

Computational Intelligence 2023/24

(Solution) Exam – February 12th, 2024

Please mind the following:

- Fill in your **personal information** in the fields below.
- You may use an unmarked dictionary during the exam. **No** additional tools (like calculators, e.g.) might be used.
- Fill in all your answers **directly into this exam sheet**. You may use the backside of the individual sheets or ask for additional paper. In any case, make sure to mark **clearly** which question you are answering. Do not use pens with green or red color or pencils for writing your answers. You may give your answers in both **German and English**.
- Points annotated for the individual tasks only serve as a preliminary guideline.
- At the end of the exam hand in your **whole exam sheet**. Please make sure you do not undo the binder clip!
- To forfeit your exam (“entwerten”), please cross out clearly this **cover page** as well as **all pages** of the exam sheet. This way the exam will not be evaluated and will not be counted as an exam attempt.
- Please place your LMU-Card or student ID as well as photo identification (“Lichtbildausweis”) clearly visible on the table next to you. Note that we need to check your data with the data you enter on this cover sheet.

Time Limit: 90 minutes

First Name:																										
Last Name:																										
Matriculation Number:																										
<table border="1"><tr><td>Topic 1</td><td>max. 10 pts.</td><td>pts.</td></tr><tr><td>Topic 2</td><td>max. 18 pts.</td><td>pts.</td></tr><tr><td>Topic 3</td><td>max. 10 pts.</td><td>pts.</td></tr><tr><td>Topic 4</td><td>max. 14 pts.</td><td>pts.</td></tr><tr><td>Topic 5</td><td>max. 10 pts.</td><td>pts.</td></tr><tr><td>Topic 6</td><td>max. 10 pts.</td><td>pts.</td></tr><tr><td>Topic 7</td><td>max. 18 pts.</td><td>pts.</td></tr><tr><td colspan="2">Total max. 90 pts.</td><td>pts.</td></tr></table>			Topic 1	max. 10 pts.	pts.	Topic 2	max. 18 pts.	pts.	Topic 3	max. 10 pts.	pts.	Topic 4	max. 14 pts.	pts.	Topic 5	max. 10 pts.	pts.	Topic 6	max. 10 pts.	pts.	Topic 7	max. 18 pts.	pts.	Total max. 90 pts.		pts.
Topic 1	max. 10 pts.	pts.																								
Topic 2	max. 18 pts.	pts.																								
Topic 3	max. 10 pts.	pts.																								
Topic 4	max. 14 pts.	pts.																								
Topic 5	max. 10 pts.	pts.																								
Topic 6	max. 10 pts.	pts.																								
Topic 7	max. 18 pts.	pts.																								
Total max. 90 pts.		pts.																								

1 General Knowledge / Single Choice

10pts

For each of the following questions select **one** correct answer ('1 of n '). Every correct answer is awarded one point. Multiple answers or incorrect answers will be marked with zero points.

(a) If you train a large language model to imitate how you text with people and then ask your fellow students to discern if they are writing with you or with the large language model, you are effectively running an instance of...

i the Turing Test	ii Prolog	iii the Singularity	iv FunSearch
----------------------	-----------	---------------------	--------------

(b) In accordance with the main focus of the programming language Prolog, its name is derived from ...

i producing log files	ii programming logic	iii probability logarithms	iv pranking old ogres
-----------------------	-------------------------	----------------------------	-----------------------

(c) Given a fixed graph structure, a neural network is fully defined by three categories of entities. Which is not one of them?

i weights	ii biases	iii curvatures	iv activation functions
-----------	-----------	-------------------	-------------------------

(d) In a standard evolutionary algorithm, *variation* operators are usually employed to increase which property of the optimization?

i selection	ii exploration	iii exploitation	iv production
-------------	-------------------	------------------	---------------

(e) Compared to the game of Go, which was only considered artificially 'mastered' as of 2016, the game of Chess was solvable much earlier. The software that was able to beat then world champion Garri Kasparow in the late 1990s was called ...

i Deep Blue	ii AlphaGo	iii DeepMind	iv ChatGPT
----------------	------------	--------------	------------

(f) A multi-agent system made up of many similar agents that have a shared interest but little communication and only locally executed algorithms is called ...

i market	ii evolutionary stable	iii goal class	iv <u>swarm</u>
----------	------------------------	----------------	-----------------

(g) One of the main difficulties for finding effective strategies in the game *Lemonade Stand* (sometimes also called *Glühwein Stand*) lies in the opponent strategies' ...

i action space	ii <u>non-stationarity</u>	iii beverage quality	iv policy model
----------------	----------------------------	----------------------	-----------------

(h) In a multi-iteration tournament of the game *Prisoner's Dilemma* where you play against each opponent for an unknown number of games, you are allowed to adopt any of these four strategies. Which one of these is expected to be the most successful over the course of such a tournament?

i always <i>Cooperate</i>	ii always <i>Defect</i>	iii give up	iv <u>tit-for-tat</u>
---------------------------	-------------------------	-------------	-----------------------

(i) The *Hawk-Dove* game was modeled as a possible explanation for the ecological balance of ...

i <u>aggression between peers</u>	ii differences in size	iii animal intelligence	iv CO2 emissions
-----------------------------------	------------------------	-------------------------	------------------

(j) The outlook of machine learning expert Richard Sutton on the eventual superiority of using search-based algorithms instead of human-experience-based knowledge was summarized in his famous 2019 publication titled ...

i <u>"The Bitter Lesson"</u>	ii "The Bitter Lemon"	iii "The Barren Line"	iv "The Better LISP"
------------------------------	-----------------------	-----------------------	----------------------

2 Agents and Goals

18pts

Scenario 1. We consider a multi-agent system made up of $n = 20$ cows and a mowing robot called *MowBot*. All these 21 agents live on an infinite meadow which we model as an infinite 2d plane of positions $l \in \mathbb{R}^2$. The agent $G^{[0]}$ is the *MowBot* and all other agents $G^{[1]}, \dots, G^{[20]}$ are cows.

