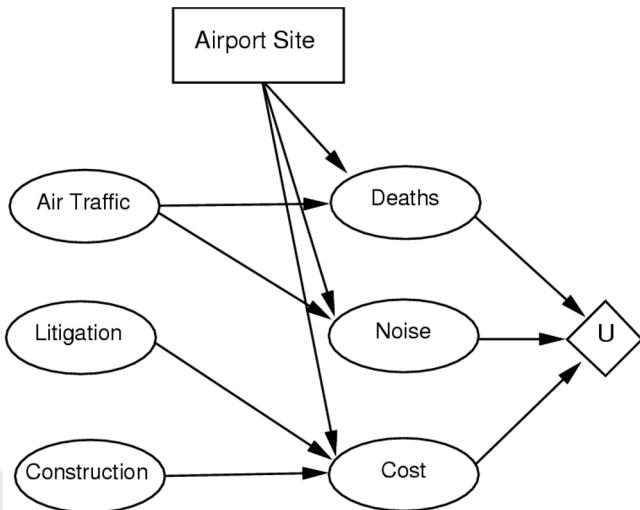
$$\text{EU}(A) = \sum_i P(\text{Result}_i(A) | \text{Do}(A)) \times U(\text{Result}_i(A))$$

The Utility of Money



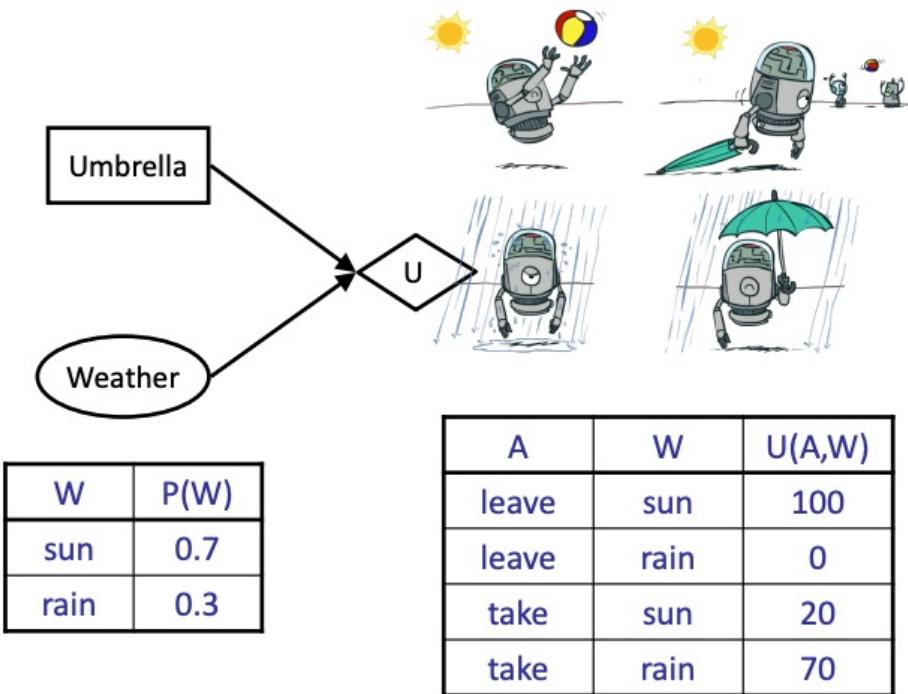
- Utility represents preferences... it is not necessarily logical
- Choices between
 - \$5
 - a coin flip for either \$10 or nothing
 - Expected value of each is \$5
 - \$1,000,000
 - A coin flip for either \$3,000,000 or nothing

Decision Networks



- Chance nodes: (ovals) represent random variables, just as in belief nets
 - Has a conditional probability table (CPT) that is indexed by the state of the parent nodes
- Decision nodes: (rectangles) represent a choice of actions
- Utility nodes: (diamonds) has as parents all those variables describing the outcome state that directly affect utility
 - Has an action-utility table that calculates the utility function based on the parents

Example



Umbrella = leave

$$\begin{aligned} \text{EU}(\text{leave}) &= \sum_w P(w)U(\text{leave}, w) \\ &= 0.7 \cdot 100 + 0.3 \cdot 0 = 70 \end{aligned}$$

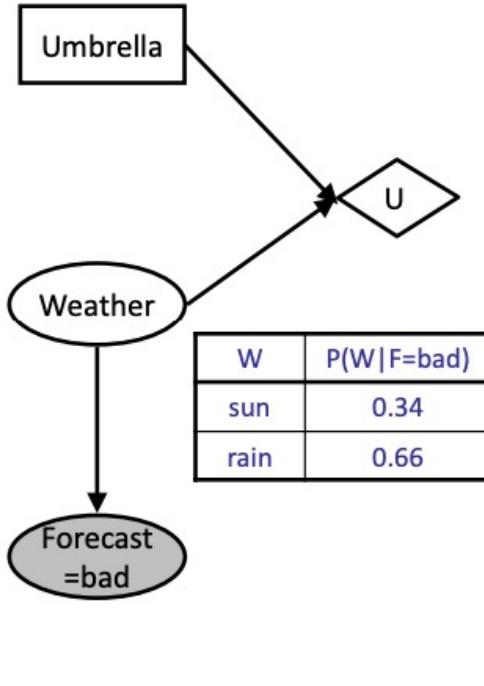
Umbrella = take

$$\begin{aligned} \text{EU}(\text{take}) &= \sum_w P(w)U(\text{take}, w) \\ &= 0.7 \cdot 20 + 0.3 \cdot 70 = 35 \end{aligned}$$

