



# The Risks of Recourse in Binary Classification

2023-10-27

# The Risk of Recourse in Binary Classification

- ▶ Preprint on ArXiv (2306.00497)
- ▶ All work presented was created in collaboration with:



Dr. Tim van Erven



Dr. Damien Garreau

# Programme of today

- Introduction to the problem
- Modelling the setting
- Optimal classifiers
- Near Optimal classifiers
- Strategising
- Conclusion

# Introduction to the problem

## With a peculiar example

# Leading example

2 parties:

► Credit Loan Applicant (A)



► Credit Loan Provider (P)



Loan application process:

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

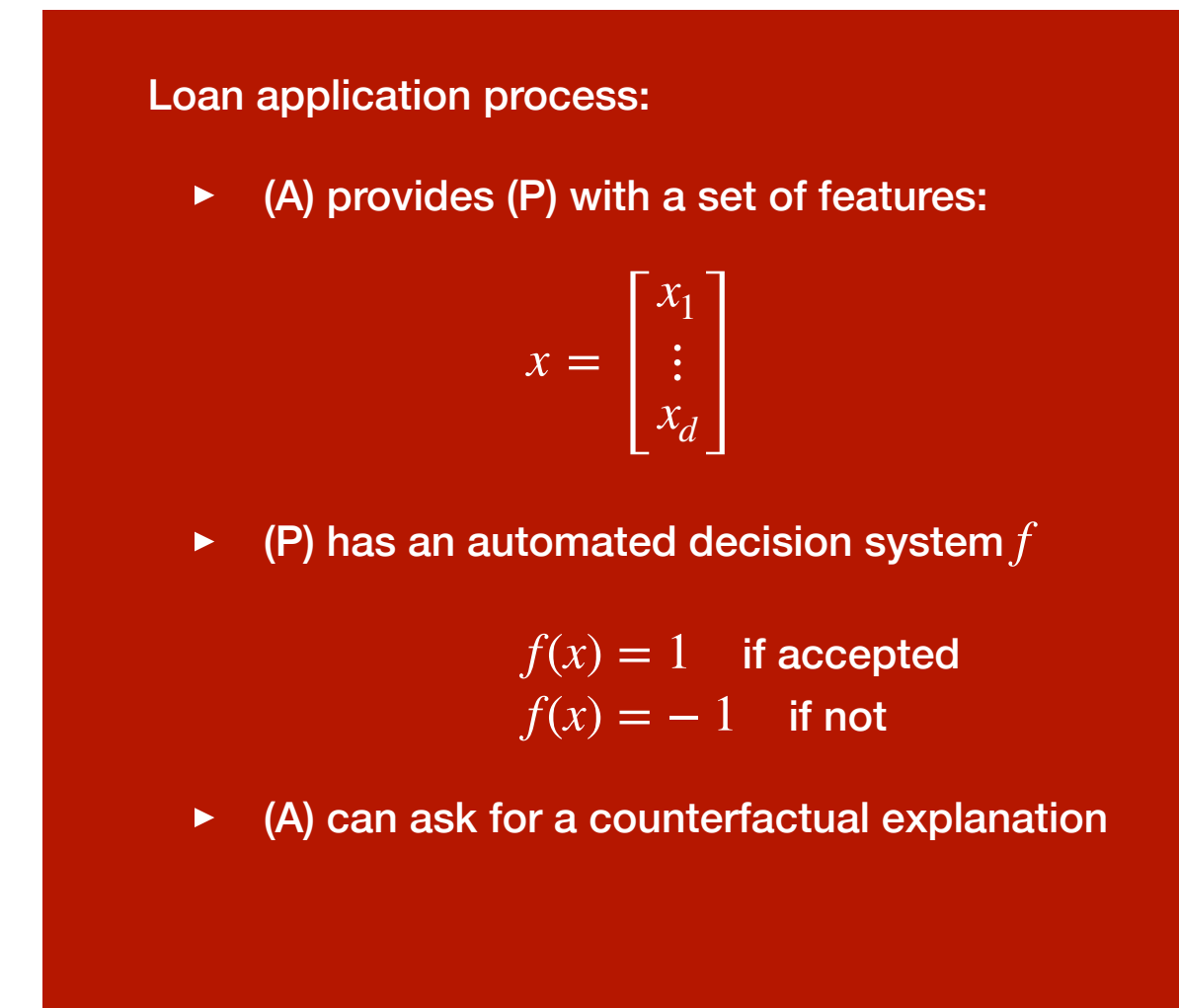
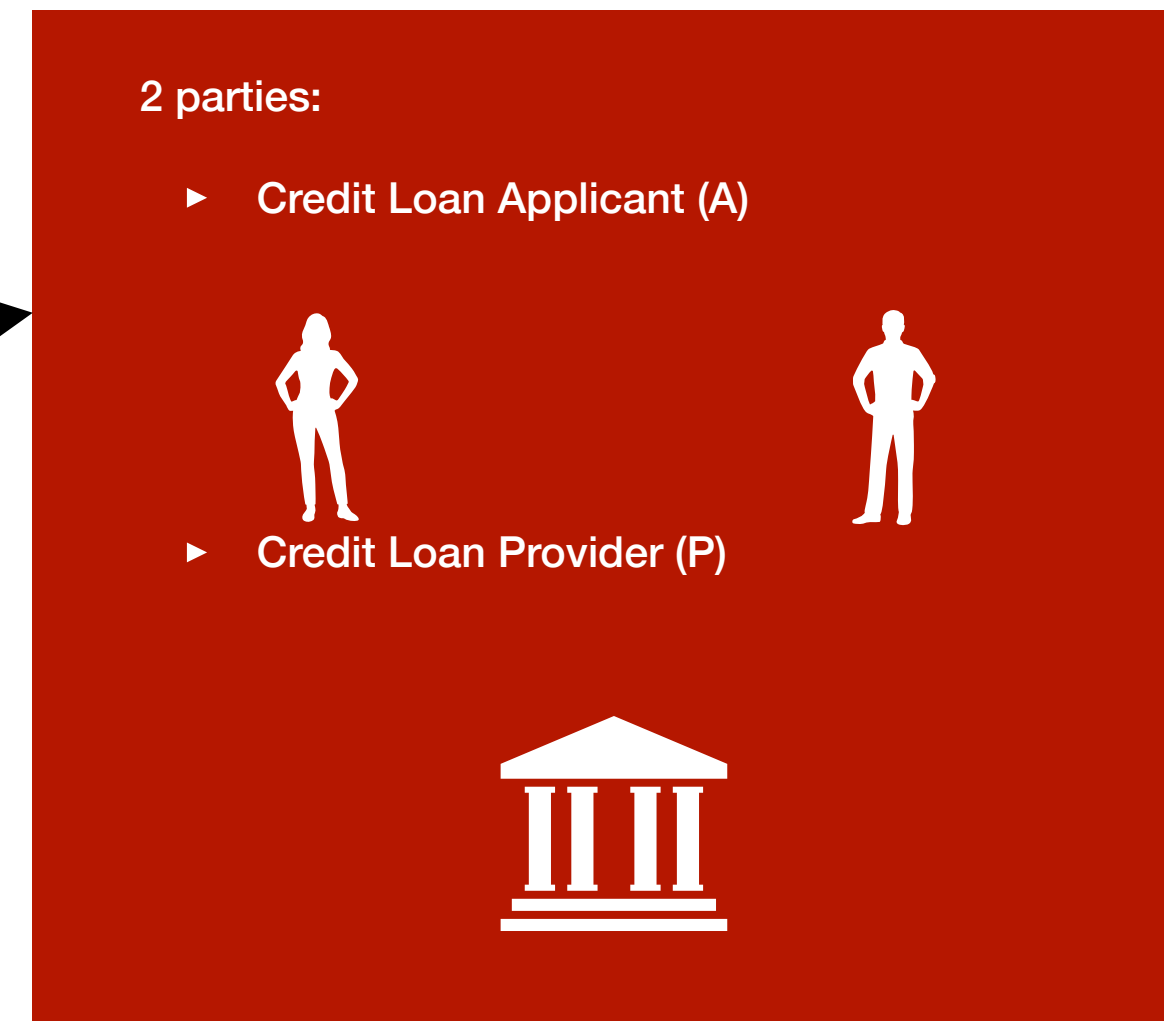
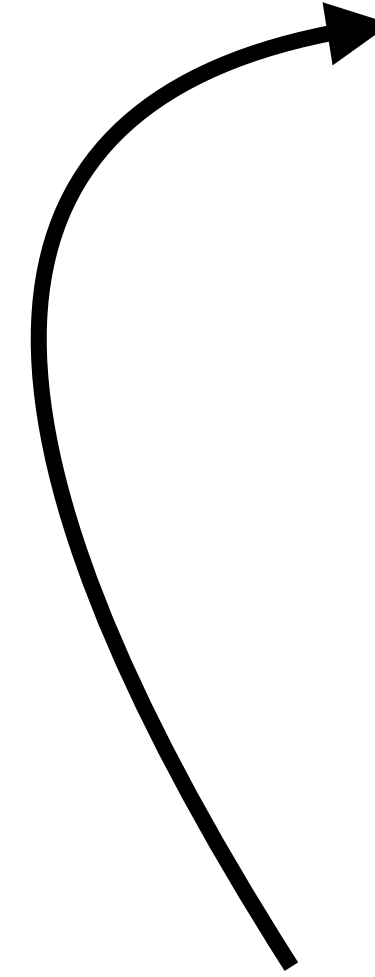
$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