The evolution of the multi-agent system is encoded as a series of states $\langle s_t \rangle_t$ of unspecified length for states $s_t \in \mathcal{S}$ at time t , where $\mathcal{S} = \mathbb{N} \times (\mathbb{R}^2)^{21}$ is the (in total 43-dimensional) state space, as well as a series of joint observations $\langle o_t \rangle_t = \langle (o_t^{[0]}, \dots, o_t^{[20]}) \rangle_t$ made by the agents and a series of joint actions $\langle a_t \rangle_t = \langle (a_t^{[0]}, \dots, a_t^{[20]}) \rangle_t$ executed by the agents.

Our system is fully observable. Thus, each agent $G^{[i]}$ makes observations

$$o_t^{[i]} = s_t = (t, l_t^{[0]}, \dots, l_t^{[20]}) \in \mathcal{S}$$

at each time step t where $l_t^{[i]} \in \mathbb{R}^2$ is the position of the agent $G^{[i]}$ at time $t \in \mathbb{N}$. Furthermore, any agent $G^{[i]}$ uses the identical action space $\mathcal{A}^{[i]} = \{\text{do_nothing}, \text{go_northeast}, \text{go_east}, \text{go_southeast}, \text{go_south}, \text{go_southwest}, \text{go_west}, \text{go_northwest}, \text{go_north}, \text{eat/mow}\}$ with $|\mathcal{A}^{[i]}| = 10$ where `do_nothing` has the agent not do anything, all actions starting with `go_` cause the agent to move for a certain length in the respective direction, and `eat/mow` causes any cow agent to eat the grass and any *MowBot* agent to mow the grass.

You can use a 2d geometric distance function $\text{dist}(l, l') = \sqrt{(x - x')^2 + (y - y')^2}$ where $l = (x, y) \in \mathbb{R}^2$ and $l' = (x', y') \in \mathbb{R}^2$.

- (i) Give a goal predicate $\gamma_1 : \langle \mathcal{S} \rangle \times \langle \mathcal{A} \rangle \rightarrow \mathbb{B}$ so that γ_1 holds iff the *MowBot* has never (at full time steps) been closer than a distance of 1 to any of the cows. (4pts)

$$\gamma_1 \left(\langle (t, l_t^{[0]}, \dots, l_t^{[20]}) \rangle_t, \langle (a_t^{[0]}, \dots, a_t^{[20]}) \rangle_t \right) \iff \forall t : \forall j \in [1; 20] \subset \mathbb{N} : \text{dist}(l_t^{[0]}, l_t^{[j]}) \geq 1$$

(ii) Let $(\mathcal{S}, \mathcal{A}, \mathbb{R}, P, R)$ be a Markov decision process (MDP, cf. Definition 3 in the appendix) with state space \mathcal{S} from Scenario 1, joint action space \mathcal{A} from Scenario 1, target space $\mathcal{T} = \mathbb{R}$, and a given transition probability function P . Give — from our *MowBot*'s perspective — a reward function R that both incentivizes to (a) fulfill γ_1 and (b) execute **mow/eat** as often as possible. (8pts)

$$R(s_t, a_t, s_{t+1}) = \begin{cases} -10 & \text{if } \exists j \in [1; 20] : \text{dist}(l_{t+1}^{[0]}, l_{t+1}^{[j]}) < 1, \\ +1 & \text{else if } a_t = \text{eat/mow}, \\ 0 & \text{else.} \end{cases}$$

(iii) Now, in any run at time step $t = 26$ all cows suddenly change their policy to a specific new (from then on again fixed) policy.

Can R still be learned as the reward function of an MDP? Briefly state your reasoning. (3pts)

Yes, because the time is part of the state space and time dependent changes can thus be recognized when learning.

(iv) Let $\gamma_2 : \langle \mathcal{S} \rangle \times \langle \mathcal{A} \rangle \rightarrow \mathbb{B}$ be goal predicate so that γ_2 holds iff the *MowBot* has never executed the action **eat/mow** within a radius of 1 from any point where any cow has ever executed **mow/eat**. Assume that the cows' policies are fixed.

Can the *MowBot* learn to fulfill γ_2 effectively for an MDP with the given state space \mathcal{S} for some reward function R' ? Briefly state your reasoning. (3pts)

No, because past eaten patches are not encoded in the state, so if the *MowBot* implements an MDP, it can never be certain where cows might have previously eaten grass.

3 Fuzzy Logic

10pts

To evaluate the impact of the *MowBot* in Scenario 1 on the lives of our cows, we performed a field study and measured a cow's happiness at various positions $(x, y) \in \mathbb{R}^2$ on the infinite meadow. For our experiments, we firmly placed the *MowBot* at position $(3, -1)$ and never let it change that position.

