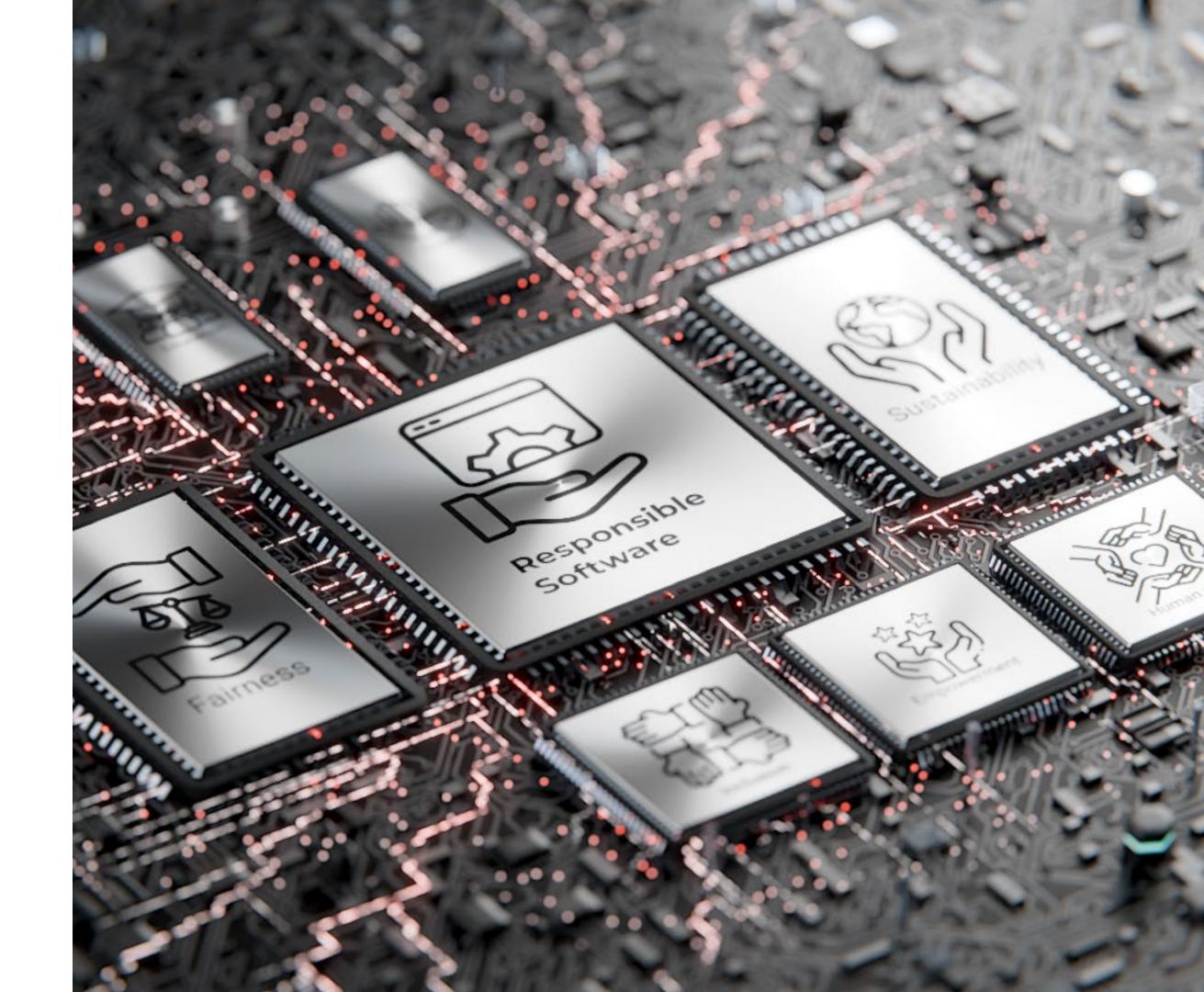
EPFL

Conclusion
Case studies
+ Q&A
10 dec.

Cécile Hardebolle

Responsible Software



Agenda for today

- 1. Review case studies:
 - a) Digital Ethics Canvas Emotion Cancelling Al
 - b) Ethics Canvas Be My Eyes
- 2. Questions & Answers

3. Some interactive review questions

Case studies

Where to find the cases?

