



UNIVERSITY OF AMSTERDAM

Korteweg de Vries Institute for Mathematics

Theoretical limitations of Explainability methods

Xomnia Data & Drinks

2023-09-14

Programme for today

- ▶ Introduce myself and my research
- ▶ Brief XAI introduction
- ▶ Attribution methods that provide recourse and are robust cannot exist
- ▶ The risk of counterfactual explanations
- ▶ Conclusion

Short introduction

Short Introduction

- ▶ PhD student at the UvA: *Formalising Explainable AI*
- ▶ Previously, MSc in Mathematics & worked at Amsterdam Data Collective
- ▶ All work presented was created in collaboration with:



Dr. Tim van Erven



Dr. Rianne de Heide



Dr. Damien Garreau

Explainable Artificial Intelligence

Setting

Leading example

2 parties:

- Credit Loan Applicant (A)



- Credit Loan Provider (P)



Loan application process:

- (A) provides (P) with a set of features:

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

- (P) has an automated decision system f

$$f(x) = 1 \quad \text{if accepted}$$

$$f(x) = -1 \quad \text{if not}$$

- (A) get decision and an explanation from (P)

Call for XAI

Some Reasons

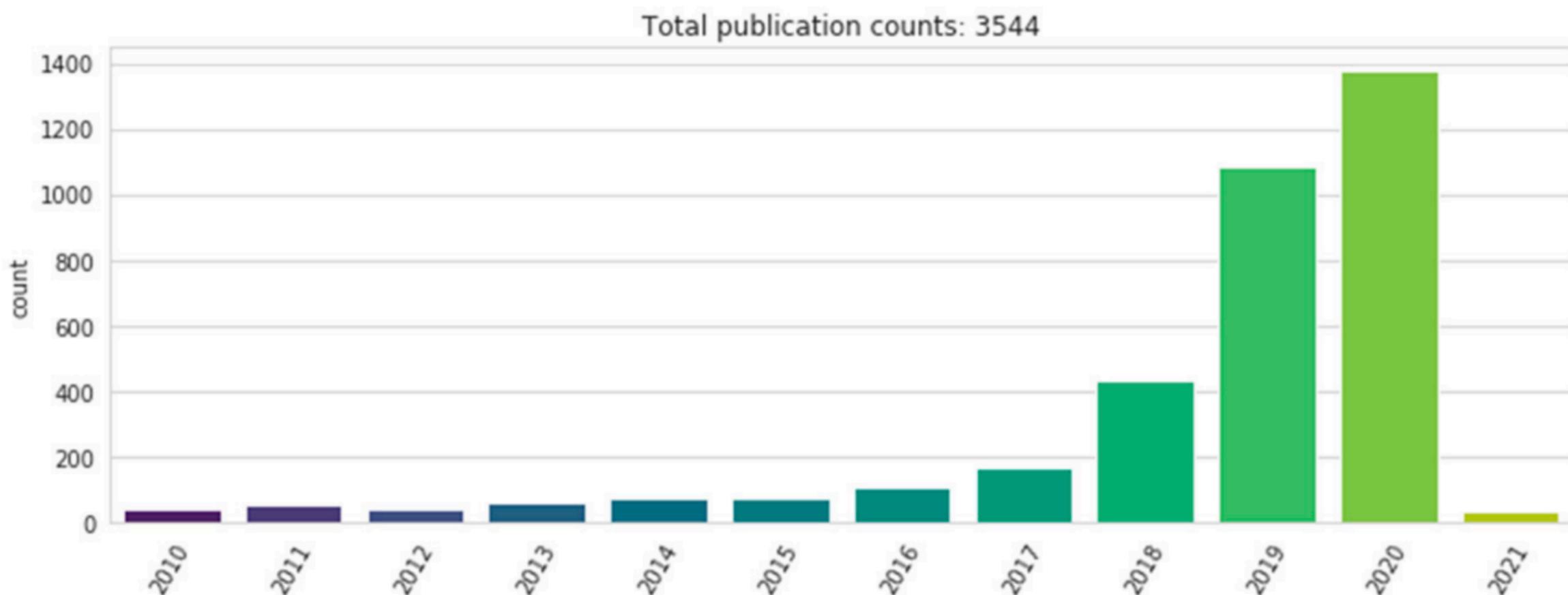
- ▶ Fairness: Biases can be detected earlier
- ▶ Trustworthiness
- ▶ Increases reliability
- ▶ Regulation



Explanations Explosion

Methods
CAM with global average pooling [42], [91]
+ Grad-CAM [43] generalizes CAM, utilizing gradient
+ Guided Grad-CAM and Feature Occlusion [68]
+ Respond CAM [44]
+ Multi-layer CAM [92]
LRP (Layer-wise Relevance Propagation) [13], [53]
+ Image classifications. PASCAL VOC 2009 etc [45]
+ Audio classification. AudioMNIST [47]
+ LRP on DeepLight. fMRI data from Human Connectome Project [48]
+ LRP on CNN and on BoW(bag of words)/SVM [49]
+ LRP on compressed domain action recognition algorithm [50]
+ LRP on video deep learning, <i>selective relevance method</i> [52]
+ BiLRP [51]
DeepLIFT [57]
Prediction Difference Analysis [58]
Slot Activation Vectors [41]
PRM (Peak Response Mapping) [59]
LIME (Local Interpretable Model-agnostic Explanations) [14]
+ MUSE with LIME [85]
+ Guidelinebased Additive eXplanation optimizes complexity, similar to LIME [93]
Also listed elsewhere: [56], [69], [71], [94]
Others. Also listed elsewhere: [95]
+ Direct output labels. Training NN via multiple instance learning [65]
+ Image corruption and testing Region of Interest statistically [66]
+ Attention map with autofocus convolutional layer [67]
DeconvNet [72]
Inverting representation with natural image prior [73]
Inversion using CNN [74]
Guided backpropagation [75], [91]
Activation maximization/optimization [38]
+ Activation maximization on DBN (Deep Belief Network) [76]
+ Activation maximization, multifaceted feature visualization [77]
Visualization via regularized optimization [78]
Semantic dictionary [39]
Network dissection [36], [37]
Decision trees
Propositional logic, rule-based [82]
Sparse decision list [83]
Decision sets, rule sets [84], [85]
Encoder-generator framework [86]
Filter Attribute Probability Density Function [87]
MUSE (Model Understanding through Subspace Explanations) [85]

Methods	H
Linear probe [101]	
Regression based on CNN [106]	
Backwards model for interpretability of linear models [107]	
GDM (Generative Discriminative Models): ridge regression + least square [100]	
GAM, GA ² M (Generative Additive Model) [82], [102], [103]	
ProtoAttend [105]	
Other content-subject-specific models:	N
+ Kinetic model for CBF (cerebral blood flow) [131]	N
+ CNN for PK (Pharmacokinetic) modelling [132]	N
+ CNN for brain midline shift detection [133]	N
+ Group-driven RL (reinforcement learning) on personalized healthcare [134]	N
+ Also see [108]–[112]	N
PCA (Principal Components Analysis), SVD (Singular Value Decomposition)	N
CCA (Canonical Correlation Analysis) [113]	
SVCCA (Singular Vector Canonical Correlation Analysis) [97] = CCA+SVD	
F-SVD (Frame Singular Value Decomposition) [114] on electromyography data	
DWT (Discrete Wavelet Transform) + Neural Network [135]	
MODWPT (Maximal Overlap Discrete Wavelet Package Transform) [136]	
GAN-based Multi-stage PCA [118]	
Estimating probability density with deep feature embedding [119]	
t-SNE (t-Distributed Stochastic Neighbour Embedding) [77]	
+ t-SNE on CNN [120]	
+ t-SNE, activation atlas on GoogleNet [121]	
+ t-SNE on latent space in meta-material design [122]	
+ t-SNE on genetic data [137]	
+ mm-t-SNE on phenotype grouping [138]	
Laplacian Eigenmaps visualization for Deep Generative Model [124]	
KNN (k-nearest neighbour) on multi-center low-rank rep. learning (MCLRR) [125]	
KNN with triplet loss and <i>query-result activation map pair</i> [139]	
Group-based Interpretable NN with RW-based Graph Convolutional Layer [123]	
TCAV (Testing with Concept Activation Vectors) [96]	
+ RCV (Regression Concept Vectors) uses TCAV with Br score [140]	
+ Concept Vectors with UBS [141]	
+ ACE (Automatic Concept-based Explanations) [56] uses TCAV	
Influence function [129] helps understand adversarial training points	
Representer theorem [130]	
SocRat (Structured-output Causal Rationalizer) [127]	
Meta-predictors [126]	
Explanation vector [128]	
# Also listed elsewhere: [14], [43], [85], [94]	N
# Also listed elsewhere: [14], [60], [85] etc	N
CNN with separable model [142]	
Information theoretic: Information Bottleneck [98], [99]	
Database of methods v.s. interpretability [10]	N
Case-Based Reasoning [143]	

