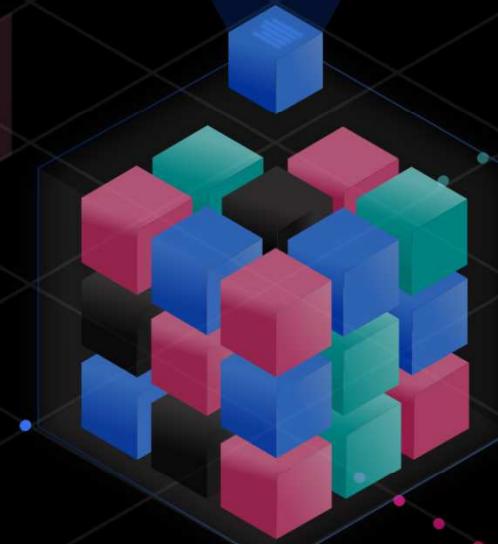


(Un-)trusted AI: Strategien um Diskriminierung durch KI entgegenzuwirken

Thomas Jirku, IBM



**where to find discrimination
(bias)**



people



data



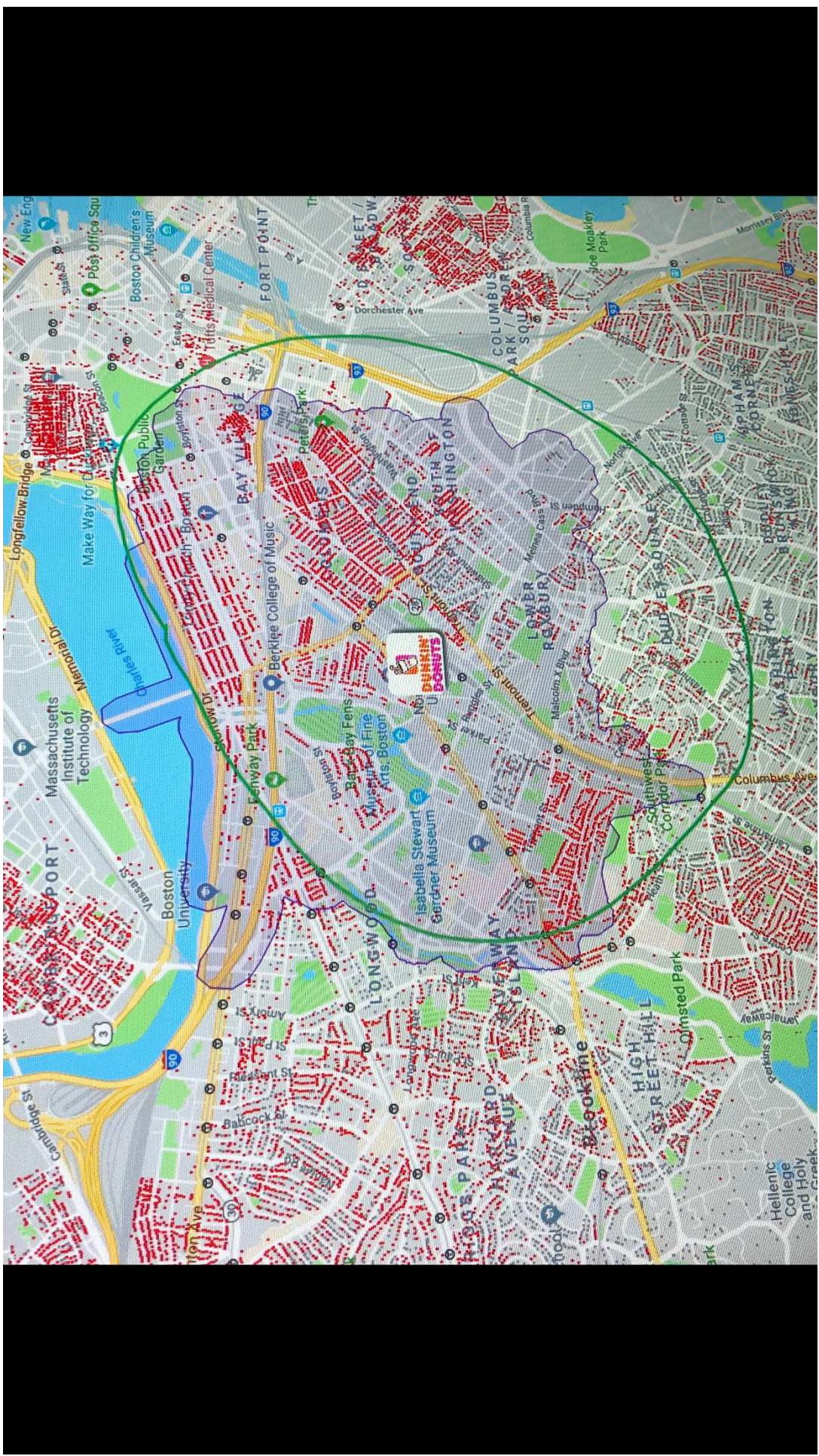
algorithms

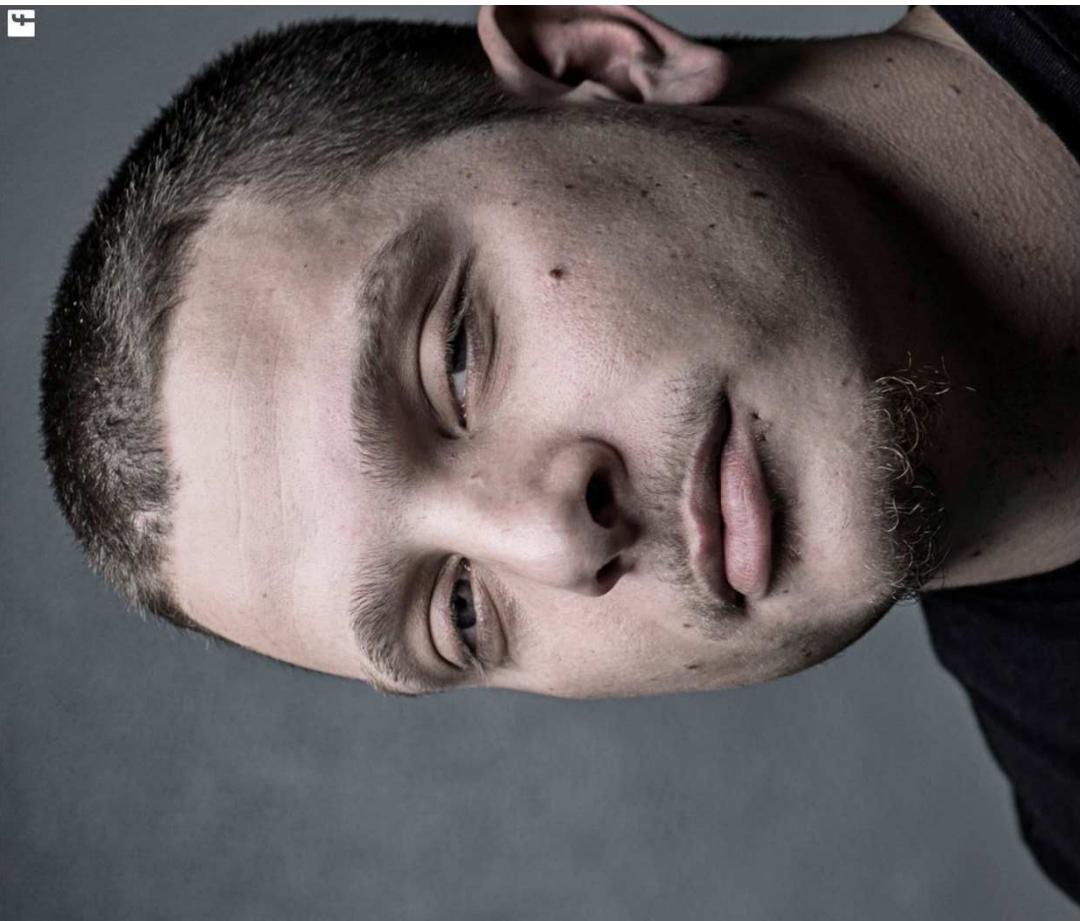