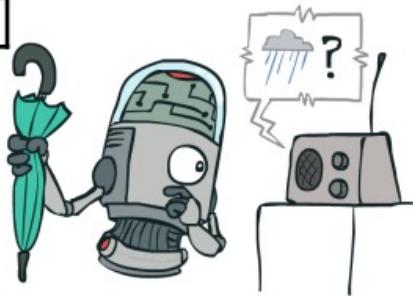
Optimal decision = leave

$$\text{MEU}(\emptyset) = \max_a \text{EU}(a) = 70$$

Example



A	W	U(A,W)
leave	sun	100
leave	rain	0
take	sun	20
take	rain	70



Umbrella = leave

$$\begin{aligned} \text{EU}(\text{leave}|\text{bad}) &= \sum_w P(w|\text{bad})U(\text{leave}, w) \\ &= 0.34 \cdot 100 + 0.66 \cdot 0 = 34 \end{aligned}$$

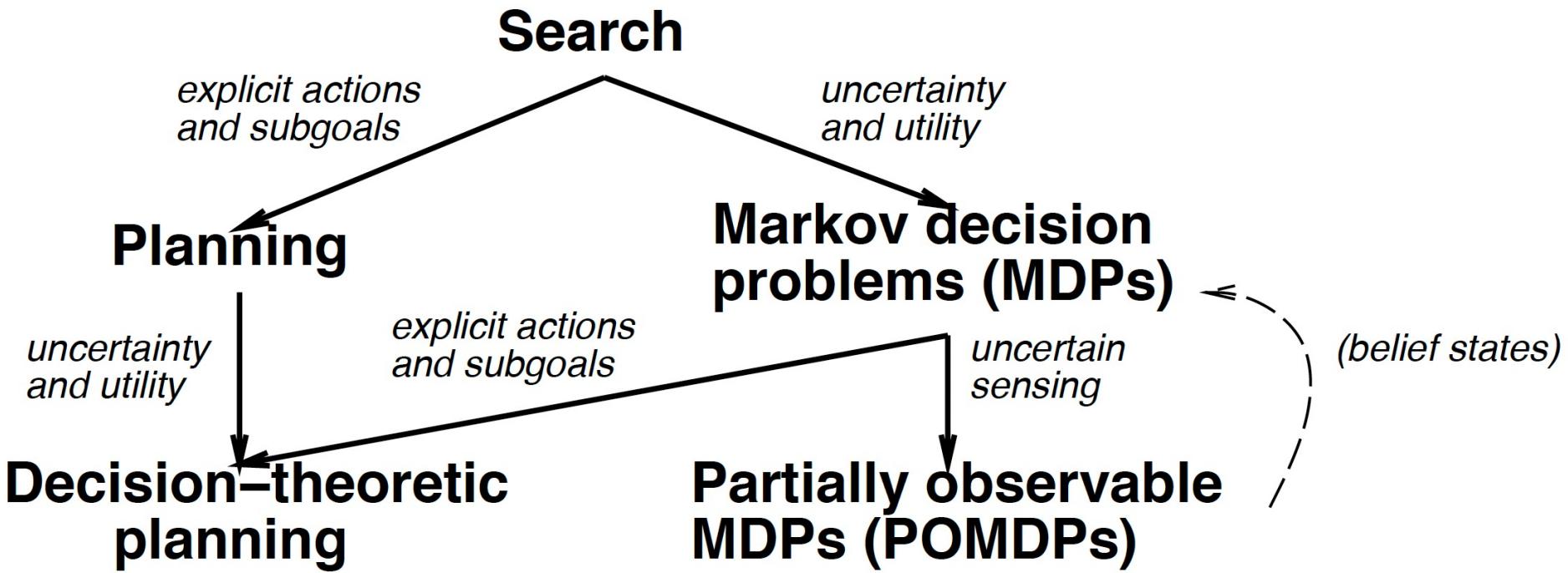
Umbrella = take

$$\begin{aligned} \text{EU}(\text{take}|\text{bad}) &= \sum_w P(w|\text{bad})U(\text{take}, w) \\ &= 0.34 \cdot 20 + 0.66 \cdot 70 = 53 \end{aligned}$$

