



# Human-Centric Machine Learning

## Feedback loops, Human-AI Collaboration and Strategic Behavior

**Manuel Gomez Rodriguez**

Includes joint work with Behzad Tabibian, Stratis Tsirtsis, Niki Kilbertus, Moein Khajehnejad, Paramita Koley, Abir De, Isabel Valera, Krikamol Muandet, Adish Singla, Niloy Ganguly and Bernhard Schölkopf

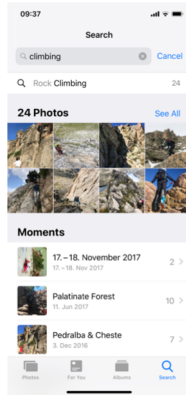


MAX PLANCK INSTITUTE  
FOR SOFTWARE SYSTEMS

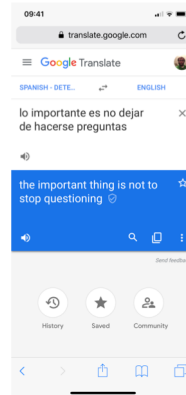
# Machine learning for automation

Machine learning (ML) has *taught* machines to...

recognize images



translate between languages



play complex games



...to name a few

**Machines achieve, or surpass, human performance  
at tasks for which intelligence is required**

# Machine learning for decision making

ML has the potential to support & enhance high-stakes decision making in a wide range of applications:



**Justice**



**Hiring**



**Information  
integrity**



**Education**



**Security**



**Health**



**Finance**

# Increasing number of missteps

Machine learning has been blamed to be one of the root causes of an increasing number of missteps

misleading people in  
social media



increasing polarization



causing car accidents



discriminating minorities

PROPUBLICA | MACHINE BIAS

## Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.



# What went wrong in these cases?

Machine learning has been mostly  
developed for **automation**



Take decisions autonomously  
on the basis of

**passively collected data**

passive setting

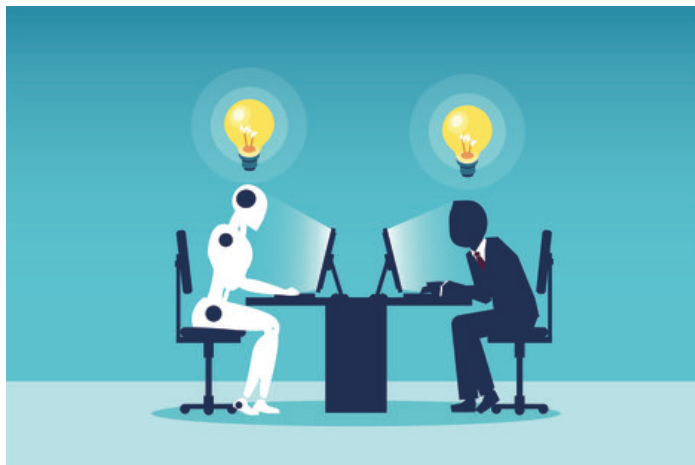
# What went wrong in these cases?

Algorithmic and human decisions **feed** and **influence**  
each other



**Sequential decision  
making process**

reactive setting



# Shortcomings of (traditional) ML models & algorithms



**Ignore feedback loop between  
algorithmic and human decisions**

# Shortcomings of (traditional) ML models & algorithms



**Ignore feedback loop between  
algorithmic and human decisions**



**Fail to anticipate how individuals will react  
to algorithmic decisions**

# Shortcomings of (traditional) ML models & algorithms



**Ignore feedback loop between  
algorithmic and human decisions**



**Do not account for strategic  
human behavior**

# Shortcomings of (traditional) ML models & algorithms



**Ignore feedback loop between  
algorithmic and human decisions**



**Do not account for strategic  
human behavior**



**Unexpected & undesirable personal,  
social and economic consequences**

# Shortcomings of (traditional) ML models & algorithms



**Ignore feedback loop between  
algorithmic and human decisions**



**Do not account for strategic  
human behavior**



**Fail to balance decisions between  
machines and humans**

# Shortcomings of (traditional) ML models & algorithms



**Ignore feedback loop between  
algorithmic and human decisions**



**Do not account for strategic  
human behavior**



**Fail to balance decisions between  
machines and humans**



**They are unable to collaborate with humans**



# Shortcomings of (traditional) ML models & algorithms



**Ignore feedback loop between  
algorithmic and human decisions**



**Do not account for strategic  
human behavior**



**Fail to balance decisions between  
machines and humans**



**Do not provide actionable  
insights**

# Shortcomings of (traditional) ML models & algorithms



**Ignore feedback loop between  
algorithmic and human decisions**



**Do not account for strategic  
human behavior**



**Fail to balance decisions between  
machines and humans**



**Do not provide actionable  
insights**



**Interpretability is necessary to use  
ML in critical domains with consequential decisions.**

# Outline of this talk

A glimpse on recent advances on human-centric ML models and algorithms. We will focus on:



Accounting for the feedback loop between algorithmic and human decisions



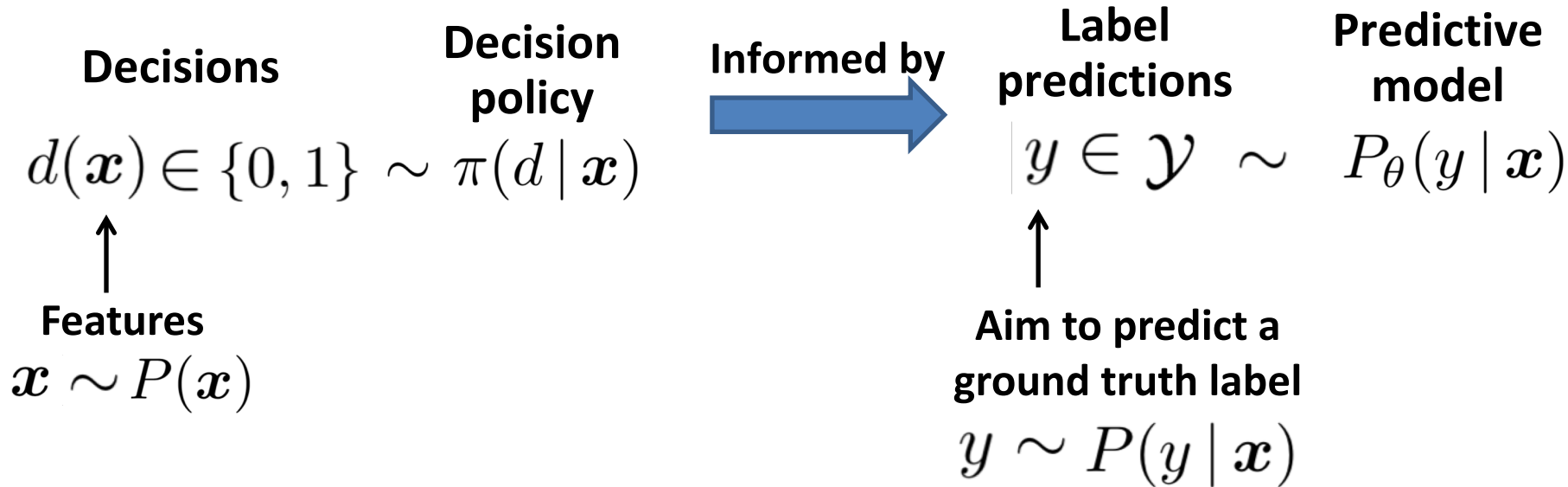
Balancing decisions between human and algorithmic decisions



Accounting for strategic human decisions

**Disclaimer.** These are emerging topics. The goal of this talk is to introduce you to a new set of problems and, for each problem, show you one solution, not *the* solution.

# A general problem setting



# Example 1: loan decisions



**Decisions**

$$d(x) \in \{0, 1\} \sim \pi(d | x)$$

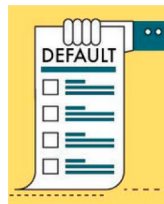
Individual  
is rejected

Individual  
receives loan



**Decision  
policy**

Informed by



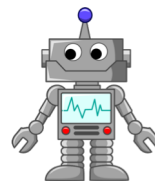
**Label  
predictions**

$$y \in \mathcal{Y} \sim P_{\theta}(y | x)$$

$$\mathcal{Y} = \{0, 1\}$$

Individual  
defaults

Individual  
pays back



**Predictive  
model**

# Example 2: bail decisions



**Decisions**

$$d(x) \in \{0, 1\} \sim \pi(d | x)$$

Individual  
remains jailed

Individual  
is released



**Decision  
policy**

Informed by



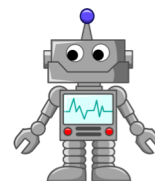
**Label  
predictions**

$$y \in \mathcal{Y} \sim P_{\theta}(y | x)$$

$$\mathcal{Y} = \{0, 1\}$$

Individual  
reoffends

Individual  
does not reoffend



**Predictive  
model**

# Example 3: medical diagnosis



**Decisions**

$$d(x) \in \{0, 1\} \sim \pi(d | x)$$

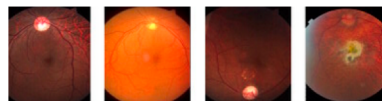
Individual  
doesn't need  
further tests

Individual  
needs further  
tests



**Decision  
policy**

Informed by



**Label  
predictions**

$$y \in \mathcal{Y} \sim P_{\theta}(y | x)$$

$$\mathcal{Y} = \{0, 1\}$$

Low severity of  
Drusen's disease

High severity of  
Drusen's disease



**Predictive  
model**

# Utility of a decision policy

The **decision maker** aims to deploy a **decision policy** that **maximizes** a very general definition of **utility**:

$$\begin{aligned} u(\pi, c) &= \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}), y \sim P(y | \mathbf{x}), d \sim \pi(d | \mathbf{x})} [y d(\mathbf{x}) - c d(\mathbf{x})] \\ &= \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}), d \sim \pi(d | \mathbf{x})} [P(y = 1 | \mathbf{x}) d(\mathbf{x}) - c d(\mathbf{x})] \end{aligned}$$



# Utility of a decision policy

The **decision maker** aims to deploy a **decision policy** that **maximizes** a very general definition of **utility**:

$$\begin{aligned} u(\pi, c) &= \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}), y \sim P(y | \mathbf{x}), d \sim \pi(d | \mathbf{x})} [y d(\mathbf{x}) - c d(\mathbf{x})] \\ &= \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}), d \sim \pi(d | \mathbf{x})} [P(y = 1 | \mathbf{x}) d(\mathbf{x}) - c d(\mathbf{x})] \end{aligned}$$

Example (loan decisions)

If a loan is granted and individual... repays:	$1 - c$	If a loan is	
... defaults:	$-c$	not granted:	$0$

The parameter  $c$  measures the cost of offering a loan in units of repaid loans