► (P) has an automated decision system  $f$

$$\begin{aligned} f(x) &= 1 && \text{if accepted} \\ f(x) &= -1 && \text{if not} \end{aligned}$$

► (A) can ask for a counterfactual explanation

# Leading example



This example is seen as a:

Counterfactual literature

Positive example

Strategic classification

Negative example

# Modelling the situation

# Model

Learning theoretic setting for classification

$$f: \mathcal{X} \subseteq \mathbb{R}^d \rightarrow \{-1, 1\}$$

We assume that

$$(X_0, Y) \sim P$$

Loss is measured by counting wrong classifications

$$\ell(f(x), y) = 1\{f(x) \neq y\}$$

Goal is to minimize **expected loss (Risk/Accuracy)**

$$R = \mathbb{E}_P[\ell(f(X_0), Y)] = P(f(X_0) \neq Y).$$

The optimal classifier is the **Bayes Classifier**

$$f_P^* = \arg \min \mathbb{E}_P[\ell(f(X_0), Y)]$$



# Model

## Adding recourse

By adding recourse in the mix,

$$X_0 \rightarrow X,$$

where  $X$  is either  $X_0$  or  $X^{\text{CF}}$ , we induce a new distribution

$$(X_0, X, Y) \sim Q.$$

Counterfactual point is defined as

$$\varphi(X_0) = X^{\text{CF}} \in \arg \min_{z: f(z)=1} c(X_0, z)$$

For simplicity, we assume that every negative  $X_0$  accepts recourse

Risk with **Recourse** is defined as

$$R_Q(f) = \mathbb{E}_Q[\ell(f(X), Y)] = Q(f(X) \neq Y)$$

Note that  $Q$  depends on  $f$  in general

# Model

## When is Recourse accepted

In general  $X_0$  does not need to change to  $X^{\text{CF}}$ ,

This is modelled by setting

$$X = BX^{\text{CF}} + (1 - B)X_0, \quad B \sim \text{Ber}(r(X_0)).$$

The function  $r(X_0)$  models how likely  $X_0$  is to accept recourse

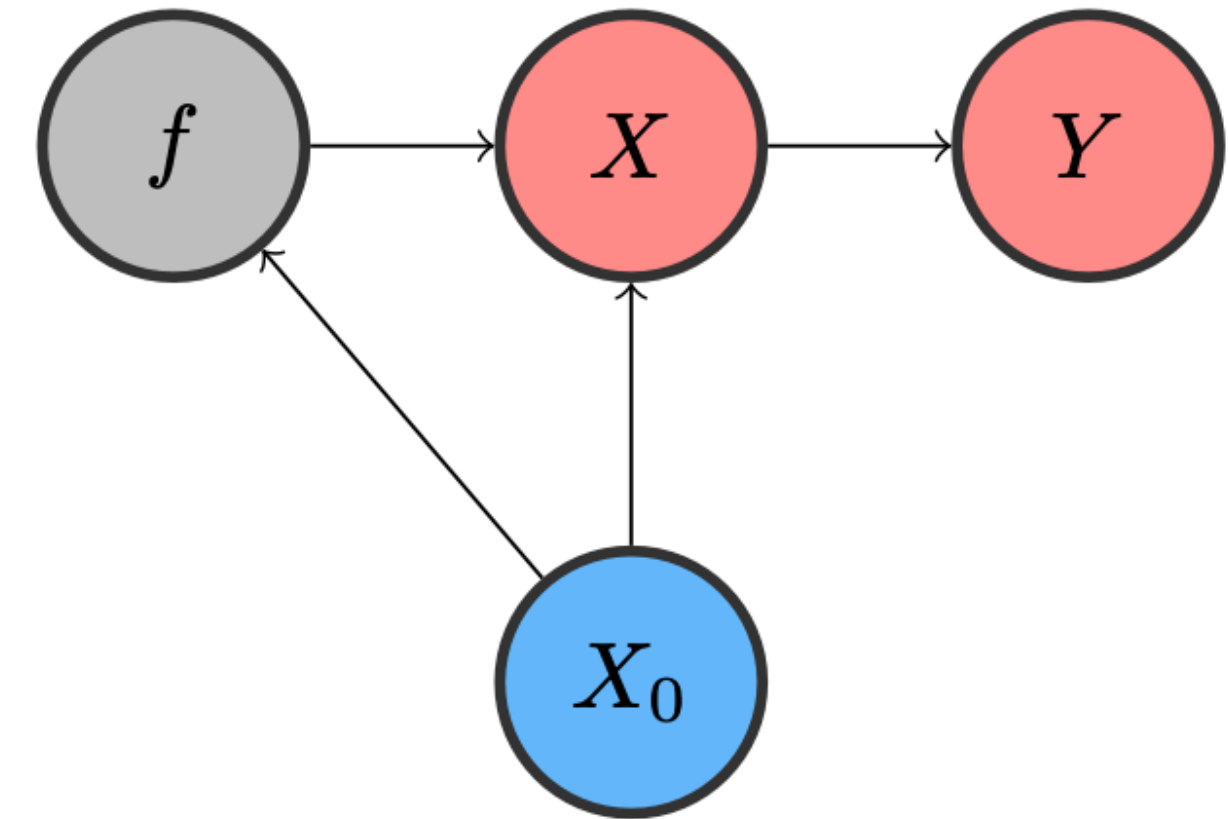
For the rest of the talk we will assume  $r(X_0) = 1$