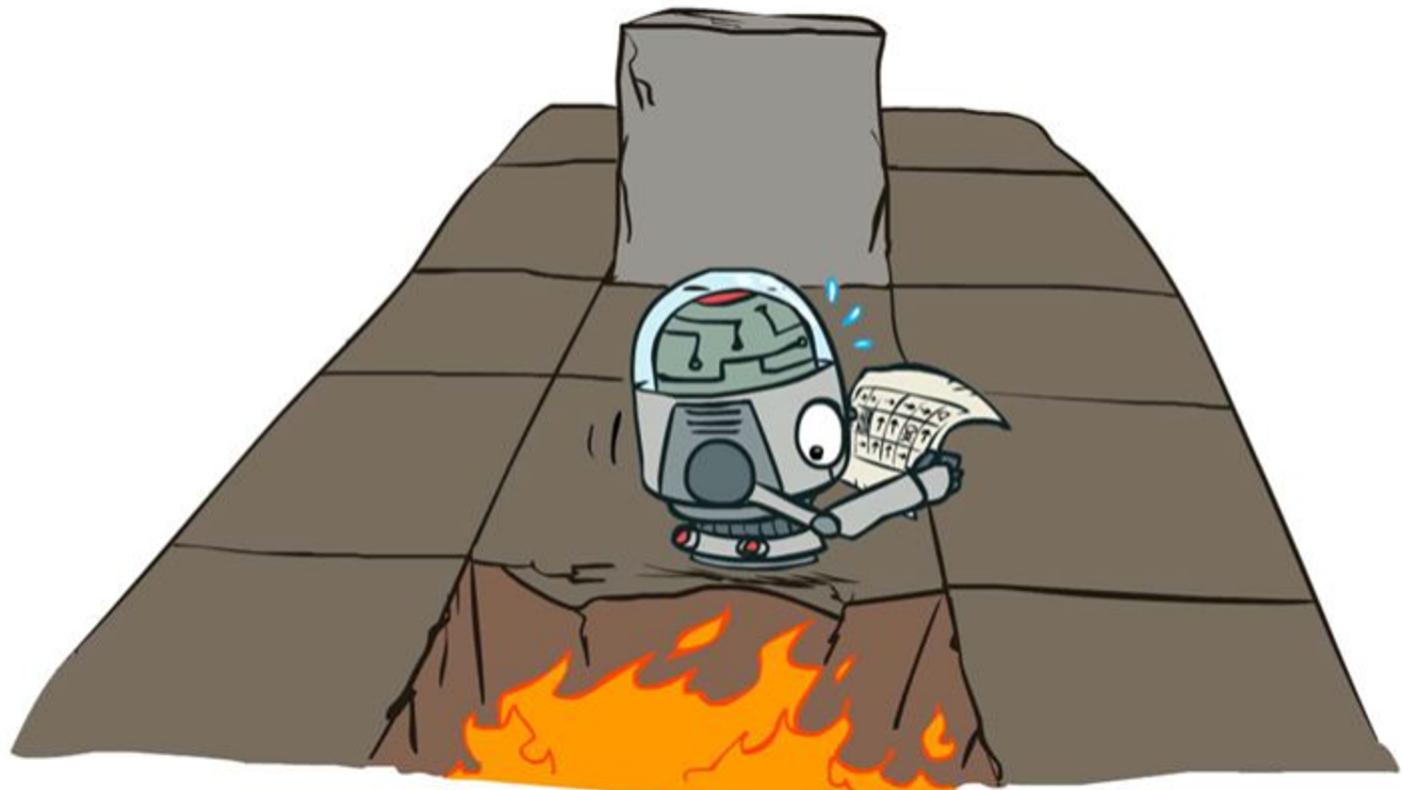
Optimal decision = take

$$\text{MEU}(F = \text{bad}) = \max_a \text{EU}(a|\text{bad}) = 53$$

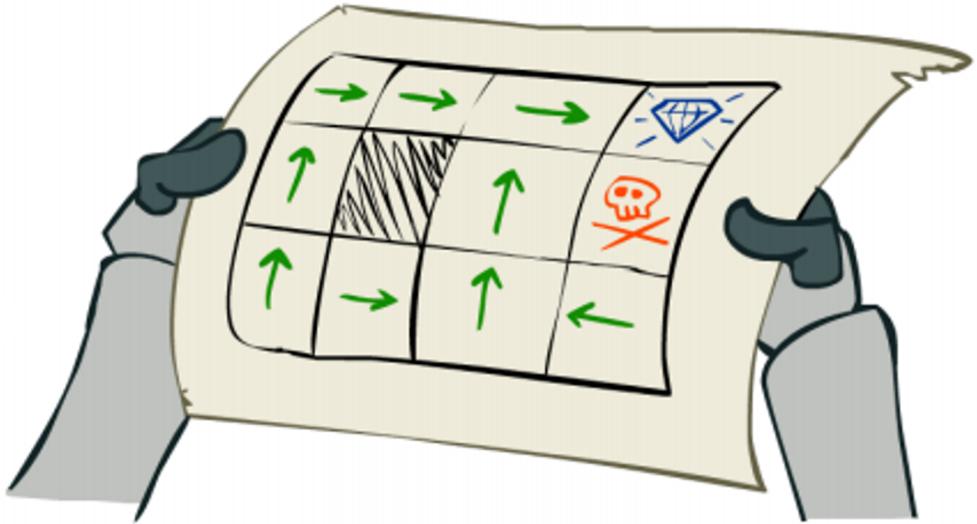
Sequential Design Problems



Sequential decision problems



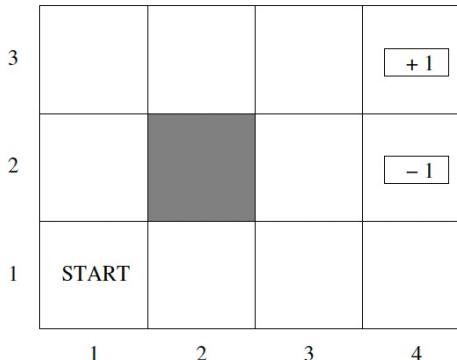
Sequential decision problems



Sequential decision problems

In **sequential decision problems** the utility of agent's actions do not depend on single decisions, expressed with the state, which the agent would have gotten into, as the result of this decision, but rather on the whole sequence of agent's action.

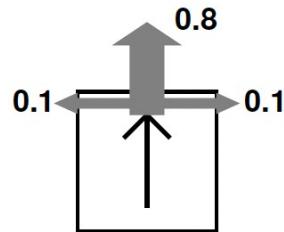
EXAMPLE: an agent is in the field START, and can move in any direction between the field. Its actions ends when it reaches one of the fields (4,2) or (4,3), with the result marked in those fields.