1. Go to moodle

- 2. Find the link to the case studies for today: Conclusion
- 3. Download the instruction sheet
- + From previous chapters, you will need:
 - Digital Ethics Canvas (7 Empowerment 1)
 - Ethics Canvas (2 Safety 2)

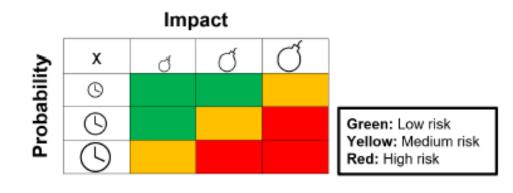
Digital Ethics Canvas

(review from Empowerment 1)

Instructions

- Read the software description
 (you can also take a look at the referenced news article)
- Fill out the Digital Ethics Canvas:
 - Context & Solution
 - Benefits: list 3 benefits (think about a range of stakeholders)
 - Risks:
 - ◆ For each of the 5 lenses identify and describe 1 risk
 - ◆ Select 1 risk and evaluate its **overall level**:
 - Severity of **impacts**
 - Probability to happen





Mitigation: for each risk, identify a corresponding mitigation measure

"Emotion Cancelling AI"



f 1 post = 1 benefit



https://speakup.epfl.ch

Room key: 57217



• 1 post = 1 risk + ethical lens

Post your ideas:

https://speakup.epfl.ch

Room key: 38467



Risks:

Make sure to explain how the risk relates to the corresponding ethical lens (e.g. if you put a risk into "Fairness", it must be clear what is unfair or biased).



Evaluating the level of risk - 1

URL: ttpoll.eu

Session ID: cs290

Consider the following Privacy risk: "Identifying customer emotions can lead to the disclosure of information the customers might consider private". How would you evaluate the level of this risk in terms of probability and severity of impacts? (select 2 options: 1 for probability, 1 for severity)

- a. Probability: low
- b. Probability: medium
- c. Probability: high
- d. Severity: low
- e. Severity: medium
- f. Severity: high

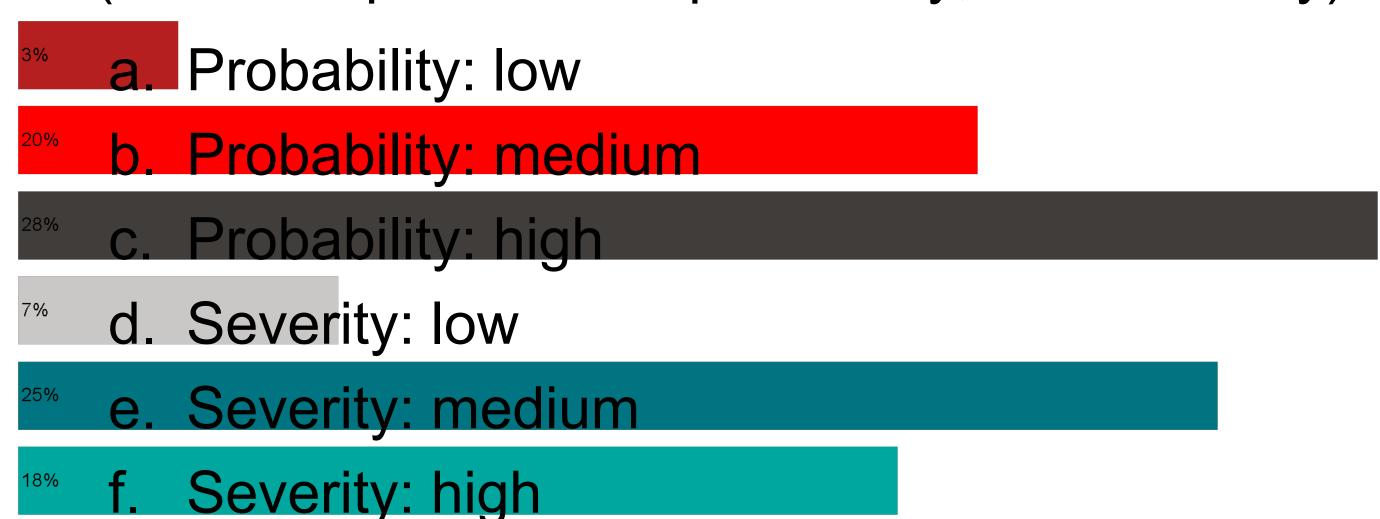
Make sure you know how to get the overall level of risk using the risk matrix

Evaluating the level of risk - 1

<u>URL:</u> ttpoll.eu

Session ID: cs290

Consider the following Autonomy risk: "The system's real-time alterations may reduce employee's ability to rely on their own judgment in emotionally charged situations". How would you evaluate the level of this risk in terms of probability and severity of impacts? (select 2 options: 1 for probability, 1 for severity)



"Emotion Cancelling AI": mitigation

Consider the following Privacy risk: "Identifying customer emotions can lead to the disclosure of information the customers might consider private" [HIGH RISK]

Which mitigation options could help reduce the risk?