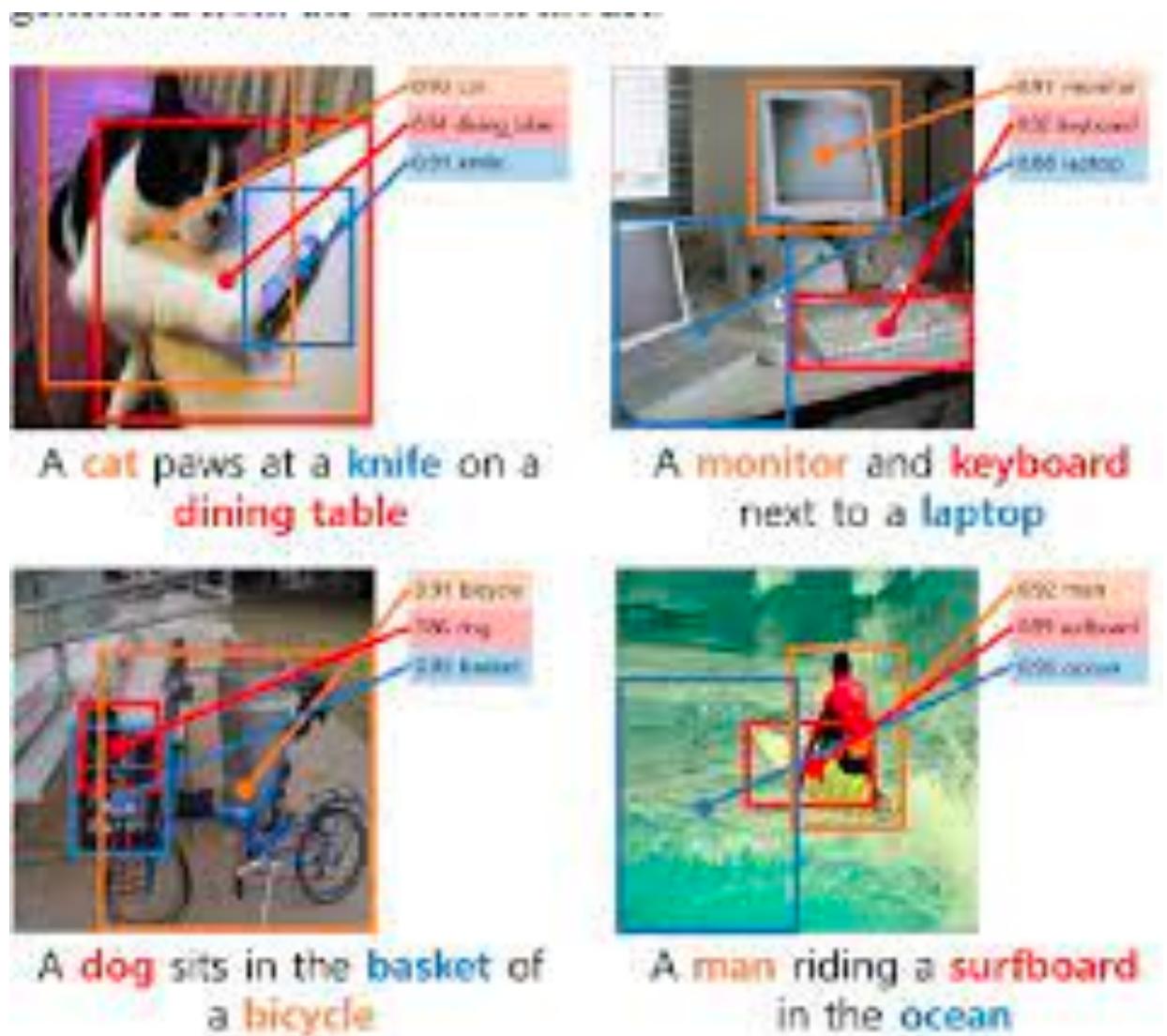


Explanations

Examples

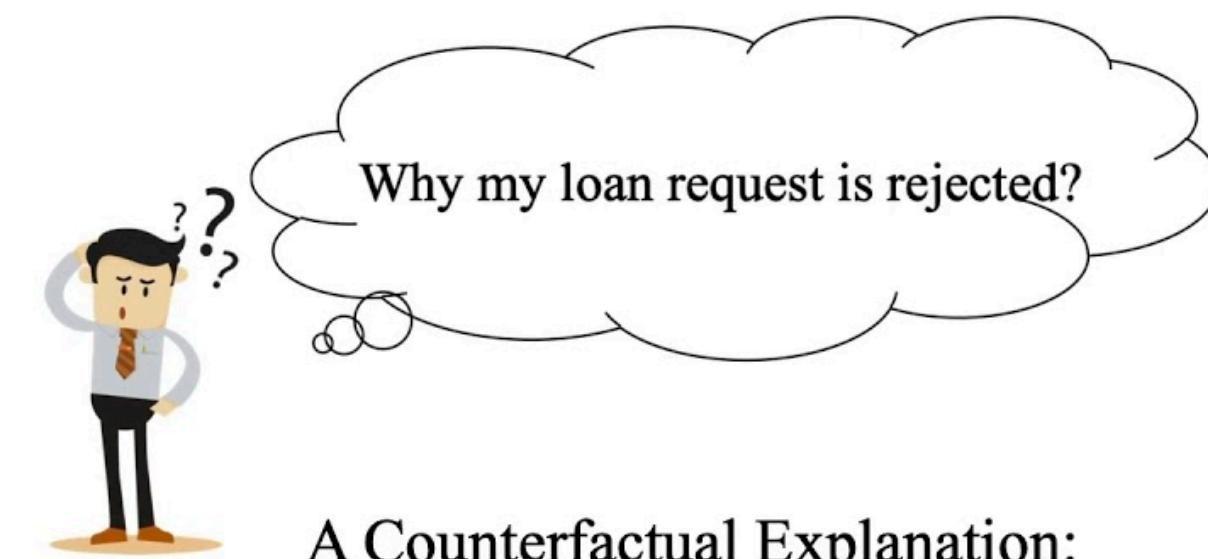
- ▶ Highlight important features
- ▶ Counterfactual explanation

- ▶ Caption generation
- ▶ Example based
- ▶ Activation Probing



Text with highlighted words

Why does the older generation think that just because they don't understand video games and technology, they feel like they have to hate them and blame every bad thing on them?



A Counterfactual Explanation:

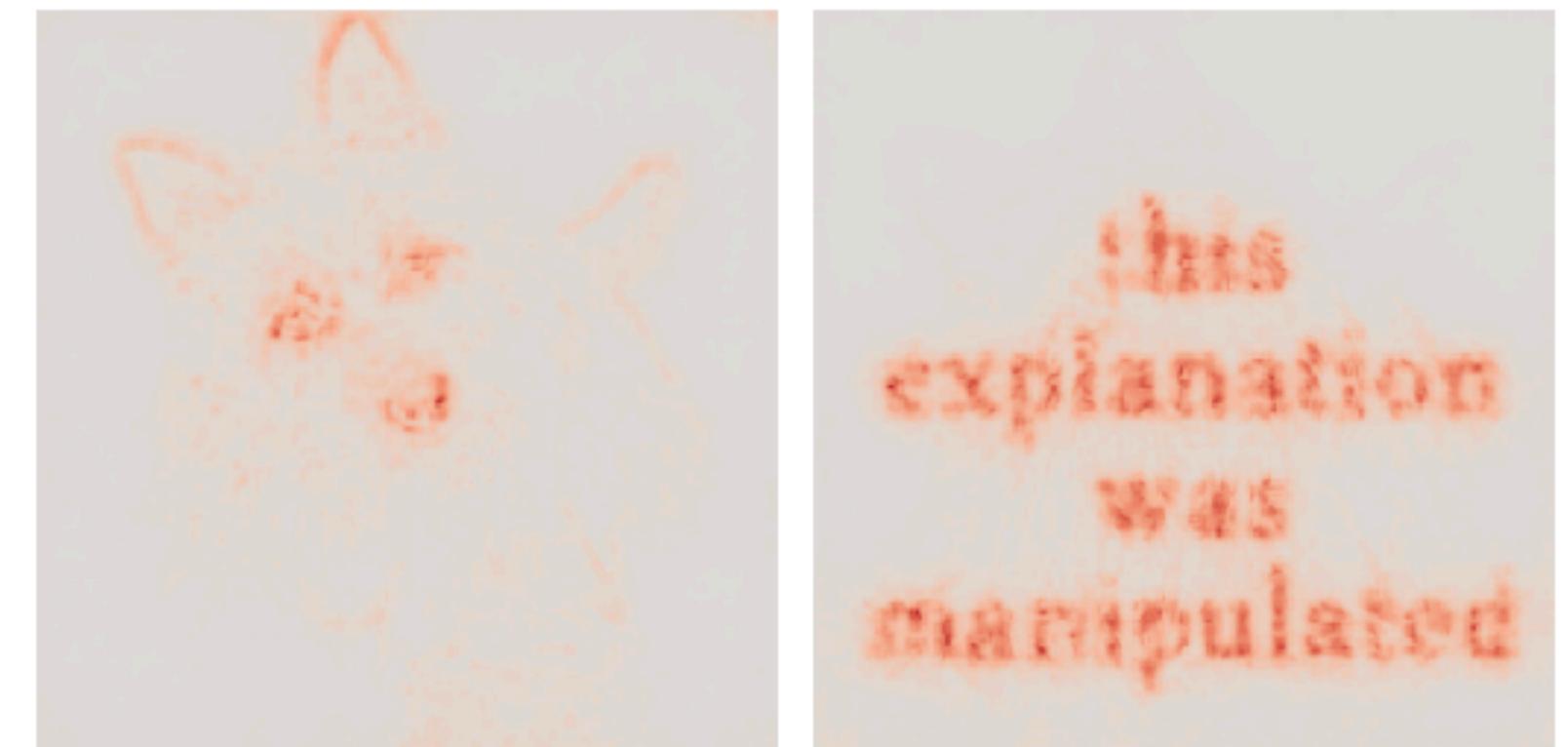
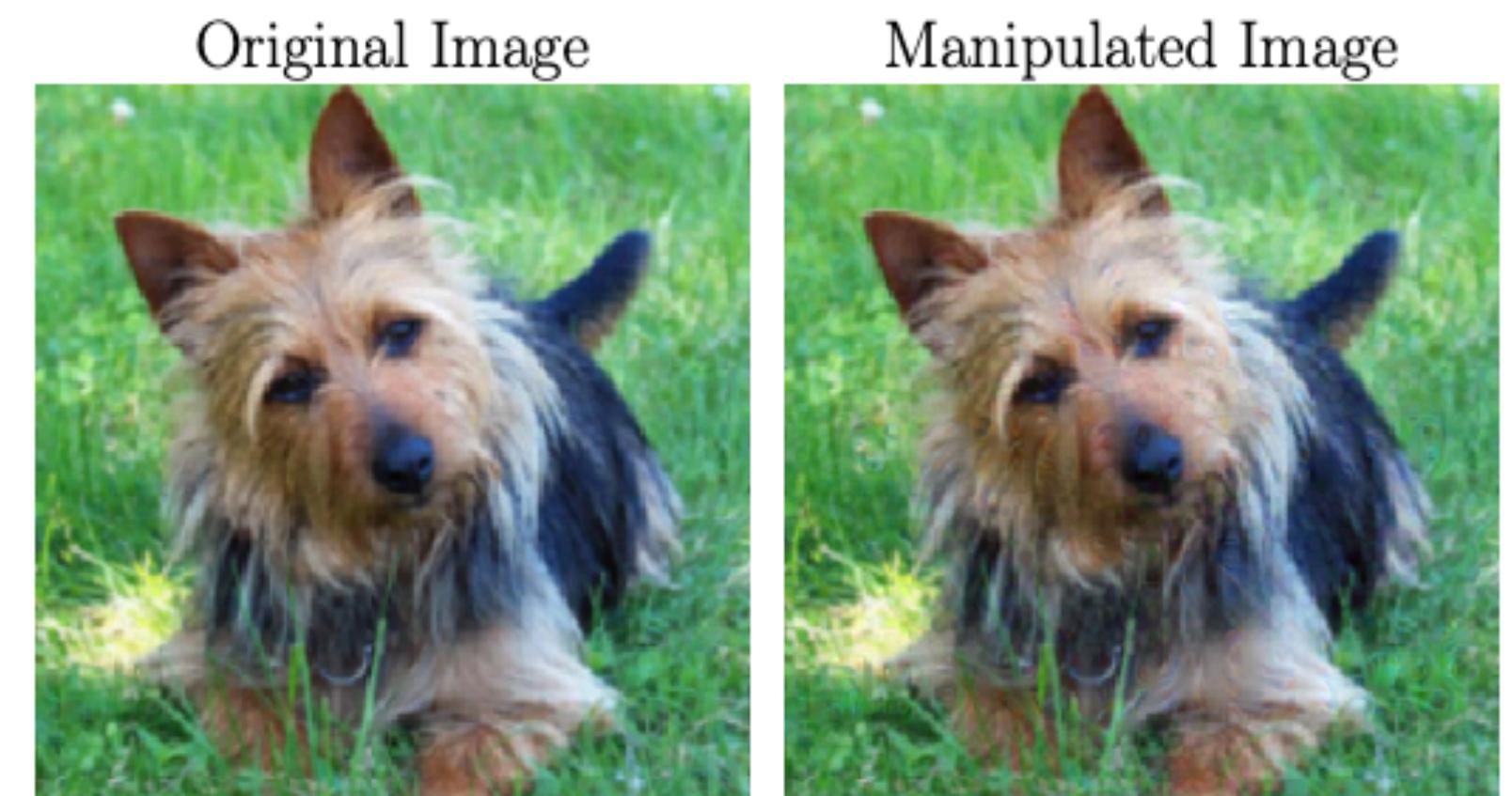
If you had an income of \$40,000 rather than \$30,000, your loan request would have been approved.