# Modelling $Q$

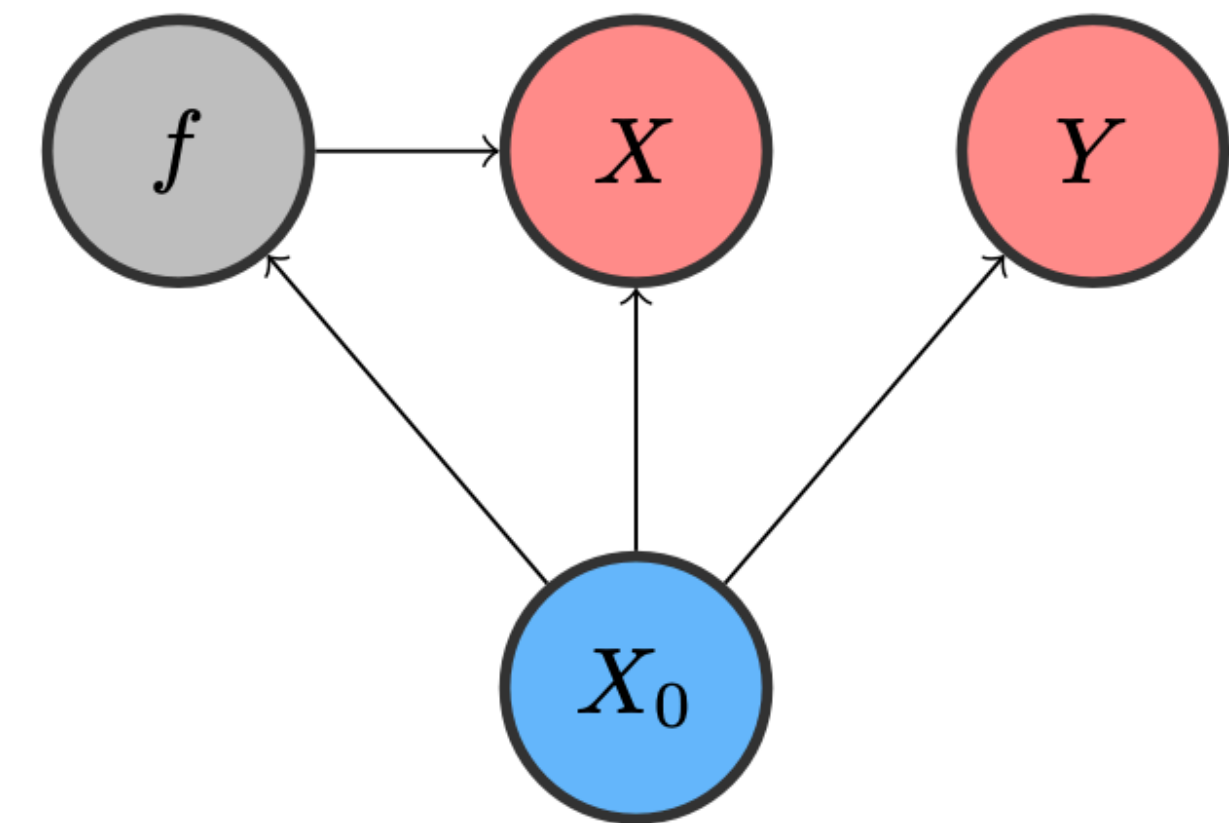
## Choice in dependency structure

We define 2 extreme cases:

- **Compliant:**  $Q(Y|X_0, X) = P(Y|X)$
- **Defiant:**  $Q(Y|X_0, X) = P(Y|X_0)$