models



actions





BLICA

SITE

News & Politics

Culture

Technology

Business

Human Interest

MONEYBOX

SEARCH

SIGN IN

FOLLOW US

G

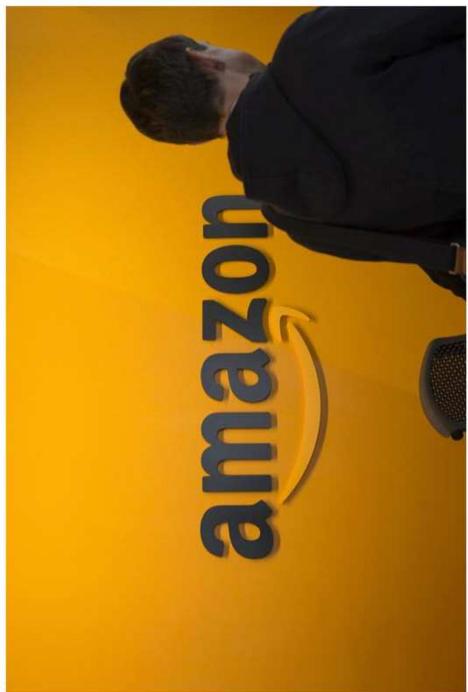
T

F

Amazon Created a Hiring Tool Using A.I. It Immediately Started Discriminating Against Women.