# Benefits of a decision policy

To ensure **fairness**, the **decision maker** may constrain the **benefits** individuals obtain:

$$b(\boldsymbol{x}, c) = \mathbb{E}_{d \sim \pi(d | \boldsymbol{x})} [\underbrace{f(d(\boldsymbol{x}), c)}_{\text{Problem dependent}}]$$


# Benefits of a decision policy

To ensure **fairness**, the **decision maker** may constrain the **benefits** individuals obtain:

$$b(\mathbf{x}, c) = \mathbb{E}_{d \sim \pi(d | \mathbf{x})} [\underbrace{f(d(\mathbf{x}), c)}_{\text{Problem dependent}}]$$

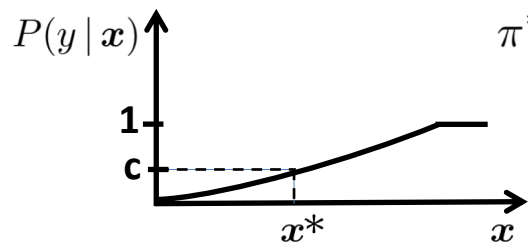
Example (loan decisions)

**Statistical parity:**  
Ensure the men and women  
have the same probability  
of receiving a loan

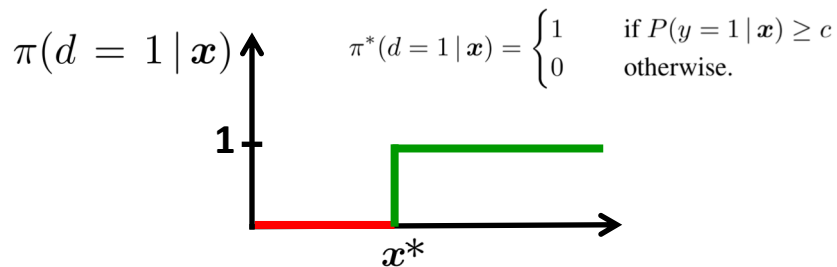

$$\left\{ \begin{array}{l} \sum_{\mathbf{x} \in \text{men}} b(\mathbf{x}, c) \approx \sum_{\mathbf{x} \in \text{women}} b(\mathbf{x}, c) \\ f(d(\mathbf{x}), c) = d(\mathbf{x}) \end{array} \right.$$

# Deterministic threshold rules

Under some *technical conditions*, deterministic threshold rules are optimal decision policies:

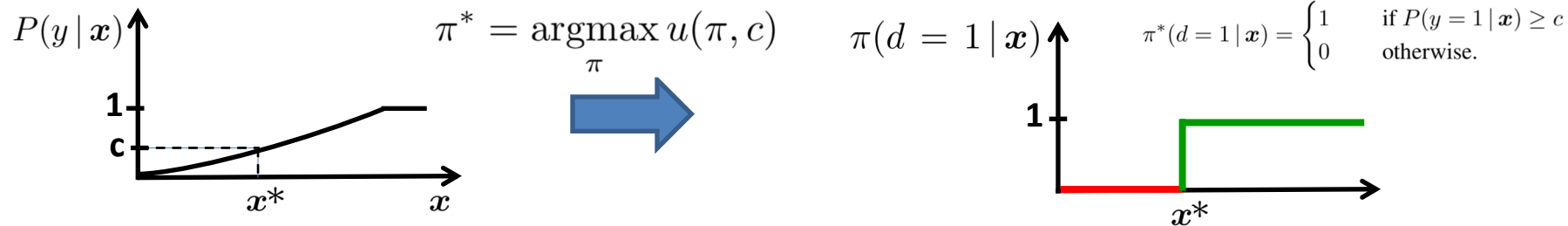


$$\pi^* = \underset{\pi}{\operatorname{argmax}} u(\pi, c)$$

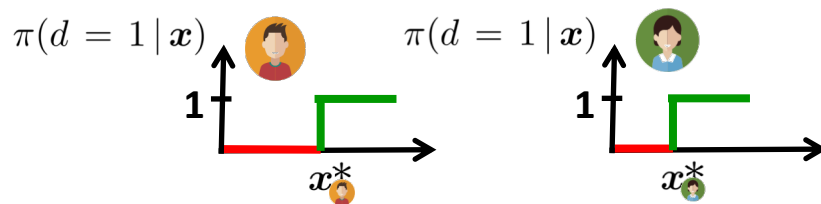


# Deterministic threshold rules

Under some **technical conditions**, deterministic threshold rules are optimal decision policies:



Under **fairness constraints**, we just need **two thresholds**:

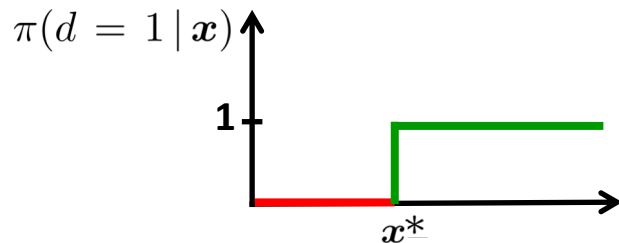


# Why are deterministic threshold rules optimal?

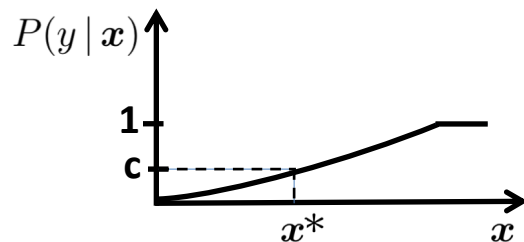
To realize **why deterministic threshold rules are optimal**, rewrite the utility of the policy as follows:

$$u(\pi, c) = \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x}), d \sim \pi(d | \mathbf{x})} [d(\mathbf{x}) \underbrace{(P(y = 1 | \mathbf{x}) - c)}_{\text{It is only positive if } P(y = 1 | \mathbf{x}) \geq c}]$$

It is only **positive** if  
 $P(y = 1 | \mathbf{x}) \geq c$



If  $P(y = 1 | x) \geq c$ ,  
make  $d = 1$ ,  
otherwise, make  $d = 0$



# So, are we done?

Let's look into the *technical conditions*

1. The **predictive model** is *perfect*

$$\rightarrow P_{\theta}(y \mid \mathbf{x}) = P(y \mid \mathbf{x})$$

2. The **label/feature distributions** and the **policy** are *independent*

$$\rightarrow P(y \mid \mathbf{x}, \pi) = P(y \mid \mathbf{x}) \quad P(\mathbf{x} \mid \pi) = P(\mathbf{x})$$

3. **Individuals** are *not strategic*

$\rightarrow$  Individuals do not seek to maximize their benefit  $b(\mathbf{x}, c)$

# So, are we done?

Let's look into the **technical conditions**

1. The predictive model is **perfect**

**In practice, these technical conditions  
are (often) violated.**

3. Individuals are **not strategic**

→ Individuals do not seek to maximize their benefit  $b(x, c)$



# Dealing with imperfect predictions

Assume the **predictive model**  $P_{\theta}(y | \mathbf{x}) = Q(y | \mathbf{x})$  trained using **historical data** is **imperfect**, i.e.,

$$Q(y | \mathbf{x}) = P(y | \mathbf{x}) + \varepsilon(y | \mathbf{x})$$

We will distinguish two different cases:

(a) **Historical data** suffers from ***selective labels***

[Lakkaraju et al., KDD 2017]

(b) **Historical data** is sampled from the **ground truth data distribution**

More common

Less common

# Historical data suffers from selective labels

**Historical data** is not **sampled** from the **ground truth distribution** but a **distribution induced** by a **previously deployed policy**

$$\underbrace{P_{\pi_0}(x, y)}_{\text{Data distribution induced by historical policy}} \propto P(y | x) \underbrace{\pi_0(d = 1 | x)}_{\text{Deployed historical policy}} P(x)$$

Data distribution **induced by**  
**historical policy**

Deployed historical policy

## Example

**Loan decisions:**


Historical data only contains individuals who received a loan in the past

The (induced) label/feature distribution & policy are **dependent!**

# Historical data suffers from selective labels

**Historical data** is not **sampled** from the **ground truth distribution** but a **distribution induced** by a **previously deployed policy**

$$P_{\pi_0}(\mathbf{x}, y) \propto P(y | \mathbf{x}) \pi_0(d = 1 | \mathbf{x}) P(\mathbf{x})$$

This creates a  feedback loop between human decisions and algorithmic decisions

Loan decisions:

Historical data only contains individuals who received a loan in the past

**dependent!**

# Are threshold rules provably suboptimal?

Take the **optimal policy** under the **original data distribution** and the **data distribution induced by the historical policy**:

$$Q^* \in \operatorname{argmax}_{Q \in \mathcal{Q}} \mathbb{E}_{\mathbf{x}, y \sim P} [\mathbf{1}[Q(y = 1 | \mathbf{x}) \geq c](y - 1)]$$

$$Q_0^* \in \operatorname{argmax}_{Q \in \mathcal{Q}} \mathbb{E}_{\mathbf{x}, y \sim P_{\pi_0}} [\mathbf{1}[Q(y = 1 | \mathbf{x}) \geq c](y - 1)]$$

# Are threshold rules provably suboptimal?

Take the **optimal policy** under the **original data distribution** and the **data distribution induced by the historical policy**:

$$Q^* \in \operatorname{argmax}_{Q \in \mathcal{Q}} \mathbb{E}_{\mathbf{x}, y \sim P} [\mathbf{1}[Q(y = 1 | \mathbf{x}) \geq c](y - 1)]$$

$$Q_0^* \in \operatorname{argmax}_{Q \in \mathcal{Q}} \mathbb{E}_{\mathbf{x}, y \sim P_{\pi_0}} [\mathbf{1}[Q(y = 1 | \mathbf{x}) \geq c](y - 1)]$$

**Proposition (negative result!).** If  $\pi_0 \neq \pi^*$  then

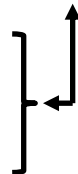
$$u(\pi_{Q_0^*}, c) < u(\pi_{Q^*}, c)$$

# In which class of policies lies the optimal decision policy?

It lies in the set of *exploring policies*.

$\pi_0(d = 1 \mid \mathbf{x}, s) > 0$  on any measurable set with positive probability under  $P$

A policy  $\pi$  is *exploring* iff the true distribution  $P$  is absolutely continuous with respect to  $P_\pi$



# In which class of policies lies the optimal decision policy?

It lies in the set of **exploring policies**.

$\pi_0(d = 1 \mid \mathbf{x}, s) > 0$  on any measurable set with positive probability under  $P$

A policy  $\pi$  is **exploring** iff the true distribution  $P$  is absolutely continuous with respect to  $P_\pi$

**Proposition (positive result!).** If  $\pi_0$  is an exploring policy,

$$u(\pi^*, c) = \sup_{\pi \in \Pi} \mathbb{E}_{\mathbf{x}, y \sim P_{\pi_0}} \left[ \frac{\pi(d = 1 \mid \mathbf{x})}{\pi_0(d = 1 \mid \mathbf{x})} (y - c) \right]$$

Set of exploring policies

Induced distribution!