From that study, we found that a cow's happiness can be described by the following formula in fuzzy logic:

$$\begin{aligned} \text{happy}((x, y)) &= \text{AND}(\text{green}((x, y)), \text{NOT}(\text{mowbot_nearby}((x, y)))) \\ \text{where } \text{green}((x, y)) &= \begin{cases} 1.0 & \text{if } |x| < 3 \text{ and } |y| < 3, \\ 0.8 & \text{otherwise,} \end{cases} \\ \text{and } \text{mowbot_nearby}((x, y)) &= 0.5 \cdot \max\{2 - \text{dist}((x, y), (3, -1)), 0\}. \end{aligned}$$

- (i) What is the happiness of a cow at position $(3, 0)$? Show your computation steps. (4pts)

$$\begin{aligned} \text{green}((3, 0)) &= 0.8 \\ \text{mowbot_nearby}((3, 0)) &= 0.5 \cdot \max(2 - 1, 0) = 0.5 \\ \text{NOT}(\text{mowbot_nearby}((3, 0))) &= 1 - 0.5 = 0.5 \\ \text{happy}((3, 0)) &= \min(0.8, 0.5) = 0.5 \end{aligned}$$

- (ii) We now consider that cows can walk one step in any direction by adding one of the eight vectors $\{(+1, +1), (+1, 0), (+1, -1), (0, +1), (0, -1), (-1, +1), (-1, 0), (-1, -1)\}$ to their position. What is the most drastic increase in happiness a cow can achieve with one step? What is one possible corresponding start position and one possible corresponding end position of such a drastic step? (6pts)

Since it was not directly specified, depending on the following assumptions:

- cows and the *MowBot* can stand on the same position: Maximum happiness increase $+0.7$, from *MowBot* $(3, -1)$ (happiness 0.0) to e.g. $(2, -2)$ 0.7 , or
- cows and the *MowBot* cannot stand on the same position: Maximum happiness increase $+0.5$ from close position e.g., position $(3, 0)$ with (happiness 0.5) moving further away to e.g., position $(2, 1)$ with (happiness 1.0)

4 Optimization

14pts

Scenario 2. The policy $\pi : \mathcal{S} \rightarrow \mathcal{A}^{[\cdot]}$ of a *math cow* is given by a real number $0 \leq u < 1$, i.e.,

$$\pi((t, -, \dots, -)) = \text{map}(\lfloor 10^t \cdot u \rfloor \bmod 10)$$

where state space \mathcal{S} and action space \mathcal{A} are given as in Scenario 1, $\lfloor w \rfloor$ for $w \in \mathbb{R}$ is the nearest integer from w rounded down, $w \bmod 10$ for positive real numbers $w \in \mathbb{R}$ is given recursively via

$$w \bmod 10 = \begin{cases} w & \text{if } w < 10, \\ w - 10 \bmod 10 & \text{otherwise,} \end{cases}$$

and $\text{map} : \{0, \dots, 9\} \rightarrow \mathcal{A}^{[\cdot]}$ is given exhaustively via

$$\begin{aligned} \text{map}(0) &= \text{do_nothing}, \\ \text{map}(1) &= \text{go_northeast}, \\ &\dots \\ \text{map}(8) &= \text{go_north}, \\ \text{map}(9) &= \text{eat/mow}. \end{aligned}$$

In short, the *math cow* policy encoding described in Scenario 2 uses a single real number u to encode a (potentially very long or even infinitely long) behavior by interpreting every single decimal digit of u as a subsequent action to be taken by a corresponding policy. We now give two examples if you need them:

- If a *math cow*'s policy π is given via the number $u = 0.199$, then that cow's action at time point $t = 1$ is given via

$$\begin{aligned} \pi((1, -, \dots, -)) &= \text{map}(\lfloor 10^1 \cdot 0.199 \rfloor \bmod 10) \\ &= \text{map}(\lfloor 1.99 \rfloor \bmod 10) \\ &= \text{map}(\lfloor 1.99 \rfloor) \\ &= \text{map}(1) \\ &= \text{go_northeast}. \end{aligned}$$

Analogously, $\pi((2, -, \dots, -)) = \text{eat/mow}$, then $\pi((3, -, \dots, -)) = \text{eat/mow}$, and then $\pi((t, -, \dots, -)) = \text{do_nothing}$ for all $t \geq 4$.

- If a *math cow*'s policy π is given via the number $u = \frac{8}{9} = 0.\bar{8} = 0.888\dots$, then that cow's action at time point $t = 1$ is given via

$$\begin{aligned} \pi((1, -, \dots, -)) &= \text{map}(\lfloor 10^1 \cdot 0.\bar{8} \rfloor \bmod 10) & (1) \\ &= \text{go_north}. & (2) \end{aligned}$$

This action is then also chosen for any other time point $t \geq 1$.

We now want to optimize one cow's policy for any given target function τ via *simoolated annealing*, a special form of simulated annealing (cf. Algorithm 1 in the appendix) with a specific acceptance probability function

$$A(Q, Q', K) = \frac{1}{K}$$

for temperature schedule

$$\kappa(t) = t.$$

(i) What is the effect of this definition for A and κ in *simoolated annealing* compared to more typical or more general variants of simulated annealing? How does this effect vary with the quality of the newly generated solution candidate? (4pts)

Here, the acceptance probability decreases with time independent of solution quality. That means that (in contrast to simulated annealing) acceptance probability does not vary with the quality of the solution candidates.

(ii) A friendly *neighbors* function, through which we generate from a policy π a new policy $\pi' \sim \text{neighbors}(\pi)$, has the following properties:

1. On average, substantial parts of the information in π remain within π' .
2. π' always differs in at least one action within the first 1000 actions from π .
3. *neighbors* is symmetric in the sense that generating π' from π has the same chance as generating π from π' would have.
4. For a given π , the value of π' is non-deterministic.

Give a friendly *neighbors* function. Give a brief explanation why your function fulfills each of the four properties enumerated above. (10pts)

Hint: You can use policies π, π' and their corresponding real numbers u, u' interchangeably, without any need to type cast from π to u or vice versa.

Note: You can write $x \sim X$ to sample a (uniformly) random element x from a finite set X . You can also write $x = \text{random}()$ to sample a random real number x with $x \in [0; 1] \subset \mathbb{R}$.

$$\begin{aligned} \text{sample } t &\sim [1; 1000] \subset \mathbb{N} \\ \text{neighbors}(u) &= \{u \pm 10^{-t}\} \end{aligned}$$

1. In most cases, only one digit of u' is changed from u , so the rest of the information remains.
2. At least one digit within the first 1000 (see t) of u' is changed from u .
3. Choosing the plus case from \pm in $\text{neighbors}(u)$ is reverted by sampling the same t and choosing the minus case starting from u' .
4. We sample t (and then again sample from the set $\text{neighbors}(u)$).

