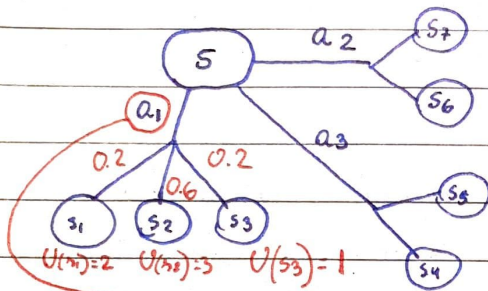Expected Utility of $X$; where $X$ is a random variables that can take multiple values $K$

$$E(X) = \sum_{K} K \, Pr(X = K)$$

→ we don't know to which state $s_a$ an action $a$ will lead us → Expected utility



$$E(u(a)) =$$

$$\sum_{s' \in Sa} u(s') \cdot Pr(s_a = s')$$

Example → expected utility of action $a_1$:

$$E(u(a_1)) = \underset{2}{u(s_1)} \times \underset{0.2}{Pr(s_a = s_1)} + \underset{3}{u(s_2)} \times \underset{0.6}{Pr(s_a = s_2)} + \dots$$

$$\dots \underset{1}{u(s_3)} \times \underset{0.2}{Pr(s_a = s_3)}$$

- We repeat the process for every $a$
And we choose the $a$ with highest expected utility....

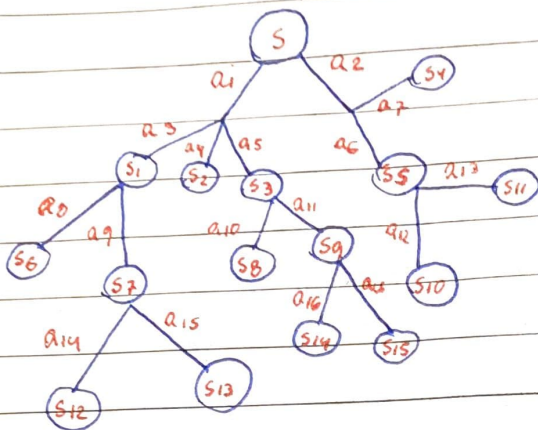Alternative methods to choose an action:

| maximin | maximax |
|---|---|
| $a^* = \underset{a \in A}{arg\ max} \left\{ \underset{s' \in Sa}{min}\ u(s') \right\}$ | $a^* = \underset{a \in A}{arg\ max} \left\{ \underset{s' \in Sa}{max}\ u(s') \right\}$ |
| Among the minimum utilities of each action, we choose the action with the highest minimum utility | Among the maximum utilities of each action, we choose the action with the highest maximum utility |

1

# Sequential decision making



We cannot use a simple expected utility for a single action ...........

As This action might be good in the short term ......
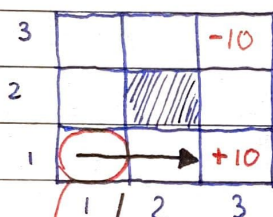
But what about in the long term with a sequence of actions?

Basic framework: Markov decision process (MDP)

- Set of states $s \in S$
- Set of Actions $A(s)$ in each state
- transition model $P(s' | s, a) \rightarrow$ current action

   $\downarrow$    $\searrow$ current state

   next state

- A reward function

- SOLUTION of MDPs $\rightarrow$ a policy $\pi$

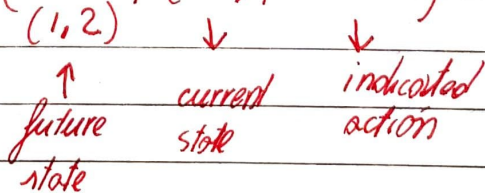   we would prefer the optimum policy $\pi^*$



States: $(1,1), (2,1), (3,1), (1,1)$,

$(1,3), (2,3), (3,3), (3,2)$

$\downarrow$ policy $\pi \rightarrow$ is the optimum policy $\rightarrow \pi^*$

$\rightarrow$ Action: Stop, North, South, East, West

Example transition model:

$P(\text{~~~~~~} | (1,1), \text{North}) = 0.8$

$(1,2)$     $\downarrow$      $\downarrow$

$\uparrow$     current    indicated

future    state     action

state

3

# SUMMARY III

Each state has a Reward assigned
To it → Reward as cost of moving
is a frequent approach



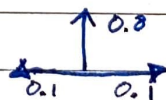R(1,1) | R(1,2) | +10

*Terminal state final
Reward*

$(1,1) \longrightarrow$ Agent moves to $(1,2) \rightarrow$ Receives a
Reward (which can be negative).