# Not all exploring policies are (equally) acceptable



Consider a  
**lending scenario**

The following **decision policies** are **exploring**:

**Give loans to everyone,  $\pi_0(d | x) = 1$  for all  $x$**

**Gives loans to every individual at random,  $\pi_0(d = 1 | x) = 0.5$  for all  $x$**

**Who thinks a bank will do well under these policies? 😊**



# Learning exploring policies

1. **Deploy** an **initial exploring policy**  $\pi_0$ , which may be far from optimal for not too long.
2. Use **data gathered with this initial exploring policy** to fit a new **parameterized exploring policy**  $\pi_\theta$  using **SGA**, i.e.,

$$\theta_{i+1} = \theta_i + \alpha_i \underbrace{\nabla_{\theta} u(\pi_{\theta}, c)|_{\theta=\theta_i}}_{\text{Log-derivative and reweighting tricks}}$$

Log-derivative and reweighting tricks

Distribution induced by initial policy

$$\mathbb{E}_{\mathbf{x}, y \sim P_{\pi_0}, d \sim \pi_{\theta}} \left[ \frac{d(y-1)}{\pi(d=1|\mathbf{x})} \nabla_{\theta} \log \pi_{\theta}(d|\mathbf{x}) \right]$$

New policy

3. **Deploy & gather data** with  $\pi_\theta$  and fit a better exploring policy. Repeat.

# Learning exploring policies

1. **Deploy** an **initial exploring policy**  $\pi_0$ , which may be far from optimal for not too long.

2. Update  
new  
 $\theta_{i+1}$

**Learning to decide rather than  
learning to predict!**

Log-derivative and  
reweighting tricks

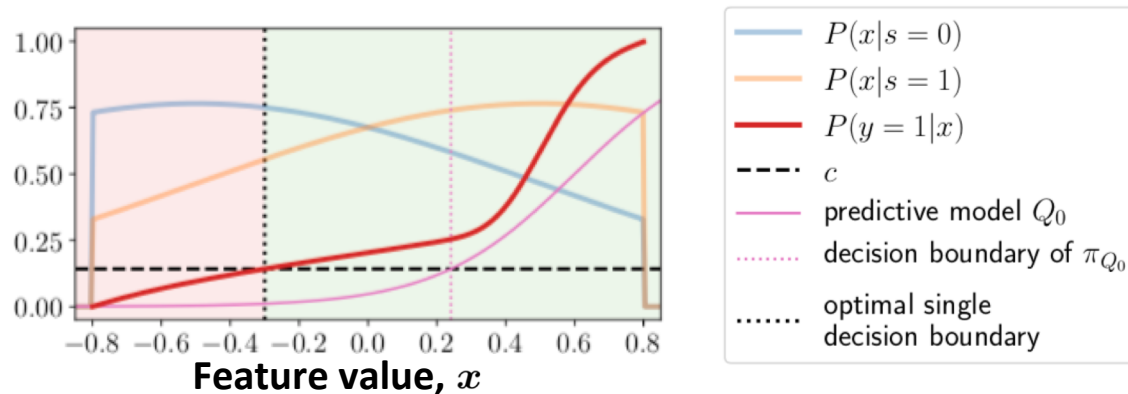
by initial policy

$$\mathbb{E}_{\mathbf{x}, y \sim P_{\pi_0}, d \sim \pi_\theta} \left[ \frac{d(y - 1)}{\pi(d = 1 | \mathbf{x})} \nabla_\theta \log \pi_\theta(d | \mathbf{x}) \right]$$

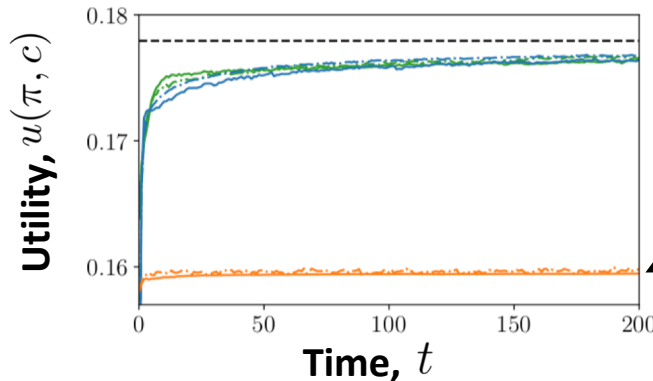
**New policy**

3. **Deploy & gather data** with  $\pi_\theta$  and fit a better exploring policy. Repeat.

# Example 1: strictly monotonic label distribution



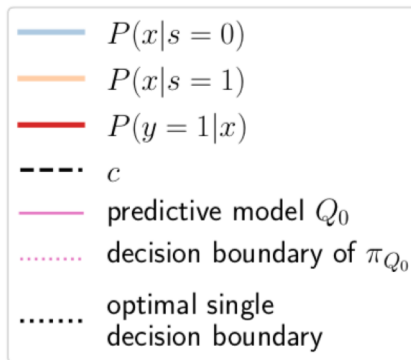
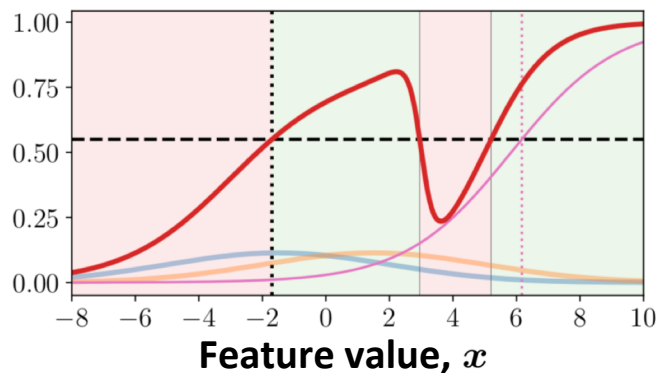
Higher is  
better



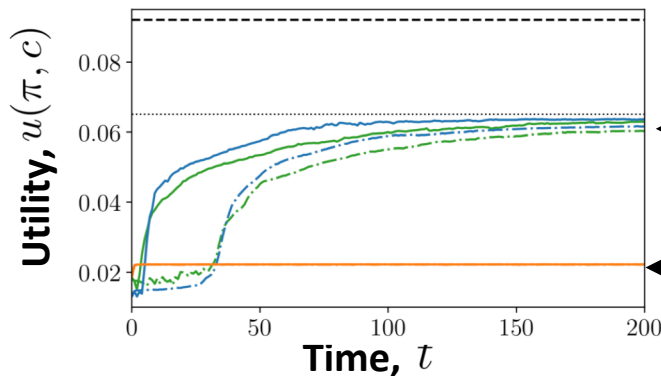
Parameterized exploring  
policies

Deterministic threshold rules

# Example 2: nonmonotonic label distribution



Higher is better

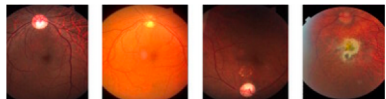


Parameterized exploring policies

Deterministic threshold rules

# Historical data is sampled from the ground truth distribution

There are **situations** where the **historical data** is sampled from the **ground truth distribution**.

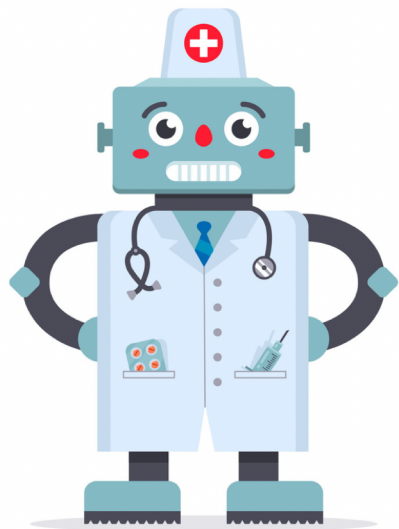


If a person has (or has not) a disease, this fact does not change after a medical diagnosis by a doctor

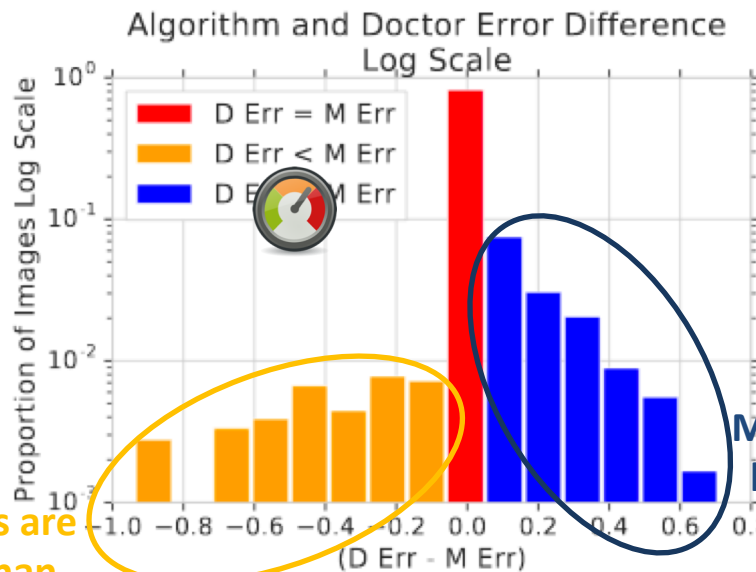


Then, given the latest deep ML model, can we just gather enough data to train a *perfect model*?