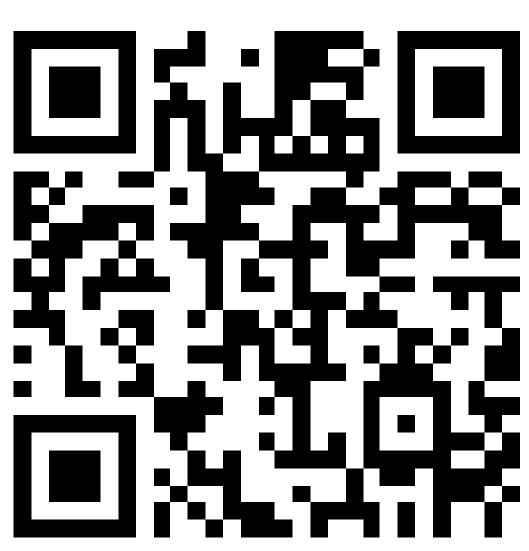


Make sure to explain how your proposal helps reduce the corresponding risk

Post your ideas:

https://speakup.epfl.ch

Room key: 02297



Ethics Canvas (review from Safety 2)

Instructions

- Read the software description
 (you can also take a look at the referenced website)
- Fill out the Ethics Canvas:
 - Stage 1: Identify relevant stakeholders
 - fill out blocks 1 and 2
 - Stage 2: Identify ethical impacts
 - fill out blocks 3, 4, 5, 6, 7 and 8
 - Stage 3: Discuss remedial actions
 - fill out block 9

"Be My Eyes"



1 post = 1 stakeholder

Post your ideas:

https://speakup.epfl.ch

Room key: **50113**





• 1 post = 1 ethical impact

Post your ideas:

https://speakup.epfl.ch

Room key: 01710



Comparison!

Ethics Canvas Ethics Canvas v1.8 - ethicscanvas.org © ADAPT Centre & Trinity College Dublin & Dublin City University, 2017. Project Title: Individuals affected Behaviour Groups affected What can we do? Worldviews Identify the types or categories of Discuss problematic changes to indi-Select the four most important Discuss how the general perception Identify the collectives or communiindividuals affected by the product vidual behaviour that may be prompt-Ethical impacts you discussed. of somebody's role in society can be ties, e.g. groups or organisations, or service, such as men/women, ed by the application e.g. differences in Identify ways of solving these affected by the project, that can be affected by your product user/non- user, age-category, etc. habits, time-schedules, choice of Impacts by changing your project's or service, such as environmental product/service design, organisaand religious groups, unions, profesactivities, people behaving more individualistic or collectivist, people tion.Or by providing recommendasional bodies, competing companies tions for its use or spelling out more behaving more or less materialistic. and government agencies, considerclearly to users the values driving ing any interest they might have in the effects of the product or service. **Group Conflicts** Relations Discuss the impact on the relation-Discuss problematic differences in ships between the groups identified, individual behaviour such as differe.g. employers and unions ences in habits, time-schedules, choice of activities, etc 9 Problematic Use of Resources Product or Service Failure Discuss the potential negative impact of your product or Discuss possible negative impacts of the consumption of service failing to operate as intended, eg technical or human resources of your project, e.g. climate impacts, privacy error, financial failure/ receivership/acquisition, security impacts, employment impacts etc. breach, data loss, etc.

DIGITAL ETHICS CANVAS			
NTEXT	SOLUTION	BENEFITS	

	WEI	FARE	
Can the solution be used in harmful ways, in particular with reg to vulnerable populations? What kind of impacts can errors from the solution have? What type of protection does the solution have against attacks.		MITIGATION	
00	000		o Ø Ø
	FAIF	RNESS	
How accessible is the solution? What kinds of biases may affect the results ? Can the outcomes of the solution be different for different user Could the solution contribute to discrimination against people of	a or grouped	MITIGATION	
00	000		ø ⊘ ⊙
AUTONOMY			
ISK Can users understand how the solution works and what its limit Are users able to make choices (e.g. concert, settings) in their is the solution and how? How does the solution affect user autoromy and agency?	os of	MITIGATION	
00	000		⊚ ⊙ ⊙
PRIVACY			
ISK What data does the solution collect Is it collecting personal or sensitive data Who has access to the data? How is the data protected? Could the solution disclose / be used to disclose private informs		MITIGATION	
00	000		ø Ø Ø
	SUSTAI	NABILITY	
ISK What is the carbon feetprint of the solution? What types of resources does it consume (e.g., water) -and prod What type of human labor is involved?	uch (e.c. wastel?	MITIGATION	
00	000		o Ø Ø
ek s iconsecuelor ICS1-56-40		Sighal Ethica Carrest (104 – E. Hardelalle, Y. Racks, Y. Annachandras, Y. Su, N. Bartisenigcoph, K. Kolac F. Jermann	

<u>**Q&A**</u>

Attributes: Sensitive or Protected?

"Are any sensitive or protected attributes defined in law? Can we have an exact definition of both? In the course we defined sensitive attributes as those that can have ethical implications, and protected that can harm."

- "Sensitive": can have ethical implications (privacy, fairness...)
- "Protected":
 - = Sensitive

 - Should not be used / should be treated specifically to prevent harm in software design, in particular ML

GDPR (EU) - Personal data

What is personal data?