As with Previous (non-sequential) decision
making processes, we need to know
the utilities. Then, calculate the
expected utility of each action. However,
This time we have to include
The future actions we might take.....

→ Image we have the right utilities of
each state. Then To decide on action
we need to calculate the expected
utility ⟶

Transition model



$$E(u(north)) = 0.8 \times 0.660 + 0.1 \times 0.655 + 0.1 \times 0.388$$

$$E(u(east)) = 0.8 \times 655 + 0.1 \times 0.660 + 0.1 \times 0.611$$

*The agents stays in the
same place*

$$E(u(west)) = 0.8 \times 0.388 + 0.1 \times 0.660 + 0.1 \times 0.611$$

Then, we choose the action with highest utility

S

How To calculate The utilies that make
the agent follow an optimal policy $\pi^*$

→ OPTIMAL
POLICY
(given the
right utilities)

$$\pi^*(s) = \arg\max_{a \in A(s)} \sum_{s'} P(s'|s,a) \, U^{\pi^*}(s')$$

We need To calculate
This

→ We can use the Bellman equation :

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

discount factor

We will have To solve as many bellman
equations as states in the map...

This is hard To solve in big scenarios

Example:

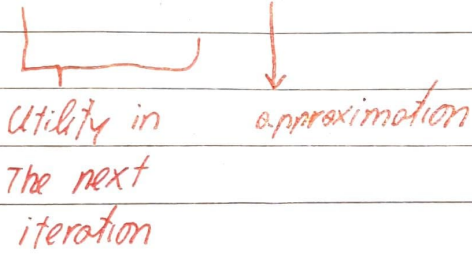$$U(1,1) = R(1,1) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s,a) \, U(s')$$
$$U(1,2) = R(1,2) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s,a) \, U(s')$$
$$U(1,3) = R(1,3) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s,a) \, U(s')$$
$$\vdots$$

Solve
This
is
Really
hard

6

We con get the Utilities with a computational approch called / labelled as

## VALUE    ITERATION

$$\rightarrow \quad U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s,a) U_i(s')$$

Utility in The next iteration

approximation

Repeat it several times until the value of $U$ does not change

# 1

a) Expected utility

Non-sequential decision making

$a_1 =$ office          $a_2 =$ going out

$P(\text{work} \mid \text{office}) = 0.5$          $P(\text{coffe} \mid \text{out}) = 0.95$
$P(\text{distracted} \mid \text{office}) = 0.3$          ~~P(coffe)~~
$P(\text{colleague} \mid \text{office}) = 0.2$          $P(\text{spill} \mid \text{out}) = 0.05$

$U(\text{work}) = 8$          $U(\text{coffe}) = 10$
$U(\text{distracted}) = 1$          $U(\text{spill}) = -20$
$U(\text{colleague}) = 5$



- $E(U(\text{office})) =$
  $0.5 \times 8 + 0.3 \times 1 + 0.2 \times 5 = 5.3$

- $E(U(\text{coffe})) =$
  $0.95 \times 10 + 0.05 \times -20 = 8.5$

$E(U(\text{coffe})) > E(U(\text{office})) \rightarrow$ we choose coffe

b) Which action should we choose

coffe expected utility is higher

c) Maximax and maximin

maximin $\rightarrow$ $a^* \; \arg\max_{a \in A} \{ \min_{s' \in S_a} u(s') \}$

→ Among the minimum utilities of each action, we choose the action with the highest minimum utility

maximax

$$a^* = \arg\max_{a \in A} \left\{ \max_{s' \in S_a} U(s') \right\}$$

Among the maximum utilities, we choose the action with the highest maximum utilities

maximin:

$$\arg\max_{s \in \{office,\ out\}} \left\{ \min^{office}(8,1,5), \min^{out}(10,-20) \right\}$$

$\downarrow$ min office  $\downarrow$ min out

$$\arg\max_{s \in \{office,\ out\}} \left\{ \boxed{\min\ office\ (1)}, \min\ out\ (-20) \right\}$$

$\downarrow$ maximin chooses office

maximax

$$\arg\max_{s \in \{office,\ out\}} \left\{ \max^{office}(8,1,5), \max^{out}(10,-20) \right\}$$