# Machines learning is sometimes worse than humans



On some instances, machine predictions are still worse than predictions made by human experts



Machines are worse than humans

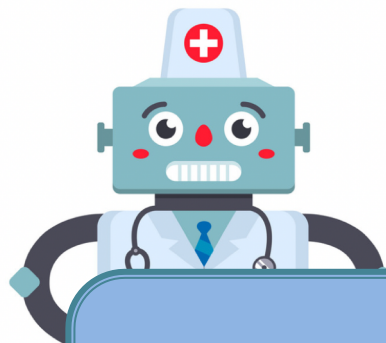


Task: estimating severity of diabetic retinopathy


Machines are better than humans

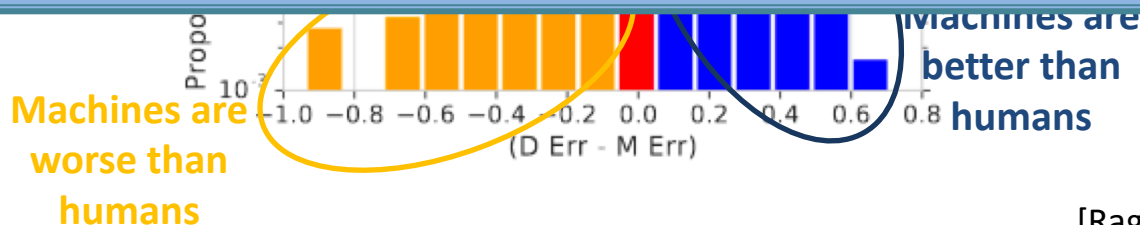
# Machines learning is sometimes worse than humans

On some instances, machine predictions are still worse than predictions made by human experts



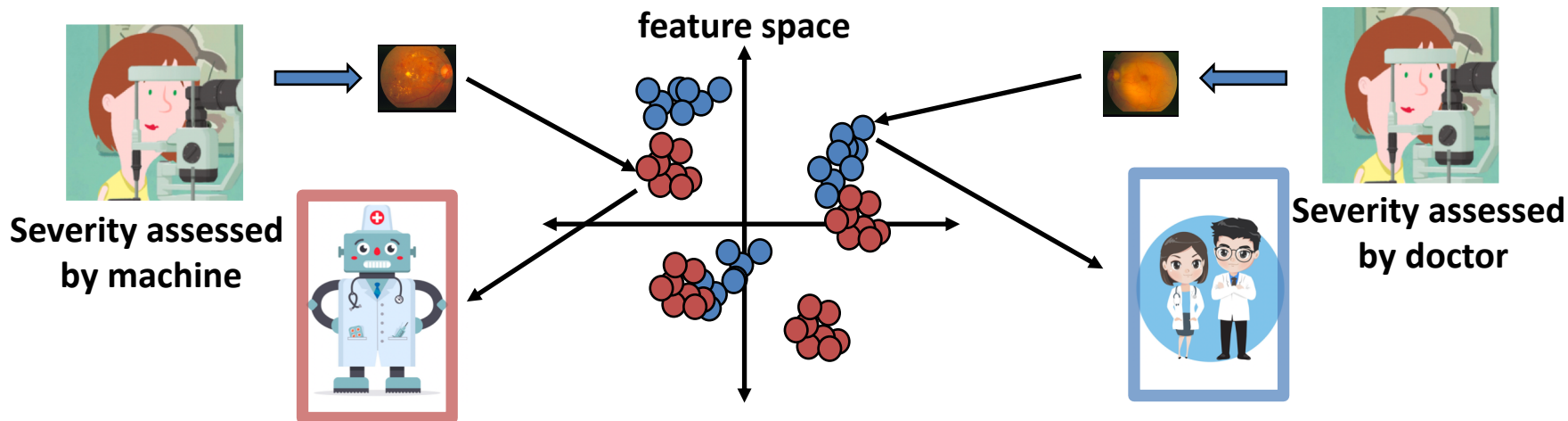
Algorithm and Doctor Error Difference  
Log Scale

Can we then  balance predictions between humans and machines?



# Machine learning for different automation levels

**Key idea:** develop machine learning models that are **optimized** to operate under **different automation levels**



They take **decisions** for a given **fraction of the instances** and leave the **remaining ones** to humans



# Optimizing the machine during training and test time

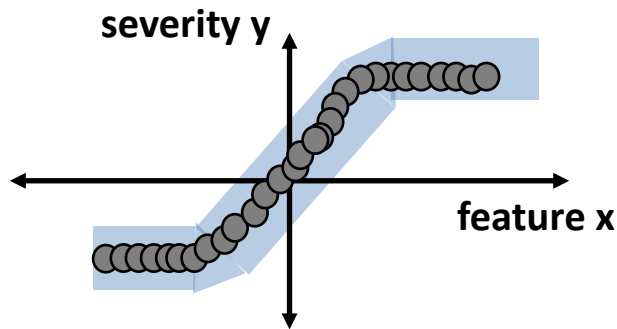
Key idea

**optimize** the **design** of the machine during training

# Optimizing the machine during training and test time

## Key idea

optimize the **design** of the machine during training

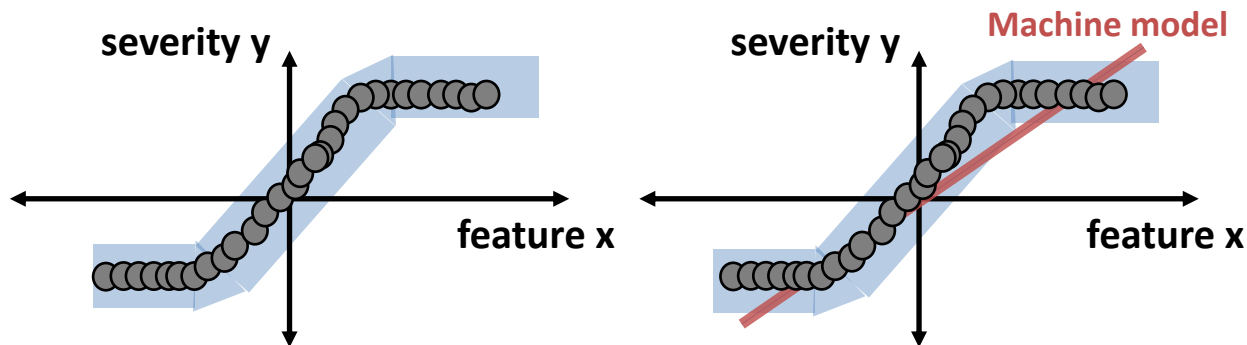


1. **The machine model** is a linear function
2. We can defer some samples to **humans**

# Optimizing the machine during training and test time

## Key idea

optimize the **design** of the machine during training



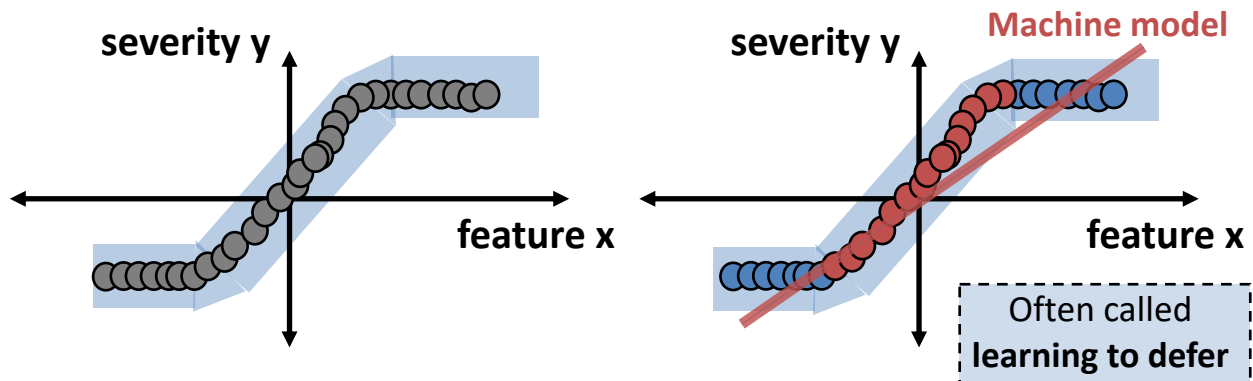
1. **The machine model** is a linear function
2. We can defer some samples to **humans**

**Machine model**  
is not optimized  
during training

# Optimizing the machine during training and test time

## Key idea

optimize the **design** of the machine during training



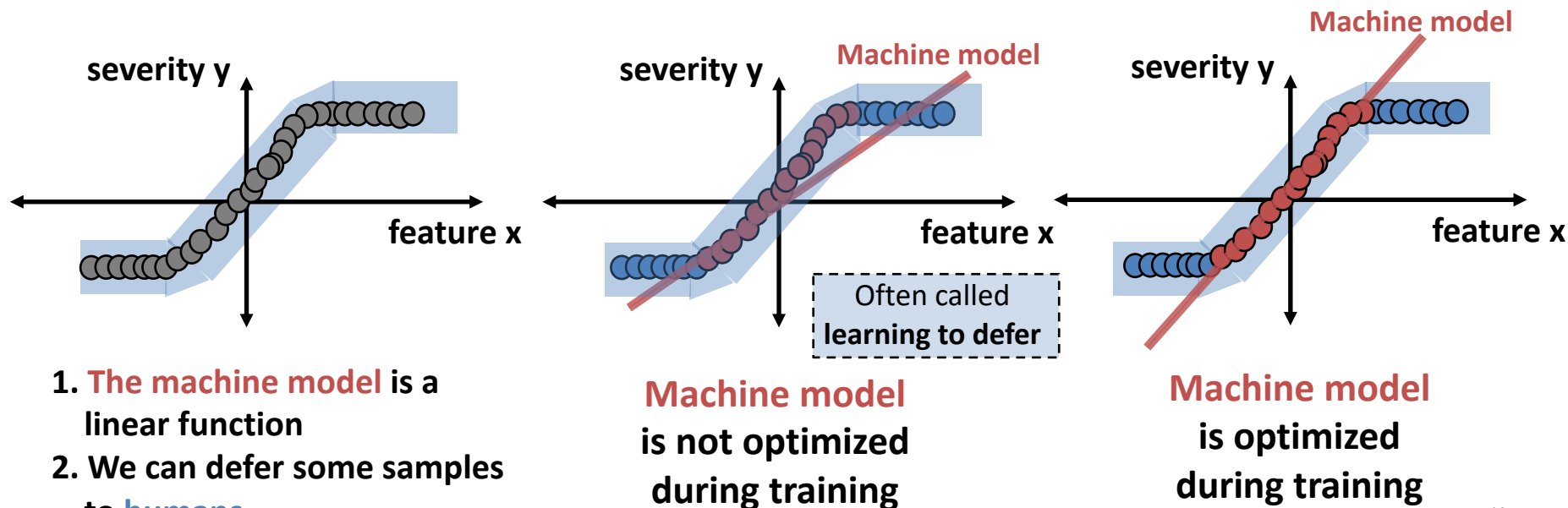
1. The **machine model** is a linear function
2. We can defer some samples to **humans**

**Machine model**  
is not optimized  
during training

# Optimizing the machine during training

## Key idea

optimize the **design** of the machine during training



1. The machine model is a linear function
2. We can defer some samples to humans

# Optimizing the machine during training

Key idea

optimize the design of the machine during training

Next, we will show how to design a ridge regression model optimized to operate under different automation levels

1. The machine model is a linear function
2. We can defer some samples to humans

Machine model  
is not optimized  
during training

Machine model  
is optimized  
during training

# Ridge regression, revisited

## Training

$$\underset{\mathbf{w}, \mathcal{S}}{\text{minimize}} \quad \sum_{i \in \mathcal{S}} \overbrace{c(\mathbf{x}_i, y_i)}^{\text{Human error per sample}} + \sum_{j \in \mathcal{S}^c} [(y_j - \mathbf{x}_j^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2]$$

subject to  $|\mathcal{S}| \leq n$

Annotations:

- $\sum_{i \in \mathcal{S}} c(\mathbf{x}_i, y_i)$ : Training samples assigned to **humans** (blue dots).
- $\sum_{j \in \mathcal{S}^c} [(y_j - \mathbf{x}_j^\top \mathbf{w})^2]$ : Training samples assigned to **machines** (red dots).
- $\mathbf{x}_j^\top \mathbf{w}$ : Machine model.
- $\lambda$ : Regularization parameter.
- $|\mathcal{S}| \leq n$ : Max. number of samples that can be assigned to humans.