Compliant case



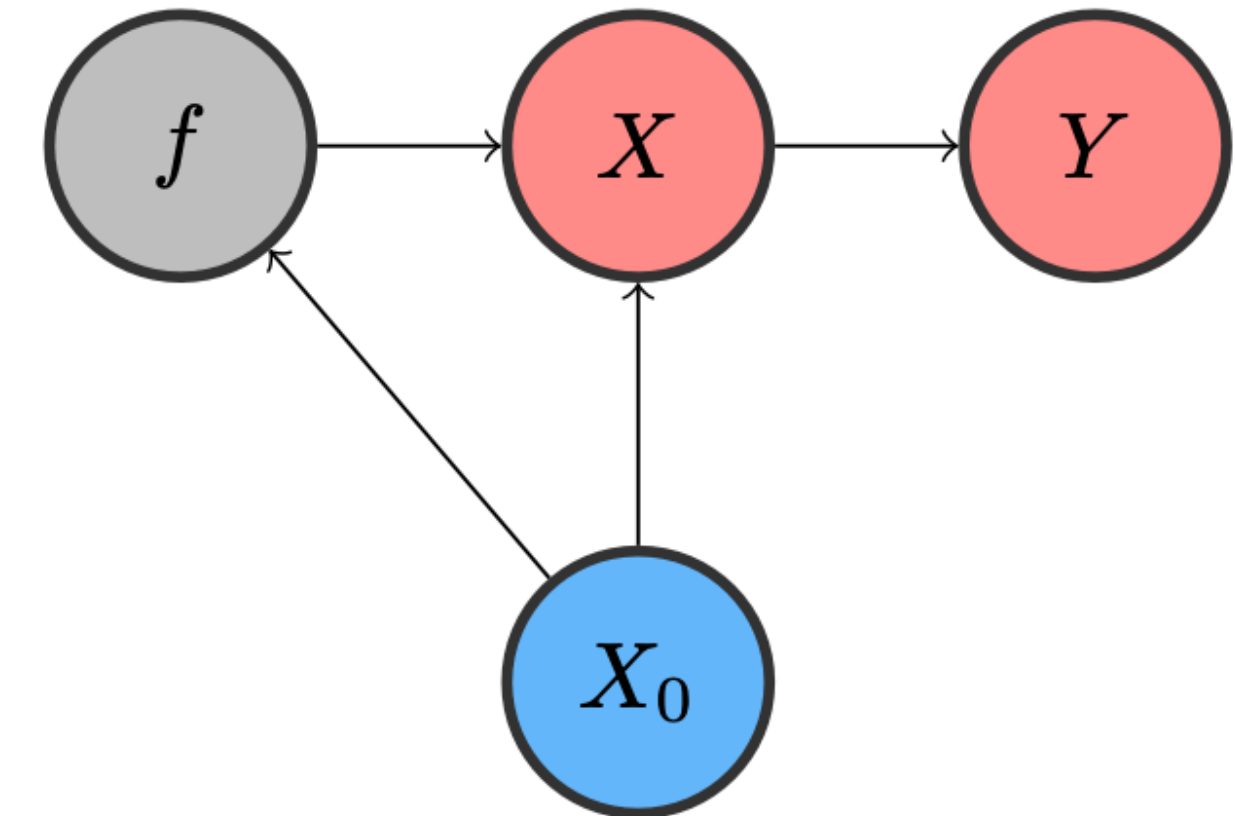
Defiant case

# Modelling $Q$

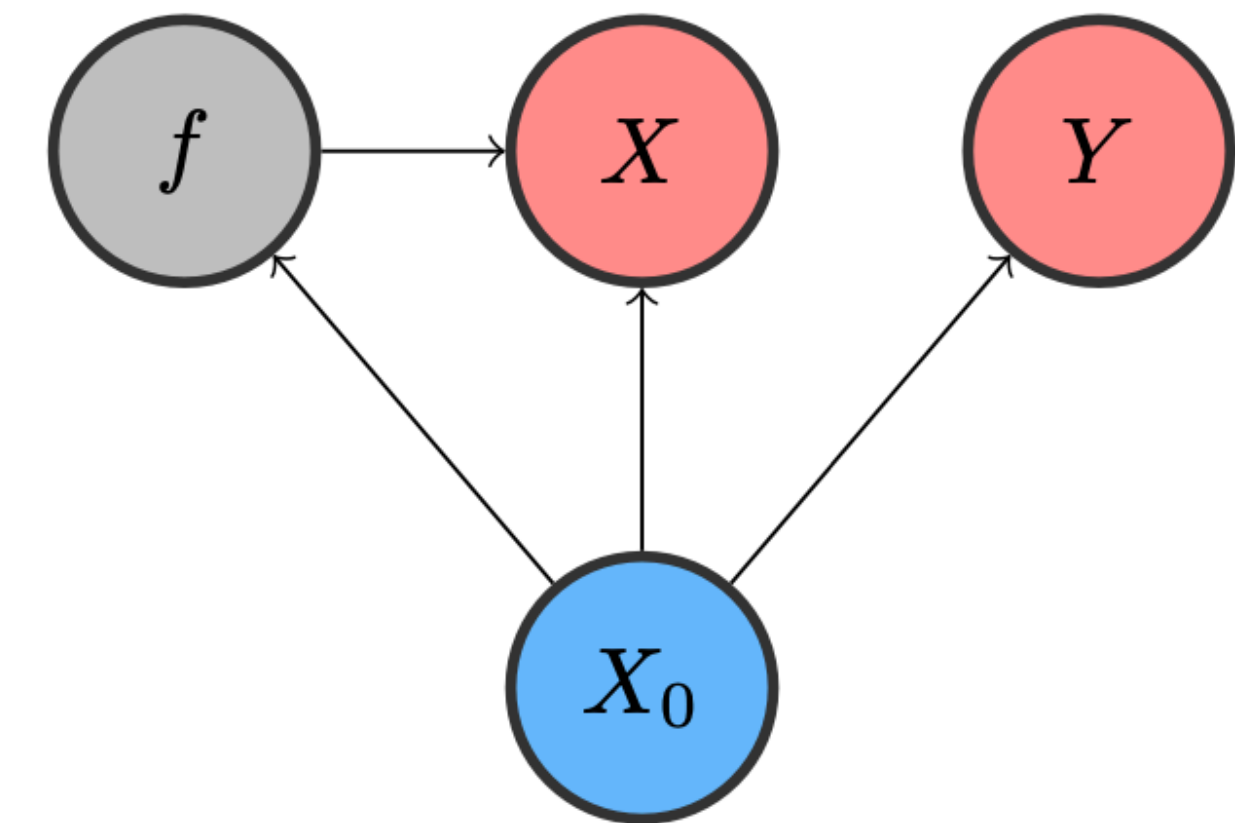
## Examples

Some examples:

- ▶ Credit loan application:
  - ▶ Compliant: Applicant improves risky behaviour
  - ▶ Defiant: Applicant tries to “game the system”
- ▶ Medical Diagnosis:
  - ▶ Compliant: Patient improves their health
  - ▶ Defiant: Patient takes medicine to reduce symptoms
- ▶ Job applications:
  - ▶ Compliant: Applicant improves their skills
  - ▶ Defiant: Applicant improves their CV



Compliant case



Defiant case

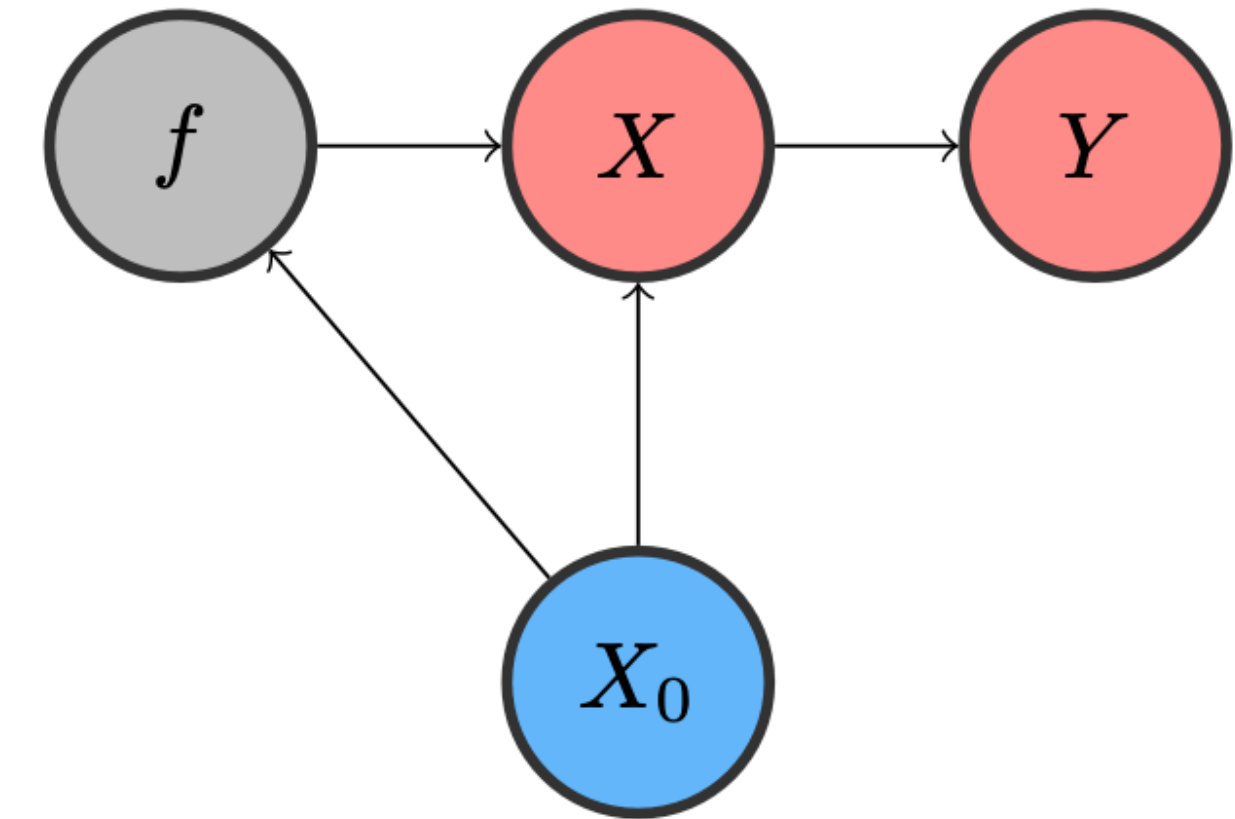
# Modelling $Q$

## Causality

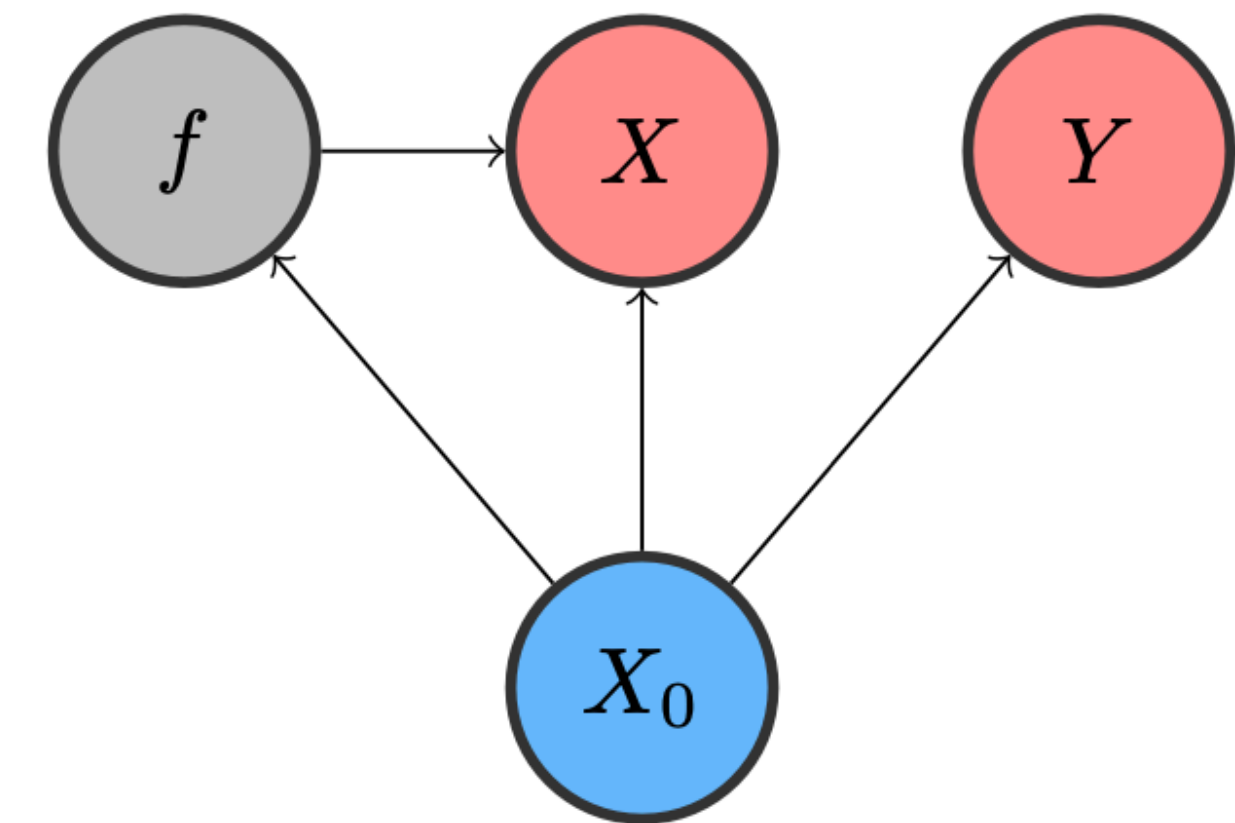
There is a very Causal interpretation to  $Q$ .