By JORDAN WEISSMANN

OCT 10, 2018 • 4:52 PM



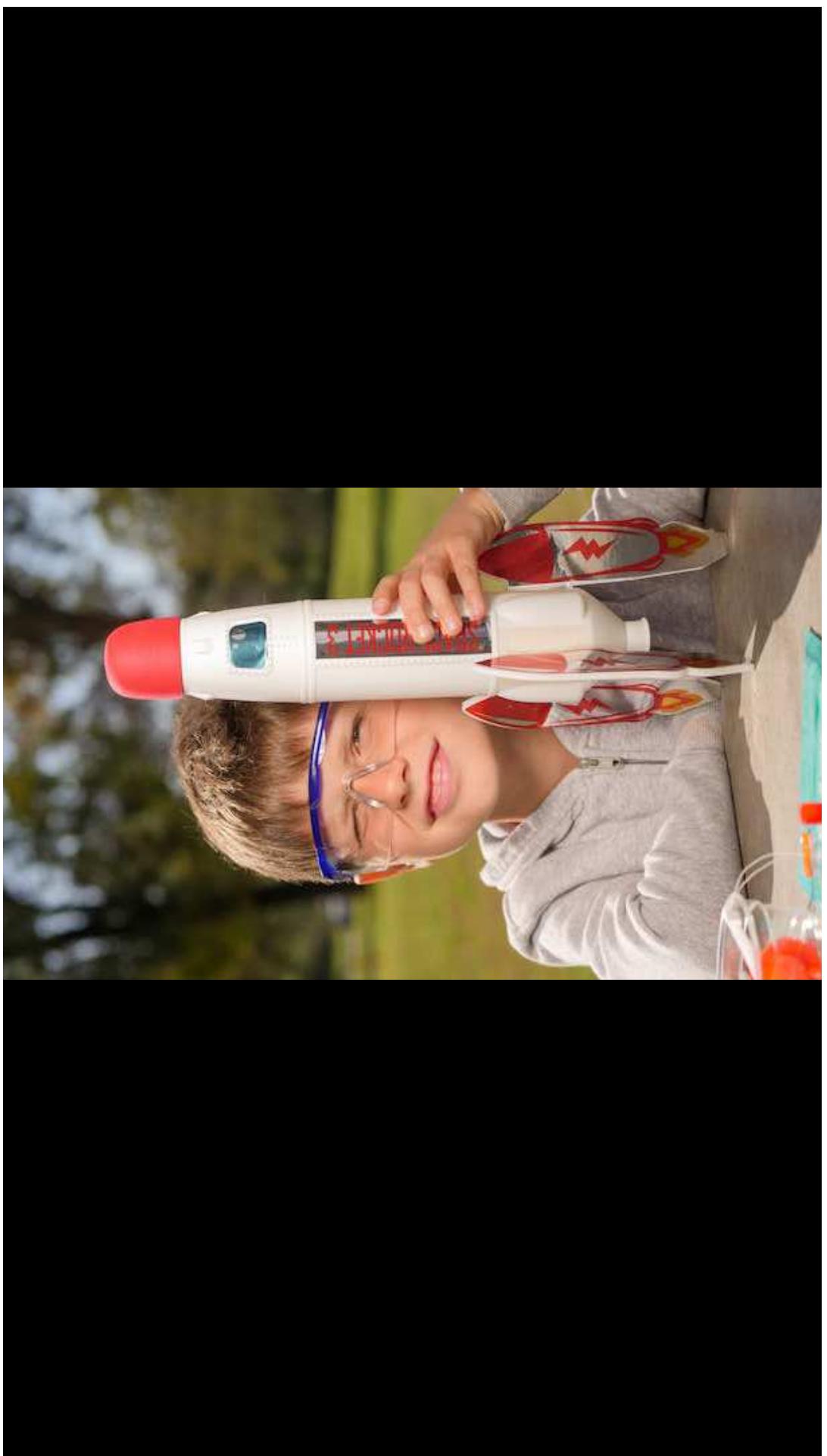
TWEET

SHARE

COMMENT

$$\begin{aligned} \text{BE_INT} \\ = f(0,10) \\ - 0,14 \times \text{GESCHLECHT_WEIBLICH} \\ - 0,13 \times \text{ALTERSGRUPPE_30_49} \\ - 0,70 \times \text{ALTERSGRUPPE_50_PLUS} \\ + 0,16 \times \text{STAATENGRUPPE_EU} \\ - 0,05 \times \text{STAATENGRUPPE_DRITT} \\ + 0,28 \times \text{AUSBILDUNG_LEHRE} \\ + 0,01 \times \text{AUSBILDUNG_MATURA_PLUS} \\ - 0,15 \times \text{BETREUUNGSPFLICHTIG} \\ - 0,34 \times \text{RGS_TYP_2} \\ - 0,18 \times \text{RGS_TYP_3} \\ - 0,83 \times \text{RGS_TYP_4} \\ - 0,82 \times \text{RGS_TYP_5} \\ - 0,67 \times \text{BEEINTRÄCHTIGT} \\ + 0,17 \times \text{BERUFSGRUPPE_PRODUKTION} \\ - 0,74 \times \text{BESCHÄFTIGUNGSTAGE_WENIG} \\ + 0,65 \times \text{FREQUENZ_GESCHÄFTSFALL_1} \\ + 1,19 \times \text{FREQUENZ_GESCHÄFTSFALL_2} \\ + 1,98 \times \text{FREQUENZ_GESCHÄFTSFALL_3_PLUS} \\ - 0,80 \times \text{GESCHÄFTSFALL_LANG} \\ - 0,57 \times \text{MN_TEILNAHME_1} \\ - 0,21 \times \text{MN_TEILNAHME_2} \\ - 0,43 \times \text{MN_TEILNAHME_3}) \end{aligned}$$





THE DATA ANALYSIS PROCESS

Step 1:

Define the question

Step 2:

Collect the data

Step 3:

Clean the data

Step 4:

Analyze the data

Step 5:

Visualize and share
your findings



It's OK!

*I removed the
[GENDER] column!*

```
ALTER TABLE "client_data_for_model" DROP "gender";
```

~~Problems~~

Solutions®

explainability

Explain a transaction



0b9b0e1d-7022-4efc-... x 688bffd7-1fe1-4d9c-a... x 01a5ea94-e77a-4bda-... x 01a5ea94-e77a-4bda-... x 01a5ea94-e77a-4bda-... x

Details ⓘ

Transaction	0b9b0e1d-7022-4efc-a74c-338bbcd41647-3
Deployment	GermanCreditRiskModel
Model Name	GermanCreditRiskModel

