



Teacher: Dr. Cécile Hardebolle
CS-290: Responsible Software
28/10/2025

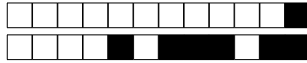
Name Firstname

SCIPER: XXXX

Do not turn the page before the start of the exam. This document is double-sided, has 16 pages, the last ones possibly blank. Do not unstaple.

- Place your student card on your table.
- No paper materials other than **one sheet of notes in A4 format, double-sided**, are allowed to be used during the exam.
- The use of a calculator or **any other electronic device** is not permitted during the exam.
- **First part: single choice questions** (12 questions, 12 points).
For the singles choice questions, we give :
 - +1 point if your answer is correct,
 - 0 points if you give no answer or your answer is incorrect.
- **Second Part: true/false questions** (4 questions, 4 points).
For the true/false questions, we give :
 - +1 point if your answer is correct,
 - 0 points if you give no answer or your answer is incorrect.
- **Third part: case studies** (3 questions, 20 points).
The number of points is noted above each question. Leave the checkboxes empty.
- Use a **black or dark blue ballpen** and clearly erase with **correction fluid** if necessary.
- If a question is wrong, the teacher may decide to nullify it.

Respectez les consignes suivantes Observe this guidelines Beachten Sie bitte die unten stehenden Richtlinien		
choisir une réponse select an answer Antwort auswählen	ne PAS choisir une réponse NOT select an answer NICHT Antwort auswählen	Corriger une réponse Correct an answer Antwort korrigieren
ce qu'il ne faut PAS faire what should NOT be done was man NICHT tun sollte		



First part: single choice questions

For each question, mark the box corresponding to the single correct answer.
1 point per question.

Question 1

The CEO of a tech company stated in the media: “In the past, we’ve invested in technology to positively impact people’s lives, and we have no intention of changing that strategy in the future - technology remains the best alternative.”

We may interpret this as:

- Representation bias
- System 1 thinking
- Implicit stereotype
- Sunk cost fallacy

Question 2

A company has developed a complex algorithm to predict whether athletes suspected of doping actually do it. A positive result means that the algorithm classifies the athlete as at risk of doping, while a negative result means no risk of doping. The system has been used for 5 years and we have access to data about athletes that were indeed caught for doping. We found that the proportion of athletes predicted to dope amongst all predictions is higher for men rather than for women.

The type of fairness metric we have used is:

- Conditional use accuracy equality
- Error rate balance
- Equal accuracy
- Demographic parity

Question 3

Imagine that you develop software for people from a single country. If you nonetheless envision cultural differences in this context, which strategy are you probably using?

- Edge cases
- STRIDE
- The people behind the data
- Bad actors

Question 4 You develop a software that analyzes the weather forecast to send the population a notification in case of upcoming extreme rain (positive result).

In this context:

- True Positive = rain is predicted and the prediction is correct
- False Positive = rain is predicted and the prediction is correct
- True Negative = no rain is predicted, and the prediction is incorrect
- False Negative = no rain is predicted, and the prediction is correct



Question 5

Fill the blanks:

If a piece of software behaves in a ___ way at first glance, but puts people of ___ at ___, then it is a case of ___ discrimination.

- neutral / several groups / a disadvantage / direct
- positive / identified groups / an advantage / inverse
- negative / several groups / an advantage / direct
- negative / specific groups / a disadvantage / indirect
- neutral / specific groups / a disadvantage / indirect

Question 6 You develop a software that analyzes the weather forecast to send the population a notification in case of upcoming extreme rain (positive result).

The False Negative Rate (FNR) is:

- The number of times no rain is predicted among all times it actually didn't rain
- The number of times rain is predicted among all times it actually didn't rain
- The number of times rain is predicted among all times it actually rained
- The number of times no rain is predicted among all times it actually rained

Question 7

A group of computer scientists with similar background, all experts in software development, are starting a new software project for healthcare. They are aware of cognitive biases and want to minimize the impact of these biases when making design decisions.

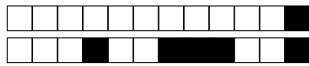
Which is the only strategy that could effectively help them in this context?

- Use a structured approach and slow down the decision-making process
- Choose one or two of them to play the devil's advocate
- Systematically include all members of their group to increase heterogeneity
- Systematically include all members of their group to apply a participatory design method

Question 8 You work on a chatbot to provide students assistance on campus questions. When evaluating the quality of the responses it provides, you identify that the responses contain hallucinations (i.e. content that is incorrect or wrong but looks perfectly plausible) with a 15% rate.