the minimal changes made to alter the decision

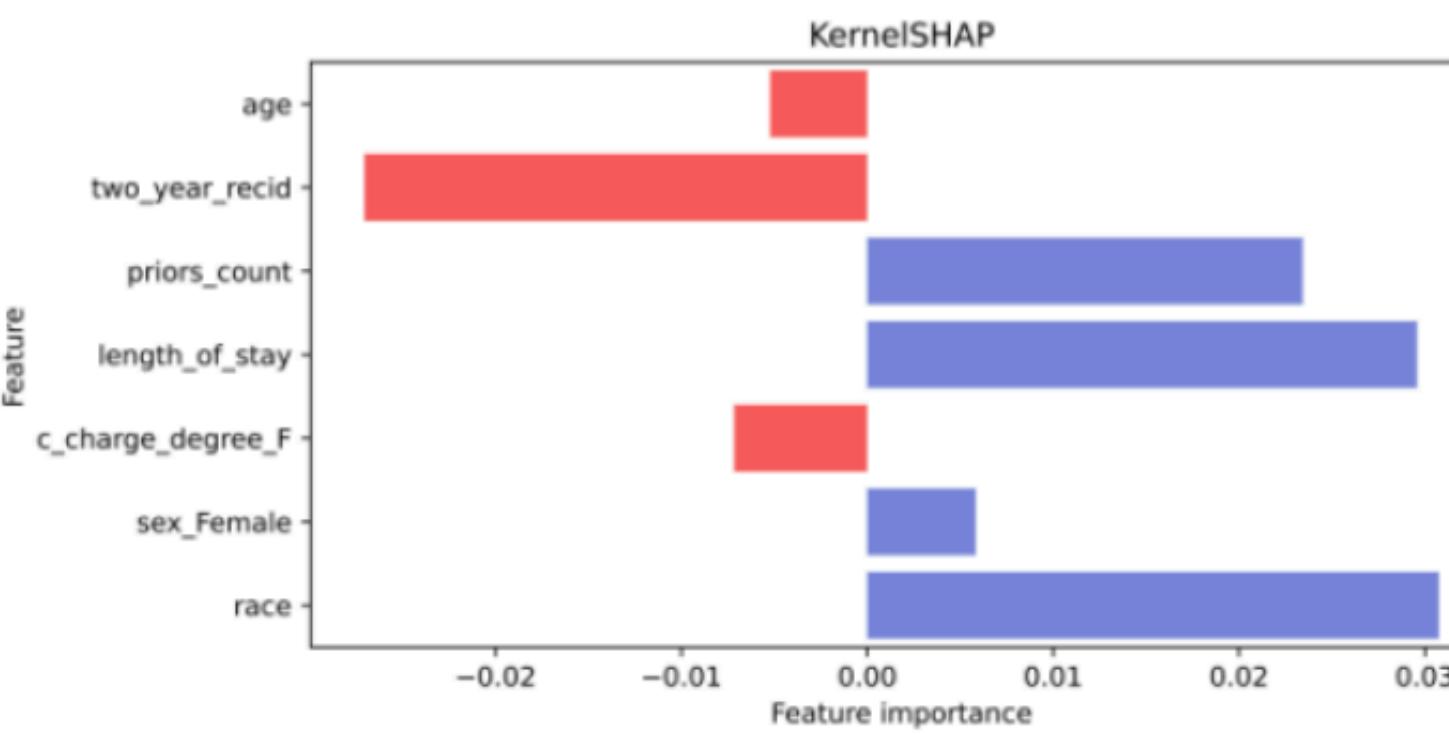
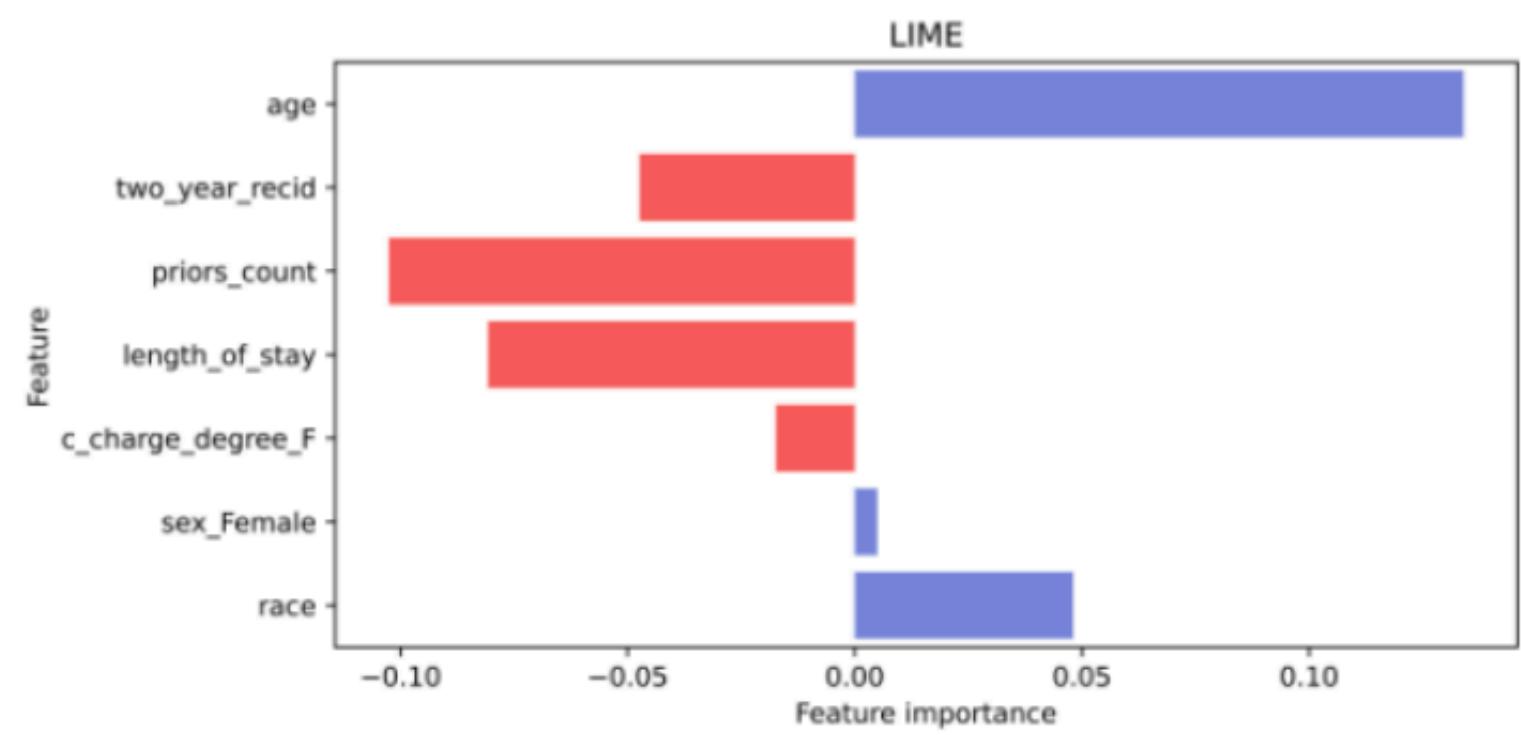
Explanations

Some issues

- ▶ Easily manipulated
- ▶ Disagreement problem
- ▶ Post-Hoc can be unfaithful



Below, you see a data point, as well as its explanation using methods **LIME** and **KernelSHAP**.



Explanations

Call for Rigor

- ▶ Mythos of model interpretability,
[Lipton, 2017]
- ▶ Towards a rigorous science of
interpretable machine learning,
[Doshi-Velez, Kim, 2017]

“Interpretability research suffers from an over-reliance on intuition-based approaches that risk — and in some cases have caused — illusory progress and misleading conclusions”,

[Leavitt, Morcos, 2020]

Attribution methods & Counterfactual methods

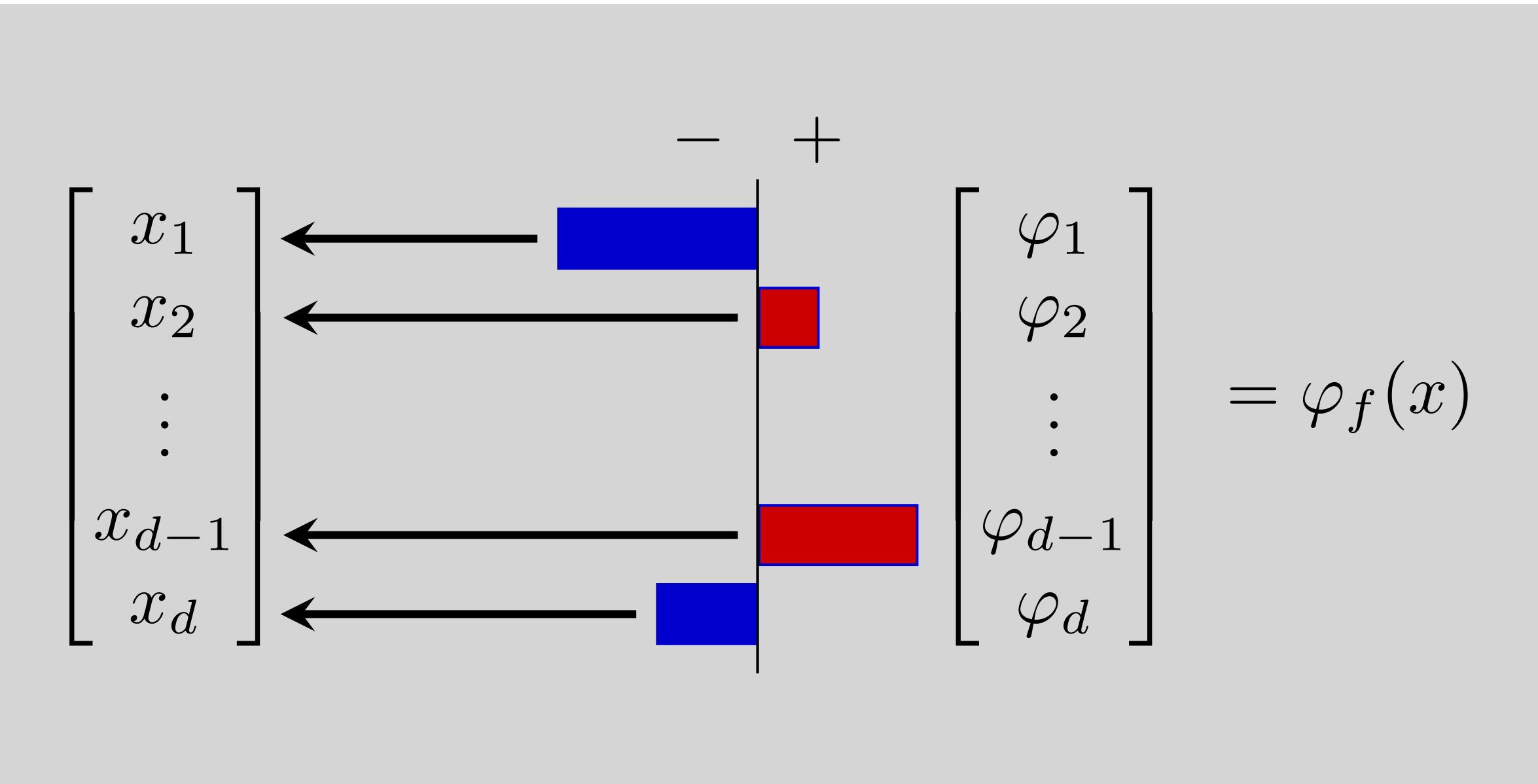
Attribution method

Attribute how important each feature was

► Positive value (typically) means that feature correlates positively with the outcome

► Negative value (typically) means that feature correlates negatively with the outcome

Not always the correct interpretation!



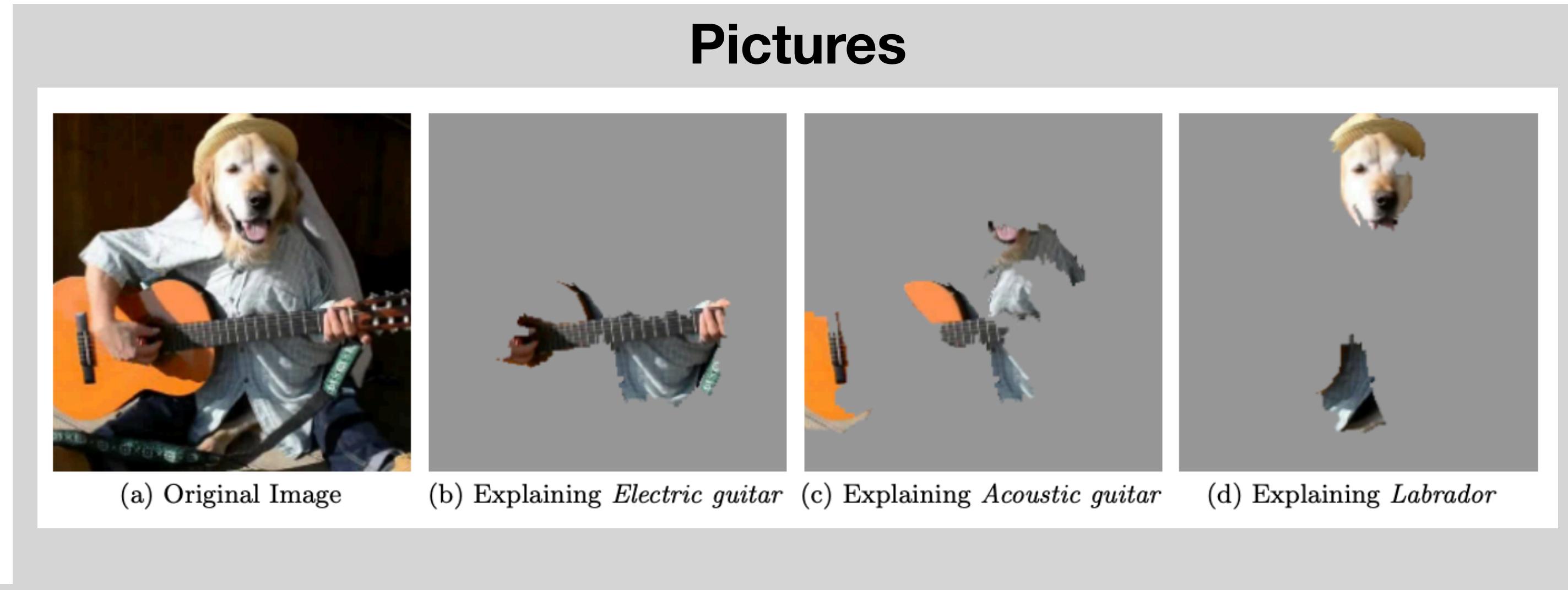
Your income of € 40 000 contributed positively.