If the problem was fully deterministic — and the agent's knowledge of its position was complete — then the problem would be reduced to action planning. For example, for the above example the correct solution would be the action plan: U-U-R-R-R. Equally good would be the plan: R-R-U-U-R. If the single actions did not cost anything (ie. only the final state did matter), then equally good would also be the plan: R-R-R-L-L-L-U-U-R-R-R, and many others.

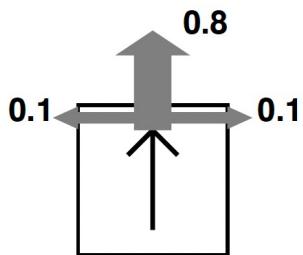
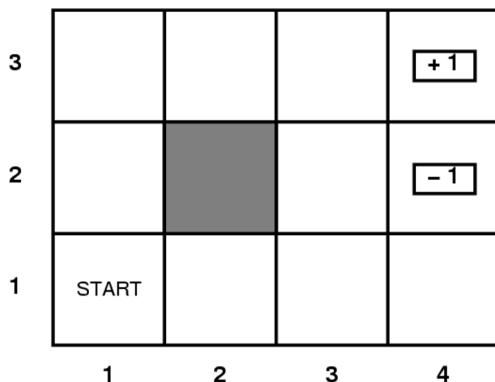
Sequential decision problems

However, considering the uncertainty, the result of the agent's actions corresponds with its intentions only with some probability. For example, we may assume that the action U (Up) transfers the agent to the desired position with probability 0.8, and with probabilities 0.1 takes the agent left or right. It is only certain that the agent will not end up in the direction opposite to the intended one. For simplicity let us also assume that the presence of the walls does not change this probability distribution, and only makes the agent stay in place, if it turned out that it should move into the wall.



We can now compute the expected values of the sequences of agent's moves. In general, there is no guarantee, that after executing any of the above sequences, the agent will indeed end up in the desired terminal state.

Sequential decision problems

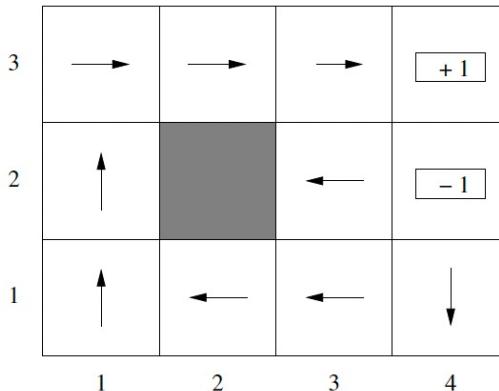


- Agent may move North, South, East, West
- Non-deterministic version
 - Moves succeed 80%
 - Deflected 90 degrees 20% (10% left and 10% right)
- Transition Model:
 - $M_{ij}^a = P(j|i, a)$ Probability that doing a in i leads to j

The agent's policy

As opposed to the action planning algorithms, here the agent should work out its strategy not as a specific sequence of actions, but as its **policy**, which is a scheme determining actions the agent should take for any specific state it would find itself in.

We can determine the optimal policy for the previous example problem. Note that at point (3,2) the policy makes the agent move left, which may seem wrong, but allows the agent to avoid ending up in state (4,2). Similarly in state (4,1).

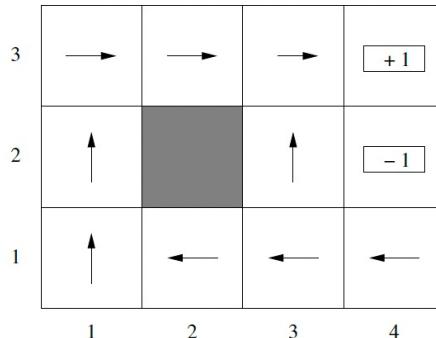


This policy obviously results from the assumption of zero cost of the moves. If the agent's outcome depended not only on the final state, but also on the number of moves, then such conservative policy would probably no longer be optimal.

Cost of the moves

Considering a positive cost of the moves, lowers the result achieved in the final states by the cumulative cost of all the moves. This certainly affects the agent's optimal policy.

For example, the following diagram shows the optimal policy when the cost of the moves equals $1/25$ of a unit. Let's note that in the states $(4,1)$ and $(3,2)$ the optimal policy dictates the move toward the state $(4,3)$, despite the risk. However, in the states $(2,1)$ and $(3,1)$ the policy still recommends going around.



Formally, the cost of the moves is introduced as a **reward** function R for states, so $R(s) = -0.04$ for all the nonterminal states.

Markov decision problems

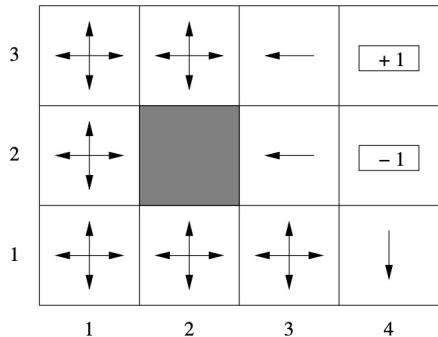
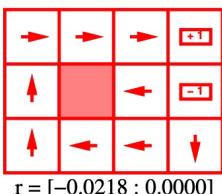
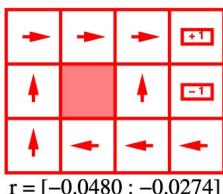
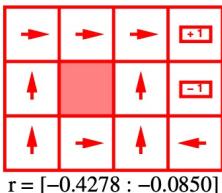
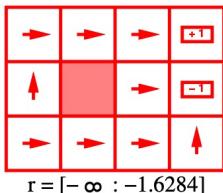
Formally, a Markov decision problem is defined by the following elements:

- the set of states with the starting state s_0 ,
- the set $\text{ACTIONS}(s)$ of actions allowed in state s ,
- transition model $P(s'|s, a)$,
- reward function $R(s)$ (also possible is: $R(s, a), R(s, a, s')$).