5 The π -Calculus

10pts

- (i) Evaluate the following π -process step by step until it cannot be reduced any further.
(6pts)

$$!(\overline{panic.beep.0} + \overline{beep}.MowBot + panic.beep.Cow + \overline{brr}.moo.Cow) \mid \overline{moo}.Cow \mid \overline{panic.moo.0}$$

$$\begin{aligned} &\rightarrow !(\overline{panic.beep.0} + \overline{beep}.MowBot + panic.beep.Cow + \overline{brr}.moo.Cow) \mid \overline{moo}.Cow \mid beep.Cow \mid moo.0 \\ &\rightarrow !(\overline{panic.beep.0} + \overline{beep}.MowBot + panic.beep.Cow + \overline{brr}.moo.Cow) \mid Cow \mid beep.Cow \mid 0 \\ &\rightarrow !(\overline{panic.beep.0} + \overline{beep}.MowBot + panic.beep.Cow + \overline{brr}.moo.Cow) \mid Cow \mid Cow \mid MowBot \mid 0 \\ &\rightarrow !(\overline{panic.beep.0} + \overline{beep}.MowBot + panic.beep.Cow + \overline{brr}.moo.Cow) \mid Cow \mid Cow \mid MowBot \end{aligned}$$

- (ii) Evaluate the following π -process step by step until it cannot be reduced any further.
(4pts)

$$!(\overline{beep}.MowBot) \mid \overline{beep}.beep.(panic.\overline{moo}.beep.0 + beep.panic.MowBot) \mid beep.0 \mid \overline{beep}.0$$

one possible path (out of several):

$$\begin{aligned} &\rightarrow !(\overline{beep}.MowBot) \mid \overline{beep}.beep.(panic.\overline{moo}.beep.0 + beep.panic.MowBot) \mid 0 \mid MowBot \mid \overline{beep}.0 \\ &\rightarrow !(\overline{beep}.MowBot) \mid \overline{beep}.beep.(panic.\overline{moo}.beep.0 + beep.panic.MowBot) \mid MowBot \mid \overline{beep}.0 \end{aligned}$$

6 Game Theory

10pts

Since cows (and their farmer) can be terribly dangerous animals, the *MowBot* has to think twice how it should behave around the meadow, depending on which behavior it encounters. Within the two-player normal-form game for our true-to-life scenario, we can observe these dynamics unfold in the following payoff matrix for the *Cowmunication Game*:

	Cool Cow	Shy Cow	Default Cow	Brave Cow	Farmer
'moo?'	2, 2	3, 0	4, 0	4, 2	2, 1
'brrrr'	4, 1	3, 1	3, 0	4, 4	2, 4
stop	1, 0	0, 1	3, 4	4, 1	2, 4
ignore	3, 0	0, 3	1, 3	0, 0	4, 4

- (i) Compute all pure Nash equilibria (cf. Definitions 6 and 9 in the appendix) for the *Cowmunication Game* by denoting all best responses. (6pts)

The best responses are....

	Cool Cow	Shy Cow	Default Cow	Brave Cow	Farmer
'moo?'	2, <u>2</u>	3, 0	4, 0	4, <u>2</u>	2, 1
'brrrr'	4, 1	3, 1	3, 0	4, <u>4</u>	2, <u>4</u>
stop	1, 0	0, 1	3, <u>4</u>	4, 1	2, <u>4</u>
ignore	3, 0	0, 3	1, 3	0, 0	4, <u>4</u>

	Cool Cow	Shy Cow	Default Cow	Brave Cow	Farmer
'moo?'	2, 2	<u>3</u> , 0	<u>4</u> , 0	<u>4</u> , 2	2, 1
'brrrr'	<u>4</u> , 1	<u>3</u> , 1	3, 0	<u>4</u> , 4	2, 4
stop	1, 0	0, 1	3, 4	<u>4</u> , 1	2, 4
ignore	3, 0	0, 3	1, 3	0, 0	<u>4</u> , 4

	Cool Cow	Shy Cow	Default Cow	Brave Cow	Farmer
'moo?'	2, <u>2</u>	<u>3</u> , 0	<u>4</u> , 0	<u>4</u> , <u>2</u>	2, 1
'brrrr'	<u>4</u> , 1	<u>3</u> , 1	3, 0	<u>4</u> , <u>4</u>	2, <u>4</u>
stop	1, 0	0, 1	3, <u>4</u>	<u>4</u> , 1	2, <u>4</u>
ignore	3, 0	0, 3	1, 3	0, 0	<u>4</u> , <u>4</u>

Thus, the pure Nash equilibria are the joint strategies: (ignore, Farmer) \mapsto (4, 4), ('moo?', Brave Cow) \mapsto (4, 2), ('brrrr?', Brave Cow) \mapsto (4, 4)

- (ii) Give all strategies on the Pareto front (cf. Definition 7 in the appendix) for the *Cowmunication Game*. Give your reasoning! (4pts)

All strategies are Pareto-dominated by the (4,4)-payoff pairs (ignore, Farmer) and ('brrrr?', Brave Cow), which are therefore the only two remaining strategies on the Pareto front.

7 Scientific Reading: Model-Based Reinforcement Learning

18pts

After some tweaking, a large language model generated the following (correct) explanation of the established technique *model-based reinforcement learning*, which we did not discuss in the lecture. Read this explanation and answer the questions afterwards!

In the context of model-based reinforcement learning (MBRL), a “model” refers to an abstract representation or approximation of the environment in which the agent operates. This model captures the dynamics of the environment, including how states transition to other states and the rewards associated with these transitions, based on actions taken by the agent.