# Ridge regression, revisited

## Training

minimize  $w, \mathcal{S}$

$$\sum_{i \in \mathcal{S}} \overbrace{c(\mathbf{x}_i, y_i)}^{\text{Human error per sample}} + \sum_{j \in \mathcal{S}^c} [(y_j - \mathbf{x}_j^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2]$$

subject to  $|\mathcal{S}| \leq n$

Annotations:

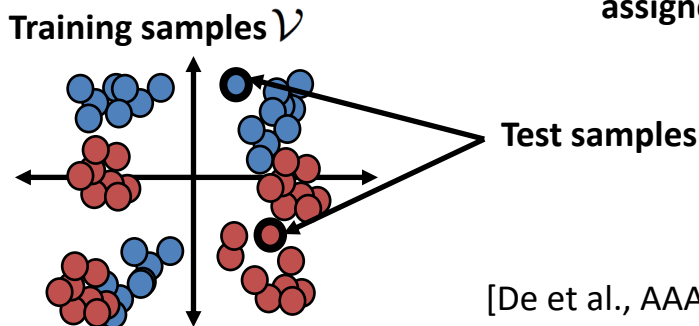
- Training samples assigned to **humans** (blue dots) →  $\sum_{i \in \mathcal{S}} c(\mathbf{x}_i, y_i)$
- Training samples assigned to **machines** (red dots) →  $\sum_{j \in \mathcal{S}^c} [(y_j - \mathbf{x}_j^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2]$
- Machine model →  $\mathbf{x}_j^\top \mathbf{w}$
- Regularization parameter →  $\lambda$
- Max. number of samples that can be assigned to humans →  $|\mathcal{S}| \leq n$

## Test

We assign a test sample with features  $\mathbf{x}$  to **humans** if

$$\operatorname{argmin}_{i \in \mathcal{V}} \|\mathbf{x}_i - \mathbf{x}\| \in \mathcal{S}^*$$

Training samples assigned to **humans** (blue dots) →  $\mathcal{S}^*$



[De et al., AAAI 2020]



# Ridge regression becomes a combinatorial problem

Given a fixed set  $\mathcal{S}$ , the **optimal machine model** is given by

$$\mathbf{w}^*(S) = (\lambda |\mathcal{S}^c| \mathbb{I} + \mathbf{X}_{\mathcal{S}^c} \mathbf{X}_{\mathcal{S}^c}^\top)^{-1} \mathbf{X}_{\mathcal{S}^c} \mathbf{y}_{\mathcal{S}^c}$$

# Ridge regression becomes a combinatorial problem

Given a fixed set  $\mathcal{S}$ , the **optimal machine model** is given by

$$\mathbf{w}^*(S) = (\lambda|\mathcal{S}^c|\mathbb{I} + \mathbf{X}_{\mathcal{S}^c}\mathbf{X}_{\mathcal{S}^c}^\top)^{-1} \mathbf{X}_{\mathcal{S}^c}\mathbf{y}_{\mathcal{S}^c}$$

Then, we can **rewrite the ridge regression problem** as a purely **combinatorial maximization problem**

$$\underset{\mathcal{S}}{\text{maximize}} \quad \underbrace{-\log \ell(\mathcal{S})}_{\downarrow} \quad \text{subject to} \quad |\mathcal{S}| \leq n$$

$$\sum_{i \in \mathcal{S}} c(\mathbf{x}_i, y_i) + \mathbf{y}_{\mathcal{S}^c}^\top \mathbf{y}_{\mathcal{S}^c} - \mathbf{y}_{\mathcal{S}^c}^\top \mathbf{X}_{\mathcal{S}^c}^\top (\lambda|\mathcal{S}^c|\mathbb{I} + \mathbf{X}_{\mathcal{S}^c}\mathbf{X}_{\mathcal{S}^c}^\top)^{-1} \mathbf{X}_{\mathcal{S}^c}\mathbf{y}_{\mathcal{S}^c}$$

# Ridge regression under human assistance is hard

Finding the solution to

$$\underset{\mathcal{S}}{\text{maximize}} \quad -\log \ell(\mathcal{S}) \quad \text{subject to} \quad |\mathcal{S}| \leq n$$

is a **NP-hard problem**

# Ridge regression under human assistance is hard

## Finding the solution to

$$\underset{\mathcal{S}}{\text{maximize}} \quad -\log \ell(\mathcal{S}) \quad \text{subject to} \quad |\mathcal{S}| \leq n$$

## is a **NP-hard** problem

### Proof sketch

Assume  $c(\mathbf{x}_i, y_i) = 0$ ,  $\lambda = 0$  and  $\mathbf{y} = \mathbf{X}^\top \mathbf{w}^* + \mathbf{b}^*$

k-sparse noise vector



Then, **the problem** can be viewed as the **robust least square (RLSR) problem**, which **has been shown to be NP-hard**:

$$\underset{\mathbf{w}, \mathcal{S}}{\text{minimize}} \quad \sum_{i \in \mathcal{S}} (y_i - \mathbf{x}_i^\top \mathbf{w})^2 \quad \text{subject to} \quad |\mathcal{S}| = |\mathcal{V}| - n$$

# A simple greedy algorithm

The **greedy algorithm** proceeds **iteratively**.

At each iteration, it assigns to a human the sample in the training set that provides the **largest marginal gain**, i.e.,

$$k^* \leftarrow \operatorname{argmax}_{k \in \underbrace{\mathcal{V} \setminus \mathcal{S}}_{\text{Points not yet assigned to humans}}} \overbrace{-\log \ell(\mathcal{S} \cup k) + \log \ell(\mathcal{S})}$$
$$\mathcal{S} \leftarrow \mathcal{S} \cup \{k^*\}$$

Does this simple greedy algorithm has approximation guarantees? 😊

# The greedy algorithm has approximation guarantees

The function  $-\log \ell(\mathcal{S})$  satisfies an **approximate notion of submodularity**

$$-\log \ell(\mathcal{S} \cup k) + \log \ell(\mathcal{S}) \geq (1 - \alpha) [-\log \ell(\mathcal{T} \cup k) + \log \ell(\mathcal{T})]$$

for all  $\mathcal{S} \subseteq \mathcal{T} \subset \mathcal{V}$

where  $\alpha \geq \alpha^*$  is the **generalized curvature**



Data dependent constant

We can conclude that the greedy algorithm will find a set  $\mathcal{S}$  such that

$$-\log \ell(\mathcal{S}) \geq \left(1 + \frac{1}{1 - \alpha}\right)^{-1} OPT$$

Optimal value

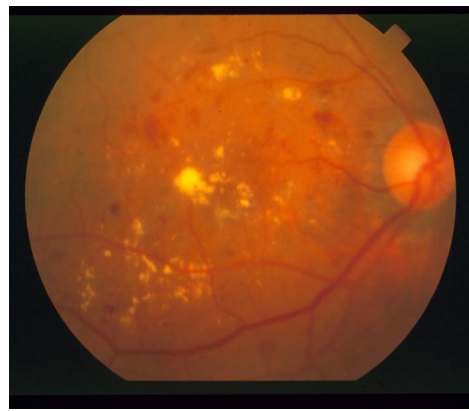
# What samples are outsourced?

**Drusen disease is characterized by pathological yellow spots...  
...however, both images are given a score of severity zero**



**Easy sample**

It is assigned to the machine

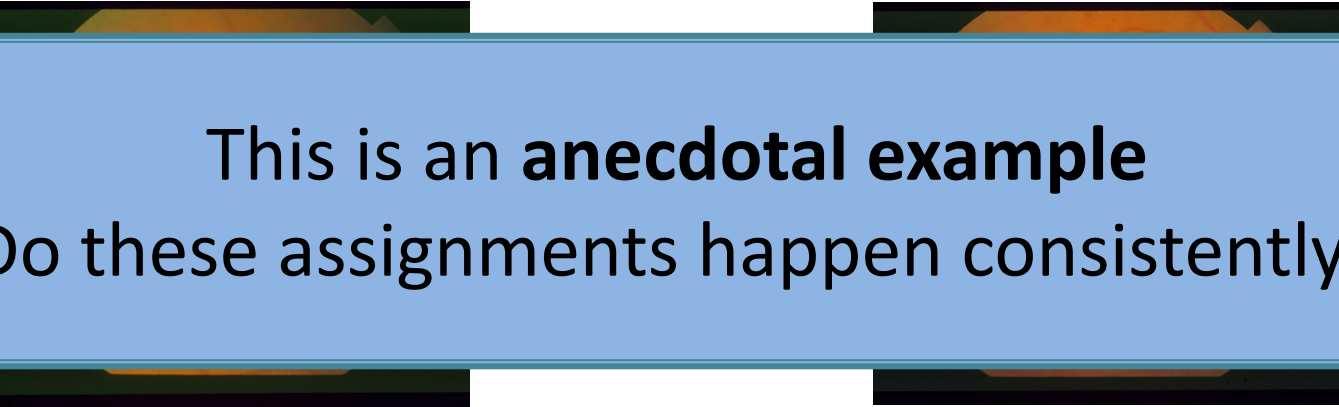


**Difficult sample**

It is assigned to the human

# What samples are outsourced?

**Drusen disease** is characterized by **pathological yellow spots**...  
...however, **both images** are given a **score of severity zero**



This is an **anecdotal example**  
Do these assignments happen consistently?