The solution to an MDP is the policy $\pi(s)$ mapping from states to actions. Let's note that under uncertainty each pass of the agent through the environment according to the policy may result in a different sequence of states, and possibly different outcome. The **optimal policy** $\pi^*(s)$ is the policy achieving the greatest expected utility.

Influence of rewards on the policy

Varying the reward values results in the changes of the optimal policy for a problem. With very high negative rewards the policy recommends going directly to the final state, regardless which one. With the rewards approaching zero the initial policy gradually returns.



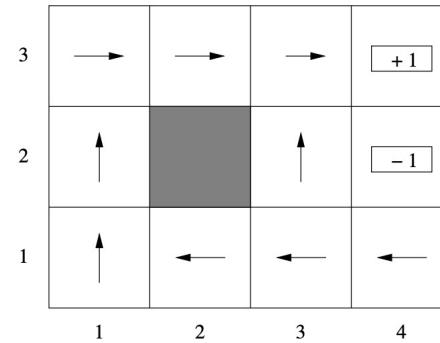
In case of positive rewards, it is no longer profitable for the agent to terminate the game, so the optimal policy tells it to avoid the terminal states. Executing actions is profitable, so the agent should keep running, and avoid termination.

Policy vs. Utility

In order to determine the optimal policy it would be useful to have state utilities, such as these marked in the diagram on the left (the question where these came from we shall postpone until later). We could then employ the MEU (Maximum Expected Utility) principle, and for each state designate the move, which maximizes the expected utility.

3	0.812	0.868	0.918	+ 1
2	0.762		0.660	- 1
1	0.705	0.655	0.611	0.388

1 2 3 4



However, in MDP problems states do not have utilities, except for the final states. The “utility” of a state (intermediate) depends on the agent’s policy, ie. what it intends to do in that state. At the same time, the agent’s policy depends on the “utilities” of the states.

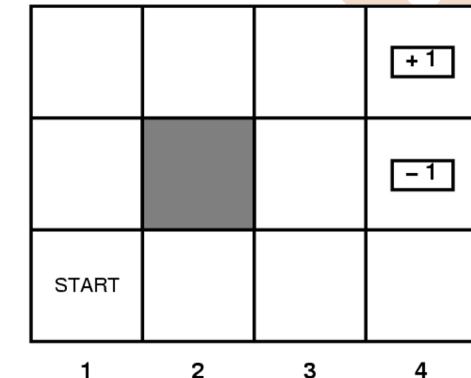
We can introduce state utilities based on policies.

Markov decision problem

Each time a given policy is executed starting from the initial state, the stochastic nature of the environment may lead to a different environment history. The quality of a policy is therefore measured by the *expected* utility of the possible environment histories generated by that policy. An **optimal policy** is a policy that yields the highest expected utility. We use π^* to denote an optimal policy. Given π^* , the agent decides what to do by consulting its current percept, which tells it the current state s , and then executing the action $\pi^*(s)$. A policy represents the agent function explicitly and is therefore a description of a simple reflex agent, computed from the information used for a utility-based agent.

The horizon problem

- The performance of the agent is measured by the sum of rewards for the transitions
 - $U_h(s_0, a_0, s_1, a_1, \dots, s_n)$
- Question: is there a finite or infinite horizon?
 - Finite: there is a fixed time N after which nothing matters
 - $U_h(s_0, a_0, s_1, a_1, \dots, s_{N+k}) = U_h(s_0, a_0, s_1, a_1, \dots, s_N)$ for all k
 - Example: $N=3$ and start at [3,1]
 - A policy that depends on time is called **nonstationary**
 - Infinite: no reason to behave differently in the same state at different time
 - Optimal action depends only on the current state → **stationary policy**



Discounting

- How to calculate the utility of state sequences?
 - $U_h(s_0, a_0, s_1, a_1, \dots, s_n) = R(s_0, a_0, s_1) + \gamma R(s_1, a_1, s_2) + \gamma^2 R(s_2, a_2, s_3) + \dots$
- λ is the **discount factor** between 0 and 1
 - γ close to 0: rewards in the distant future are viewed as insignificant
 - γ close to 1: an agent is willing to wait for long-term rewards
 - $\gamma = 1$: purely additive rewards
- For $\gamma < 1$ and $R \leq R_{max}$ the utilities are always finite

Practicality of discount factor

- Human appear to value near-term rewards more highly than rewards in the distant future
- Economic: if the rewards are monetary, then it really is better to get them sooner rather than later because early rewards can be invested and produce returns while you're waiting for the later rewards
 - Discount factor $\gamma \Leftrightarrow$ interest rate of $\frac{1}{\gamma} - 1$
 - $\gamma = 0.9 \Leftrightarrow$ interest rate of 11.1%
- True rewards may never arrive for all sorts of reasons that are not taken into account in the transition model
 - A discount factor of $\gamma \Leftrightarrow$ adding a probability $1 - \gamma$ of accidental termination at every step, independent of the action taken