Model-based reinforcement learning is a sophisticated approach in the field of artificial intelligence, where the system explicitly constructs this model of the environment. Unlike model-free methods, which learn solely from trial-and-error interactions with the environment, MBRL leverages its internal model to simulate outcomes, enabling it to anticipate future states and make more informed decisions. The internal model can be either learned from interactions with the real environment or provided *a priori*.

In MBRL, the model interacts with other components of the learning process in a pivotal way. It serves as a predictive engine that forecasts the next state and potential rewards for given actions. The agent uses this model for planning by simulating actions in a “mental rehearsal” before executing them in the real environment. This ability to predict and evaluate potential future scenarios allows the agent to optimize its policy with fewer real-world interactions, which is particularly valuable in scenarios where experimentation is costly or dangerous.

The core components of an MBRL system include the model of the environment, the policy that dictates the agent’s behavior, and the value function, which estimates the expected return from states or state-action pairs. The environment model is central, providing a bridge between the policy and value function by predicting the outcomes of actions. Learning can occur in two key areas: improving the environment model to better predict future states and rewards, and optimizing the policy and value function based on predictions from the model.

The distinction between model-based RL and simpler value or Q-value approaches (model-free methods) lies in their strategies for improving the reward. Model-free methods learn a direct mapping from states or state-action pairs to values or Q-values, respectively, relying heavily on extensive interaction with the environment to iteratively improve their estimates. They excel in environments where the model is hard to learn, but they often require a large amount of experience to converge to an optimal policy. On the other hand, MBRL seeks to understand the dynamics of the environment through its model, which allows it to anticipate and plan for future states. This foresight enables MBRL to potentially achieve higher rewards with fewer interactions by using its model to explore outcomes of actions without physically executing them, making it more efficient in environments where building a reasonably accurate model is feasible.

(i) In Scenario 1, a cow agent $G^{[7]}$ might want to mirror all actions performed by the much cooler cow agent $G^{[6]}$ except that $G^{[7]}$ wants to execute `eat/mow` whenever $G^{[6]}$ executes `do_nothing` and vice versa. Give a policy $\pi^{[7]}$ for $G^{[7]}$ that uses a *model* $\mu^{[6]}$ of $G^{[6]}$'s behavior. Explain how the usage of a *model* matters in your example. (12pts)

$$\pi^{[7]}(\dots) = \begin{cases} \text{eat/mow} & \text{if } \mu^{[6]}(t) = \text{do_nothing}, \\ \text{do_nothing} & \text{if } \mu^{[6]}(t) = \text{eat/mow}, \\ \mu^{[6]}(t) & \text{otherwise} \end{cases}$$

with $\mu^{[6]} : \mathbb{N} \rightarrow \mathcal{A}^{[6]}$

The model $\mu^{[6]}$ is what allows $\pi^{[7]}$ to predict agent $G^{[6]}$'s actions, so it is essential to adapt specifically to $G^{[6]}$ in a given situation instead of just learning a policy that “accidentally” matches the desired behavior. This is important, when training data on $G^{[6]}$ is easy to come by.

(ii) In Scenario 2, imagine an “irrational math cow” agent whose policy is given by the irrational number $u = \frac{\pi}{10} = 0.314\dots$, i.e., π (the policy) is given by $\frac{\pi}{10}$ (the number). Recall that, for any irrational number, there is no repeating pattern in its decimal digits. Will model-based reinforcement learning be helpful for the other cows to model the behavior of the “irrational math cow” agent? Explain your reasoning. (6pts)

Since decimals of π (the number) cannot be predicted from previous decimals, let alone from the position of the digit (as it is done with $\mu^{[6]}$ above), we will have a hard time modeling the crazy cow with traditional models (assuming finite memory, e.g.). Thus, we can assume that model-based reinforcement learning will have a hard time modeling an “irrational math cow” and thus will probably not be helpful.

Appendix: Definitions

Notation. \mathbb{B} is the set of truth values or Booleans, i.e., $\mathbb{B} = \{0, 1\}$. \mathbb{N} is the set of natural numbers starting from zero, i.e., $\mathbb{N} = \{0, 1, 2, \dots\}$ so that it holds that $\mathbb{B} \subset \mathbb{N} \subset \mathbb{Z} \subset \mathbb{R} \subset \mathbb{C}$. \mathbb{P} denotes the space of probabilities and \mathbb{F} denotes the space of fuzzy values with $\mathbb{P} = \mathbb{F} = [0; 1] \subset \mathbb{R}$ being only discerned for semantic but not mathematical reasons. $\wp(X)$ denotes the power set of X . \mathbb{E} denotes the expected value. $\#$ denotes vector or sequence concatenation, i.e., given two vectors $\mathbf{x} = \langle x_1, \dots, x_{|\mathbf{x}|} \rangle$ and $\mathbf{y} = \langle y_1, \dots, y_{|\mathbf{y}|} \rangle$, $\mathbf{x} \# \mathbf{y} = \langle x_1, \dots, x_{|\mathbf{x}|}, y_1, \dots, y_{|\mathbf{y}|} \rangle$. A vector $\langle x_0, \dots, x_{n-1} \rangle$ with length $n \in \mathbb{N}$ can also be written as $\langle x_i \rangle_{0 \leq i \leq n-1}$ for a new iteration variable i . $_$ denotes unspecified function arguments ($f(_) = 0$ is the constant function that always returns zero, e.g.) or other variables that are not used. For any finite set $X = \{x_0, \dots, x_n\}$, $|X| = n$ denotes the number of elements in X . For infinite sets, $|X| = \infty$.