See *Improvement-Focused Causal Recourse (ICR)* [König et al.]  
for a more extensive treatment of this view

We view everything on a distributional level



Compliant case



Defiant case

# Optimal classifier

# Optimal Classifier

## Example (Compliant)

We assume that

$$X|Y = +1 \sim N(\mu, \Sigma)$$

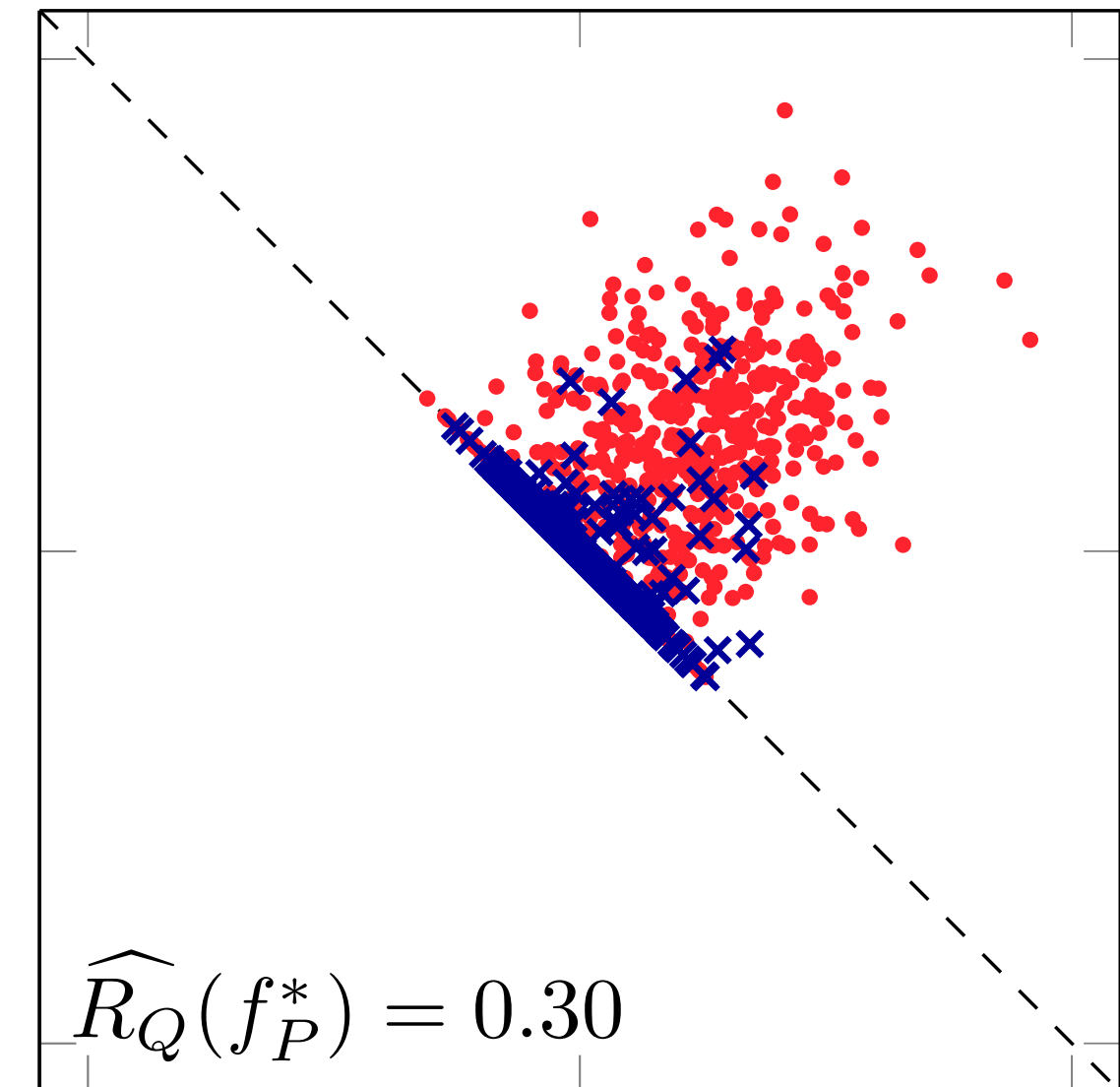
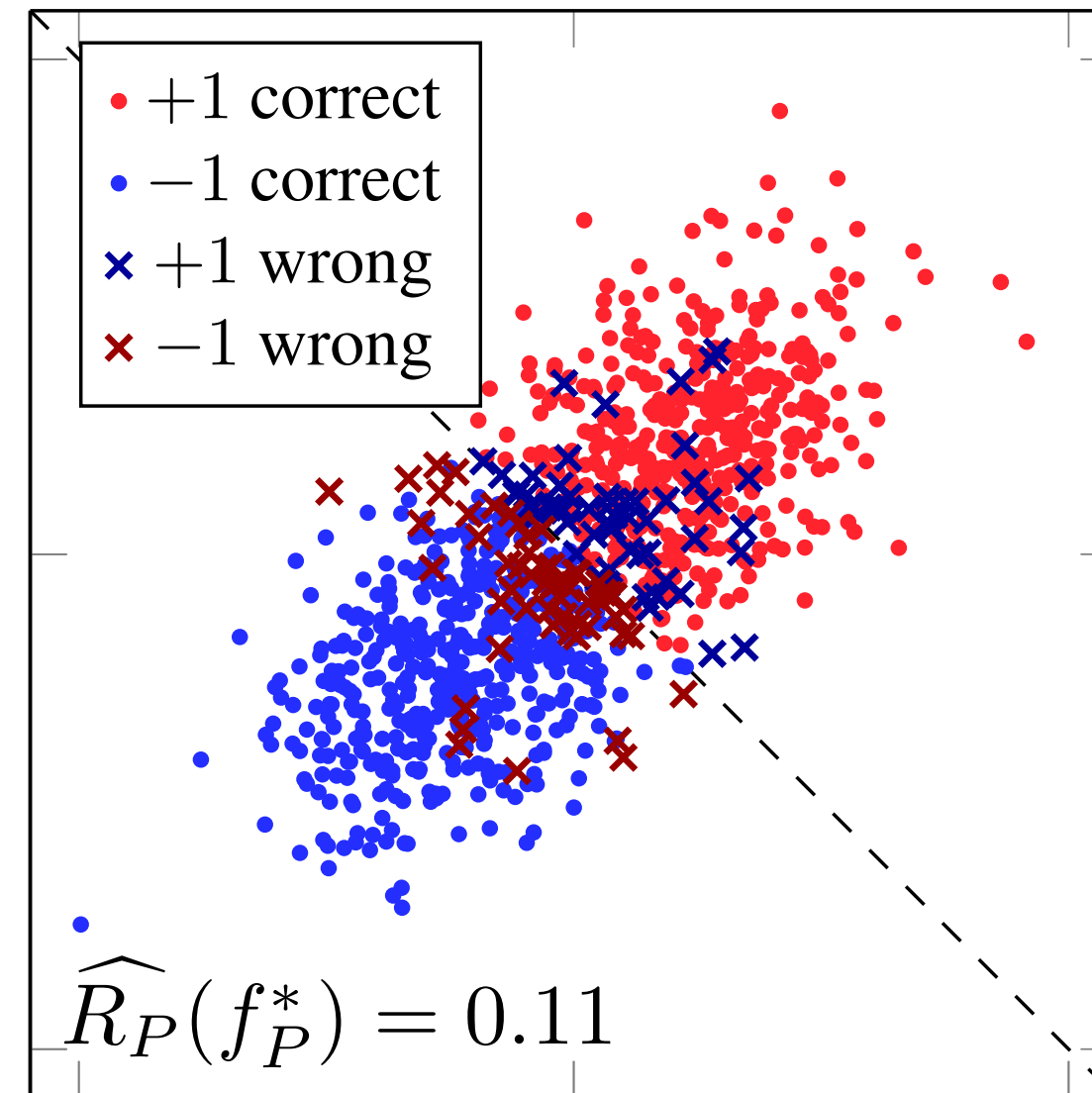
$$X|Y = -1 \sim N(\nu, \Sigma)$$