Utilities of state sequences

A utility function for the sequences of states is called **separable**, if:

$$U([s_0, s_1, \dots, s_n]) = f(s_0, U([s_1, \dots, s_n]))$$

In our example 4×3 problem the utility function is separable since we can compute it using the formula:

$$U([s_0, s_1, \dots, s_n]) = R(s_0) + R(s_1) + \dots + R(s_n)$$

We call a utility function **additive** if it has the following property:

$$U([s_0, s_1, \dots, s_n]) = R(s_0) + U([s_1, \dots, s_n])$$

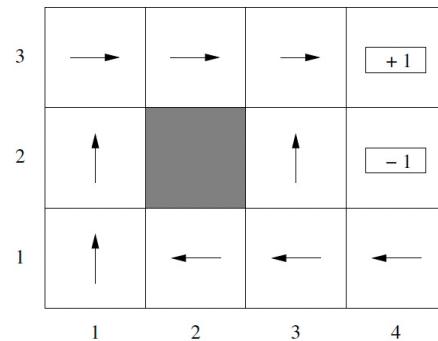
It turns out that for many practical cases the utility functions are additive. For example, when considering the cost function in the state space search, we implicitly assumed that they are additive. The incurred cost in a state was simply the cost in the previous state, plus the cost of the move.

Optimal policies

In order to determine the optimal policy it would be useful to have state utilities, such as these marked in the diagram on the left (the question where these came from we shall postpone until later). We could then employ the MEU (Maximum Expected Utility) principle, and for each state designate the move, which maximizes the expected utility.

3	0.812	0.868	0.918	+ 1
2	0.762		0.660	- 1
1	0.705	0.655	0.611	0.388

1 2 3 4



However, in MDP problems states do not have utilities, except for the final states. The “utility” of a state (intermediate) depends on the agent’s policy, ie. what it intends to do in that state. At the same time, the agent’s policy depends on the “utilities” of the states.

We can introduce state utilities based on policies.

Optimal policies

The utility of a state with respect to a policy can be defined as the expected value of the rewards obtained by initiating actions from this state:

$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \right]$$

By S_t we denote here the random variable signifying the state the agent will be at step t after starting from state s and executing the policy π .

It turns out that, even though theoretically the optimal policy $\pi^* = \operatorname{argmax}_\pi U^\pi(s)$ may depend on the choice of the starting state, for the decision processes with the Markov property, for infinite sequences with discounting, there is no such dependence. The agent's optimal policy, determining all her actions, is independent on the starting point.

For the utility of a state $U(s)$ we will take its utility computed with respect to the optimal policy $U^{\pi^*}(s)$.

Dynamic programming

The optimal policy π^* , being a function defined on the set of states, may then be associated with the (yet unknown) state utility function:

$$\pi^*(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a)U(s')$$

where $P(s'|s, a)$ is the probability that the agent will reach the state s' if she executes the action a in the state s .

Since we want to express the utility of a state as the expected value of a discounted sum of rewards of a sequence of states, then it can be tied to the utilities of its descendant states with the following equation (Bellman, 1957):

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

For the n states we thus obtain n equations — unfortunately nonlinear, due to the max term — with n unknowns. Solving these equations is called **dynamic programming**.

n-step decision problems

If in some problem the final states were achieved with known utilities in exactly n steps, then we could solve the Bellman equation by first determining the utilities for the states at step $n - 1$, then at step $n - 2$, etc., until reaching the start state. Problems of such type are called **n-step decision problems**, and solving them is relatively easy.

Unfortunately, in most practical cases we may not assume to have constant, n-step sequences, due, for example, to looping.

The value iteration algorithm

For problems, which cannot be stated as the above n-step decision problems, we can compute approximate values of the state utilities in an iterative procedure called the **value iteration**:

$$U_{t+1}(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) U_t(s')$$

At step ($t = 0$) we assume arbitrary state utility values, and at the subsequent steps of the algorithm we compute their subsequent approximations.

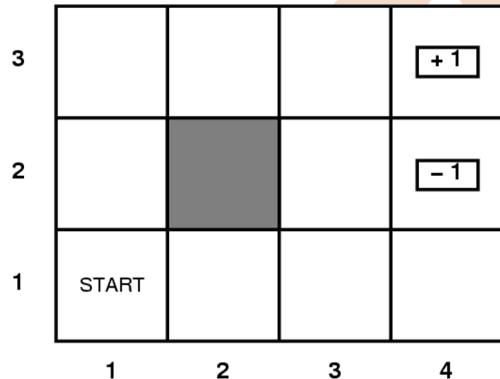
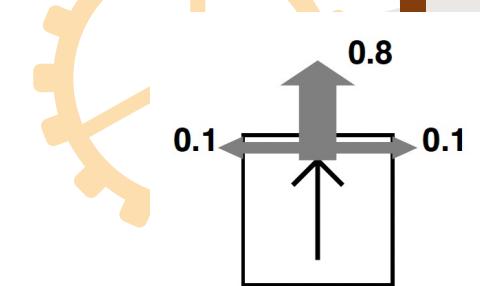
The algorithm may be terminated by comparing the obtained utilities and estimating the error. Note that the optimal agent's policy may be precisely determined by approximate utility values, even before they converge.