$\uparrow$ max office  $\uparrow$ max out
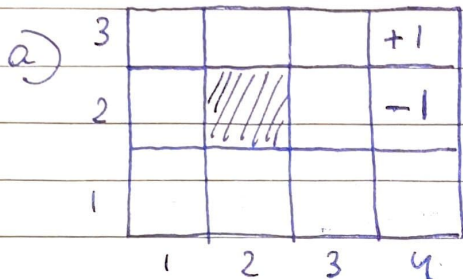
$$\arg\max_{s \in \{office,\ out\}} \left\{ \max\ office\ (8), \boxed{\max\ out\ (10)} \right\}$$

$\downarrow$ maximax chooses out

# Exercise 2

a)



States: $(1,1)$, $(1,2)$,
$(1,3)$, ~~(1,4)~~, $(2,1)$,
$(2,3)$, ~~(2,4)~~, $(3,1)$,
$(3,2)$, $(3,3)$, $(4,1)$
$(4,2)$, $(4,3)$

Initial state: $(1,1)$    Actions: up, dawn, ~~left~~,
right

$$R(s)\begin{cases} 1 & \text{for } s = (4,3) \rightarrow \text{terminal state} \\ -1 & \text{for } s = (4,2) \rightarrow \text{terminal state} \\ -0.04 & \text{for } s \neq (4,3) \text{ and } (4,2) \end{cases}$$

## Transition model



indicated direction
0.8
0.1      0.1

Example:


$P((1,2) \mid (1,1), Up) = 0.8$


$P((1,1) \mid (1,1), Up) = 0.1$


$P((2,1) \mid (1,1), Up) = 0.1$

b)   Normal Bellman equation $\rightarrow U(s) =$
$$\rightarrow U(s) = R(s) + Y \max_{a \in A(s)} \sum_{s'} P(s' \mid s, a) U(s')$$

Deterministic Bellman (No Probability)
$$\rightarrow U(s) = R(s) + Y \max_{a \in A(s)} U(s')$$

We know which state is going to be the result of
every of the agents actions

# c) Deterministic Value Iteration

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s,a) \times U(s')$$

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} U(s')$$

deterministic

→ $\gamma = 1$ and $U(s) = 0$

<mark>A 1ˢᵗ iteration:</mark>

| 3 | -0.04 | -0.04 | 0.96 | +1 |
|---|---|---|---|---|
| 2 | -0.04 | ▨ | -0.04 | -1 |
| 1 | -0.04 | -0.04 | -0.04 | -0.04 |
|   | 1 | 2 | 3 | 4 |

Examples

$U_{i+1}(1,1) = -0.04 + 1 \times \max$ {
- Up: **0**
- East: 0
- Sout: 0
- West: 0
}

$\underbrace{\qquad}_{-0.04}$

$U_{i+1}(3,3) = -0.04 + 1 \times \max$ {
- Up: 0
- **East: 1**
- South: 0
- West: 0
}

$\underbrace{\qquad}_{0.96}$

<mark>2ⁿᵈ Round / Iteration</mark>

| 3 | -0.08 | 0.92 | 0.96 | +1 |
|---|---|---|---|---|
| 2 | -0.08 | ▨ | 0.92 | -1 |
| 1 | -0.08 | -0.08 | -0.08 | -0.08 |
|   | 1 | 2 | 3 | 4 |

Examples

$U_{i+1}(2,3) = -0.04 + 1 \times \max$ {
- Up: 0.04
- **East: 0.96**
- West: -0.04
- South: -0.04
}

$\underbrace{\qquad}_{0.92}$

2nd Continuation

5th iteration

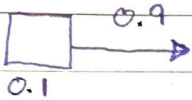| 0.88 | 0.92 | 0.96 | |
|------|------|------|--|
| 0.84 | ///// | 0.92 | |
| 0.8 | 0.84 | 0.8 | |

d) Calculate the optimal policy

→ Remember: we have to calculate the expected $u$ to know the best policy

However, we are in a deterministic environment. So, is just the utility

| 0.88 → | 0.92 → | 0.96 → | +1 |
|--------|--------|--------|----|
| 0.84 ↑ | ///// | 0.92 ↑ | 0.92 |
| 0.8 → | 0.84 → | 0.8 ↑ | 0.84 ← |