$$P(Y = +1) = P(Y = -1) = \frac{1}{2}$$

Optimal classifier is (assuming  $\|\mu\|_{\Sigma^{-1}} = \|\nu\|_{\Sigma^{-1}}$ )

$$f_P^*(x) = \text{sign}(\theta^\top x),$$

$$\theta = \Sigma^{-1}(\mu - \nu)$$



$$\triangleright R_P(f_P^*) = \Phi(\|\mu - \nu\|_{\Sigma^{-1}})$$

$$\triangleright R_Q(f_P^*) = \frac{1}{4} + \frac{1}{2}\Phi(\|\mu - \nu\|_{\Sigma^{-1}})$$

$$R_Q(f_P^*) > R_P(f_P^*) \text{ if } R_P(f_P^*) < \frac{1}{2}$$

# Optimal Classifier

## Formal result

### *Theorem*

Let  $\ell$  be the 0/1 loss and suppose that  $P(Y = 1 \mid X_0 = x) = \frac{1}{2}$  for all  $x$  on the decision boundary of  $f_P^*$ , then:

A. For the Compliant case,

$$R_Q(f_P^*) = \frac{1}{2}P(f_P(X_0) = -1) + P(f_P(X_0) = 1, Y = -1) > R_P(f_P^*)$$

B. For the Defiant case,

$$R_Q(f_P^*) = P(Y = -1) > R_P(f_P^*)$$



# Optimal Classifier

## Proof sketch (Compliant)

$$R_Q(f_P^*) = \frac{1}{2}P(f_P(X_0) = -1) + P(f_P(X_0) = 1, Y = -1) > R_P(f_P^*)$$

- ➡ Every point is now classified as +1
- ➡ The mistakes you make are
  - ➡ Original  $f_P^*(X_0) = +1$  but  $Y = -1$ ,
  - ➡ Half of the original  $f_P^*(X_0) = -1$ ,
  - ➡ because  $P(Y = +1 | X) = P(Y = -1 | X) = \frac{1}{2}$  on the decision boundary

# Optimal Classifier

## Proof sketch (Defiant)

$$R_Q(f_P^*) = P(Y = -1) > R_P(f_P^*)$$

- ➡ Every point is now classified as +1
- ➡ The mistakes you make are
  - ➡ Original  $f_P^*(X_0) = +1$  but  $Y = -1$ ,
  - ➡ Original  $f_P^*(X_0) = -1$ , but  $Y = -1$ , because the label does not change in this case

# Non-Optimal classifier

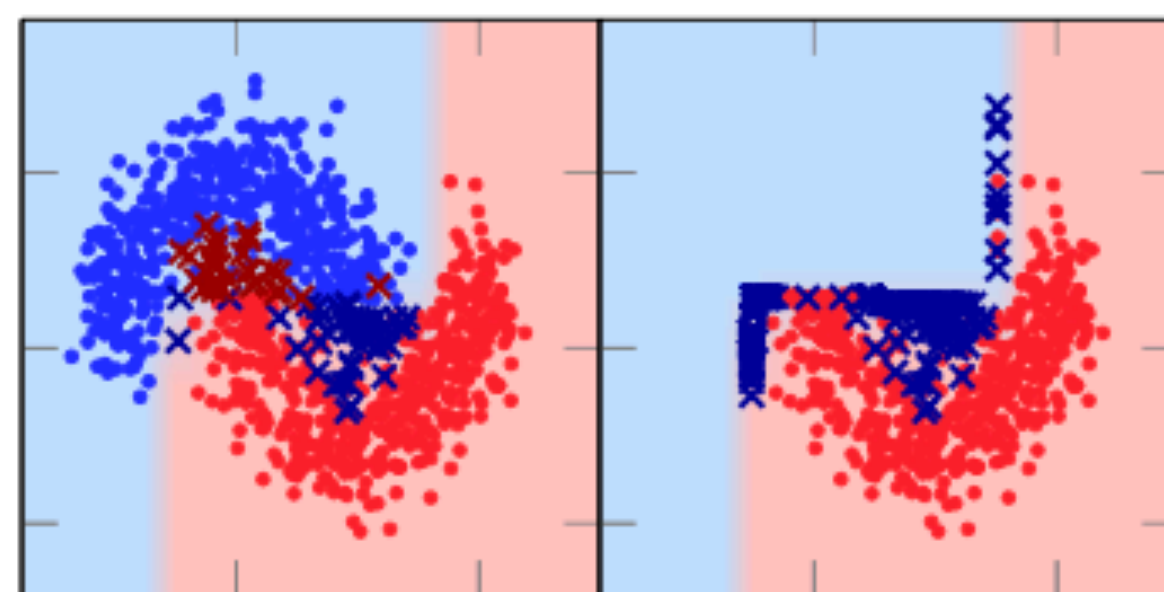
# Non-Optimal Classifier

## More general

What if we consider non-optimal classifiers?

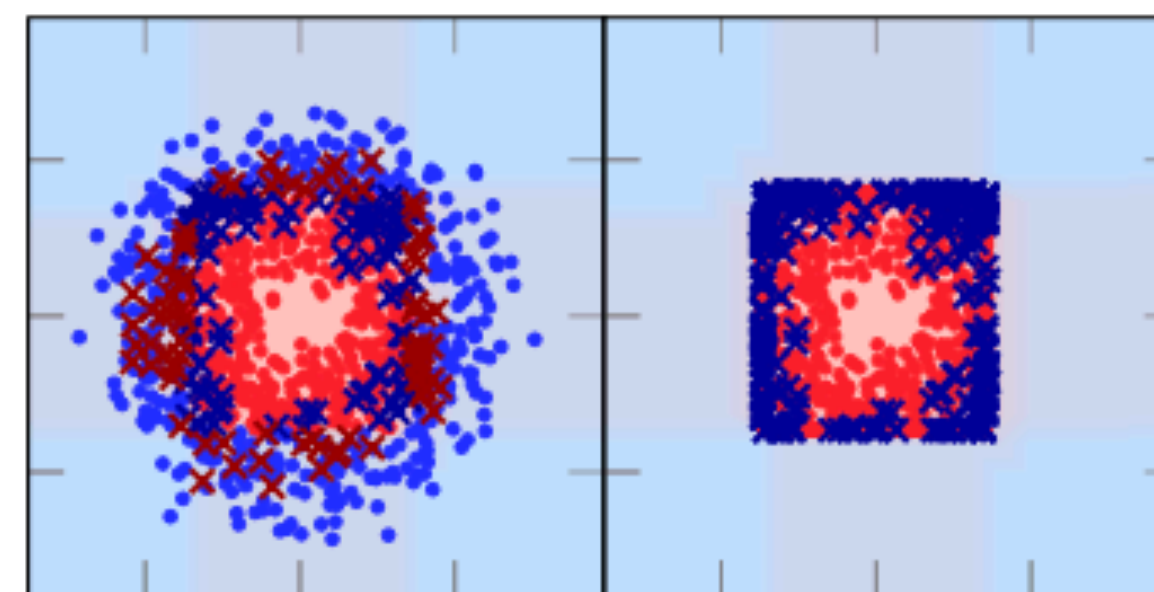
Need to make some extra assumptions:

- ▶ Assume  $f(x) = \text{sign}(g(x) - \frac{1}{2})$  for some probabilistic classifier  $g(x): \mathcal{X} \rightarrow [0,1]$
- ▶ The function  $g$  is “ $\epsilon$ -close” to  $P(Y = 1 | X = x)$  along the decision boundary



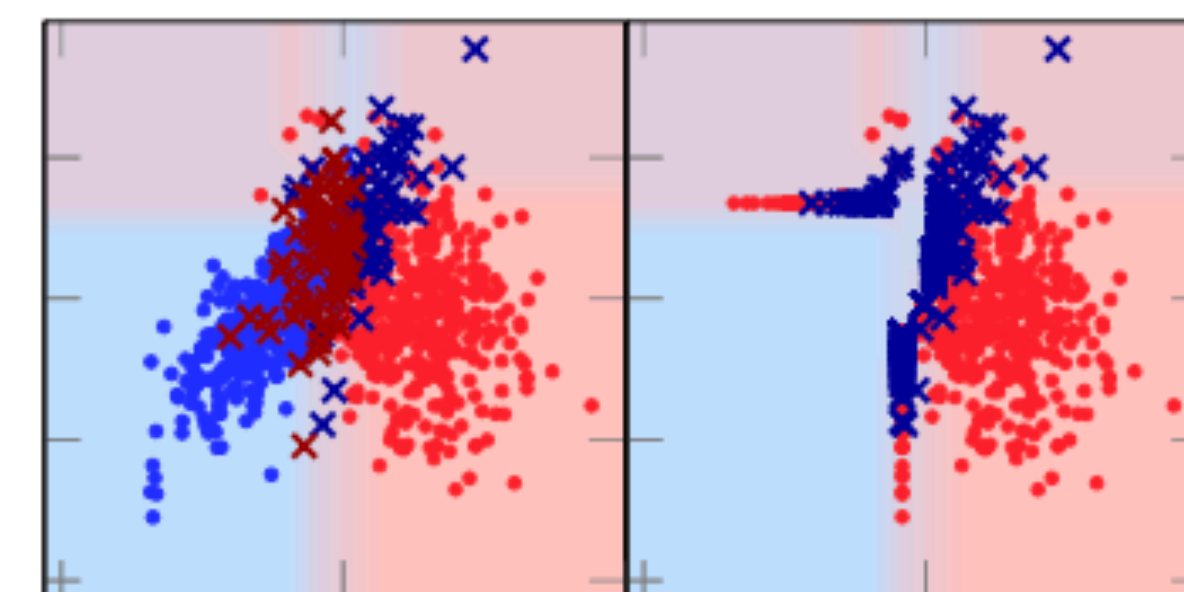
$$\widehat{R}_P(f) = 0.09$$

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



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

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



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

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

# $\epsilon$ -close assumption

## More general

What is the assumption?

Informally:

- ▶ The function  $g$  is “ $\epsilon$ -close” to  $P(Y = 1 \mid X = x)$  along the decision boundary

Formally:

$$\int_{\{x_0: g(x_0) < \frac{1}{2}\}} \left| \frac{1}{2} - P(Y = 1 \mid X = \varphi(x_0)) \right| P(dx_0) < \epsilon$$

Implied by uniform bound,

$$\left| \frac{1}{2} - P(Y = 1 \mid X_0 = x) \right| < \epsilon \text{ for all } x \text{ such that } g(x) = \frac{1}{2}$$

# Non-Optimal Classifier

## Formal result

### *Theorem*

Let  $\ell$  be the 0/1 loss,  $g: \mathcal{X} \rightarrow [0,1]$  a continuous probabilistic classifier and assume the  $\epsilon$ -condition:

A. For the Compliant case,  $R_Q(f)$  is lower and upper bounded by

$$(\frac{1}{2} \pm \epsilon)P(f(X_0) = -1) + P(f(X_0) = +1, Y = -1)$$

B. For the Defiant case,

$$R_Q(f) = P(Y = -1)$$

# Non-Optimal Classifier

## Implications

### *Theorem*

Let  $\ell$  be the 0/1 loss,  $g: \mathcal{X} \rightarrow [0,1]$  a continuous probabilistic classifier and assume the  $\epsilon$ -condition:

A. For the Compliant case,  $R_Q(f) \geq R_P(f)$  if

$$P(Y = -1 \mid f(X_0) = -1) \geq \frac{1}{2} + \epsilon$$

B. For the Defiant case,  $R_Q(f) \geq R_P(f)$  if and only if

$$P(Y = -1 \mid f(X_0) = -1) \geq \frac{1}{2}$$

# Non-Optimal Classifier

## Interpretation

*Recourse will harm the risk if*

- A. For the Compliant case, if  $f$  approximates the true conditional distribution and  $f$  performs  $\epsilon$  better on the negative class
- B. For the Defiant case, if  $f$  performs better than random on the negative class



# Non-Optimal Classifier

## Experimental results

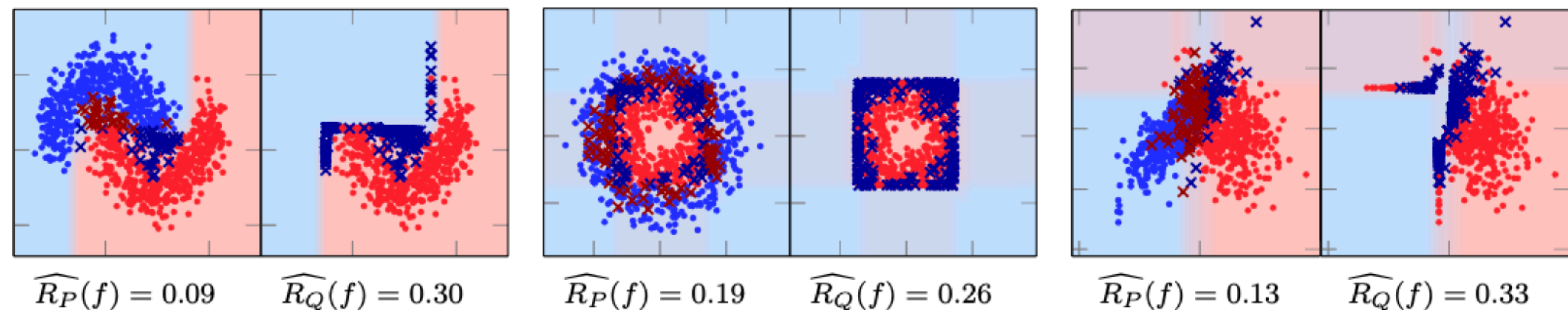


Table 1: Estimated risks on synthetic data sets. Lower risk is bold.

	Moons data		Circles data		Gaussians data	
	$R_P$	$R_Q$	$R_P$	$R_Q$	$R_P$	$R_Q$
Logistic Regression (LR)	<b>0.13</b>	0.33	0.51	<b>0.34</b>	<b>0.14</b>	0.32
GradientBoostedTrees (GBT)	<b>0.08</b>	0.30	<b>0.19</b>	0.26	<b>0.13</b>	0.33
Decision Tree (DT)	<b>0.08</b>	0.29	<b>0.19</b>	0.23	<b>0.13</b>	0.34
Naive Bayes (NB)	<b>0.13</b>	0.33	<b>0.17</b>	<b>0.16</b>	<b>0.15</b>	0.28
QuadraticDiscriminantAnalysis (QDA)	<b>0.13</b>	0.33	<b>0.17</b>	<b>0.16</b>	<b>0.12</b>	0.33
Neural Network(4)	<b>0.12</b>	0.32	<b>0.23</b>	0.30	<b>0.13</b>	0.36
Neural Network(4, 4)	<b>0.04</b>	0.26	<b>0.17</b>	0.22	<b>0.12</b>	0.40
Neural Network(8)	<b>0.04</b>	0.23	<b>0.16</b>	0.20	<b>0.12</b>	0.36
Neural Network(8, 16)	<b>0.04</b>	0.26	<b>0.16</b>	0.18	<b>0.11</b>	0.35
Neural Netowrk(8, 16, 8)	<b>0.04</b>	0.26	<b>0.16</b>	0.18	<b>0.11</b>	0.35

Table 2: Estimated risks on real data sets. Lower risk is bold.