Personal data is any information that relates to an identified or identifiable living individual (data subject). Different pieces of information, which together can lead to the identification of a particular person, may also be considered personal data.

Personal data that has been de-identified, encrypted or **pseudonymised** but can be used to re-identify a person remains personal data and falls within the scope of the General Data Protection Regulation (GDPR), the EU's main data protection law.

Examples of personal data

- a name and surname
- a home address
- an email address such as 'name.surname@company.com []'
- an Internet Protocol (IP) address
- an identification card number
- a cookie ID
- the advertising identifier of your phone
- data held by a hospital or doctor, which could be a symbol that uniquely identifies a person

Data protection rights

Under the GDPR, individuals have several rights over their personal data.

The rights of individuals

- Right to be informed
- Right of access
- Right to rectification
- Right to erasure
- Right to restriction of processing
- Right to data portability
- Right to object
- Rights in relation to automated decision-making and profiling

https://commission.europa.eu/law/law-topic/data-protection_en

GDPR (EU) - Sensitive data

What personal data is considered sensitive?

Answer

The following personal data is considered 'sensitive' and is subject to specific processing conditions:

- personal data revealing racial or ethnic origin, political opinions, religious or philosophical beliefs;
- trade-union membership;
- genetic data, biometric data processed solely to identify a human being;
- health-related data;
- data concerning a person's sex life or sexual orientation.

References

- Article 4(13), (14) and (15) and Article 9 and Recitals (51) to (56) of the GDPR

https://commission.europa.eu/law/law-topic/data-protection/rules-business-and-organisations/legal-grounds-processing-data/sensitive-data_en

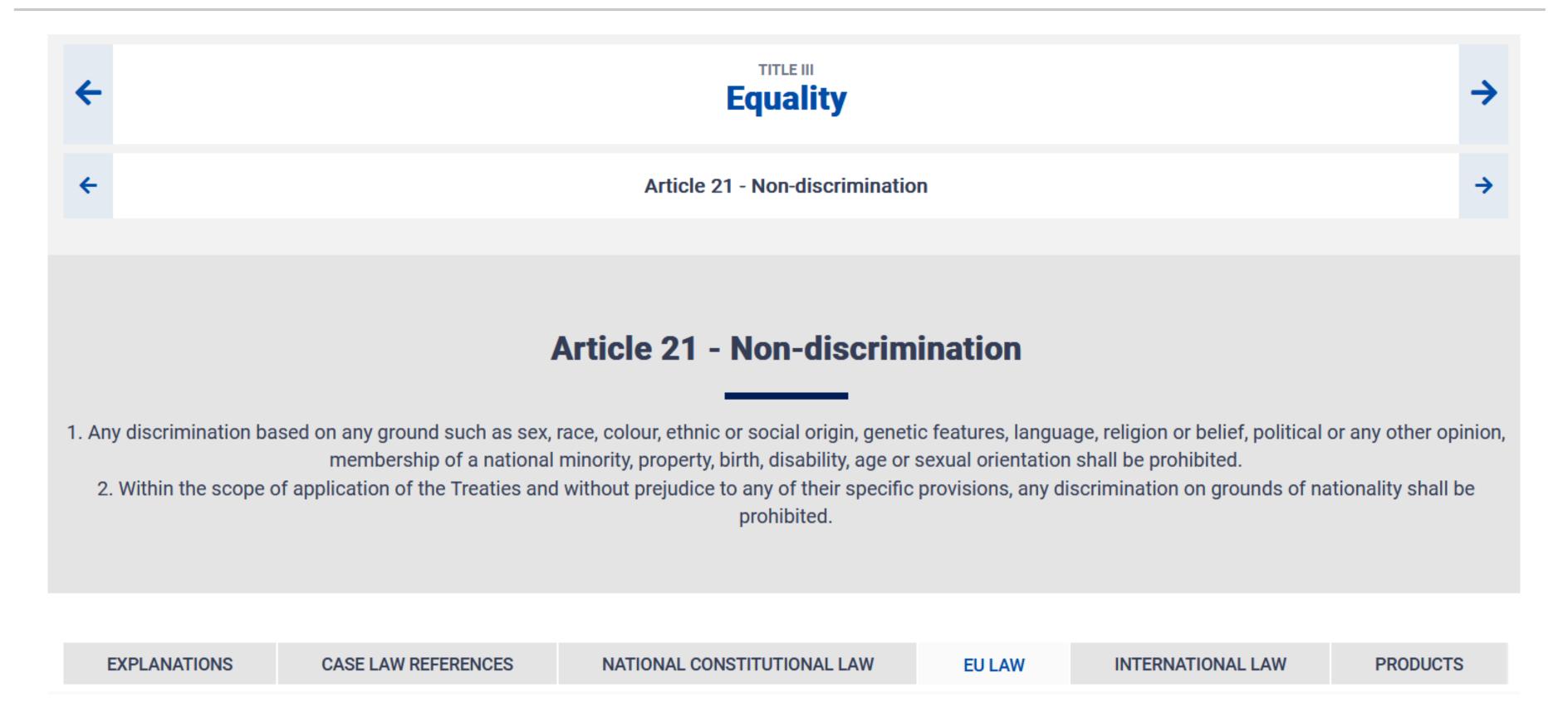
Under what conditions can my company/organisation process sensitive data?

