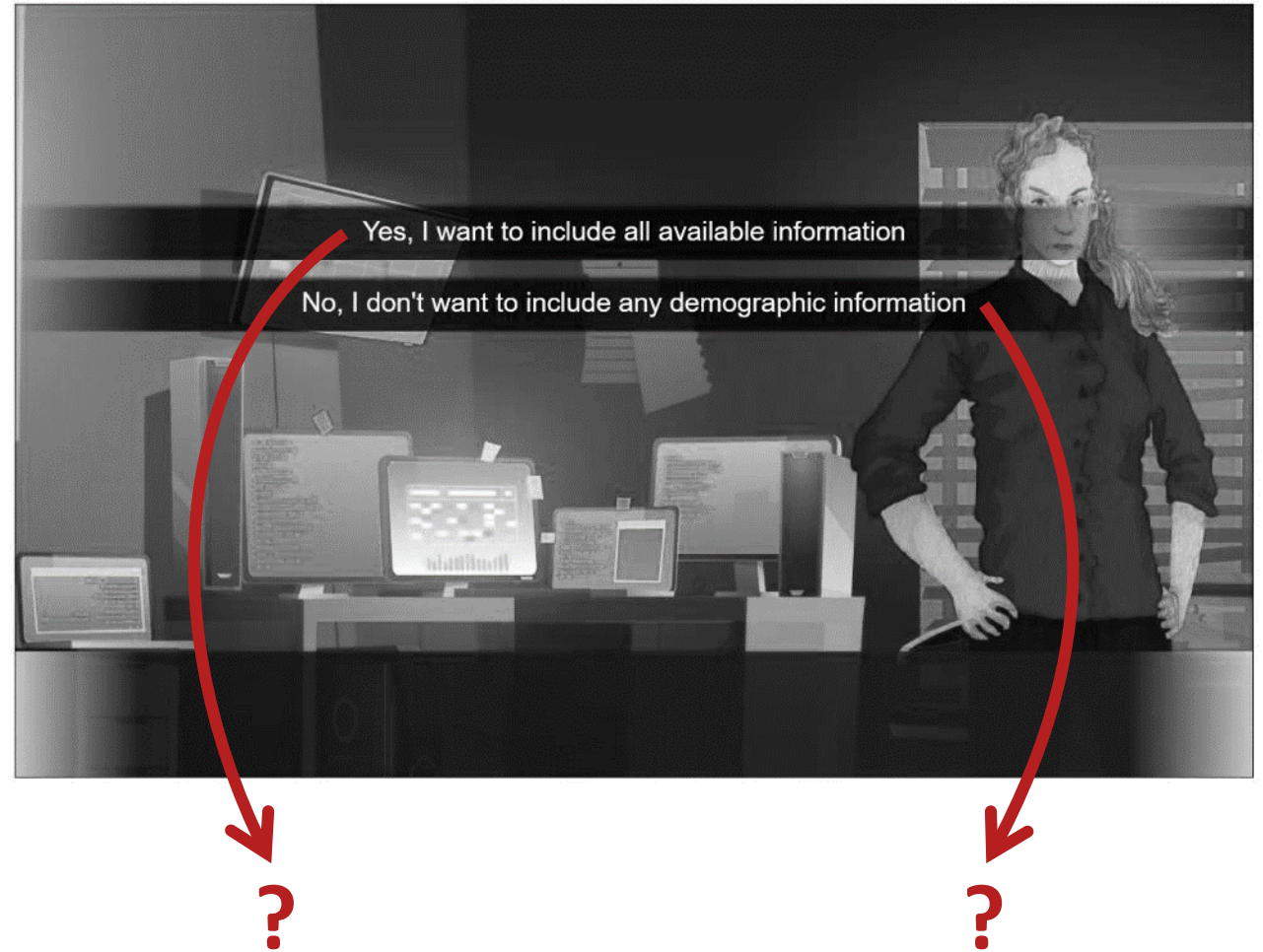# Decision making game: debriefing

## CS-433

### 21 November 2023

Cécile Hardebolle
EPFL Center for Digital Education
cecile.hardebolle@epfl.ch

# Let's refresh our memory

In the prisoner case:

- ▸ What decision did you take: include the demographic information or not ?
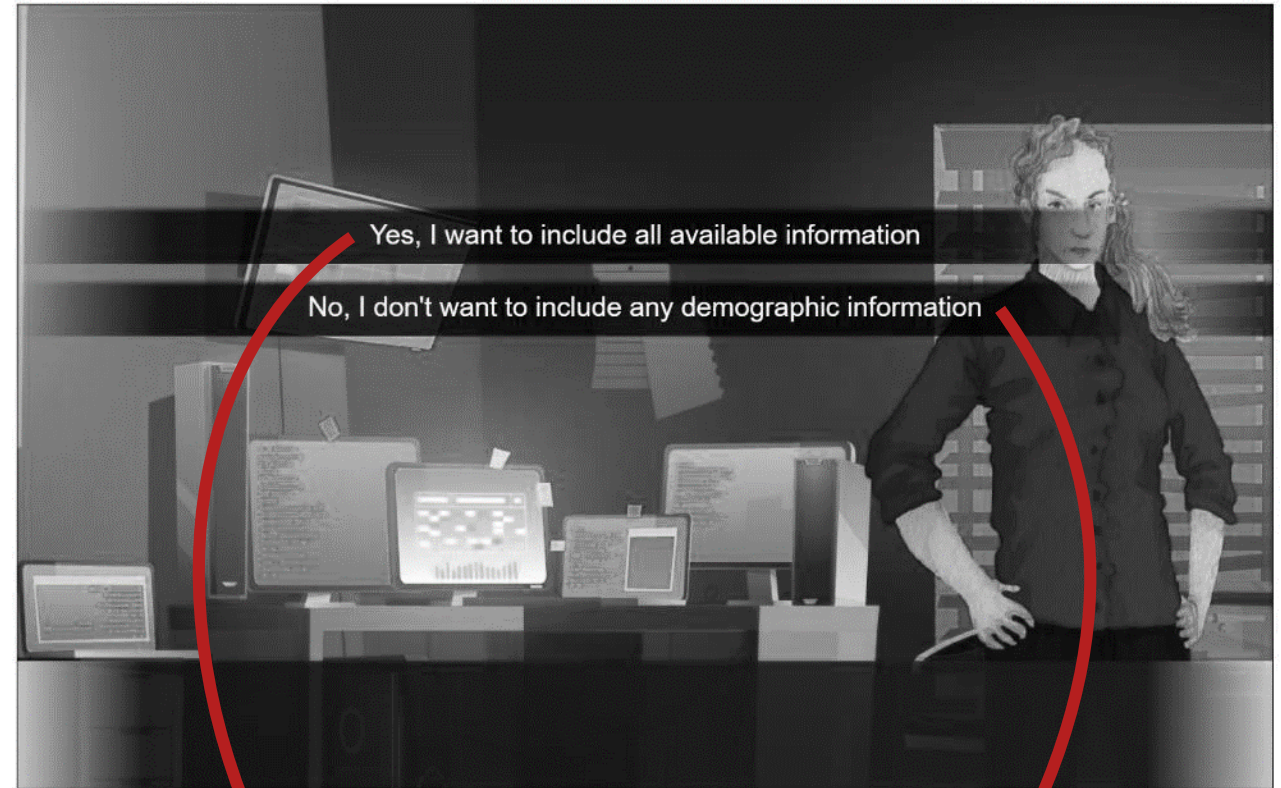- ▸ What was the consequence?

# Fairness vs. Accuracy

In the prisoner case:

▸ What decision did you take: include the demographic information or not ?

▸ What was the consequence?



Yes, I want to include all available information

No, I don't want to include any demographic information

**Racist model**
▶ **Friend not released**

**Low accuracy**
▶ **Neighbor robbed**

# Based on real cases



VERNON PRATER — LOW RISK 3
BRISHA BORDEN — HIGH RISK 8

VERNON PRATER
Prior Offenses
2 armed robberies, 1 attempted armed robbery
Subsequent Offenses
1 grand theft
LOW RISK 3

BRISHA BORDEN
Prior Offenses
4 juvenile misdemeanors
Subsequent Offenses
None
HIGH RISK 8

Angwin, J., Larson, J., Mattu, S., & Kirchner, L. (2016, May 23). Machine Bias : There's software used across the country to predict future criminals. And it's biased against blacks. ProPublica.
https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing
Dressel, J., & Farid, H. (2018). The accuracy, fairness, and limits of predicting recidivism. Science Advances, 4(1), eaao5580.
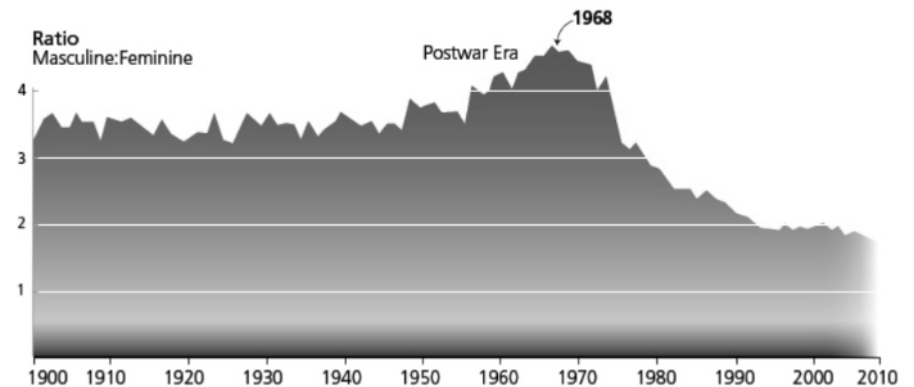https://doi.org/10.1126/sciadv.aao5580