What type of situation are you facing?

- Ethical blindness
- Ethical issue
- Ethical dilemma
- Ethical sensitivity



Question 9

Here are three variables:

- Disinformation spread
- Public trust in information
- Development of disinformation software

We know that :

- As the spread of disinformation increases, the public trust in information decreases
- As the public trust in information decreases, bad actors see a growing opportunity to develop disinformation software exploiting this mistrust

In a causal loop diagram representing the dynamics between these variables, which arrows would we have (select only one answer)?

- The arrow between “Public trust in information” and “Development of disinformation software” has a negative sign and the arrow between “Development of disinformation software” and “Disinformation spread” has a positive sign.
- The arrow between “Public trust in information” and “Development of disinformation software” has a negative sign and the arrow between “Development of disinformation software” and “Disinformation spread” has a negative sign.
- The arrow between “Public trust in information” and “Development of disinformation software” has a positive sign and the arrow between “Development of disinformation software” and “Disinformation spread” has a negative sign.
- The arrow between “Public trust in information” and “Development of disinformation software” has a positive sign and the arrow between “Development of disinformation software” and “Disinformation spread” has a positive sign.

Question 10 A bad actor launched a phishing attack on employees of Swiss public institutions to steal their login credentials. An online media outlet reported on it, with the most upvoted comments on the article criticizing the institutions for their inability to counter online threats, harming their reputation. what type of impact is the harm to reputation as a result of the attack?

- Indirect
- Direct
- Both direct and indirect
- Neither direct nor indirect

Question 11 What is a dilemma?

- A situation in which you have to weigh the pros and cons of each decision (and their consequences) and choose the one with the higher number of pros.
- A situation in which you have to weigh the pros and cons of each decision (and their consequences), with no decision 100% perfect or 100% imperfect
- A situation in which you have to decide between two alternatives using a coin flip (better to leave things to chance)
- A situation in which you should escalate the decision to your management line.

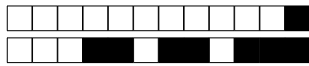


Question 12

A start-up developed a machine learning model designed to connect people based on their personal interests. A big company has then bought the start-up and is currently using the algorithm to connect jobseekers with employers.

What type of bias is likely to appear?

- Aggregation bias
- Measurement bias
- Intersectional bias
- Deployment bias



Second part: true/false questions

For each question, mark either the box TRUE if the statement is true or the box FALSE if the statement is false.

1 point per question.

You found a dataset with 5 variables, all self-reported by participants: eye-color, extraversion and 3 health-related variables. When analyzing the data you identify that:

- there are positive and substantial correlations among the 3 health variables
- there is a positive and substantial correlation between eye-color and extraversion
- there is no correlation between eye-color and the health variables

Question 13 Eye-color is a latent variable

TRUE FALSE

Question 14 Eye-color is a sensitive attribute

TRUE FALSE

Question 15 Eye-color is a proxy for health

TRUE FALSE

Question 16 Extraversion is a latent variable

TRUE FALSE



Third part: case Studies

Answer in the empty space provided. Use the extra pages at the end if you need more space. Your answer should be concise but make your reasoning clear and your argument should be explicitly justified. Leave the check-boxes empty, they are used for grading.

Question 17: Case 1: Harms modeling - Social assistant chatbot *This question is worth 5 points.*

<input type="checkbox"/>	0	<input type="checkbox"/>	.5	<input type="checkbox"/>	1	<input type="checkbox"/>	.5	<input type="checkbox"/>	2	<input type="checkbox"/>	.5	<input type="checkbox"/>	3	<input type="checkbox"/>	.5	<input type="checkbox"/>	4	<input type="checkbox"/>	.5	<input type="checkbox"/>	5
--------------------------	---	--------------------------	----	--------------------------	---	--------------------------	----	--------------------------	---	--------------------------	----	--------------------------	---	--------------------------	----	--------------------------	---	--------------------------	----	--------------------------	---

Scenario:

In the realm of technological innovation, a revolutionary social-assistant chatbot emerges, designed to offer guidance on relationships. This cutting-edge human-centered AI, inspired by Snapchat’s AI chatbot, aims to become an indispensable part of people’s lives. Sarah and James are two individuals with contrasting lives. James, a young artist, craves genuine connections with like-minded people, while Sarah, a young consultant, struggles to balance her career with her personal life. Sarah and James turn to this chatbot for relationship advice. Its sophisticated algorithm analyzes their preferences, communication styles, and social behaviors to offer tailored suggestions for interactions. In addition to helping them identify others’ emotions, it provides them with conversation starters and even helps plan memorable dates. As the chatbot gains traction and spreads throughout society, it becomes an integral part of society’s social, economic, and political landscape. It reshapes how people approach dating and relationships, influencing not only their personal lives but also impacting the dating industry, advertising strategies, and even political campaign tactics. In addition, companies rely on the chatbot to predict the emotions of their staff and their clients to maximize their benefits. Yet, there are those who remain skeptical of the chatbot’s far-reaching influence. Some individuals, wary of data privacy concerns and the potential for manipulation, opt to abstain from using the technology. They seek more traditional avenues for forming connections, believing in the value of genuine human interactions and the potential risks that come with relying on AI for personal advice.

Task:

Considering the following extract of the harms modeling table, describe what should go in the different cells:

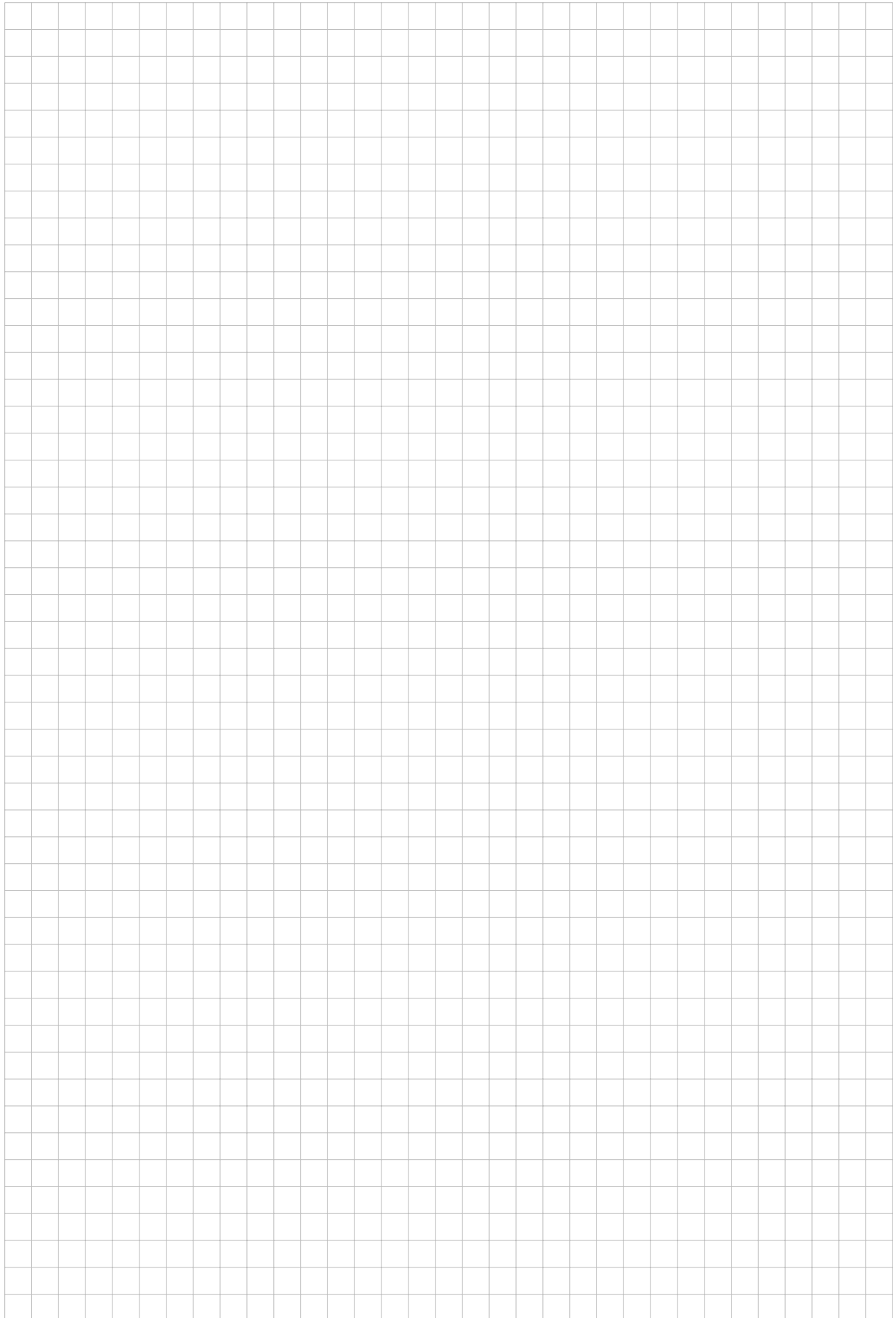
- [4 x 1 point] For cells A, B, C and E: describe 1 harm that corresponds to the category (1-2 sentences for each harm)
- [1 point] For cell D: indicate the corresponding harm category

Make sure to identify your answers with the corresponding letters (no need to reproduce the table).