**Easy sample**

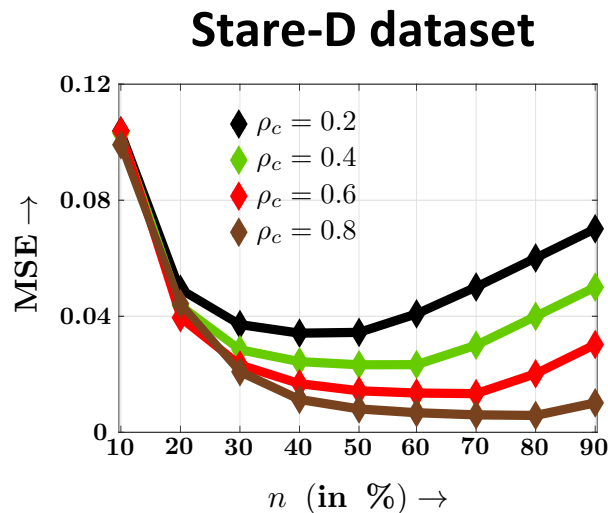
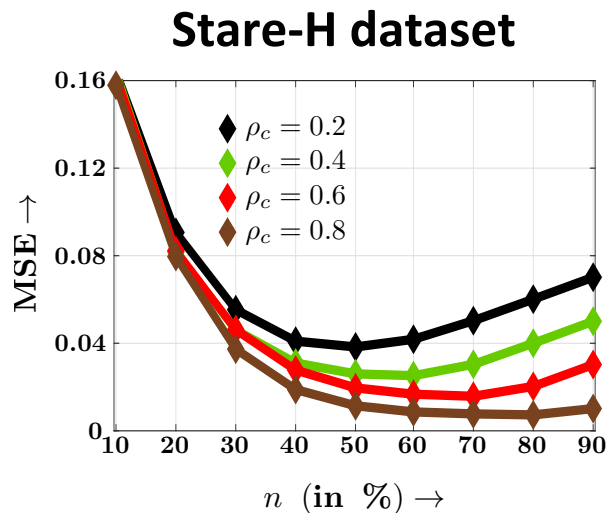
our algorithm assigns it to machine

**Difficult sample**

our algorithm assigns it to human



# The greedy algorithm spots samples where humans are accurate



**$\rho_c$ : fraction of samples with low human error**

As long as there are **samples that humans can predict with low error**, the **greedy algorithm outsources them to humans** and the **performance improves**

# Strategic behavior and transparency

Until now, we have **assumed individuals are not strategic**:

→ Individuals do not seek to maximize their **benefit**

$$b(\mathbf{x}, c) = \mathbb{E}_{d \sim \pi(d | \mathbf{x})} [f(d(\mathbf{x}), c)]$$

However, this **assumption** is **in conflict** with the **increasing pressure to be transparent** about **policies, models and features**



→ Individuals may use **knowledge**, gained by **transparency**, to **invest effort strategically** to **maximize** their **chances** of **receiving a beneficial decision**.

# Strategic behavior and transparency

Until now, we have **assumed individuals are not strategic**:

→ **Individuals do not seek to maximize their benefit**

$$U_i(x) = \pi_i(x) \quad [U_i(x) = \pi_i(x)]$$

How  
the  
ab

Can we **design transparent ML models**  
that **account for strategic behavior**?

→ **Individuals may use knowledge, gained by transparency, to invest effort strategically to maximize their chances of receiving a beneficial decision.**

# Transparency on predictive models and policies

We can be **transparent** about:

## → **Predictive models**

**Goal:** develop **accurate predictive models** under strategic behavior.  
Most work view strategic behavior as **gaming**.

[Brückner et al., JMLR 2012; Hardt et al., NIPS 2016; Dong et al., EC 2018; Hu et al., WWW 2019]

## → **Policies**

**Goal:** design **policies that maximize utility** under strategic behavior.  
Most work view strategic behavior as **self-improvement**.

[Kleinberg & Raghavan, EC 2019; Khajehnejad et al., Arxiv 2019]

# Transparency on predictive models and policies

We can be **transparent** about:

## → **Predictive models**

**Goal:** develop **accurate predictive models** under strategic behavior.  
Most work view strategic behavior as **gaming**.

[Brückner et al., JMLR 2012; Hardt et al., NIPS 2016; Dong et al., EC 2018; Hu et al., WWW 2019]

## → **Policies**

**Goal:** design **policies that maximize utility** under strategic behavior.  
Most work view strategic behavior as **self-improvement**.

[Kleinberg & Raghavan, EC 2019; Khajehnejad et al., Arxiv 2019]

# Causal vs non causal features

Individual's strategic behavior as **gaming** or **self-improvement**?

Individuals **invest effort** on  
changing **noncausal features**

$P(x)$  changes

$P(y | x)$  changes



**Goodhart's law**

**“When a measure becomes a target, it ceases  
to be a good measure”**

Individuals **invest effort** on  
changing **causal features**

$P(x)$  changes

$P(y | x)$  does not change

# Causal vs non causal features

Individual's strategic behavior as **gaming** or **self-improvement**?

Individuals **invest effort** on  
changing **noncausal features**

$P(x)$  changes

$P(y | x)$  changes



**Goodhart's law**

**“When a measure becomes a target, it ceases  
to be a good measure”**

Individuals **invest effort** on  
changing **causal features**

$P(x)$  changes

$P(y | x)$  does not change

# Example 1: car insurance decisions



**1. Insurance company reveals it uses the number of speeding tickets to decide the insurance premium**

**2. Drivers may drive more carefully to pay a lower price**



**3. This will likely make them better drivers**

**causal feature  
self-improvement**



# Example 2: loan decisions

1. A bank reveals it uses **credit card debt** to decide **loans interest rates**



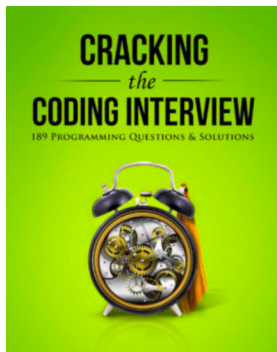
2. Applicants may **avoid credit card debt** overall to **pay less interest**

3. This will **improve** their **financial situation**

**causal feature  
self-improvement**

# Example 3: hiring decisions

1. A software company publishes the coding exercises it uses during recruiting



2. Applicants just practice only those coding exercises

3. This will not necessarily make them better employees

**noncausal feature  
gaming**

# Full transparency: a Stackelberg game

## Stackelberg game-theoretic formulation

The **decision maker publishes the decision policy  $\pi$  before individuals (best-)respond.**

For each **individual** with initial set of **features  $\mathbf{x}_i$** , her **best response** is:

$$\mathbf{x}_j = \operatorname{argmax}_{k \in [m]} \underbrace{b(\mathbf{x}_k, c)}_{\text{Benefit individual obtains for having features } \mathbf{x}_k} - \underbrace{c(\mathbf{x}_i, \mathbf{x}_k)}_{\text{Cost individual pays for changing from } \mathbf{x}_i \text{ to } \mathbf{x}_k}$$

We assume

$$\begin{aligned} b(\mathbf{x}, c) &= \mathbb{E}_{d \sim \pi(d | \mathbf{x})} [f(d(\mathbf{x}), c)] \\ &= \pi(\mathbf{x}) \end{aligned}$$



**Benefit individual obtains for having features  $\mathbf{x}_k$**

**Cost individual pays for changing from  $\mathbf{x}_i$  to  $\mathbf{x}_k$**

# From individual to population best response

## Individual best response

$$\mathbf{x}_j = \operatorname{argmax}_{k \in [m]} b(\mathbf{x}_k, c) - c(\mathbf{x}_i, \mathbf{x}_k)$$