| DETECT LANGUAGE | TURKISH | ENGLISH | SPANI... | | ENGLISH | SPANISH | ARABIC | |
|---|---|---|---|---|---|---|---|---|

O bir aşçı
o bir mühendis
o bir hemşire
o bir doktor

She is a cook
he is an engineer
she is a nurse
he is a doctor

52/5000

Send feedback

Machine Translation | Gendered Innovations
http://genderedinnovations.stanford.edu/case-studies/nlp.html



**Ratio of Masculine to Feminine Pronouns in U.S. Books, 1900-2008**
Changes parallel increases in women's labor force participation, education, age at first marriage, etc.

Ratio Masculine:Feminine

1968 — Postwar Era

The ratio of masculine pronouns ("he," "him," "his," "himself") to feminine pronouns ("she," "her," "hers," "herself") peaked at over 4:1 in 1968. By 2000 the ratio dropped dramatically to 2:1 (Twenge et al., 2012).

Data from American English corpus of the Google Books database (~1.2 million books).
Reproduced from Twenge et al., 2012.

4

# Bias and fairness in Machine Learning

# Bias and fairness in Machine Learning



(a) Data Generation

(b) Model Building and Implementation

Suresh, H., & Guttag, J. V. (2021). A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle. *Equity and Access in Algorithms, Mechanisms, and Optimization*, 1–9. https://doi.org/10.1145/3465416.3483305

# Most Machine Learning systems will affect people



"Crowd" by Amy West on Flickr, CC BY 2.0
https://www.flickr.com/photos/amy_elizabeth_west/3876549126/

# Risk: harm at scale

## Dutch scandal serves as a warning for Europe over risks of using algorithms

The Dutch tax authority ruined thousands of lives after using an algorithm to spot suspected benefits fraud – and critics say there is little stopping it from happening again.

As the world turns to AI to automate their systems, the Dutch scandal shows how devastating they can be

NEWS | ARTIFICIAL INTELLIGENCE

## The Dutch Tax Authority Was Felled by AI—What Comes Next? › European regulation hopes to rein in ill-behaving algorithms

BY RAHUL RAO | 09 MAY 2022 | 4 MIN READ

AMNESTY INTERNATIONAL

WHO WE ARE   WHAT WE DO   COUNTRIES   GET INVOLVED   DONATE NOW   LATEST   SEARCH

SHARE

NEWS                                                        October 25, 2021

## Dutch childcare benefit scandal an urgent wake-up call to ban racist algorithms

Recently added

Pakistan: Ban on film Joyland showcasing transgender character

8

# Your decisions can make a difference!



"It is important to acknowledge that **not all problems should be blamed on the data**. The ML pipeline involves a series of choices and practices, from model definition to user interfaces used upon deployment. **Each stage involves decisions that can lead to undesirable effects.**"

(Suresh & Guttag, 2021)

# Responsible design decisions

→ Evaluate **who could be affected** and **how**

> **2 strategies to practice in your project!**
>
> ↳ **Required "Ethical risks" section in your report**

# Responsible design decisions

→ Evaluate **who could be affected** and **how**

<span style="color:darkred">**Stakeholder Analysis**</span>

**Indirect**
Impacted by the system, though they never/rarely interact directly with it

**Direct**

Interact directly with the system or with the system's output

# Responsible design decisions

→ Evaluate **who could be affected** and **how**

**Stakeholder Analysis**

Prisoners
Neighbor

**Indirect**

Environment

Government

Public(s)

Neighborhood

NGOs

**Direct**

Friends

End-users

Admins

Contractors

Developers

# Responsible design decisions

→ Evaluate **who could be affected** and **how**

**Risk Analysis**

1. Speculative scenarios
2. Research + measure

**Negative impact**

+ Likelihood of occurrence

+ Severity

# Case: ML model that predicts user emotions based on smartphone touch data

Imagine it is **deployed in a social media platform**:

- ▸ Included in the user interface to display predicted emotions
- ▸ Accessible to third parties (e.g. for ads)
- ▸ Used for internal functions (e.g. content recommendation or moderation)

**What could go wrong? Generate as many ideas as possible:**

- ☐ Can the solution be used in **harmful ways**?
- ☐ What kind of impacts can **errors** from the solution have?
- ☐ Could the solution disclose / be used to **disclose private information**?
- ☐ Could the solution contribute to **discrimination** against people or groups?
- ☐ What **choices** are users able to make in their use of the solution and how?

| Context | Beneficence | Non-maleficence | | Solution |
|---|---|---|---|---|
| □ In which context is the solution evaluated? | □ What are the expected benefits of the solution in this context? | **Risks** | **Mitigation** | □ What are the characteristics of the solution under evaluation? |
| | | □ Can the solution be used in harmful ways, in particular with regards to vulnerable populations?<br><br>□ What kind of impacts can errors from the solution have? | | |
| | | □ What type of protections does the solution have against attacks or misuse? | | |

| Privacy | | Fairness | |
|---|---|---|---|
| **Risks** | **Mitigation** | **Risks** | **Mitigation** |
| □ What data does the solution collect?<br>□ Is it collecting personal or sensitive data?<br>□ Who has access to the data?<br>□ How is the data protected? | | □ How accessible is the solution?<br>□ What kinds of biases may affect the results?<br>□ Can the outcomes of the solution be different for different users or groups? | |
| □ Could the solution disclose / be used to disclose private information? | | □ Could the solution contribute to discrimination against people or groups? | |

| Sustainability | | Empowerment | |
|---|---|---|---|
| **Risks** | **Mitigation** | **Risks** | **Mitigation** |
| □ What is the carbon footprint of the solution?<br>□ What types of resources does it consume (e.g. water) and produce (e.g. waste)?<br>□ What type of human labor is involved? | | □ Can users understand how the solution works and what its limits are?<br>□ Are users able to make choices (e.g. consent, settings) in their use of the solution and how?<br>□ How does the solution affect user autonomy and agency? | |

**Have you identified a range of different types of risks?**

# Context

□ In which context is the solution evaluated?

# Beneficence

□ What are the expected benefits of the solution in this context?

# Non-maleficence

| | Mitigation |
|---|---|
| … be used in harmful … with regards to …ations? | |
| …pacts can errors from the | |
| □ …tections does the so… …nst attacks or misuse? | |

**Sticky notes:** Generate mis-communication and conflict · Use for scams · Use for bullying

# Solution

□ What are the characteristics of the solution under evaluation?

# Privacy

| | Mitigation |
|---|---|
| □ … solution collect? …onal or sensitive | |
| □ …ss to the data? | |
| □ H… …otected? | |
| □ C… …disclose / be used to dis… …rmation? | |

**Sticky notes:** People lose control over private emotions · Mental health data · Possible to identify political opinions

# Fairness

| Risks | Mitigation |
|---|---|
| … is the solution? …ases may affect the | |
| □ Can the outcomes of the solution be … t users or groups? | |
| … contribute to …nst people or groups? | |

**Sticky notes:** Bias wrt cultural differences · Bias wrt people with handicap (physical / mental)

# Sustainability

| Risks | Mitigation |
|---|---|
| □ What is the carbon footprint of the | |
| …esources does it consume (e.g. water) and produce (e.g. | |
| …uman labor is involved? | |

**Sticky notes:** CO2 & water for training + inference · Emotion labeling for training

# Empowerment

| | Mitigation |
|---|---|
| …tand how the solution …limits are? | |
| □ A… …e choices (e.g. con… …r use of the solu… | |
| □ … …lution affect user a… …ncy? | |

**Sticky notes:** People lose control over private emotions · People stop relying on their own interpretation · Attention capture, manipulation

**User guide and examples on https://go.epfl.ch/canvas**

# You can make a difference!

ML systems can have **harmful consequences for** <span style="color:red">**people**</span>

↳ Evaluate **who** could be affected and **how:**

- ▸ Stakeholder analysis
- ▸ Ethical risk analysis

**https://go.epfl.ch/canvas**

↳ Document the **ethical risk analysis** you have performed

**In your project report:**

- ▸ Describe 1 type of risk (or justify the absence of risks)
- ▸ Explain how you evaluated it
- ▸ Describe how you took it into account (or the barriers to do it)

**Ask your questions on Ed (tag me)**