Definition 1 (agent). Let \mathcal{A} be a set of actions. Let \mathcal{O} be a set of observations. An agent A can be given via a policy function $\pi : \mathcal{O} \rightarrow \mathcal{A}$. Given a time series of observations $\langle o_t \rangle_{t \in \mathcal{Z}}$ for some time space \mathcal{Z} the agent can thus generate a time series of actions $\langle a_t \rangle_{t \in \mathcal{Z}}$ by applying $a_t = \pi(o_t)$.

Definition 2 (optimization). Let \mathcal{X} be an arbitrary state space. Let \mathcal{T} be an arbitrary set called target space and let \leq be a total order on \mathcal{T} . A total function $\tau : \mathcal{X} \rightarrow \mathcal{T}$ is called target function. Optimization (minimization/maximization) is the procedure of searching for an $x \in \mathcal{X}$ so that $\tau(x)$ is optimal (minimal/maximal). Unless stated otherwise, we assume minimization.

An optimization run of length $g+1$ is a sequence of states $\langle x_t \rangle_{0 \leq t \leq g}$ with $x_t \in \mathcal{X}$ for all t .

Let $e : \langle \mathcal{X} \rangle \times (\mathcal{X} \rightarrow \mathcal{T}) \rightarrow \mathcal{X}$ be a possibly randomized or non-deterministic function so that the optimization run $\langle x_t \rangle_{0 \leq t \leq g}$ is produced by calling e repeatedly, i.e., $x_{t+1} = e(\langle x_u \rangle_{0 \leq u \leq t}, \tau)$ for all t , $1 \leq t \leq g$, where x_0 is given externally (e.g., $x_0 =_{\text{def}} 42$) or chosen randomly (e.g., $x_0 \sim \mathcal{X}$). An optimization process is a tuple $(\mathcal{X}, \mathcal{T}, \tau, e, \langle x_t \rangle_{0 \leq t \leq g})$.

Algorithm 1 (simulated annealing). Let $\mathcal{D} = (\mathcal{X}, \mathcal{T}, \tau, e, \langle x_u \rangle_{0 \leq u \leq t})$ be an optimization process. Let $neighbors : \mathcal{X} \rightarrow \wp(\mathcal{X})$ be a function that returns a set of neighbors of a given state $x \in \mathcal{X}$. Let $\kappa : \mathbb{N} \rightarrow \mathbb{R}$ be a temperature schedule, i.e., a function that returns a temperature value for each time step. Let $A : \mathcal{T} \times \mathcal{T} \times \mathbb{R} \rightarrow \mathbb{P}$ with $\mathbb{P} = [0; 1] \subset \mathbb{R}$ be a function that returns an acceptance probability given two target values and a temperature. Commonly, we use

$$A(Q, Q', K) = e^{\frac{-(Q' - Q)}{K}} \text{ with Euler's number } e$$

for $\mathcal{T} \subseteq \mathbb{R}$. The process \mathcal{D} continues via simulated annealing if e is of the form

$$e(\langle x_u \rangle_{0 \leq u \leq t}, \tau) = x_{t+1} = \begin{cases} x'_t & \text{if } \tau(x'_t) \leq \tau(x_t) \text{ or } r \leq A(\tau(x_t), \tau(x'_t), \kappa(t)), \\ x_t & \text{otherwise,} \end{cases}$$

where $x'_t \sim neighbors(x_t)$ and $r \sim \mathbb{P}$ are drawn at random for each call to e .

Definition 3 (Markov decision process (MDP)). Let \mathcal{A} be a set of actions. Let \mathcal{S} be a set of states. Let A be an agent given via a policy function $\pi : \mathcal{S} \rightarrow \mathcal{A}$. Let $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathcal{T}$ be a possibly randomized *cost (reward)* function. A Markov decision process is a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{T}, P, R)$ where $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{P}$ is the *transition probability function* often written as $P(s'|s, a) = \Pr(s_{t+1} = s' \mid s_t = s \wedge a_t = a)$.

A Markov decision process run for policy π is a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \tau, \pi, \langle s_t \rangle_{t \in \mathbb{Z}}, \langle a_t \rangle_{t \in \mathbb{Z}})$ where

$$s_{t+1} \sim P(s_{t+1} | s_t, a_t)$$

and where τ is to be minimized (maximized) and usually has a form similar to

$$\begin{aligned} \text{accumulated cost (reward)} \tau(\pi) &=_{\text{def}} \sum_{t \in \mathbb{Z}} R(s_t, a_t, s_{t+1}) \\ \text{or discounted expected cost (reward)} \tau(\pi) &=_{\text{def}} \mathbb{E} \left[\sum_{t \in \mathbb{Z}} \gamma^t \cdot R(s_t, a_t, s_{t+1}) \right] \end{aligned}$$

where $\gamma \in [0; 1] \subset \mathbb{R}$ is called a discount factor.

A decision process generates a (possibly infinite) series of rewards $\langle r_t \rangle_{t \in \mathbb{Z}}$ with $r_t = R(s_t, a_t, s_{t+1})$.

Definition 4 (multi-agent system). Let $G = \{G^{[1]}, \dots, G^{[N]}\}$ be a set of $|G| = N$ agents with observation spaces $\mathcal{O}^{[i]}$ and action spaces $\mathcal{A}^{[i]}$ controlled by policies $\pi^{[i]}$ for all $i = 1, \dots, N$, respectively. The multi-agent system G then takes a joint action $a \in \mathcal{A}$ with $\mathcal{A} = \mathcal{A}^{[1]} \times \dots \times \mathcal{A}^{[N]}$ after making a joint observation $o \in \mathcal{O}$ with $\mathcal{O} = \mathcal{O}^{[1]} \times \dots \times \mathcal{O}^{[N]}$ based on the joint policy $\pi(o^{[1]}, \dots, o^{[N]}) = (a^{[1]}, \dots, a^{[N]})$ where $a^{[i]} = \pi^{[i]}(o^{[i]})$ for all i .