} Best policy

3



Non deterministic Bellman:

$$U_{i+1}(s) = R(s) + Y \max_{a \in A(s)} \sum_{s'} P(s'/s,a)\, U(s')$$

0.9

| 3 | $\circ$ | $\circ$ | $\circ$ | 1 |
|---|---|---|---|---|
| 2 | $\circ$ | ///// | $\circ$ | $-1$ |
| 1 | $\circ$ | $\circ$ | $\circ$ | $\circ$ |
| | 1 | 2 | 3 | 4 |

Example

$U_{i+1}(3,3) \Rightarrow$

$-0.04 + 0.9 \max$

Up: $0 \times 0.9 + 0 \times 0.1 = 0$
Left: $0 \times 0.9 + 0 \times 0.1 = 0$
Right: $1 \times 0.9 + 0 \times 0.1 = 0.9$
Down: $0 \times 0.9 + 0 \times 0.1 = 0$

0.77

Repeat this for all states. Then, in the next iteration
repeat this for every state. Repeat this process until
U does not change much.

| 3 | 0.87 | 0.91 | 0.96 | +1 |
|---|---|---|---|---|
| 2 | 0.82 | ///// | 0.91 | $-1$ |
| 1 | 0.78 | 0.72 | 0.87 | 0.82 |
| | 1 | 2 | 3 | 4 |

final Utilities

To know the best policy/action......
We need to calculate the expected utility
and choose the maximum one

$$\pi^*(s) = \arg\max_{a \in A} \sum_{s'} P(s'/s,a)\, U^*(s')$$
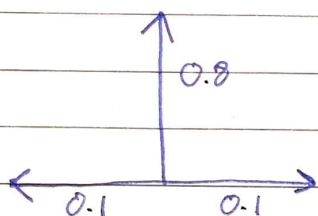
3 Example:

$\pi^*(3,2) = \arg\max$
$\begin{cases} Up: 0.96 \times 0.9 + 0.1 \times 0.91 = 0.955 \\ Down: 0.87 \times 0.9 + 0.1 \times 0.91 = 0.874 \\ Left: 0.91 \times 0.9 + 0.1 \times 0.1 = 0.91 \\ Right: -1 \times 0.9 + 0.1 \times 0.91 = -0.809 \end{cases}$

4 → New transition model


0.8
0.1     0.1

a) 1 Value iteration for (3,1)

$$V_{i+1}(s) = R(s) + \gamma \arg\max_{s \in A(s)} \sum_{s'} P(s'/a_{,s}) \cdot U(s')$$

→ $V_{i+1}(3,3) = -0.04 + \gamma \max \begin{cases} 0.8 \times 0.66 \ldots \end{cases}$

| 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 |

$\max \begin{cases} Up: 0.8 \times 0.660 + 0.1 \times 0.388 + 0.655 \times 0.1 = 0.6323 \\ Down: 0.8 \times 0.611 + 0.1 \times 0.655 + 0.1 \times 0.388 = 0.593 \\ Left: 0.8 \times 0.655 + 0.1 \times 0.611 + 0.1 \times 0.660 = 0.6511 \\ Right: 0.8 \times 0.388 + 0.1 \times 0.611 + 0.1 \times 0.660 = 0.4375 \end{cases}$

$V_{i+1}(3,3) = -0.04 + 1 \times 0.6511 = \boxed{0.6111}$

The value has not changed → convergence

4

## b) Expected utility in $(1,3)$

Up: $0.8 \times 0.660 + 0.1 \times 0.655 + 0.1 \times 0.388 = 0.6323$

Down: $0.8 \times 0.611 + 0.1 \times 0.655 + 0.1 \times 0.388 = 0.5931$

Left: $0.8 \times 0.655 + 0.1 \times 0.660 + 0.1 \times 0.611 = \boxed{0.6511}$

Right: $0.8 \times 0.388 + 0.1 \times 0.660 + 0.1 \times 0.611 = 0.4375$

## c) maximax and maximin

- Example in $(3,1)$

maximax

Up: max is 0.660

Down: max is 0.655   } among these ones the

Left: max is 0.660     max is 0.660. We will

Right: max is 0.660   } need to romdomly choose
                        one

maximin

Up: min is 0.388

Down: min is 0.388   } among these ones the

Left: min is 0.611     max is 0.611. therefore

Right: min is 0.388   } the agent will go left