However, your 5 different credit cards contributed negatively

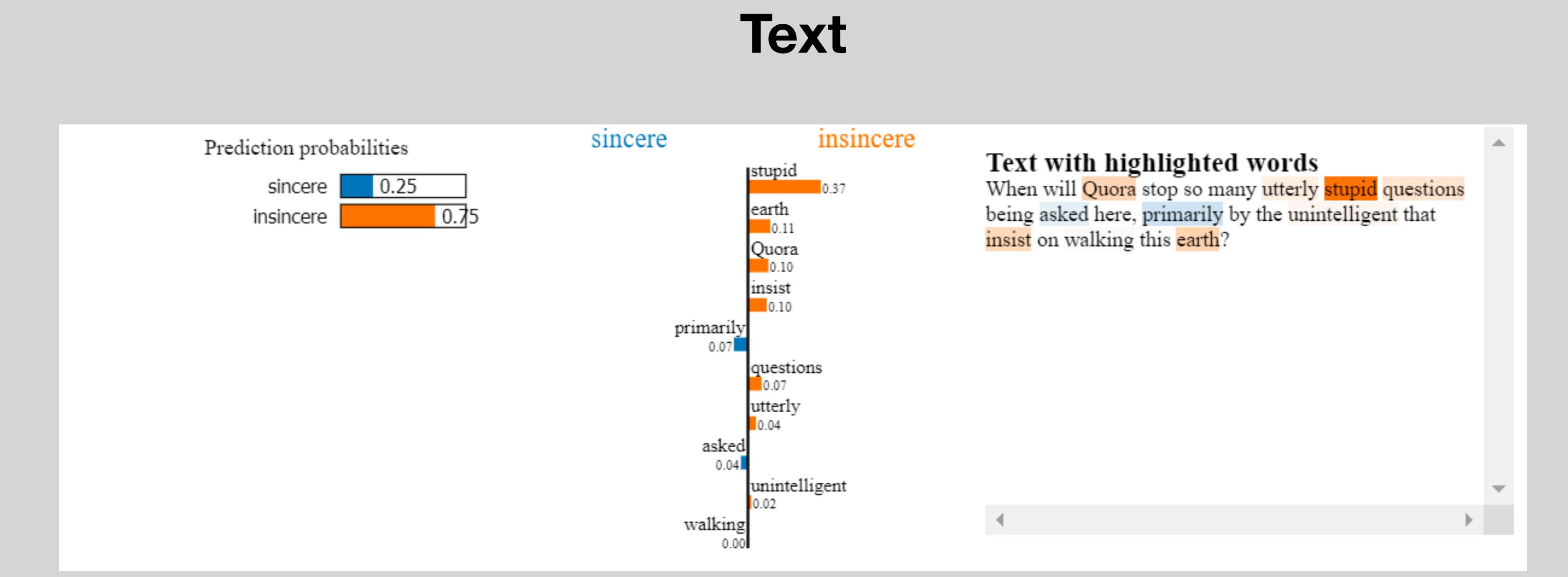
Attribution methods

Examples

- ▶ LIME
- ▶ SHAP
- ▶ ...



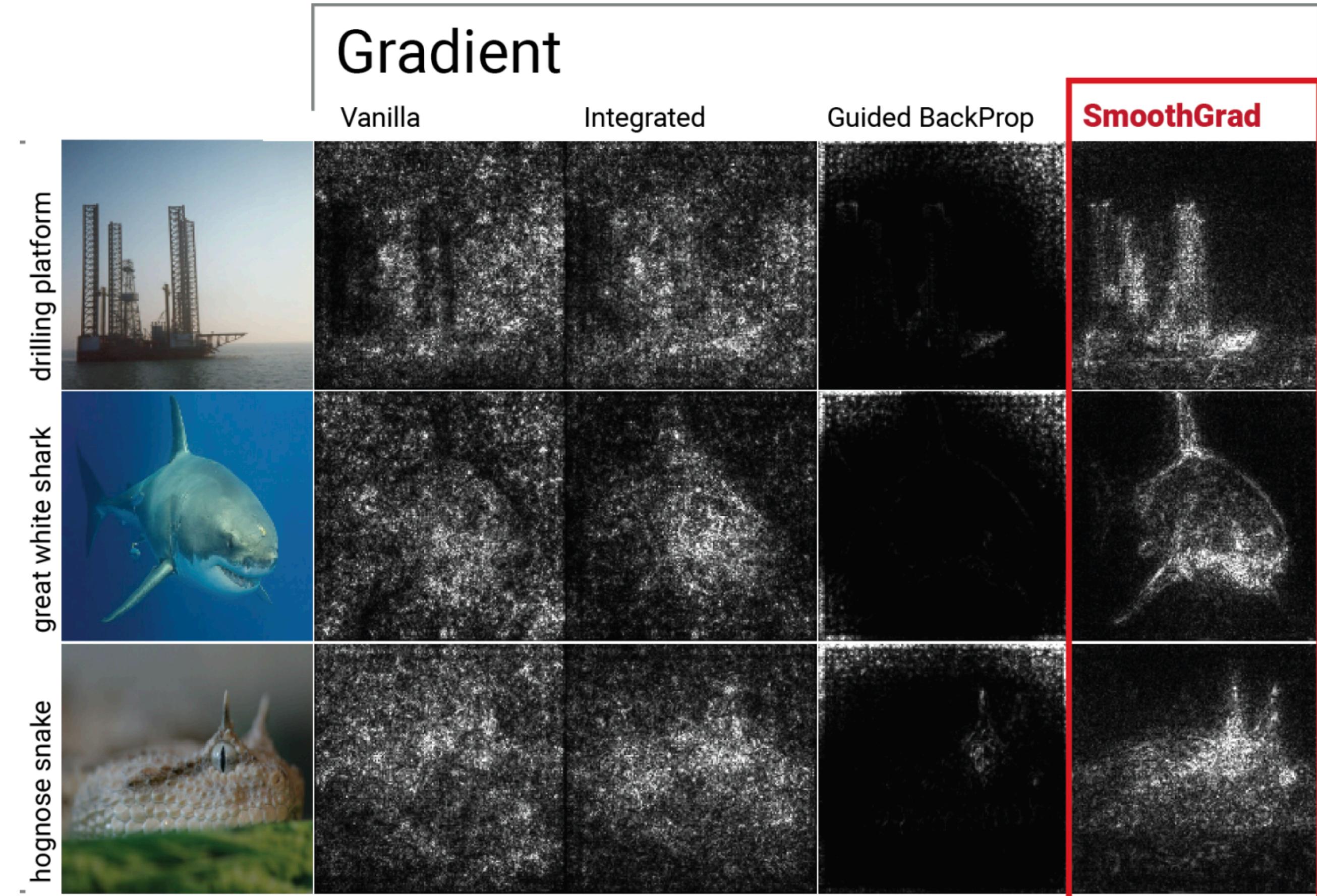
Text



Attribution methods

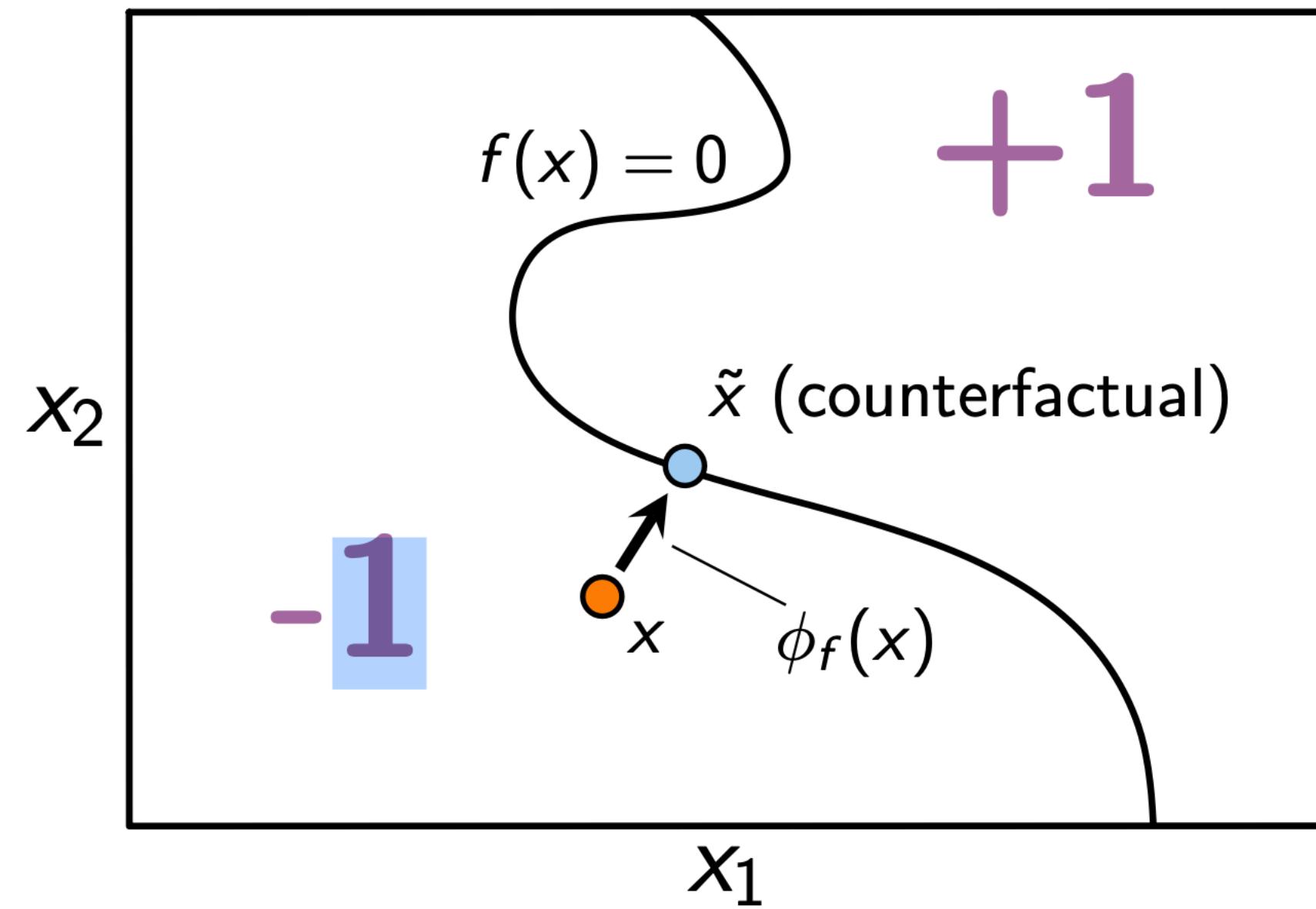
Examples

- ▶ Grad-Cam
- ▶ SmoothGrad
- ▶ Integrated Gradients
- ▶ ...



Counterfactual method

- ▶ Tell (A) how to change the decision from -1 to +1
- ▶ Minimal cost for (A)
- ▶ Provide *Recourse*



If you would have had an income of € 40 000 instead of €35 000, your loan request would have been approved.