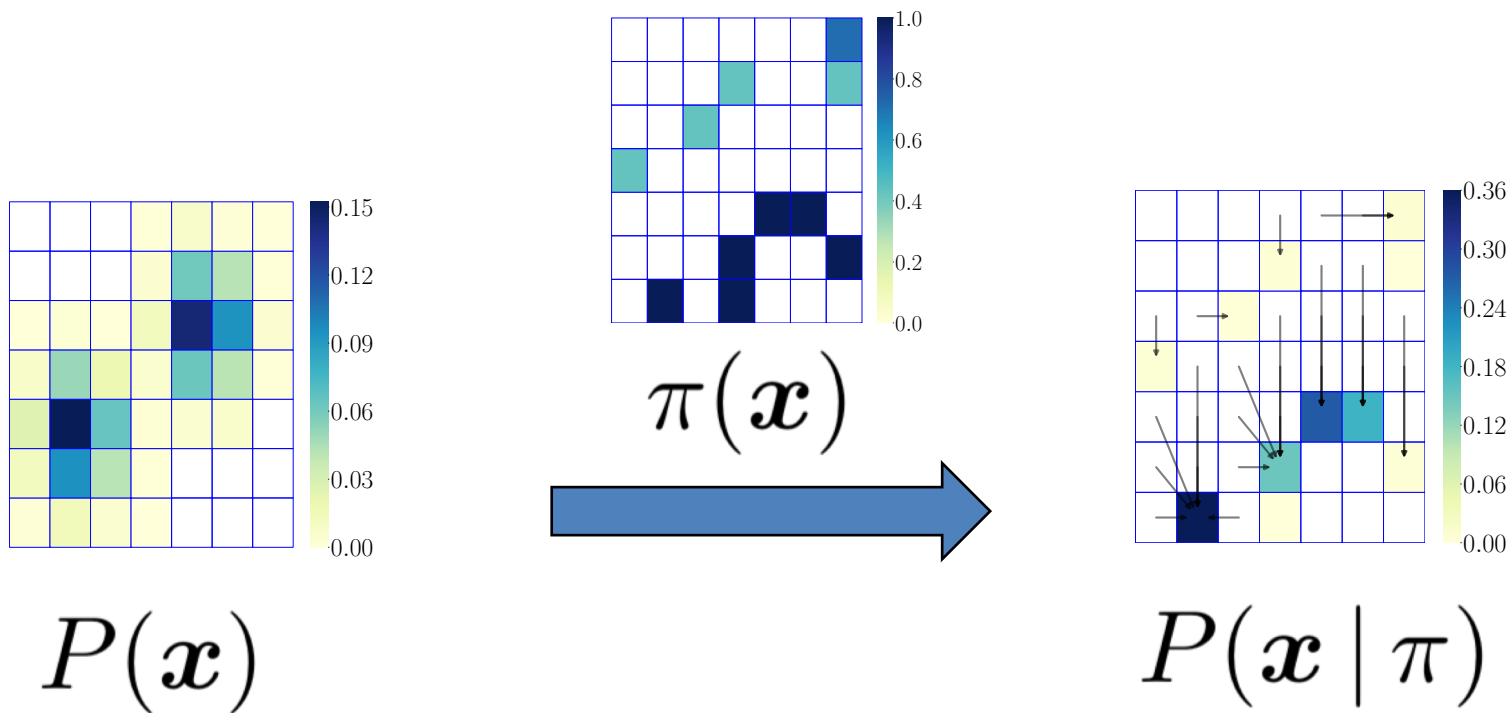
(Optimal) transportation of mass

## Induced distribution

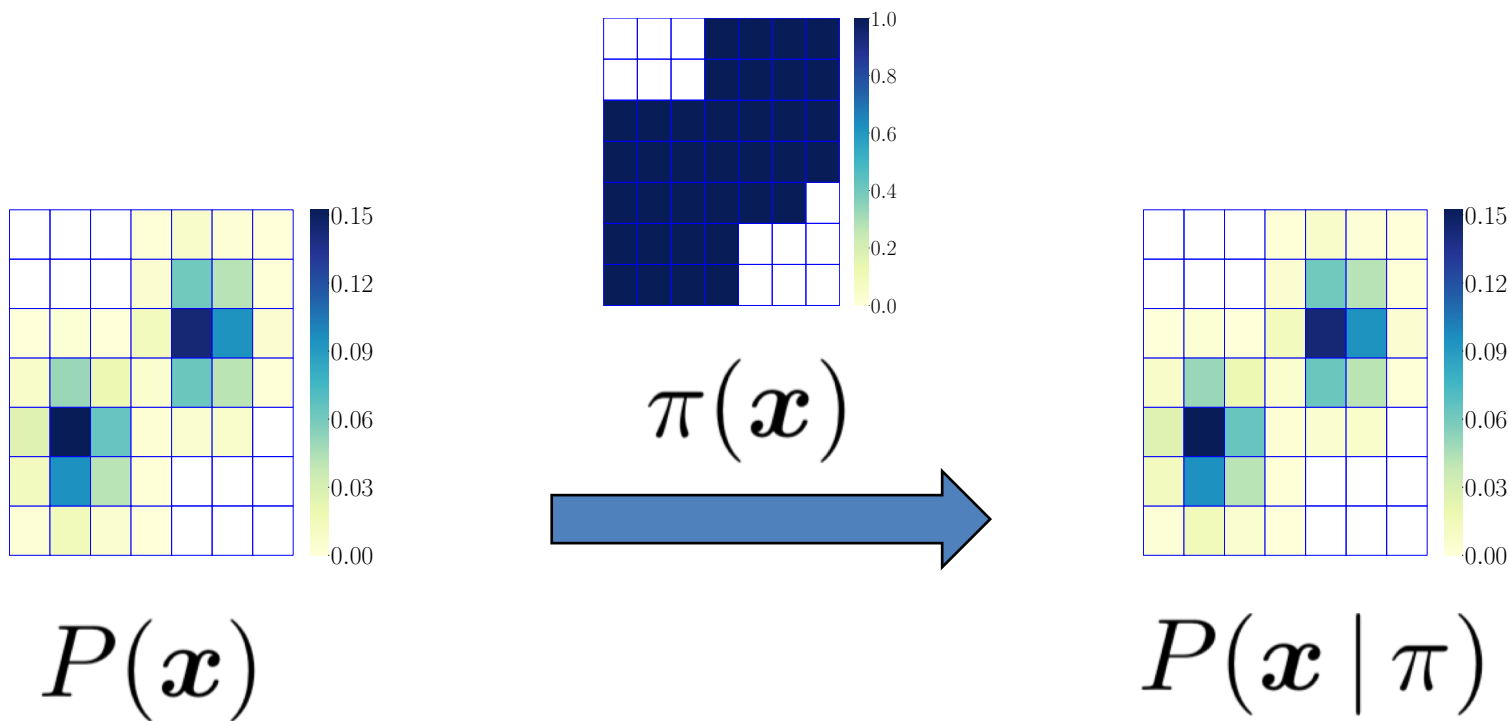
$$P(\mathbf{x}_j | \pi) = \sum_{i \in [m]} P(\mathbf{x}_i) \mathbb{I}(\mathbf{x}_j = \operatorname{argmax}_{k \in [m]} b(\mathbf{x}_k, c) - c(\mathbf{x}_i, \mathbf{x}_k))$$

└──┘  
Flow between  $P(\mathbf{x}_i)$  and  $P(\mathbf{x}_j | \pi)$

# Example 1: original and induced distributions



## Example 2: original and induced distributions




# Finding optimal decisions is hard

## Finding the solution to

$$\begin{aligned}\pi^* &= \operatorname{argmax}_{\pi} u(\pi, c) \\ &= \operatorname{argmax}_{\pi} \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x} \mid \pi), d \sim \pi(d \mid \mathbf{x})} [P(y = 1 \mid \mathbf{x}) d(\mathbf{x}) - c d(\mathbf{x})]\end{aligned}$$

That makes it hard




is a **NP-hard problem**

## Proof idea

Using a **reduction** to the **Boolean satisfiability (SAT) problem** [Karp, 1972]

# Optimal decisions may be stochastic

The **NP-hardness** result implies that **threshold rules** are **not always optimal**.


$$\pi^*(d = 1 | \mathbf{x}) = \begin{cases} 1 & \text{if } P(y = 1 | \mathbf{x}) \geq c \\ 0 & \text{otherwise.} \end{cases}$$

There are many **scenarios** in which the **optimal decision policies** are **not deterministic**. For example:

$$P(\mathbf{x}) = 0.1 \mathbb{I}(x = 1) + 0.4 \mathbb{I}(x = 2) + 0.5 \mathbb{I}(x = 3)$$

$$P(y = 1 | \mathbf{x}) = 1.0 \mathbb{I}(x = 1) + 0.7 \mathbb{I}(x = 2) + 0.4 \mathbb{I}(x = 3)$$

$$c(\mathbf{x}_i, \mathbf{x}_j) = \begin{bmatrix} 0.0 & 0.0 & 0.0 \\ 0.3 & 0.0 & 0.0 \\ 1.2 & 0.3 & 0.0 \end{bmatrix}$$

$$c = 0.1$$

**Non-strategic:**  $\pi^*(d | \mathbf{x}) = 1$  for all  $\mathbf{x}$

**Strategic:**  $\pi^*(d = 1 | \mathbf{x} = 1) = 1$   
 $\pi^*(d = 1 | \mathbf{x} = 2) = 0.7$

$$\pi^*(d = 1 | \mathbf{x} = 3) = 0$$



# Highest outcome and negative outcomes

Given any instance of the **utility maximization problem under strategic behavior**, it easy to realize that:

$\pi^*(\mathbf{x}_1) = 1$  where  $\mathbf{x}_1$  is the feature value with highest outcome  $P(y | \mathbf{x}_1) \geq c$

└→ The ***best individuals*** always receive a **beneficial decision**

$\pi^*(\mathbf{x}_i) = 0$  for all  $\mathbf{x}_i$  such that  $P(y | \mathbf{x}_i) < c$

└→ Always **decide negatively** about **individuals** providing ***negative utility***

# Highest outcome and negative outcomes

Given any instance of the **utility maximization problem under strategic behavior**, it easy to realize that:

$\pi^*(x_1) = 1$  where  $x_1$  is the feature value with highest

What about **individuals** in the **middle range** providing **positive utility**?

└ Always **decide negatively** about **individuals** providing ***negative utility***

# Outcome monotonic costs

We can **further characterize a family of optimal policies** if the **cost individuals pay to change features** satisfies **natural property, outcome monotonicity, i.e.,**

[Improving an individual's outcome requires increasing amount of effort]

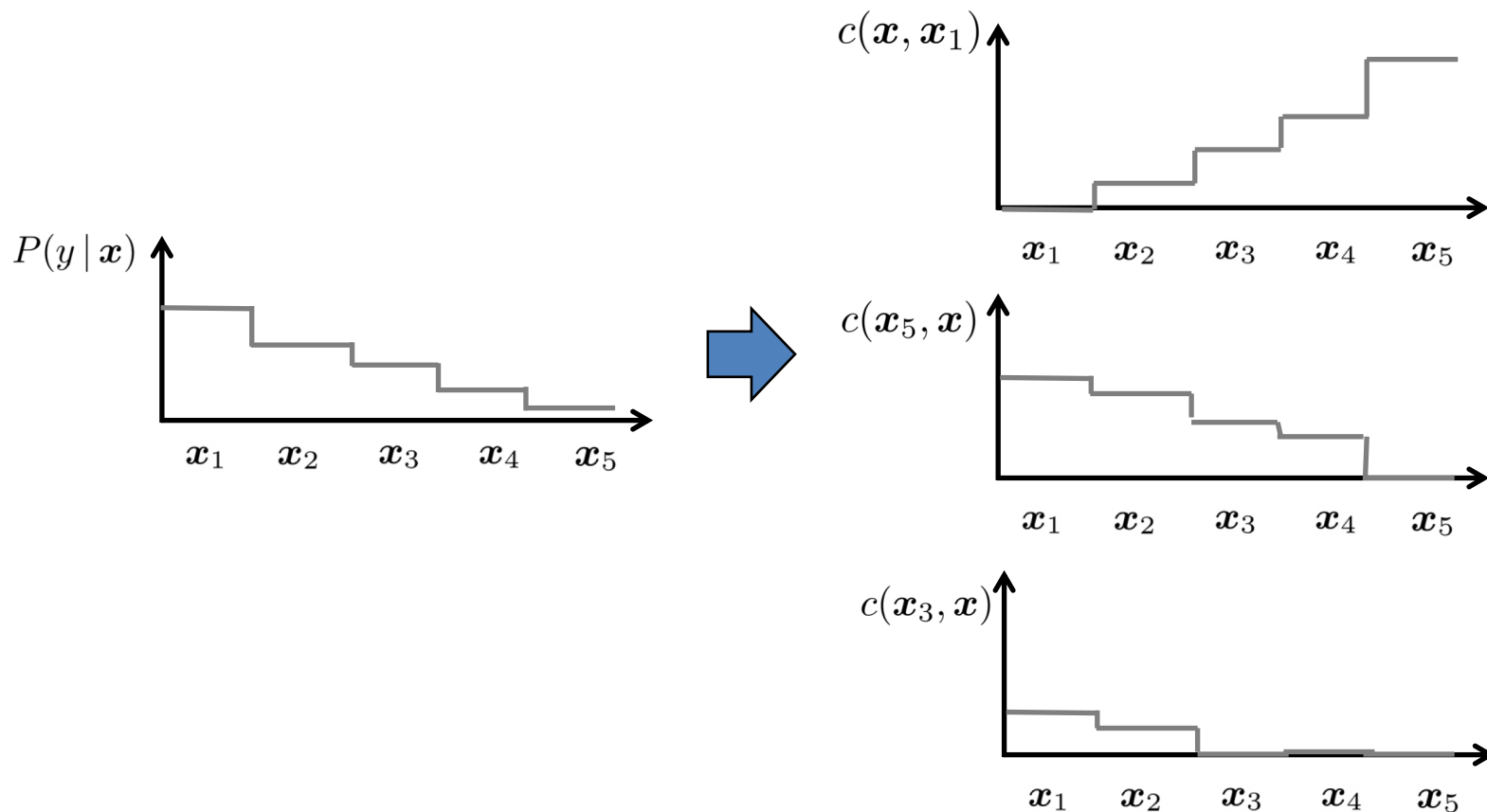
$$P(y = 1 | \mathbf{x}_i) < P(y = 1 | \mathbf{x}_j) < P(y = 1 | \mathbf{x}_k) \Leftrightarrow c(\mathbf{x}_i, \mathbf{x}_j) < c(\mathbf{x}_i, \mathbf{x}_k)$$

$$P(y = 1 | \mathbf{x}_i) > P(y = 1 | \mathbf{x}_j) > P(y = 1 | \mathbf{x}_k) \Leftrightarrow c(\mathbf{x}_j, \mathbf{x}_i) < c(\mathbf{x}_k, \mathbf{x}_i)$$

[Worsening an individual's outcome requires no effort]

$$P(y = 1 | \mathbf{x}_i) > P(y = 1 | \mathbf{x}_j) \Leftrightarrow c(\mathbf{x}_i, \mathbf{x}_j) = 0$$

# Example: outcome monotonic costs



# Outcome monotonic policies

**Proposition (positive result!).** If costs are **outcome monotonic**, there exists an **outcome monotonic policy** that is **optimal in terms of utility**.