Answer

Your company/organisation can only process sensitive data if one of the following conditions is met:

- the **explicit consent** of the individual was obtained (a law may rule out this option in certain cases);
- an EU or national law or a collective agreement, requires your company/organisation to process
 the data to comply with its obligations and rights, and those of the individuals, in the fields of
 employment, social security and social protection law;
- the vital interests of the person, or of a person physically or legally incapable of giving consent, are at stake
- you are a foundation, association or other not-for-profit body with a political, philosophical, religious or trade union aim, processing data about its members or about people in regular contact with the organisation;
- the personal data was manifestly made public by the individual;
- the data is required for the establishment, exercise or defence of legal claims
- the data is processed for reasons of substantial public interest on the basis of EU or national law;
- the data is processed for the purposes of preventive or occupational medicine, assessment of the working capacity of the employee, medical diagnosis, the provision of health or social care or treatment, or the management of health or social care systems and services on the basis of EU or national law, or on the basis of a contract as a health professional;
- the data is processed for reasons of public interest in the field of public health on the basis of EU or national law;
- the data is processed for archiving, scientific or historical research purposes or statistical purposes on the basis of EU or national law.

EU Charter of Fundamental Rights



Review questions "Whole Course"

URL: ttpoll.eu

Session ID: cs290

Ethical sensitivity refers to a person's... (select 1 answer)

- a. ... willingness to ackowledge their mistakes
- b. ... capability to understand the impact of a situation on others
 - c. ... ability to act to benefit others even at their own expense
 - d. ... ability to enforce ethical values in a team

URL: ttpoll.eu

Session ID: cs290

An ethical dilemma is a situation where... (select 1 answer)

- a. ... there is a solution that is better than the others
- b. ... it is possible to find an ethical compromise
- c. ... a good trade-off between solutions is possible
- d. ... all solutions conflict with an ethical value

Autonomous car software

URL: ttpoll.eu

Session ID: cs290

The software of an autonomous car fails to recognize traffic signs correctly.

We are in the presence of (select all that apply):

- a. A safety threat
- b. A security threat
- c. A safety hazard
 - d. A security hazard

The "confusing" matrix

URL: ttpoll.eu

Session ID: cs290

We use software to detect fissures in concrete walls before they become visible to the naked eye.

A positive result means presence of fissure.

Select all the correct statements:

- a. TN = actual absence of fissure, correct prediction
 b. TP = actual absence of fissure, correct prediction
- c. FN = actual presence of fissure, incorrect prediction
 - d. FP = actual presence of fissure, incorrect prediction

Disinformation

URL: ttpoll.eu

Session ID: cs290

A piece of information is false but created without intention to harm. It is (select all that apply):

- **'**
 - ^{59%} a. Misinformation
 - b. Disinformation
 - ° c. Malinformation
- **/**
- ^{20%} d. Fake news

False beliefs

URL: ttpoll.eu

<u>Session ID:</u> cs290

If you are exposed to a dis/mis-information post by Melon Husk, you are more likely to believe it because of (select 1 answer):

- System 2
- b. False consensus
- ^{83%} c. Source cues
 - d. Illusory truth

Attributes

Session ID: cs290

URL: ttpoll.eu

Hair color as an attribute to represent people is: (select all that apply)

- % a.
 - a. A sensitive attribute
- 0%
- b. An observed variable
- ° c. A latent variable
- d. An objective representation of people
- 0%
- e. A subjective representation of people

Bias

URL: ttpoll.eu

Session ID: cs290

The city of Lozhann decides to deploy a smartphone app that allows residents to report potholes throughout the city to help with the identification of repair needs.

The data collected by the app will probably exhibit: (select all that apply)

- a. Preexisting bias
 - b. Confirmation bias
- c. Representation bias
- d. Measurement bias
 - e. Automation bias



Biases in the ML lifecycle

URL: ttpoll.eu

Session ID: cs290

The society RetailProtect develops a ML model to identify instances of shoplifting in retail shops. They evaluate their model on a benchmark in which actors from diverse ethnicities simulate a range of shoplifting actions.