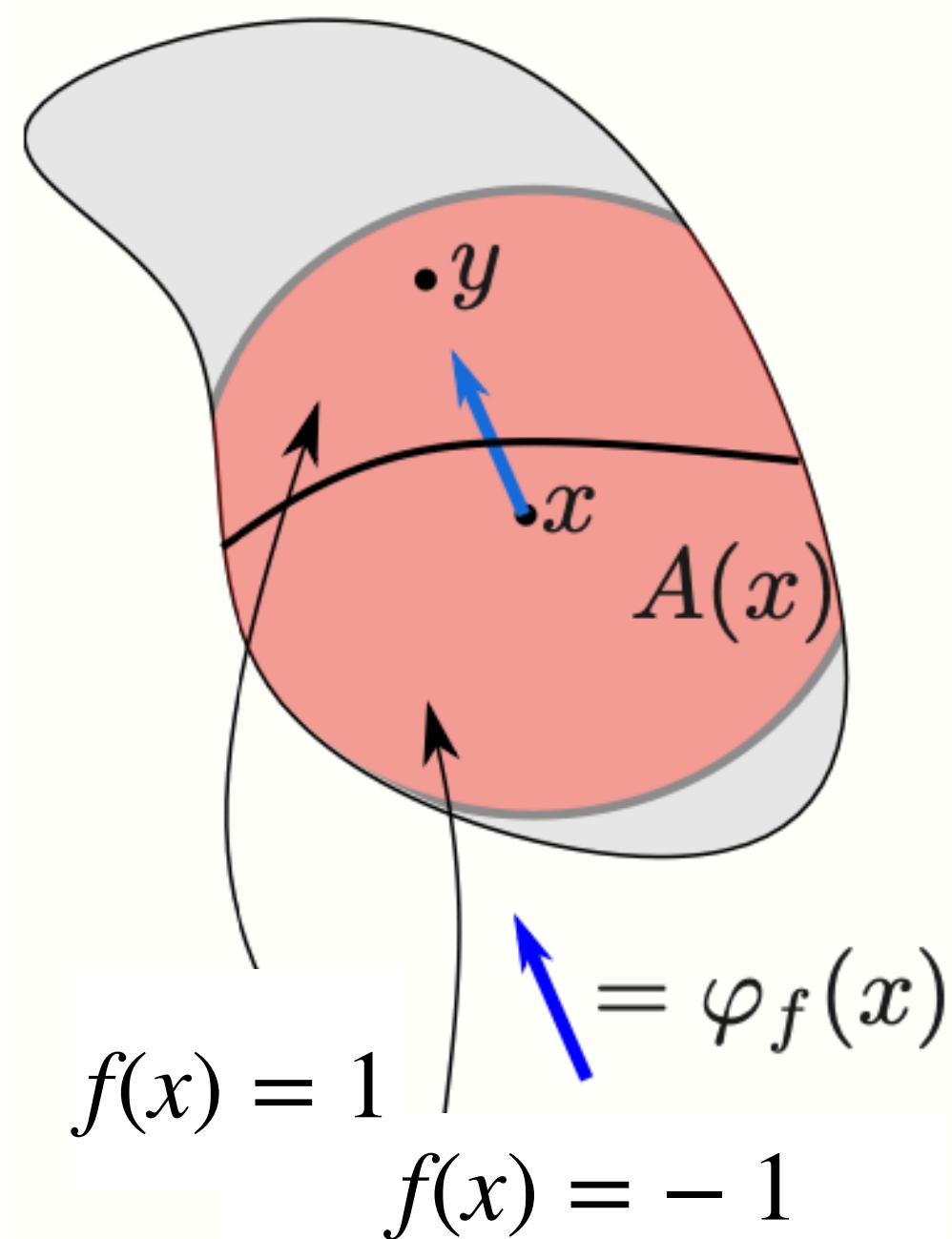
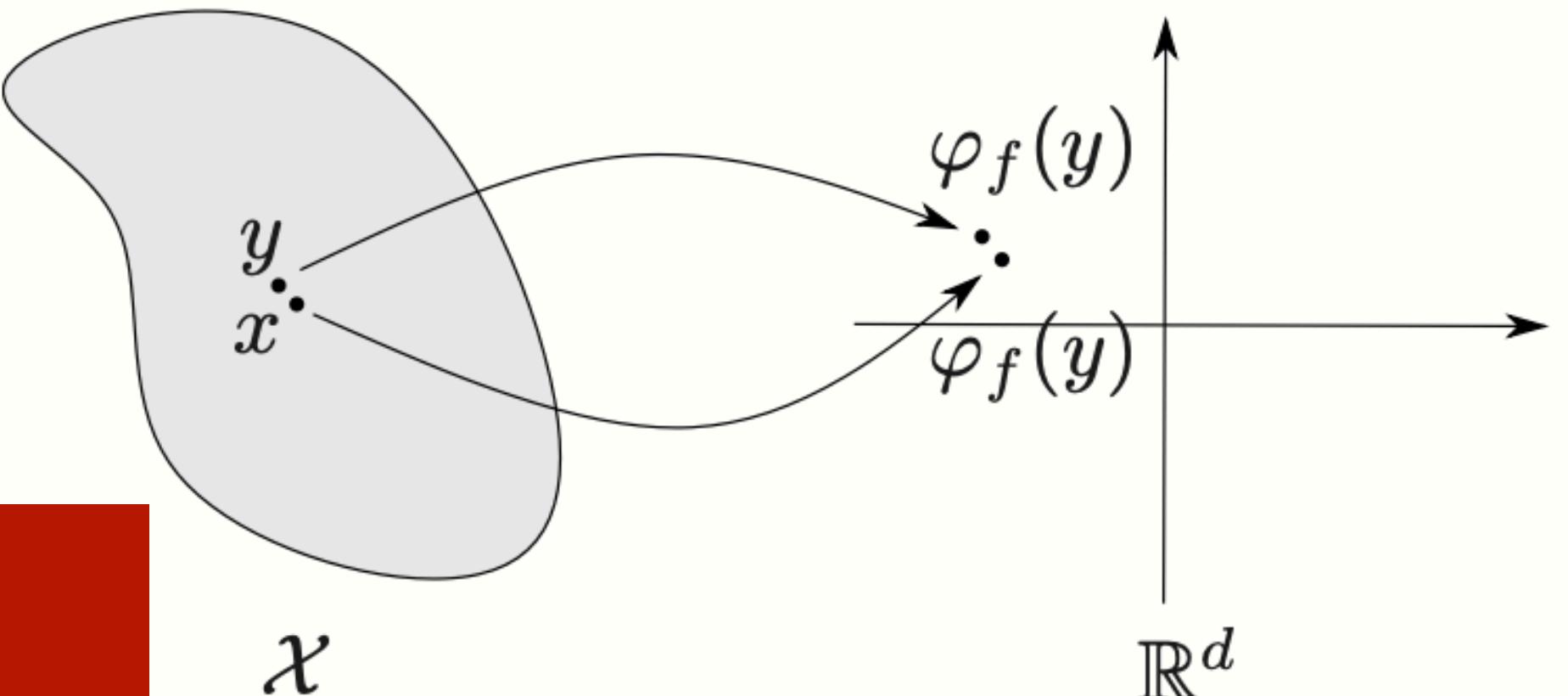
**Attribution methods that provide
recourse cannot be robust**

Desirable properties

Attribution methods

There are 2 desirable properties:

- ▶ Explanations need to be *robust*:
 - ▶ If features are similar, explanations should be similar
- ▶ Explanations should act as counterfactuals (**Recourse sensitive**):
 - ▶ Positive attribution -> increase that value
 - ▶ Negative attribution -> decrease that value
- ▶ Goal: change the decision from -1 to 1



Loan Applicant Example

Robustness



Features:

- ▶ Income: 40.000
- ▶ Creditcards: 5
- ▶ ...
- ▶ Gender: Female
- ▶



Features:

- ▶ Income: 40.000
- ▶ Creditcards: 5
- ▶ ...
- ▶ Gender: Male
- ▶



“Your income of € 40 000 contributed positively.

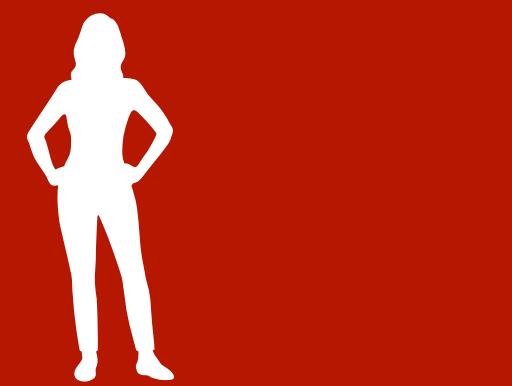
However, your 5 different credit cards contributed negatively”

Loan Applicant Example

Recourse sensitivity



***“Your income of € 40 000 contributed positively.
However, your 5 different credit cards
contributed negatively”***



***“If I decrease my number of credit cards,
I will get the loan!”***

Impossibility result

Attribution methods cannot always
Be Robust
Provide Recourse

Impossibility result

Specifically,

- ▶ Someone develops an attribution method
- ▶ I can construct a classifier, for which
 - ▶ Either **robustness** fails,
 - ▶ Or **recourse sensitivity**,

Attribution methods cannot always

Provide recourse

Be robust

A taste of hope

There are models f , that

- ▶ Allow attributions
- ▶ Which are
 - ▶ **robustness**,
 - ▶ **recourse sensitivity**,
- ▶ Linear models
- ▶ Monotone models

In general,

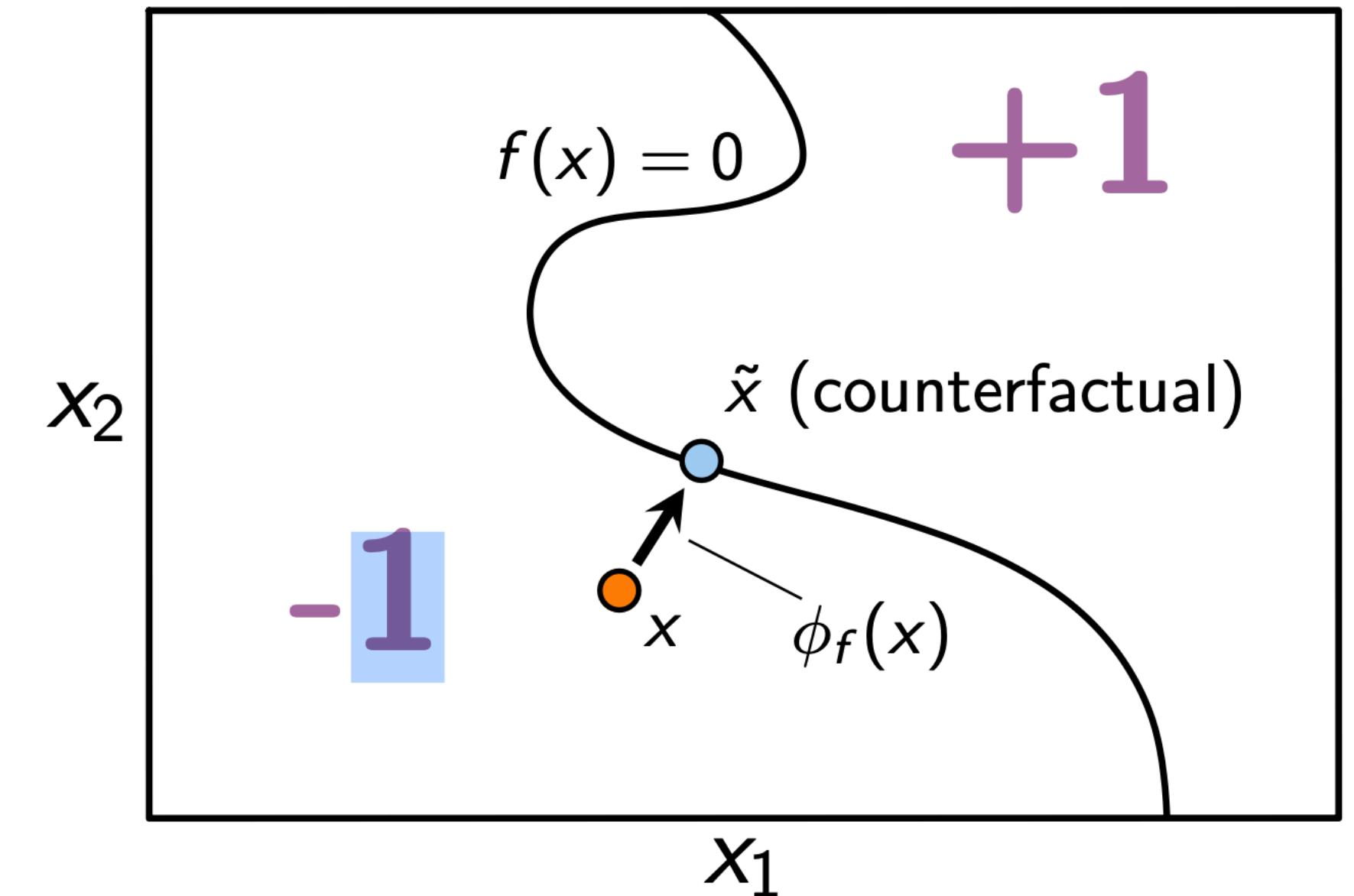
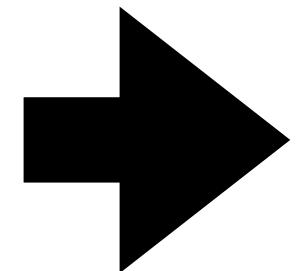
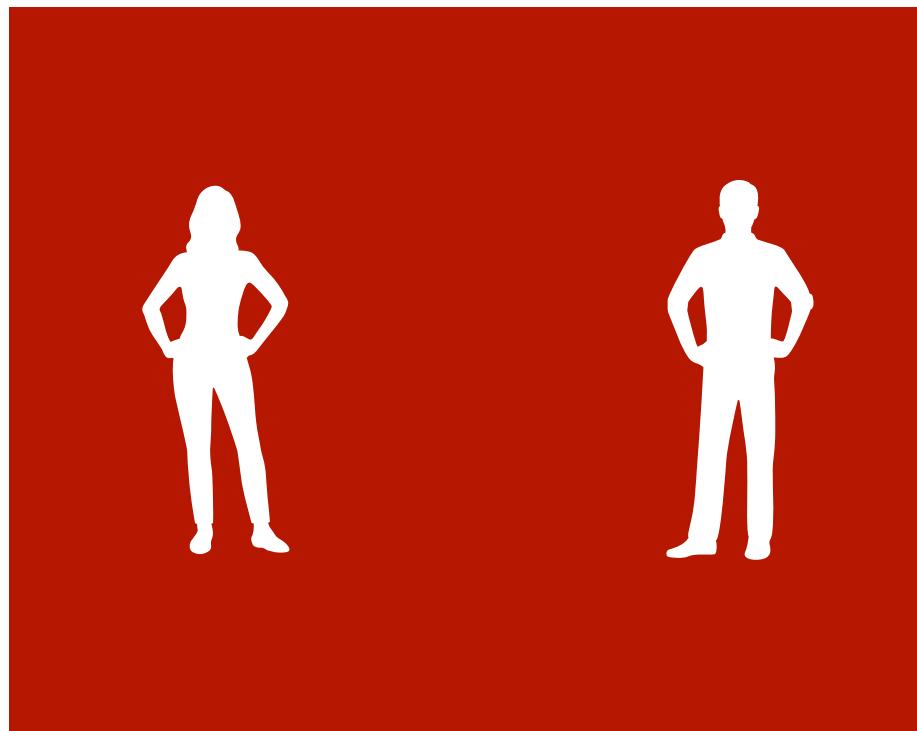
- ▶ Very difficult to check,
- ▶ Pretty easy to break