An **outcome monotonic policy** satisfies that:

$$P(y = 1 \mid \mathbf{x}_i) < P(y = 1 \mid \mathbf{x}_j) \Leftrightarrow \pi(\mathbf{x}_i) < \pi(\mathbf{x}_j)$$

***Better individuals*** are more likely to receive a **beneficial decision**

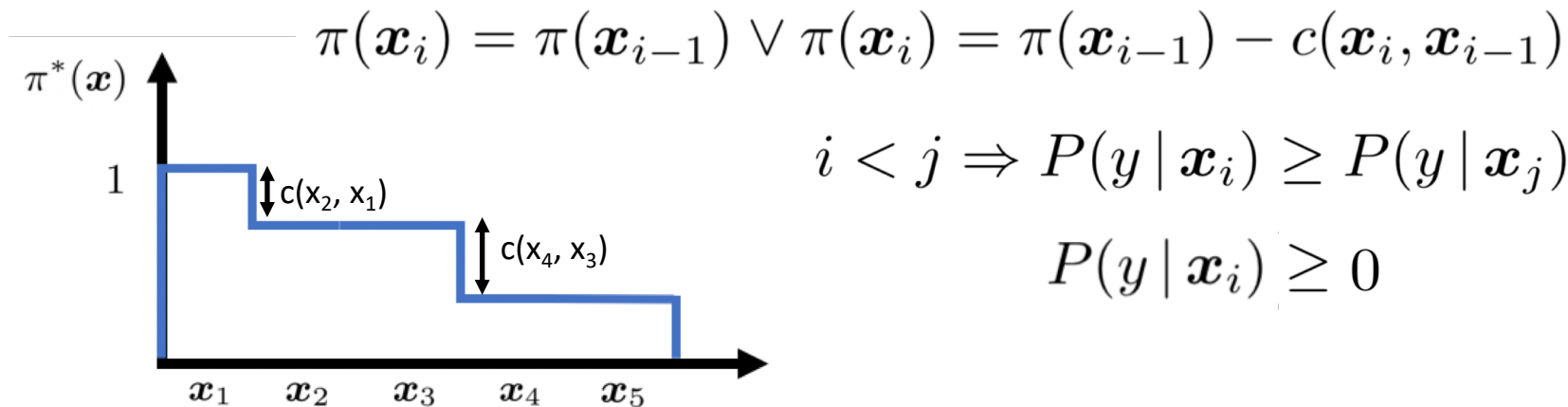


**Fair  
property**

# Outcome monotonic binary policies (I)

$$c(\mathbf{x}_i, \mathbf{x}_j) = c(\mathbf{x}_i, \mathbf{x}_k) + c(\mathbf{x}_k, \mathbf{x}_j) \longleftrightarrow \text{Unidimensional features}$$

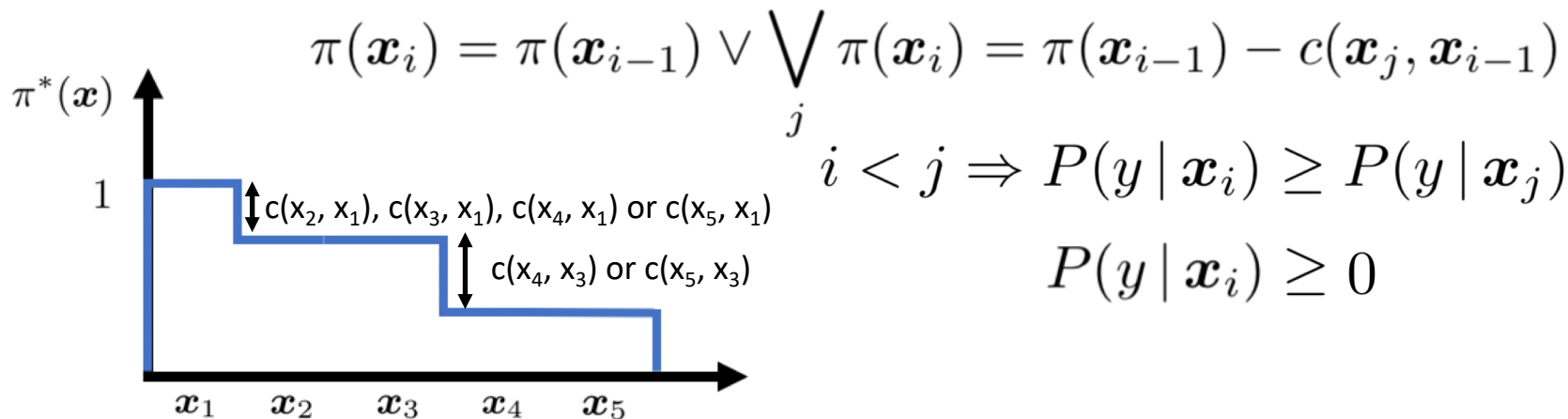
If costs are **additive**, there exists an optimal outcome monotonic “binary” policy that satisfies that:



# Outcome monotonic binary policies (II)

$$c(\mathbf{x}_i, \mathbf{x}_j) + c(\mathbf{x}_j, \mathbf{x}_k) \geq c(\mathbf{x}_i, \mathbf{x}_k) \longleftrightarrow \text{Multidimensional features}$$

If costs are **subadditive**, there exists an optimal outcome monotonic “binary” policy that satisfies that:



# An iterative algorithm for general costs

```
1:  $\pi \leftarrow \text{INITIALIZEPOLICY}$ 
2: repeat
3:    $\pi' \leftarrow \pi$ 
4:   for  $i = 1, \dots, m$  do
5:      $\pi(x_i) \leftarrow \text{SOLVE}(i, \pi, C, P, Q) \longrightarrow$ 
6:   end for
7: until  $\pi = \pi'$ 
8: Return  $\pi', u(\pi', c)$ 
```

$[c(x_i, x_j)]$      $[P(y | x_i)]$

$\uparrow$              $\uparrow$

$\downarrow$

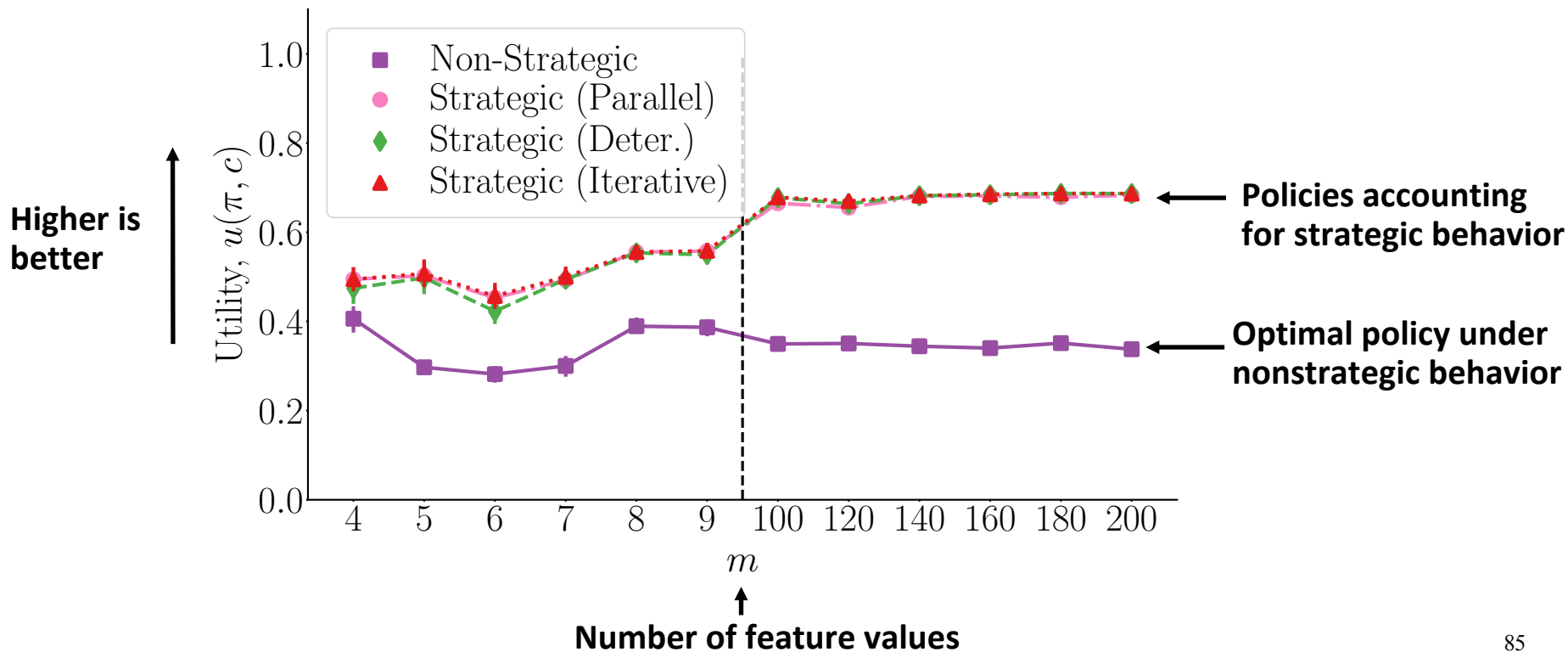
$[P(x_i)]$

**Tractable problem**  
**Fix all  $\pi(x_k)$  with  $x_k \neq x_i$**   
**and find best  $\pi(x_i)$**

The iterative algorithm is guaranteed to terminate in polynomial time and find a locally optimal policy



# Example: utility under strategic behavior



# Beyond full transparency

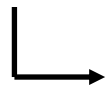
We have **assumed the decision maker publishes the entire policy.**

**In practice, it will reveal only part of the policy to each individual. For example:**



## **Counterfactual explanations**

[Wachter et al., Harvard JL and Tech 2017; Ustun et al., FAT\* 2019]



**Reveal one example of feature value “close” to the original feature value, which would lead to a beneficial decision**

# Beyond full transparency

We have assumed the decision maker publishes the entire policy.

In practice, it will reveal only part of the policy to each individual.

**Open problem: find counterfactual explanations under strategic behavior**

- Reveal one example of feature value close to the original feature value, which would lead to a beneficial decision

# Summary of this talk

Under a simple problem setting, we have learned about a few machine learning models and methods to:



Account for the feedback loop between algorithmic and human decisions



Balance decisions between human and algorithmic decisions



Account for strategic human decisions

**Disclaimer.** This was a *biased view*. These are emerging topics and there is still a lot of open problems and research directions! Join us 😊



# Thanks!

**more at [learning.mpi-sws.org](https://learning.mpi-sws.org)**