This is a case of (select 1 answer):



a. Evaluation bias

b. Aggregation bias

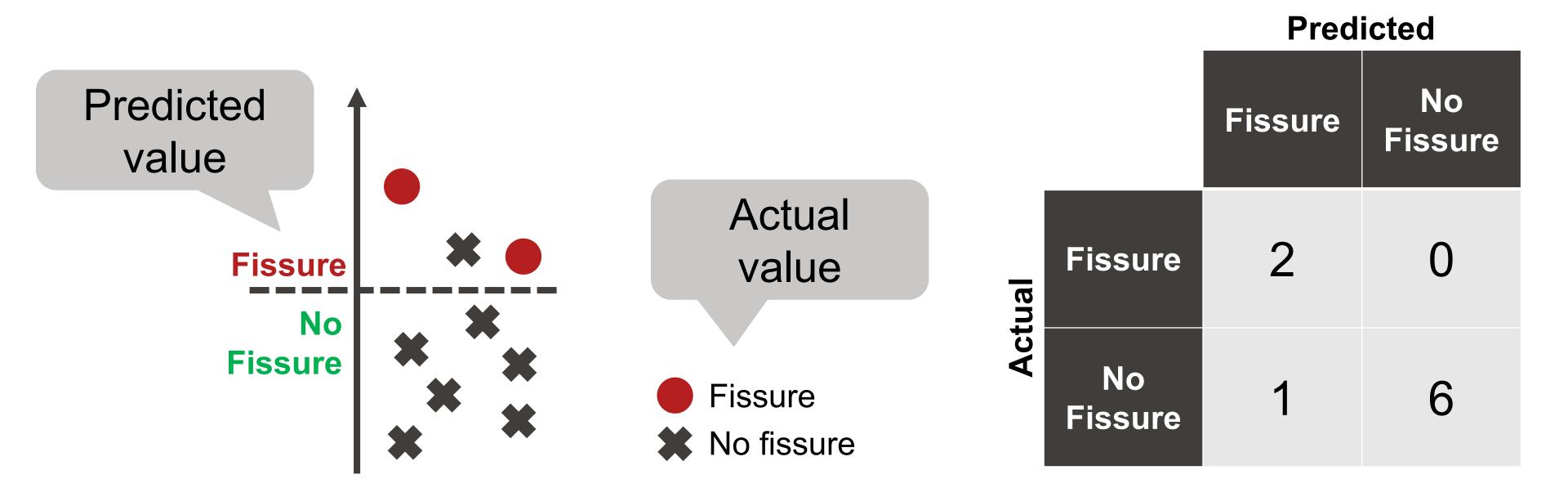
° c. Optimization bias

d. Deployment bias

Evaluation: a benchmark with **actors** does not correspond to the target application context (real shops)

Fairness metrics

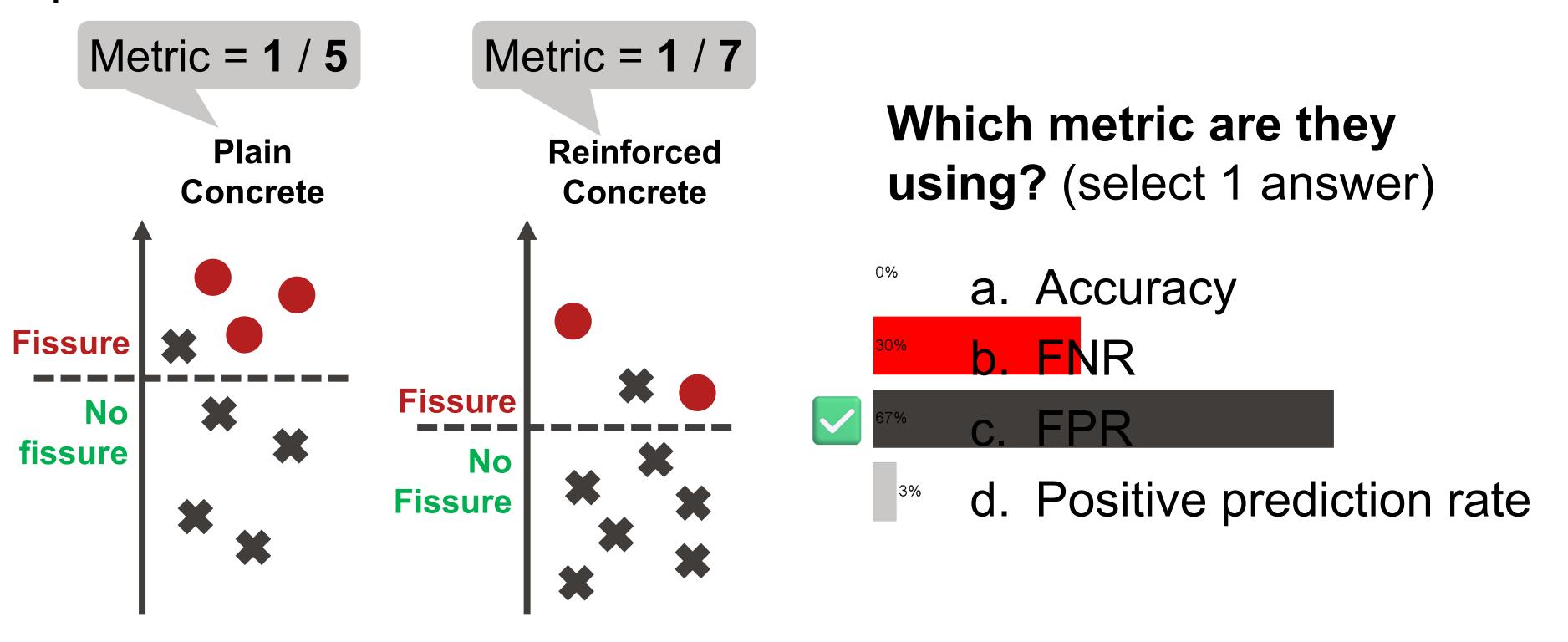
The company SuperCrack has developed a model to detect fissures in concrete before they become visible. They evaluate their model against a benchmark. The results look like this:



URL: ttpoll.eu

Session ID: cs290

They want to know whether their model performs equally well for plain concrete and for reinforced concrete. Here are the results:



Power Usage Effectiveness

URL: ttpoll.eu

Session ID: cs290

The GreenDC datacenter consumes an average of 1 MW. This means annually a total of 8 760 MWh of electricity. 50% of this electricity is used to power the IT equipment. What is the PUE of GreenDC?

°° a. 0.5

°° b. 1

°° c. 1.5

∞ d. 2

URL: ttpoll.eu

Session ID: cs290

The causes of water consumption from a datacenter which are included in the metric WUE_{source} are (select all that apply):

- ^{0%} a. Material mining
- b. Hardware manufacturing
- ° c. Concrete production
- d. Electricity production
- e. Cooling

Nudges

URL: ttpoll.eu

<u>Session ID:</u> cs290

Which of the following are examples of digital (software) nudges? (select all that apply)

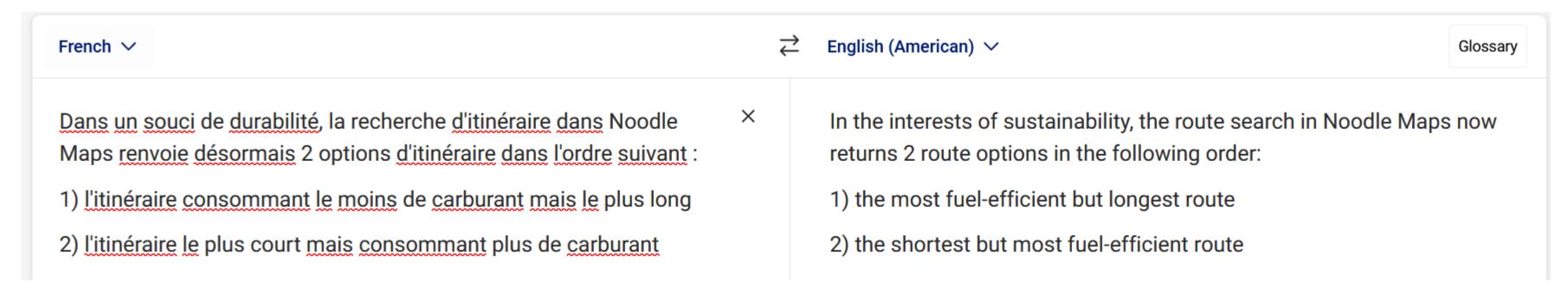
- ^o a. Automatic redirection to another website.
- b. Automatic newsletter subscription as stated in usage policy.
- c. Default value in online form
 - of d. Notice about strictly necessary cookies
- e. Notice about the behavior of other people

Translation

URL: ttpoll.eu

Session ID: cs290

Consider the following translation. What is the issue here?



- a. Parity error
- b. Factuality error
- ° c. Measurement error
- d. Faithfulness error

Conclusion

Responsible engineering of software

"The way a technology is designed determines its possibilities, which can, for better or for worse, have consequences for human wellbeing."