Maximum changes allowed for the same outcome ⓘ

Age	19
LoanDuration	11
CheckingStatus	no_checking

Risk

CONFIDENCE

No Risk

77.64%

22.36%

Factors contributing to Risk confidence level

LoanAmount: 6020

OwnsProperty: car_other

Age: 19

LoanDuration: 11

CheckingStatus: 0_to_200

Sex: female

OthersOnLoan: none

EmploymentDuration: 1_to_4

InstallmentPlans: none

Telephone: none

Factors contributing to No Risk confidence level

Age: 19

LoanDuration: 11

CheckingStatus: 0_to_200

EmploymentDuration: 1_to_4

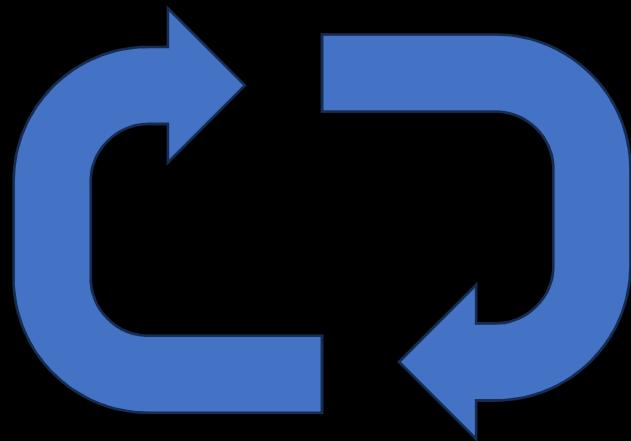
Telephone: none

fairness



robustness

continuous



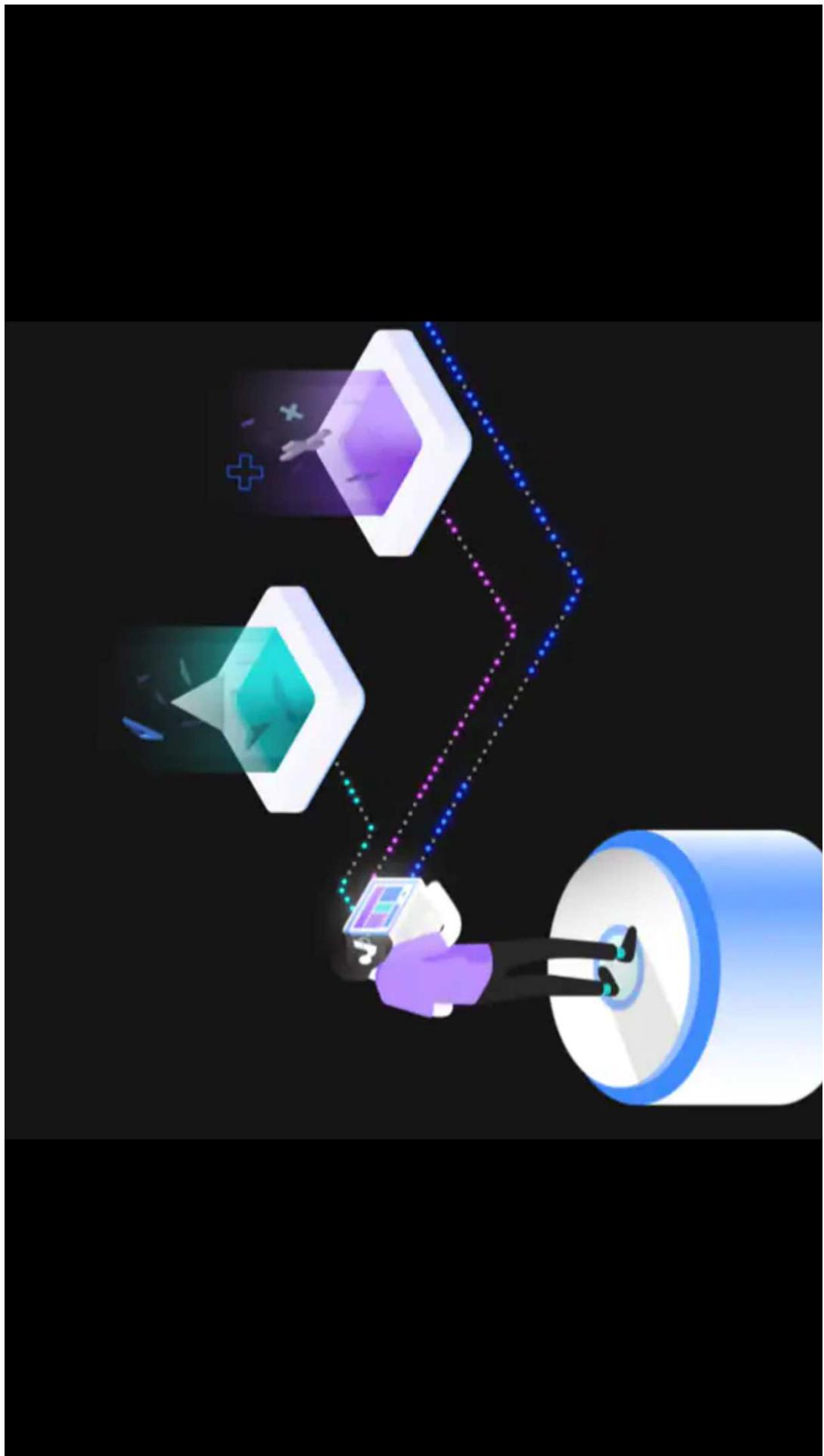
evaluation

transparency

AI FACTSHEET

Model Name	Mortgage Evaluator Privacy
Target Audience	Model Risk Officer
Overview	This evaluation tests the vulnerability of the model to select inference attacks and measures certain risk factors that are known to be strongly correlated with privacy risk.
Evaluation Details	<ul style="list-style-type: none">Membership inference is a type of attack where, given a trained model and a data sample, one can deduce whether or not that sample was part of the model's training.Attribute inference is an attack where certain sensitive features may be inferred about individuals whose data was included in training a model.
Results	Evaluation
Training set size	650,877
Overfitting	0.0011
Black-box membership inference (BBMI)	0.55
Black-box attribute inference (BBAI) - Age	0.046
Black-box attribute inference - Debt-to-income ratio	0.069
Feature influence - Age	0.003
Feature influence - Debt-to-income ratio	0.678