Risk of Recourse

Counterfactual methods

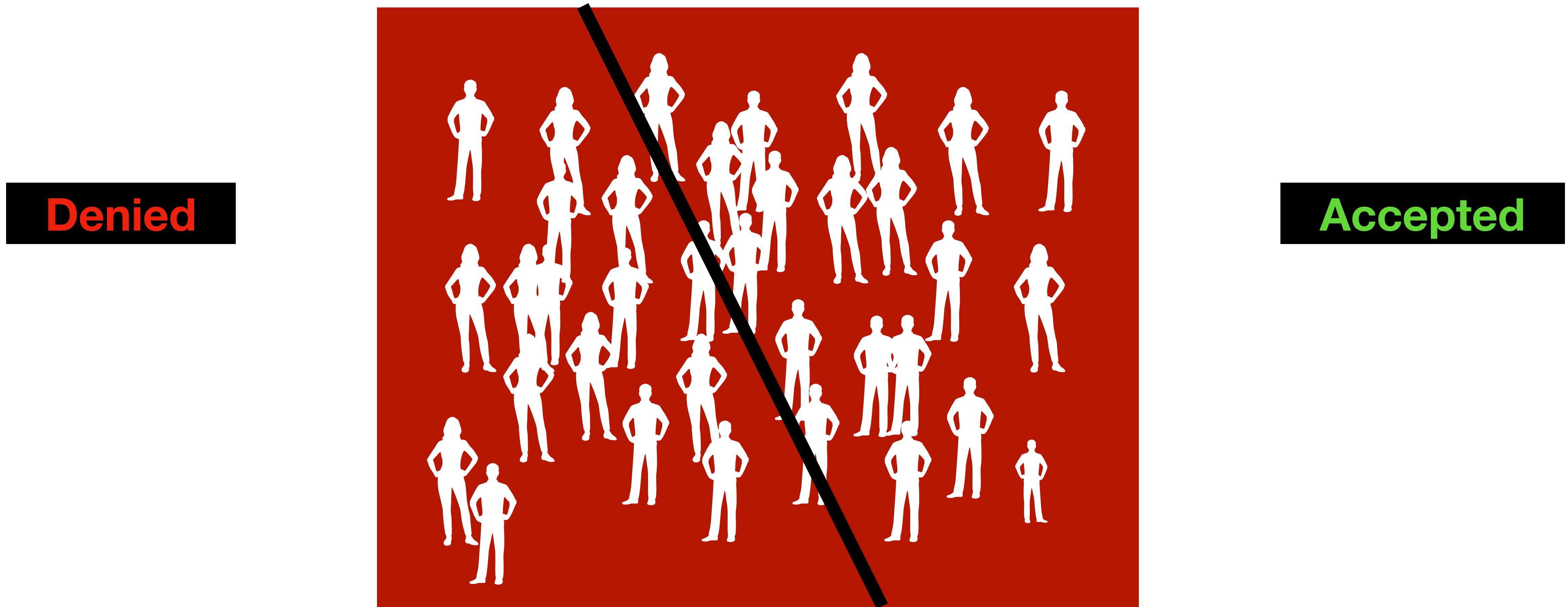
No Attribution methods

- ▶ Forget about attribution methods for now
- ▶ Consider Counterfactual Explanations
- ▶ What happens when we look at the population?