Definition 5 (normal-form game). Let $G = \{G^{[1]}, \dots, G^{[N]}\}$ be a set of $|G| = N$ agents. Let $\mathcal{A} = \mathcal{A}^{[1]} \times \dots \times \mathcal{A}^{[N]}$ be the space of joint actions where $\mathcal{A}^{[i]}$ is the set of actions available to agent $G^{[i]}$ for all i . Let $\chi : \mathcal{A} \rightarrow \mathcal{T}$ be a utility function for the joint action space \mathcal{A} and the joint target space $\mathcal{T} = \mathcal{T}^{[1]} \times \dots \times \mathcal{T}^{[N]}$ where $\mathcal{T}^{[i]}$ is the target space of agent $G^{[i]}$ for all i . Unless stated otherwise, the utility χ is to be maximized. From χ we can derive a set of single-agent utility functions $\chi^{[i]} : \mathcal{A} \rightarrow \mathcal{T}^{[i]}$ for all i . A tuple $(G, \mathcal{A}, \mathcal{T}, \chi)$ is called a *normal-form game*.

Definition 6 (strategy). Let $(G, \mathcal{A}, \mathcal{T}, \chi)$ be a normal-form game. In a single iteration of the game, an agent $G^{[i]}$'s behavior is given by a (possibly randomized) policy $\pi^{[i]} : () \rightarrow \mathcal{A}^{[i]}$. Then, $\pi^{[i]}$ is also called $G^{[i]}$'s *strategy*. If multiple iterations of the game are played, an agent $G^{[i]}$'s behavior can be given by a (possibly randomized) policy $\pi^{[i]} : \mathcal{A}^n \rightarrow \mathcal{A}^{[i]}$ where n is a number of previous iterations. The agent's strategy is then conditioned on a list of previous joint actions.

An agent $G^{[i]}$ whose actions are given via a policy $\pi^{[i]}$ of the form $\pi^{[i]}(.) = a^{[i]}$ for some action $a^{[i]} \in \mathcal{A}^{[i]}$ is playing a pure strategy.

An agent $G^{[i]}$ whose actions are given via a policy $\pi^{[i]}$ of the form $\pi^{[i]}(.) \sim A^{[i]}$ according to some distribution over $A^{[i]} \subseteq \mathcal{A}^{[i]}$ is playing a mixed strategy.

If a mixed strategy is based on a uniform distribution over actions $A^{[i]} \subseteq \mathcal{A}^{[i]}$, we write $\pi^{[i]} = A^{[i]}$ as a shorthand.

Definition 7 (Pareto front for strategies). Let $(G, \mathcal{A}, \mathcal{T}, \chi)$ be a normal-form game. A joint strategy $\pi(-) = (\pi^{[1]}(-), \dots, \pi^{[|G|]}(-))$ Pareto-dominates another joint strategy $\pi'(-) = (\pi'^{[1]}(-), \dots, \pi'^{[|G|]}(-))$ iff for all agents $G^{[i]}$ it holds that

$$\chi^{[i]}(\pi(-)) \geq \chi^{[i]}(\pi'(-))$$

and there exists some agent $G^{[j]}$ so that

$$\chi^{[j]}(\pi(-)) > \chi^{[j]}(\pi'(-)).$$

A joint strategy π is Pareto-optimal iff there is no other strategy π' so that π' Pareto-dominates π .

The set of all Pareto-optimal strategies is called the Pareto front.

Definition 8 (best response). Let $(G, \mathcal{A}, \mathcal{T}, \chi)$ be a normal-form game. Let $\pi^{[-i]}$ be a joint strategy of agents $G^{[1]}, \dots, G^{[i-1]}, G^{[i+1]}, \dots, G^{[N]}$, i.e., all agents except $G^{[i]}$. Let $\pi^{[i]} \oplus \pi^{[-i]}$ be a joint strategy of all agents then. Given a strategy $\pi^{[-i]}$ for all agents except $G^{[i]}$, $G^{[i]}$'s *best response* is the strategy $\pi^{*[i]}$ so that for all strategies $\pi'^{[i]}$ it holds that

$$\chi^{[i]}((\pi^{*[i]} \oplus \pi^{[-i]})(-)) \geq \chi^{[i]}((\pi'^{[i]} \oplus \pi^{[-i]})(-)).$$

Definition 9 (Nash equilibrium). Let $(G, \mathcal{A}, \mathcal{T}, \chi)$ be a normal-form game played for a single iteration. A joint strategy π is a *Nash equilibrium* iff for all agents $G^{[i]}$ it holds that $\pi^{[i]}$ is the best response to $\pi^{[-i]}$.

Definition 10 (π -process). Let N be a set of names ($N = \{\text{“a”}, \text{“b”}, \text{“c”}, \dots\}$, e.g.). \mathbb{L}_π is the set of valid processes in the π -calculus given inductively via:

- (null process)** $0 \in \mathbb{L}_\pi$,
- (τ prefix)** if $P \in \mathbb{L}_\pi$, then $\tau.P \in \mathbb{L}_\pi$,
- (receiving prefix)** if $a, x \in N$ and $P \in \mathbb{L}_\pi$, then $a(x).P \in \mathbb{L}_\pi$,
- (sending prefix)** if $a, x \in N$ and $P \in \mathbb{L}_\pi$, then $\bar{a}\langle x \rangle.P \in \mathbb{L}_\pi$,
- (choice)** if $P, Q \in \mathbb{L}_\pi$, then $(P + Q) \in \mathbb{L}_\pi$,
- (concurrency)** if $P, Q \in \mathbb{L}_\pi$, then $(P \mid Q) \in \mathbb{L}_\pi$,
- (scoping)** if $x \in N$, $P \in \mathbb{L}_\pi$, then $(\nu x) P \in \mathbb{L}_\pi$,
- (replication)** if $P \in \mathbb{L}_\pi$, then $!P \in \mathbb{L}_\pi$.