Bellman equation example

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

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

$$\begin{aligned}
 U(1, 1) &= -0.04 + \gamma \max[& 0.8U(1, 2) + 0.1U(2, 1) + 0.1U(1, 1), & (Up) \\
 & 0.9U(1, 1) + 0.1U(1, 2), & (Left) \\
 & 0.9U(1, 1) + 0.1U(2, 1), & (Down) \\
 & 0.8U(2, 1) + 0.1U(1, 2) + 0.1U(1, 1)]. & (Right)
 \end{aligned}$$



Value of an action

- Q-value is the value of an action

function Q-VALUE (*mdp, s, a, U*) **returns** a utility value

return $\sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma U[s']]$

- V-value is another name of utility

The value iteration algorithm

function VALUE-ITERATION(mdp, ϵ) **returns** a utility function

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

ϵ , the maximum error allowed in the utility of any state

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

δ , 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

The value iteration algorithm

function VALUE-ITERATION(mdp, ϵ) **returns** a utility function

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

ϵ , the maximum error allowed in the utility of any state

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

δ , the maximum relative change in the utility of any state

repeat

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

for each state s in S do

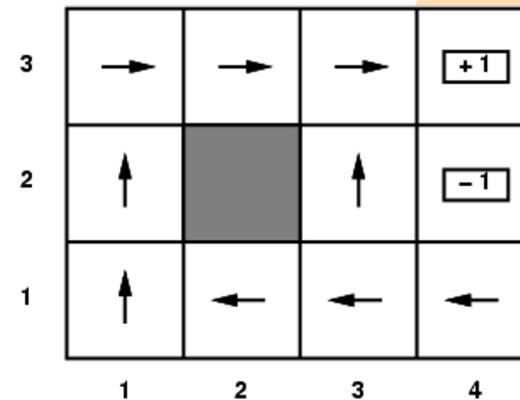
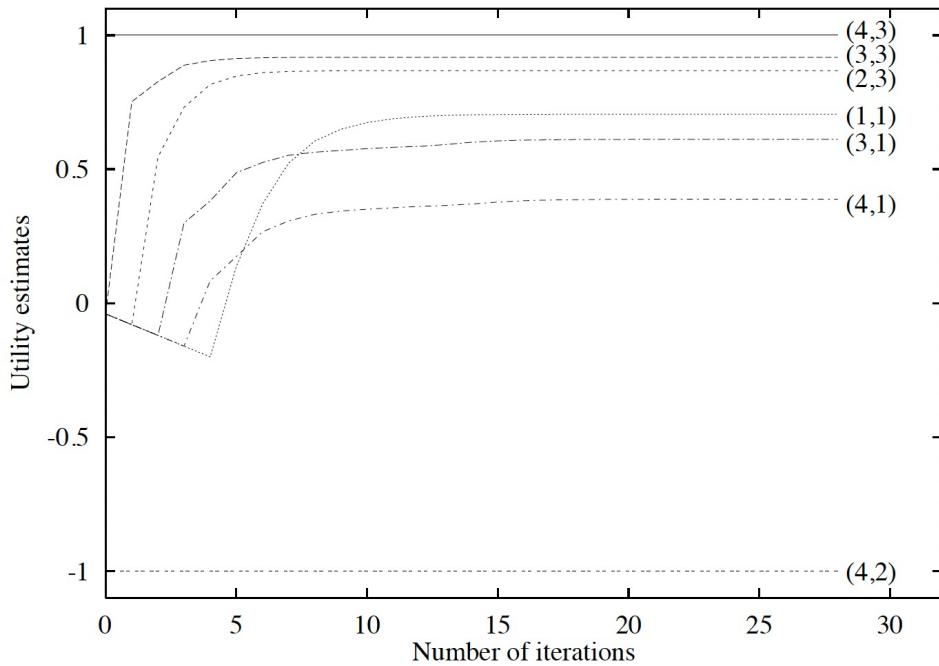
$U'[s] \leftarrow \max_{a \in A(s)} Q\text{-VALUE}(mdp, s, a, U)$

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

until $\delta \leq \epsilon(1 - \gamma)/\gamma$

return U

The value iteration algorithm



A table showing utility values for the 4x3 grid world after 30 iterations. The grid has columns labeled 1, 2, 3, 4 and rows labeled 1, 2, 3. The utility values are:

Row\Col	1	2	3	4
3	0.812	0.868	0.912	+1
2	0.762		0.660	-1
1	0.705	0.655	0.611	0.388

Convergence of value iteration

As we saw in the example the value iteration procedure converged nicely in all states. The question is, is it always this way?

It turns out that it is. The value iteration algorithm always leads to stable values of state utilities, which are the sole solutions of the Bellman equation. The number of iterations of the algorithm needed to reach an arbitrary error level ϵ is given by the following formula, where R_{\max} is the upper bound on the reward values:

$$N = \lceil \log(2R_{\max}/\epsilon(1 - \gamma))/\log(1/\gamma) \rceil$$

Convergence of value iteration

- In practice, the following termination condition can be used for the value iteration:
$$\|U_{i+1} - U_i\| < \epsilon(1 - \gamma)/\gamma$$
- In practice, the optimal policy can be reached much further than the utility values stabilize with desired small errors.
- N grows unboundedly when γ approaches one. The convergence can be sped up by decreasing γ , although this implies shortening the agent's horizon, and neglecting the long-term effects.
- For $\gamma = 1$, if there exist terminal states, similar convergence criteria and error estimates can be derived.

Policy iteration algorithm

Just because the optimal policy is often relatively insensitive to the specific values of the utilities, it can be computed by a similar iterative process, called the **policy iteration**. It works by selecting an arbitrary initial policy π_0 , and initial utilities, and then updating the utilities determined by the policy, according to the following formula:

$$U_{t+1}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_t(s)) U_t(s')$$

alternating it with subsequent update of the policy

$$\pi_{t+1}(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) U_t(s')$$

In the above formulas $\pi_t(s)$ denotes the action designated by the current policy for the state s . The first formula gives a set of linear equations, which can be solved exactly for U_{t+1} in $O(n^3)$ time (they are the exact utilities for the current approximate policy).