(Roeser, 2012)

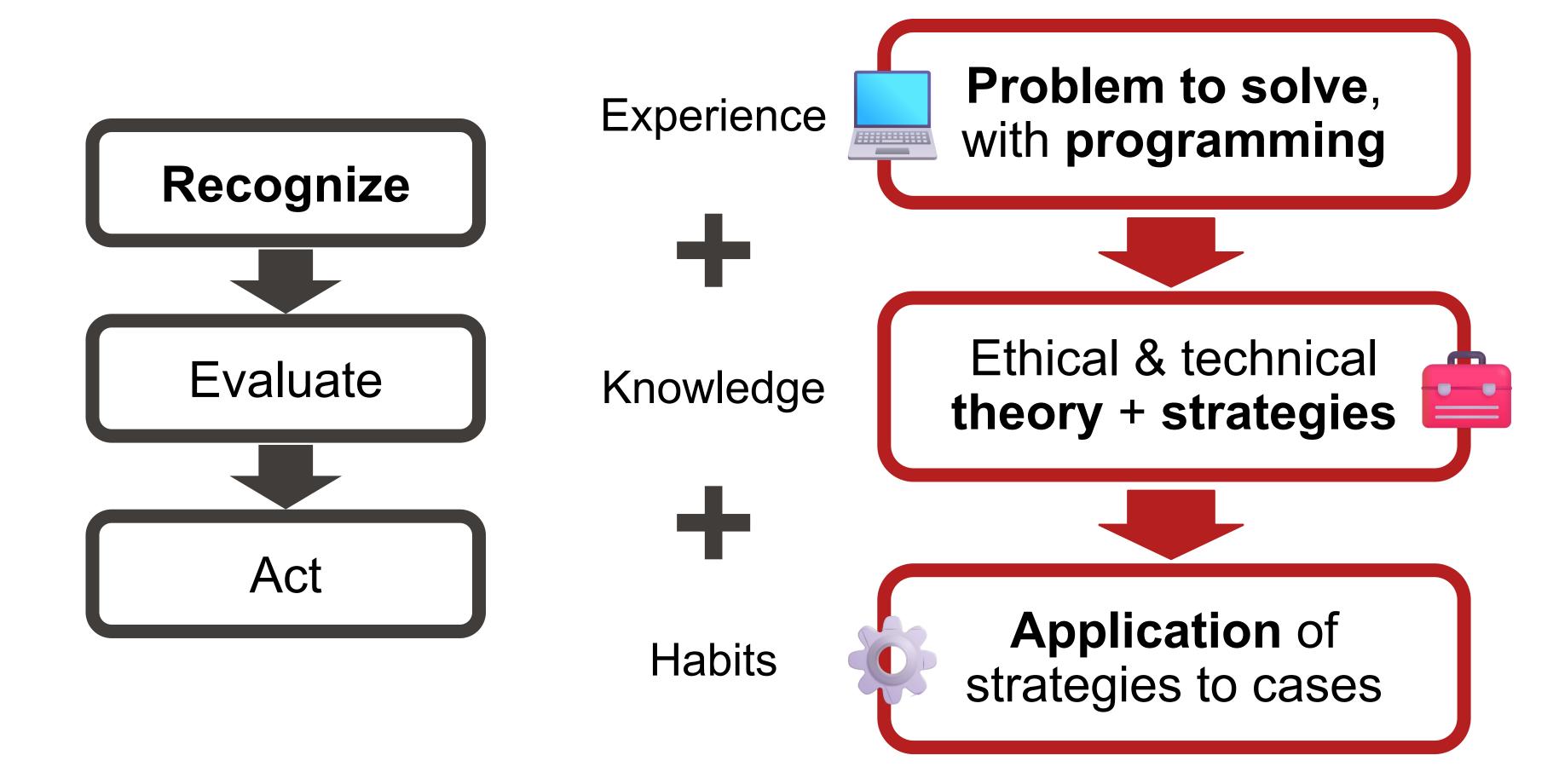
"Computing professionals' actions change the world. To act responsibly, they should reflect upon the wider impacts of their work, consistently supporting the public good."

(ACM, 2018)

Making engineering design decisions responsibly:

- 1. With a goal to do good
- 2. While preventing avoidable negative impacts
- 3. Taking **people**, other systems, **social structures** and our **planet** into account

(adapted and simplified from: Schwartz, 2016; Rest, 1986)



We can do better than that



Meta and OpenAI have spawned a wave of AI sex companions—and some of them are children

The uncensored AI economy is booming, giving rise to hard legal and ethical questions.

BY BEN WEISS AND ALEXANDRA STERNLICHT

January 8, 2024 at 3:00 PM GMT+1







Universal credit

Revealed: bias found in AI system used to detect UK benefits fraud

Exclusive: Age, disability, marital status and nationality influence decisions to investigate claims, prompting fears of 'hurt first, fix later' approach

Robert Booth UK technology editor

Fri 6 Dec 2024 06.00 CET

And now what?

- You have the power to make some change!
- Don't get fooled by the hype and the shiny useless/harmful software trends...
- A lot of questions seen in the course need more research!

Thank you for attending this course, good luck for the exam and all the best for all your projects!



References

- https://fortune.com/longform/meta-openai-uncensored-aicompanions-child-pornography/
- https://arstechnica.com/ai/2024/12/your-ai-clone-could-target-your-family-but-theres-a-simple-defense/
- https://www.theguardian.com/society/2024/dec/06/revealed-biasfound-in-ai-system-used-to-detect-uk-benefits