Any element $P \in \mathbb{L}_\pi$ is called a π -process. If the binding order is clear, we leave out parentheses.

The free names of a π -process $P \in \mathbb{L}_\pi$, written $\mathfrak{F}(P)$ with $\mathfrak{F} : \mathbb{L}_\pi \rightarrow \wp(N)$ are given inductively via:

- $\mathfrak{F}(0) = \emptyset$,
- $\mathfrak{F}(\tau.P) = \mathfrak{F}(P)$,
- $\mathfrak{F}(a(x).P) = \{a\} \cup (\mathfrak{F}(P) \setminus \{x\})$,
- $\mathfrak{F}(\bar{a}\langle x \rangle.P) = \{a, x\} \cup \mathfrak{F}(P)$,
- $\mathfrak{F}(P + Q) = \mathfrak{F}(P) \cup \mathfrak{F}(Q)$,
- $\mathfrak{F}(P \mid Q) = \mathfrak{F}(P) \cup \mathfrak{F}(Q)$,
- $\mathfrak{F}((\nu x) P) = \mathfrak{F}(P) \setminus \{x\}$,
- $\mathfrak{F}(!P) = \mathfrak{F}(P)$,

for any $a, x \in N$ and any $P, Q \in \mathbb{L}_\pi$.

Definition 11 (π -congruence). Two π -processes $P, Q \in \mathbb{L}_\pi$ are structurally congruent, written $P \equiv Q$, if they fulfill the predicate $\equiv: \mathbb{L}_\pi \times \mathbb{L}_\pi \rightarrow \mathbb{B}$. We define inductively:

(α -conversion) $P \equiv Q$ if both only differ by the choice of bound names,

(choice rules) $P + Q \equiv Q + P$, and $(P + Q) + R \equiv P + (Q + R)$, and $P \equiv P + P$

(concurrency rules) $P \mid Q \equiv Q \mid P$, and $(P \mid Q) \mid R \equiv P \mid (Q \mid R)$, and $P \equiv P \mid 0$,

(scoping rules) $(\nu x) (\nu y) P \equiv (\nu y) (\nu x) P$, and $(\nu x) (P \mid Q) \equiv P \mid ((\nu x) Q)$ if $x \notin \mathfrak{F}(P)$, and $(\nu x) 0 \equiv 0$,

(replication rules) $!P \equiv P \mid !P$,

for any names $x, y \in N$ and processes $P, Q, R \in \mathbb{L}_\pi$.

Definition 12 (π -substitution). For a π -process P we write $P[y := z]$ for the π -process where every free occurrence of name y is replaced by name z . Formally, we define:

- $0[y := z] = 0$,
- $(\tau.P)[y := z] = \tau.(P[y := z])$,
- $(a(x).P)[y := z] = z(x).P$ for $a = y$ and $x = y$,
- $(a(x).P)[y := z] = a(x).P$ for $a \neq y$ and $x = y$,
- $(a(x).P)[y := z] = z(x).(P[y := z])$ for $a = y$ and $x \neq y$,
- $(a(x).P)[y := z] = a(x).(P[y := z])$ for $a \neq y$ and $x \neq y$,
- $(\bar{a}\langle x \rangle.P)[y := z] = \bar{z}\langle z \rangle.(P[y := z])$ for $a = y$ and $x = y$,
- $(\bar{a}\langle x \rangle.P)[y := z] = \bar{a}\langle z \rangle.(P[y := z])$ for $a \neq y$ and $x = y$,
- $(\bar{a}\langle x \rangle.P)[y := z] = \bar{z}\langle x \rangle.(P[y := z])$ for $a = y$ and $x \neq y$,
- $(\bar{a}\langle x \rangle.P)[y := z] = \bar{a}\langle x \rangle.(P[y := z])$ for $a \neq y$ and $x \neq y$,
- $(P + Q)[y := z] = (P[y := z]) + (Q[y := z])$,
- $(P \mid Q)[y := z] = (P[y := z]) \mid (Q[y := z])$,
- $(\nu x) P[y := z] = ((\nu x) P)$ for $x = y$,
- $(\nu x) P[y := z] = (\nu x) (P[y := z])$ for $x \neq y$,
- $(!P)[y := z] = !(P[y := z])$,

for any names $a, x, y, z \in N$ and processes $P, Q \in \mathbb{L}_\pi$.

Definition 13 (π -evaluation). An evaluation of a π -process P is a sequence of π -processes $P_0 \succ \dots \succ P_n$ where $P_0 = P$ and P_{i+1} is generated from P_i via the application of an evaluation rule $\succ: \mathbb{L}_\pi \rightarrow \mathbb{L}_\pi$ to any sub-term of P_i . We define the following evaluation rules:

(reaction) $(a(x).P + P') \mid (\bar{a}\langle y \rangle.Q + Q') \succ_{\text{REACT}} (P[x := y]) \mid Q,$

(τ transition) $\tau.P + P' \succ_{\text{TAU}} P,$

(parallel execution) $P \mid R \succ_{\text{PAR}} Q \mid R$ if it holds that $P \succ Q,$

(restricted execution) $(\nu x) P \succ_{\text{RES}} (\nu x) Q \mid R$ if it holds that $P \succ Q,$

(structural congruence) $P' \succ_{\text{STRUCT}} Q'$ if it holds that $P \succ Q$ and $P \equiv P'$ and $Q \equiv Q',$

for any names $a, x, y \in N$ and processes $P, P', Q, Q' \in \mathbb{L}_\pi$ where $\equiv: \mathbb{L}_\pi \times \mathbb{L}_\pi \rightarrow \mathbb{B}$ is the predicate for structural congruence.