Policy iteration algorithm

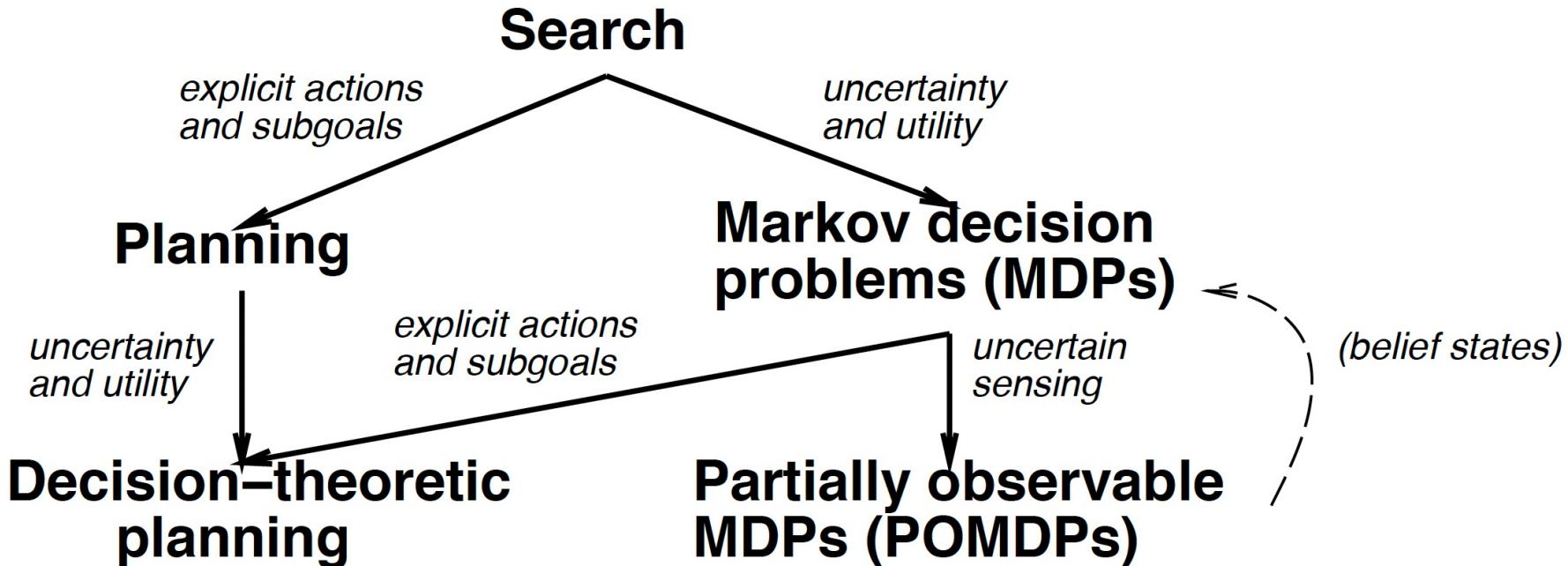
The policy iteration algorithm terminates when the policy update step make no change. Since for a finite space there exist a finite number of policies, the algorithm is certain to terminate.

For small state spaces ($n \in O(n^3)$) the above procedure is often most efficient. For large spaces, however, the $O(n^3)$ causes it to run very slowly. In these cases a **modified policy iteration** can be used, which works by iteratively updating the utilities — instead of computing it exactly each time — using a simplified Bellman updating given by the formula:

$$U_{t+1}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_t(s)) U_t(s')$$

Compared with the original Bellman equation the determination of the optimal action has been dropped here, since the actions are determined here by the current policy. Thus this procedure is simpler, and even several such update steps can be made before the next policy iteration step (updating the policy).

Sequential Design Problems



Uncertain state information

In a general case the agent may not be able to determine the state it ended up in after taking an action, or rather can only determine it with a certain probability. Such problems are called **partially observable Markov decision problems (POMDP)**.

In these problems the agent must compute the expected utility of its actions, taking into account both their various possible outcomes, and various possible new information (still incomplete), that it may acquire, depending on the state it ends up in.

The solution to the problem can be obtained by computing the probability distribution over all possible states that the agent can possibly be in, by considering the uncertain information about its environment that it was able to accumulate. In the general case, however, this computation is made difficult by the fact, that undertaking an action causes the agent to acquire new information, which may change its knowledge of the environment in ways difficult to consider. In practice, the agent must take into account new information it may acquire, along with the states it may transfer to. This may involve computing the value of information, which was covered earlier.

POMDP

A POMDP problem is defined by the following four elements:

- the set of states, but with no starting state s_0 ,
- the set $\text{ACTIONS}(s)$ of actions allowed in state s ,
- the transition function: $P(s'|s, a)$, which is a probability distribution of reaching the state s' after executing action a in state s ,
- reward function: $R(s)$,
- sensor model: $P(e|s)$, which is a probability distribution of receiving the evidence e (partially erroneous) in state s ,
- initial belief state: b_0 .

In POMDP problems there is no assumption about knowing the initial state. Instead, the agent has to work with the **belief state** $b(s)$, which is a probability distribution over all possible states s . Initially, we only know the initial belief state b_0 .

The task is to compute such policy, that would allow the agent to reach the goal with the highest probability. During the course of action the agent will change her belief state, due both to the newly received information, and to executing the actions by itself.