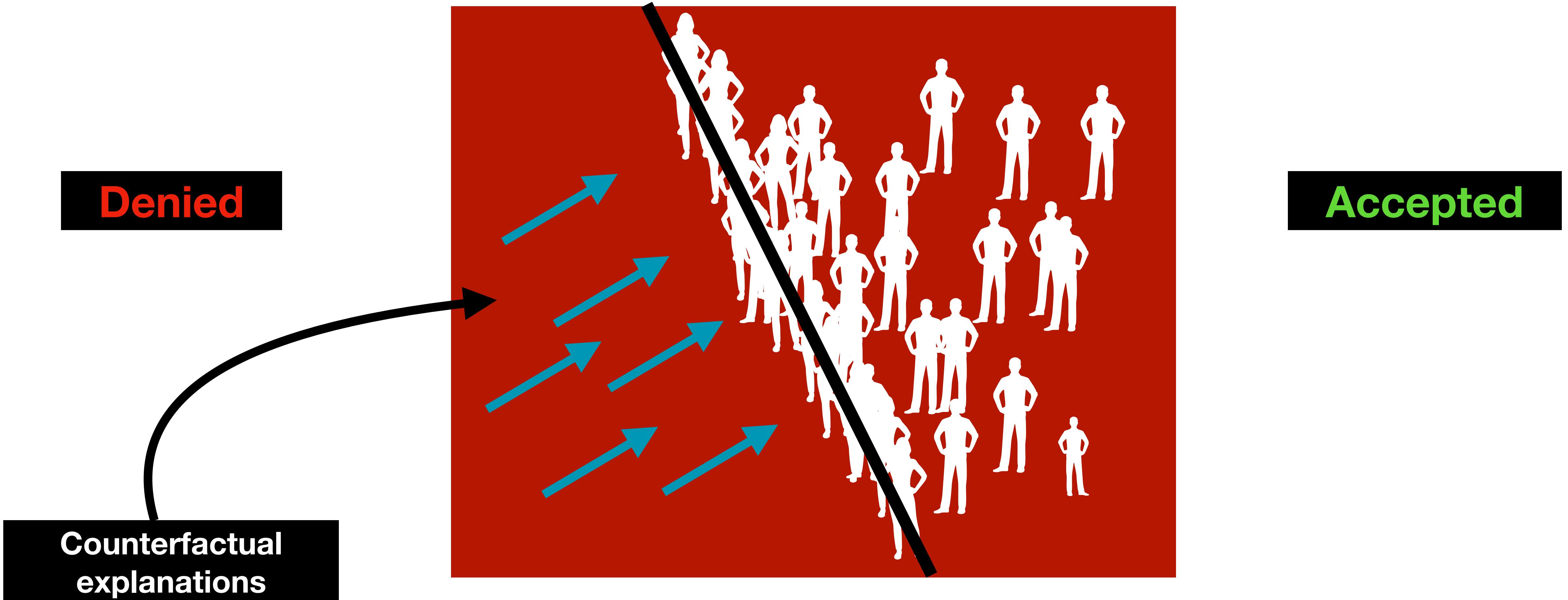
Counterfactual methods

No Attribution methods



Counterfactual methods

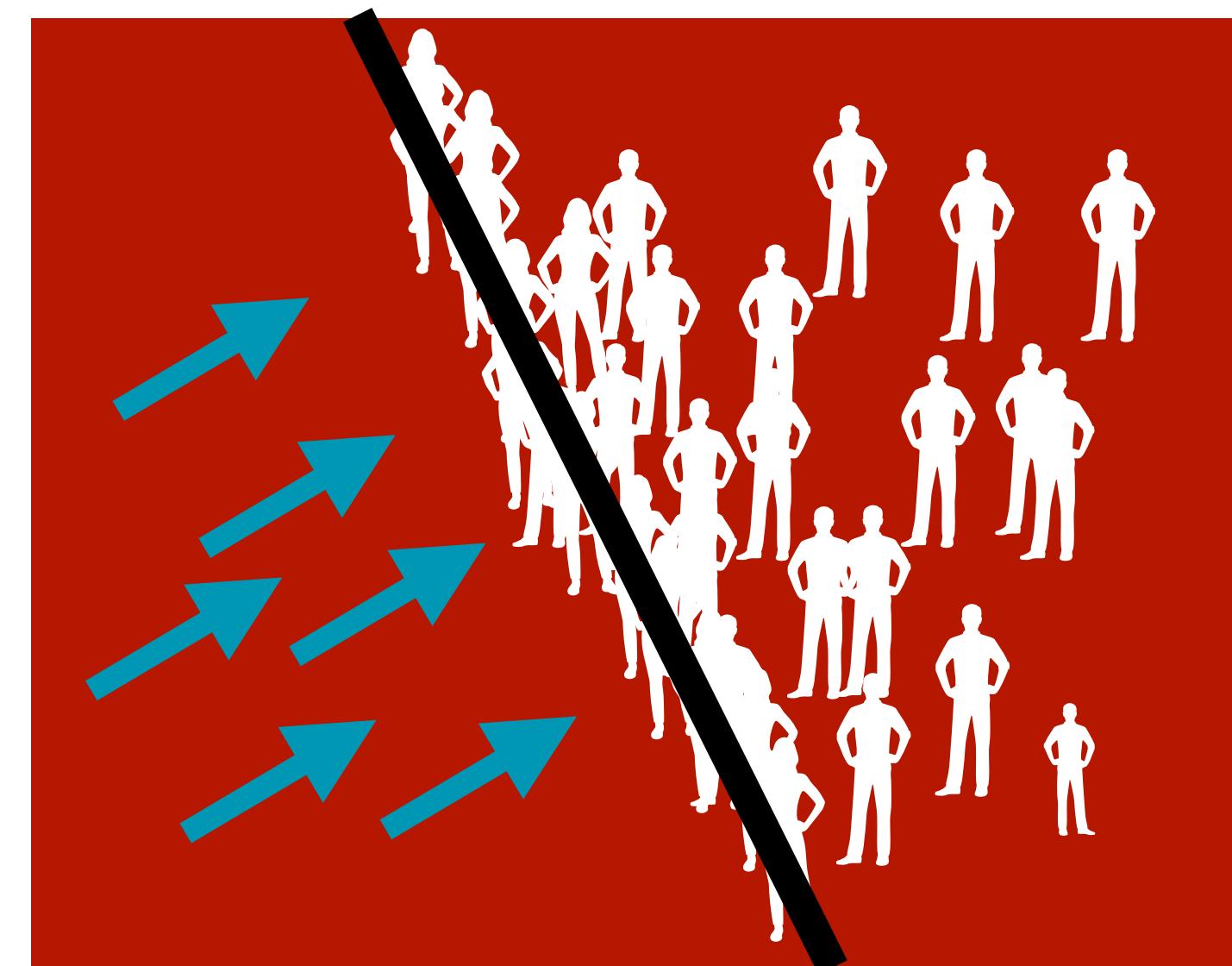
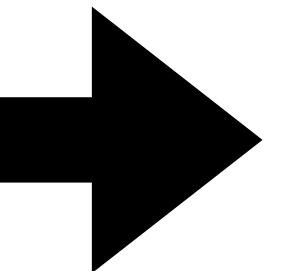
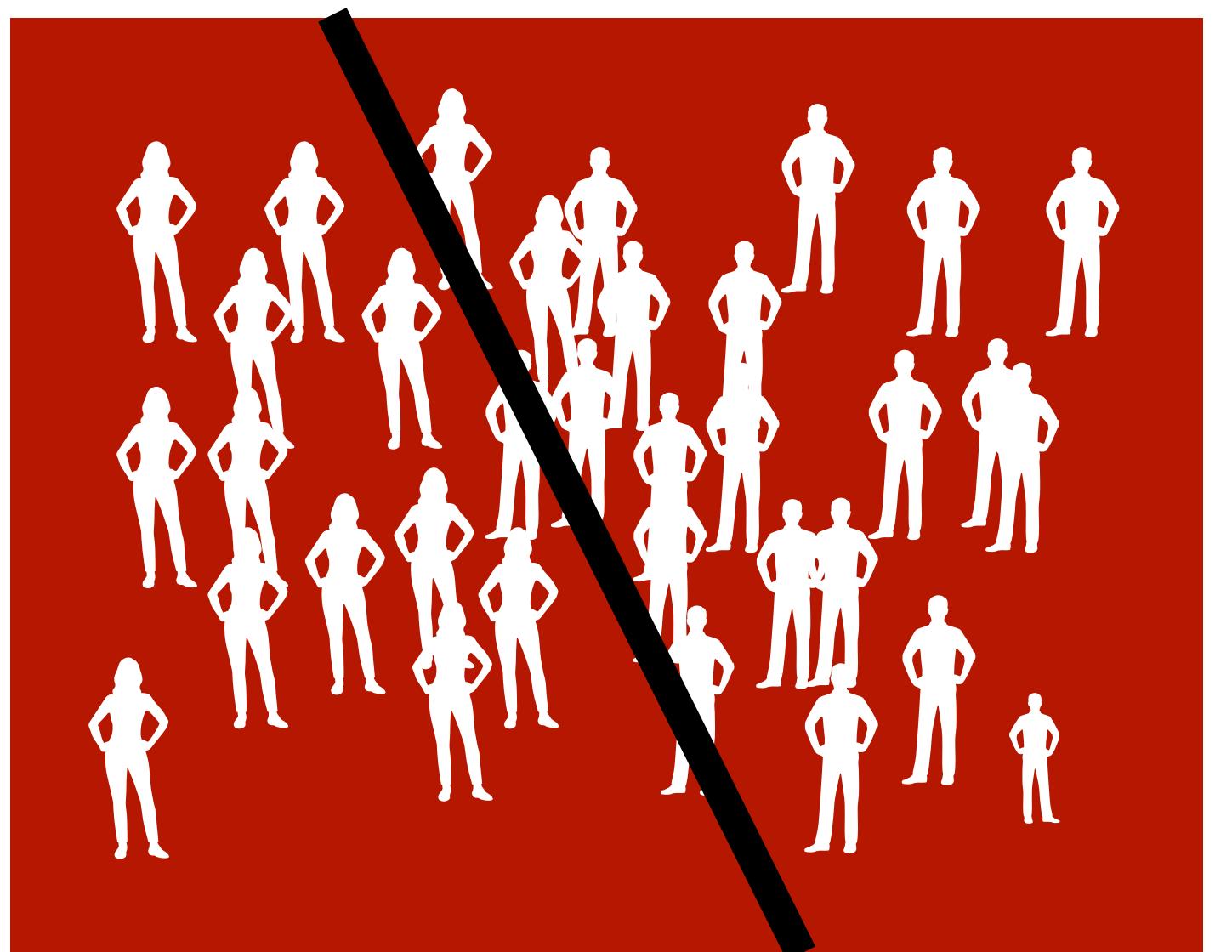
No Attribution methods



The risk of recourse

- What happens when we look at the population?
 - Underlying data distribution changes!
 - Modelling assumptions are violated!

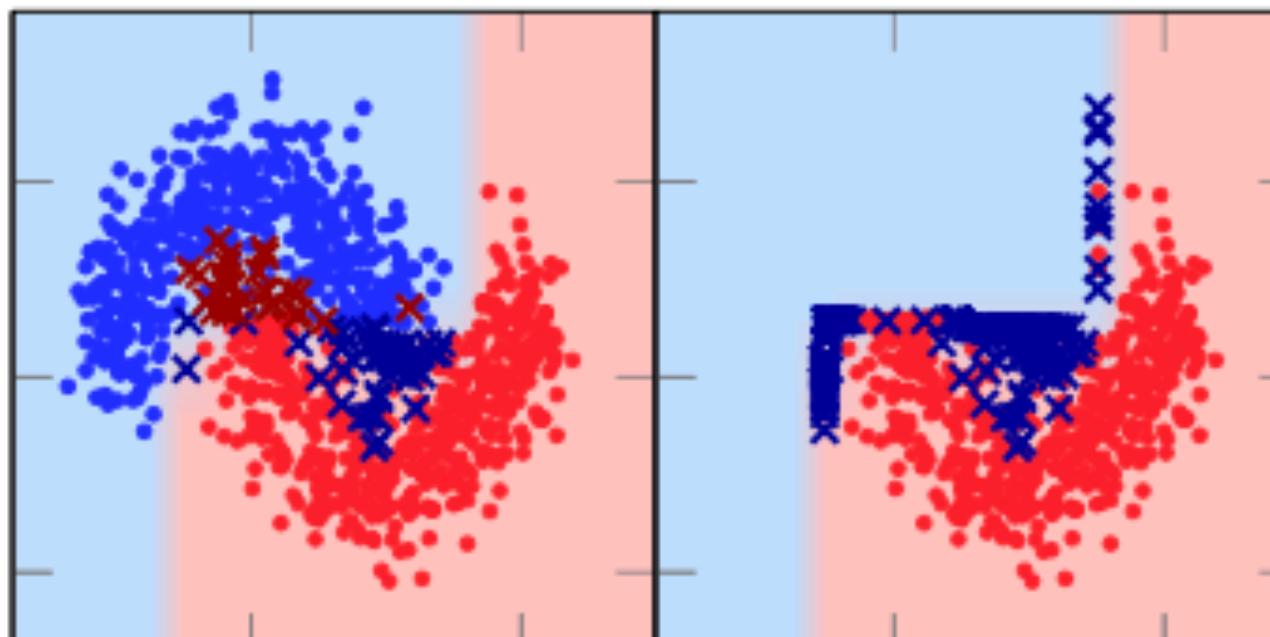
Accuracy will drop in most cases!



The risk of recourse

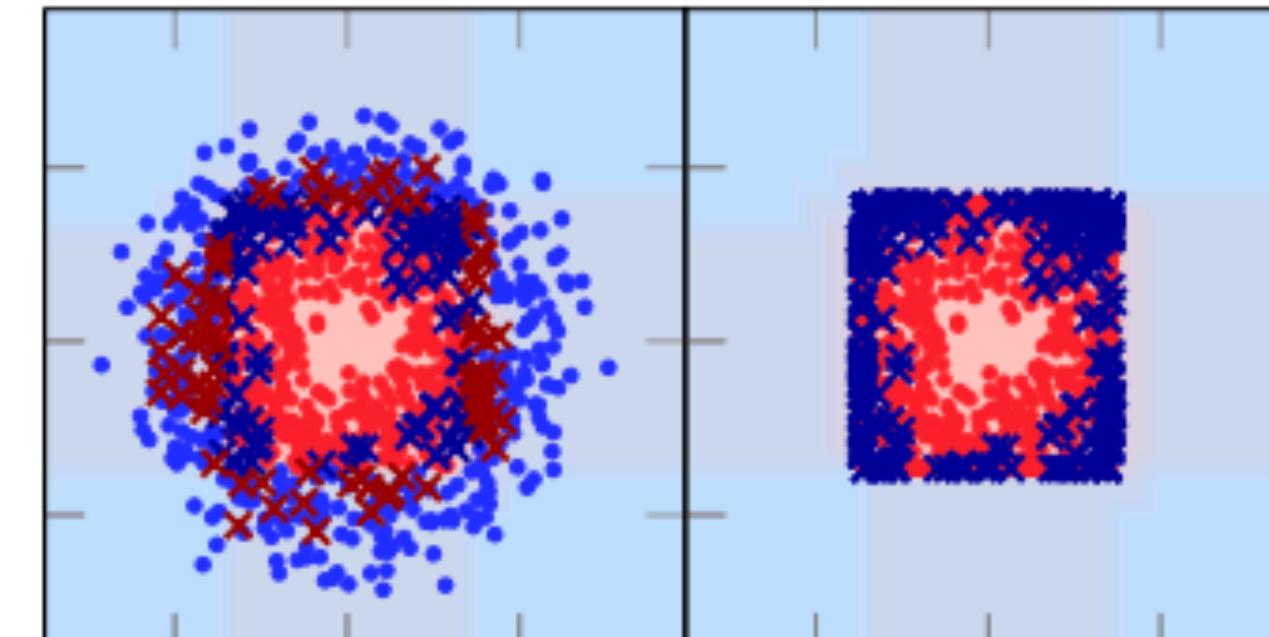
- What happens when we look at the population?
 - Underlying data distribution changes!
 - Modelling assumptions are violated!

Accuracy will drop in most cases!



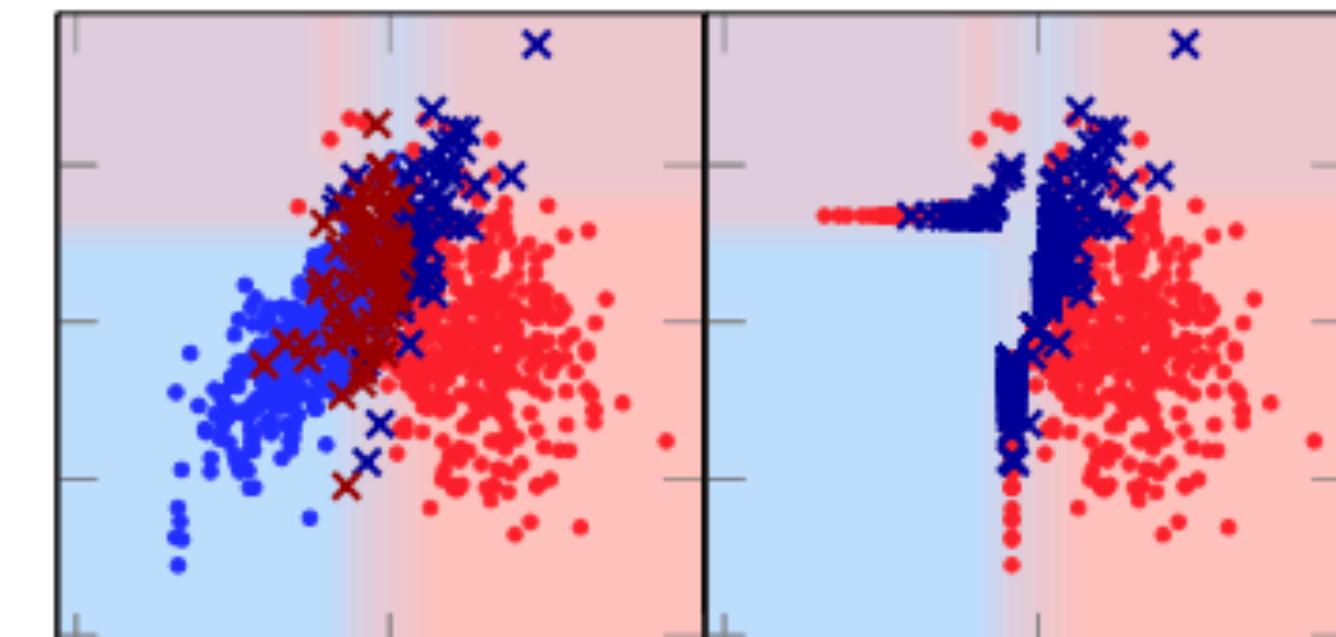
$$\widehat{R}_P(f) = 0.09$$

$$\widehat{R}_Q(f) = 0.30$$



$$\widehat{R}_P(f) = 0.19$$

$$\widehat{R}_Q(f) = 0.26$$



$$\widehat{R}_P(f) = 0.13$$

$$\widehat{R}_Q(f) = 0.33$$

The risk of recourse

Strategising against

Can (P) strategise against this accuracy drop?