Category	Type of harm	Description of harms
Humans	Physical injury	A)
Allocation of Resources	Opportunity loss	B)
Human Rights	Liberty loss	C)
	D)	Most intimate feelings are now “public”
Social System Harms	Social detriment	E)



+1/8/53+





Question 18: Case 2: Values analysis - Personalized deals *This question is worth 7 points.*

<input type="checkbox"/>	0	<input type="checkbox"/>	.5	<input type="checkbox"/>	1	<input type="checkbox"/>	.5	<input type="checkbox"/>	2	<input type="checkbox"/>	.5	<input type="checkbox"/>	3	<input type="checkbox"/>	.5
<input type="checkbox"/>	4	<input type="checkbox"/>	.5	<input type="checkbox"/>	5	<input type="checkbox"/>	.5	<input type="checkbox"/>	6	<input type="checkbox"/>	.5	<input type="checkbox"/>	7	<input type="checkbox"/>	

Scenario:

A webshop manager wants to offer interesting deals to the shop’s customers, and thinks that it would be best to offer personalized deals to each one of them. As the customers provide their email address when registering, the manager creates the following script: for each user, the script finds some account linked to the mail address (Facebook, YouTube, Amazon, Retail stores, etc.) and buys the data related to that user. With that data, a personalized offer containing deals adapted to the centers of interest of the user is sent directly by email.

Task:

Your overall task is to perform an analysis of the values and value tensions involved for the different stakeholders in the case.

We provide the following stakeholders:

- Lydia, the webstore manager.
- Hari, a customer who browses and purchases products on the webshop.

Follow the 2 steps below:

1 [5 points] Identify 2 values from stakeholders that are supported by the software (= 2 value-based benefits) and 2 values from stakeholders that are opposed by the software (= 2 value-based harms), i.e., 4 values in total.

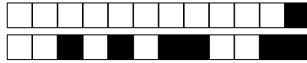
Consider the value-based benefit/harm table template below and describe what would go in each cell for each of the values you identified:

- (A) [not graded] Name the stakeholder, who must be one of the 2 stakeholders mentioned above
- (B) [4 x 0.5 points] Name the value (you should use the names in Appendix 3.1) and explain in your own words what the value means for this stakeholder
- (C) [4 x 0.25 points] Indicate if the value is supported (value-based benefit) or harmed (value-based harm) for this stakeholder
- (D) [4 x 0.5 points] Justify why it is supported / harmed by the software

Make sure to identify your answers with the corresponding letters (no need to reproduce the table). A list of Schwartz’s values is provided in appendix 3.1.

2 [2 points] Draw a value-based tension map showing at least 1 value tension and provide an explanation of the tension.

Stakeholder	Key Value	Benefits	Harms	Justification
Stakeholder: (A)	Value name and description: (B)	Benefit or Harm: (C)		It’s a value-based benefit/harm for this stakeholder because: (D)



Appendix 3.1

Table 1: Schwartz' Value Table - Source: Schwartz et al. (2012).

Self-enhancement	Power Resources	Power through control of material and social resources
	Power Dominance	Power through exercising control over people
	Achievement	Personal success through demonstrating competence according to social standards
	Hedonism	Pleasure and sensuous gratification for oneself
Openness to change	Stimulation	Excitement, novelty, and challenge in life
	Self-direction Action	The freedom to determine one's own actions
	Self-direction Thought	The freedom to cultivate one's own ideas and abilities
Self-transcendence	Universalism Tolerance	Acceptance and understanding of those who are different from oneself
	Universalism Concern	Commitment to equality, justice, and protection for all people
	Universalism Nature	Preservation of the natural environment
	Humility	Recognizing one's insignificance in the larger scheme of things
	Benevolence Dependability	Being a reliable and trustworthy member of the in-group
	Benevolence Caring	Devotion to the welfare of in-group members
Conservation	Tradition	Maintaining and preserving cultural, family, or religious traditions
	Conformity Interpersonal	Avoidance of upsetting or harming other people
	Conformity Rules	Compliance with rules, laws, and formal obligations
	Security Societal	Safety and stability in the wider society
	Security Personal	Safety in one's immediate environment
	Face	Security and power through maintaining one's public image and avoiding humiliation