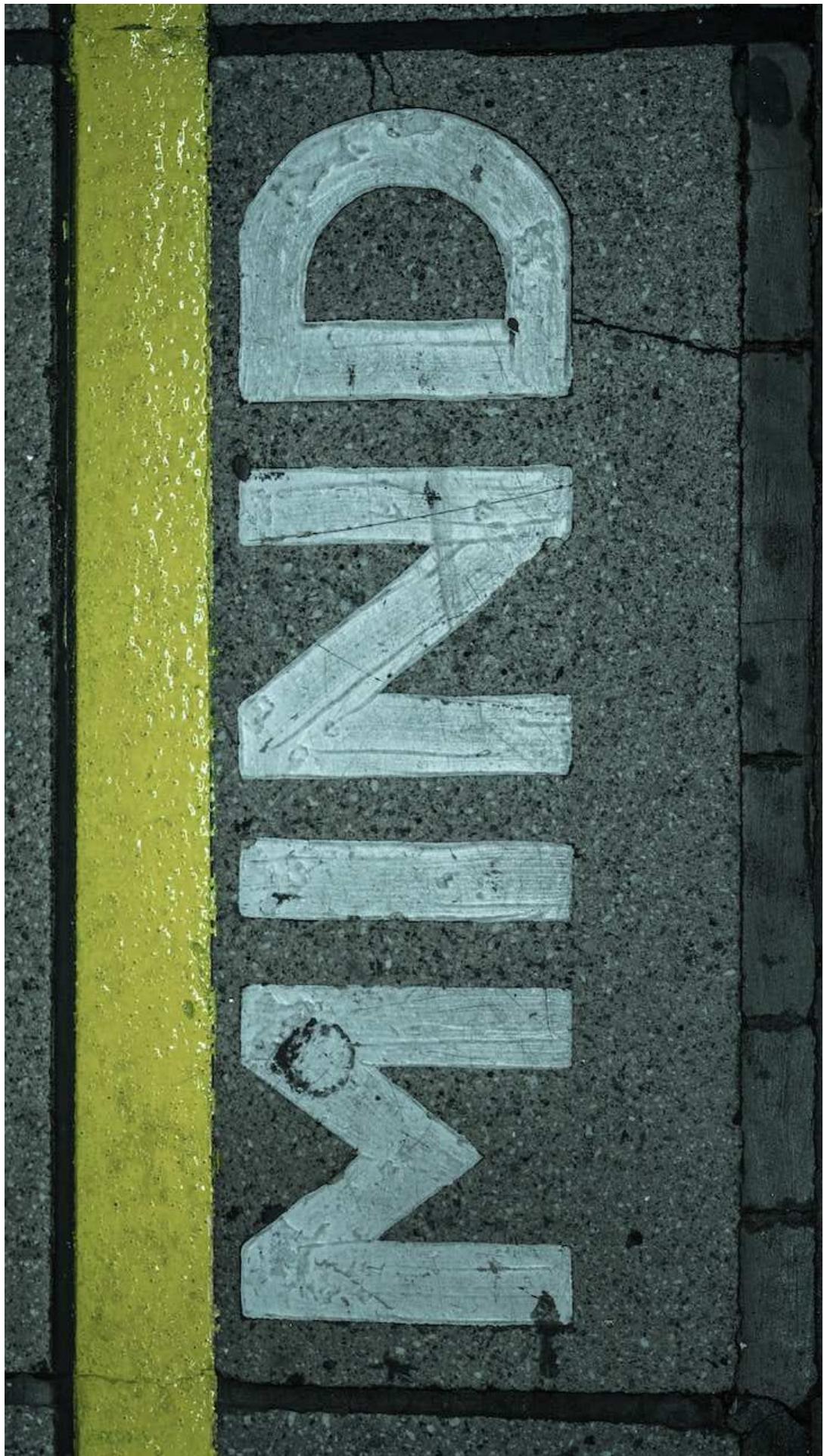
data privacy



3

2

1



literature:

- **angriff der algorithmen**, cathy o`neil
- **weapons of math destruction**, cathy o`neil
- anleitung zum unglücklich sein, paul watzlawick

Thomas Jirku, IBM

Thank
you!



LinkedIn Kontakt