- ▶ Need to assume that not everyone gets an explanation

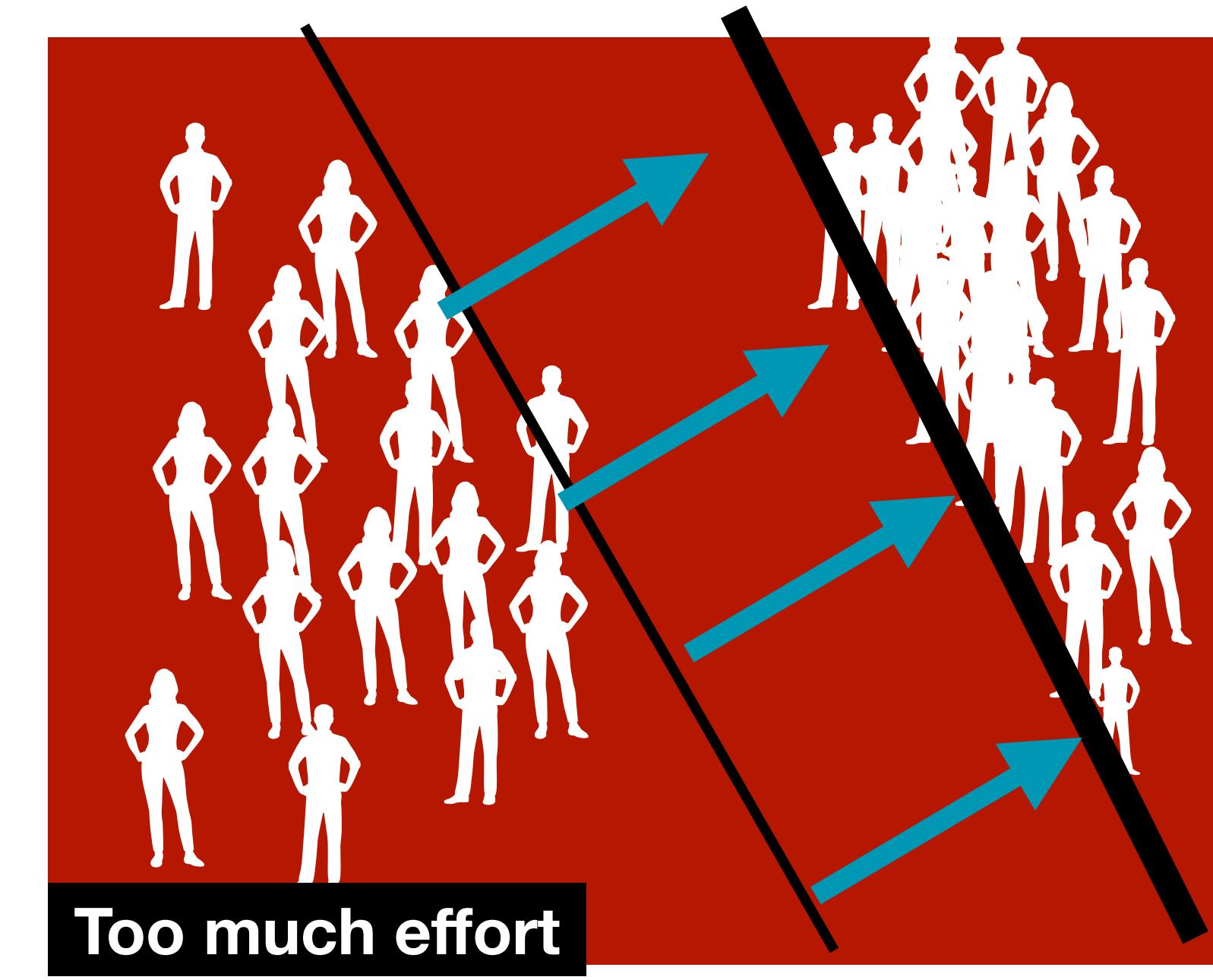
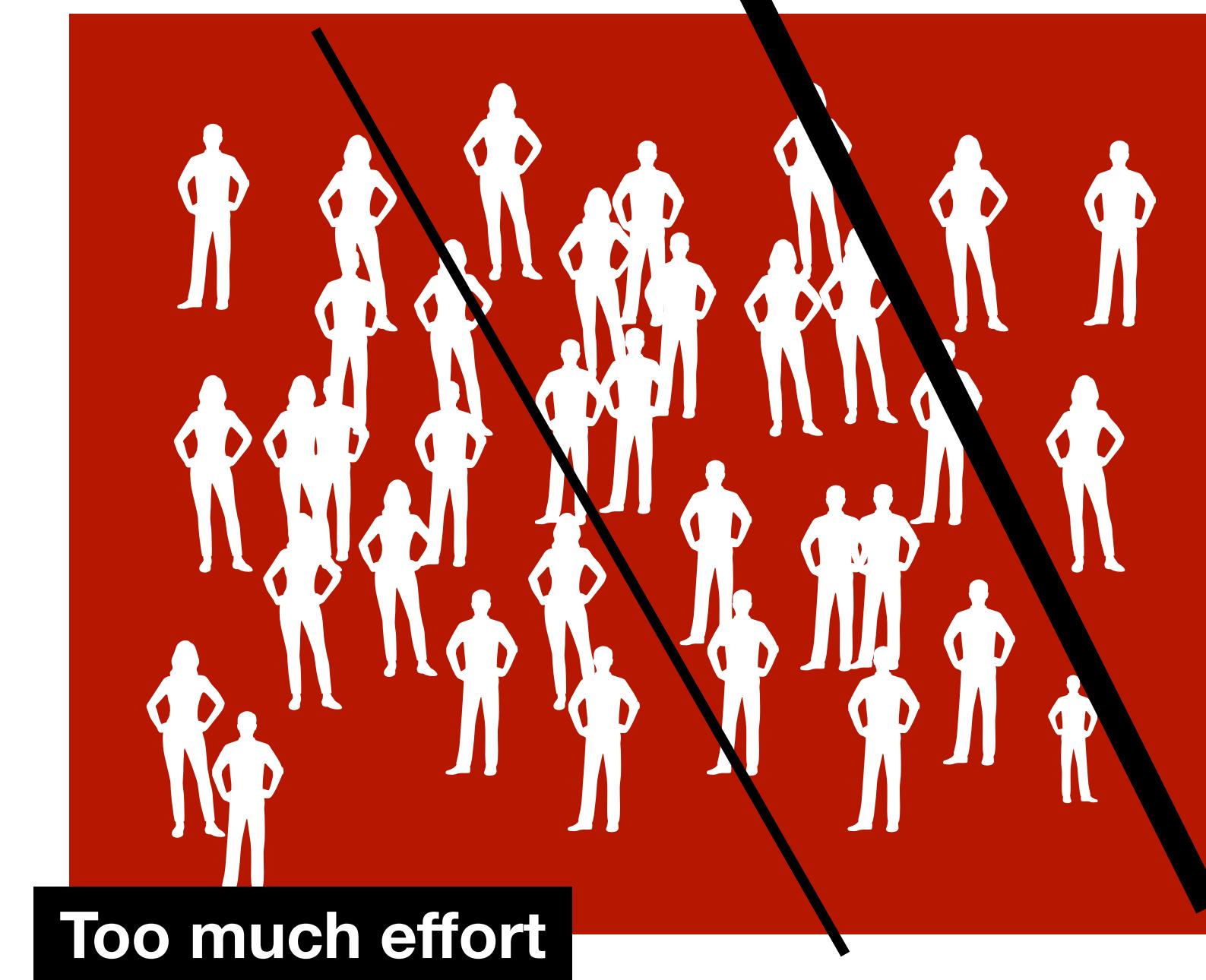
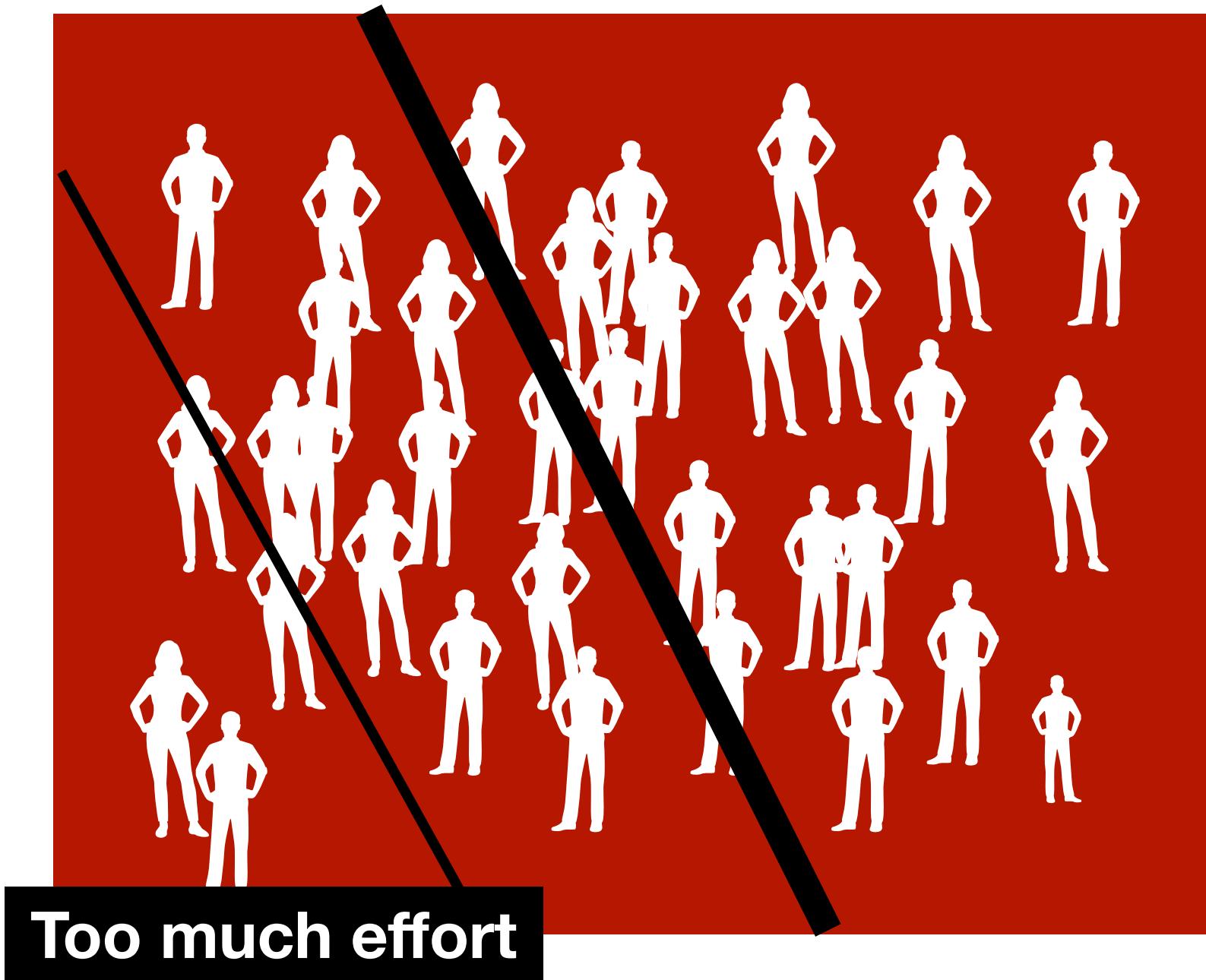
Too much effort



Yes and No

The risk of recourse

Strategising against



Exactly the same people get a loan

More effort

Practical implications

Practical implications

- ▶ Define the goal of your explanation, then pick a method
- ▶ When using attribution methods, be careful interpreting the scores
- ▶ When using counterfactual methods, ask the question
 - ▶ “How harmful is misclassification?”
- ▶ Whenever possible build simple models, increase complexity when necessary!

Some references

- ▶ *Mythos of model interpretability*, [Lipton, 2017]
- ▶ *Towards a rigorous science of interpretable machine learning*, [Doshi-Velez, Kim, 2017]
- ▶ *Towards falsifiable interpretability research* [Leavitt, Morcos, 2020]
- ▶ *Attribution-based Explanations that Provide Recourse Cannot be Robust* [F, De Heide, van Erven, 2022]
- ▶ *The Risks of Recourse in Binary Classification* [F, Garreau, van Erven, 2023]

Thank you for your attention!