	Credit data						Census data						HELOC data					
	Wachter		GS		CoGS		Wachter		GS		CoGS		Wachter		GS		CoGS	
	$R_P$	$R_Q$	$R_P$	$R_Q$	$R_P$	$R_Q$	$R_P$	$R_Q$	$R_P$	$R_Q$	$R_P$	$R_Q$	$R_P$	$R_Q$	$R_P$	$R_Q$	$R_P$	$R_Q$
LR	0.17	<b>0.05</b>	0.17	<b>0.05</b>	0.17	<b>0.04</b>	<b>0.21</b>	0.29	<b>0.21</b>	0.33	<b>0.21</b>	0.32	<b>0.29</b>	0.41	<b>0.29</b>	0.41	<b>0.29</b>	0.44
GBT	<b>0.06</b>	<b>0.06</b>	<b>0.06</b>	<b>0.07</b>	<b>0.06</b>	<b>0.07</b>	0.15	<b>0.04</b>	<b>0.15</b>	0.23	<b>0.15</b>	0.33	<b>0.20</b>	<b>0.21</b>	<b>0.20</b>	0.25	<b>0.20</b>	0.37
DT	0.29	<b>0.12</b>	0.29	<b>0.05</b>	0.29	<b>0.05</b>	<b>0.23</b>	<b>0.21</b>	<b>0.23</b>	0.43	<b>0.23</b>	0.45	<b>0.19</b>	0.25	<b>0.19</b>	0.21	<b>0.19</b>	0.31
NB	0.11	<b>0.06</b>	0.11	<b>0.06</b>	0.11	<b>0.07</b>	<b>0.19</b>	0.78	<b>0.19</b>	0.76	<b>0.19</b>	0.81	<b>0.29</b>	0.44	<b>0.29</b>	0.43	<b>0.29</b>	0.48
QDA	0.12	<b>0.06</b>	0.12	<b>0.06</b>	0.12	<b>0.07</b>	<b>0.20</b>	0.78	<b>0.20</b>	0.75	<b>0.20</b>	0.82	<b>0.32</b>	0.46	<b>0.32</b>	0.47	<b>0.32</b>	0.52
NN(4)	<b>0.06</b>	<b>0.06</b>	<b>0.06</b>	<b>0.07</b>	<b>0.06</b>	<b>0.06</b>	<b>0.16</b>	0.26	<b>0.16</b>	0.25	<b>0.16</b>	0.26	<b>0.29</b>	0.47	<b>0.29</b>	0.46	<b>0.29</b>	0.50
NN(4, 4)	<b>0.06</b>	<b>0.06</b>	<b>0.06</b>	<b>0.07</b>	<b>0.06</b>	<b>0.07</b>	<b>0.15</b>	0.30	<b>0.15</b>	0.27	<b>0.15</b>	0.30	<b>0.29</b>	0.47	<b>0.29</b>	0.47	<b>0.29</b>	0.51
NN(8)	<b>0.06</b>	<b>0.06</b>	<b>0.06</b>	<b>0.06</b>	<b>0.06</b>	<b>0.07</b>	<b>0.16</b>	0.34	<b>0.16</b>	0.33	<b>0.16</b>	0.33	<b>0.28</b>	0.44	<b>0.28</b>	0.46	<b>0.28</b>	0.51
NN(8, 16)	<b>0.06</b>	<b>0.06</b>	<b>0.06</b>	<b>0.07</b>	<b>0.06</b>	<b>0.07</b>	<b>0.15</b>	0.36	<b>0.15</b>	0.34	<b>0.15</b>	0.36	<b>0.27</b>	0.42	<b>0.27</b>	0.45	<b>0.27</b>	0.46
NN(8, 16, 8)	<b>0.06</b>	<b>0.06</b>	<b>0.06</b>	<b>0.07</b>	<b>0.06</b>	<b>0.07</b>	<b>0.15</b>	0.36	<b>0.15</b>	0.34	<b>0.15</b>	0.36	<b>0.27</b>	0.42	<b>0.27</b>	0.45	<b>0.27</b>	0.46

# Non-Optimal Classifier

## Proof sketch (Compliant)

Upper/lower:  $(\frac{1}{2} \pm \epsilon)P(f(X_0) = -1) + P(f(X_0) = +1, Y = -1)$

- ➡ Every point is now classified as +1
- ➡ The mistakes you make are
  - ➡ Original  $f(X_0) = +1$  but  $Y = -1$ ,
  - ➡ Within  $\epsilon$ -distance of half the original  $f(X_0) = -1$ ,
- ➡ Simplify  $R_P(f) \leq (\frac{1}{2} - \epsilon)P(f(X_0) = -1)$

# Non-Optimal Classifier

## Proof sketch (Defiant)

$$R_Q(f) = P(Y = -1) > R_P(f)$$

- ➡ Every point is now classified as +1
- ➡ The mistakes you make are
  - ➡ Original  $f(X_0) = +1$  but  $Y = -1$ ,
  - ➡ Original  $f(X_0) = -1$ , but  $Y = -1$ , because the label does not change in this case



# Non-Optimal Classifier

## Some more examples

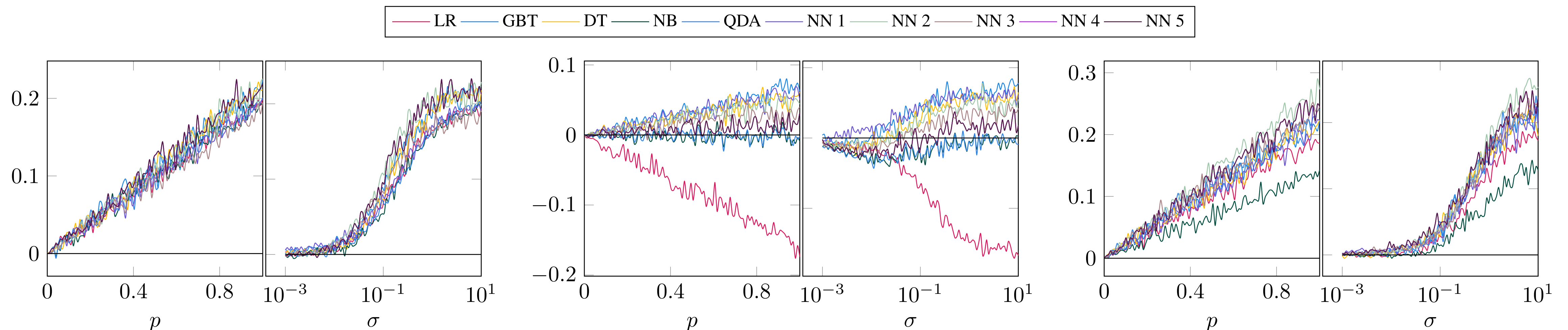
What if  $r(x)$  is not constant:

- ▶  $r(x) = p \in [0,1]$

- ▶  $r(x) = e^{-\frac{1}{2\sigma}\|x-\varphi(x)\|}$

- ▶ On y-axis:  $R_Q - R_P$

- ▶ From L to R: Moons, Circles, Gaussians



# Strategising

# Strategising

## Example

*Can (P) strategise against this accuracy drop?*

- Need to assume that not everyone gets an explanation, i.e.

$$r(x_0) = 1 \{ \|\varphi(x_0) - x_0\| < D \} \text{ for some } D > 0$$

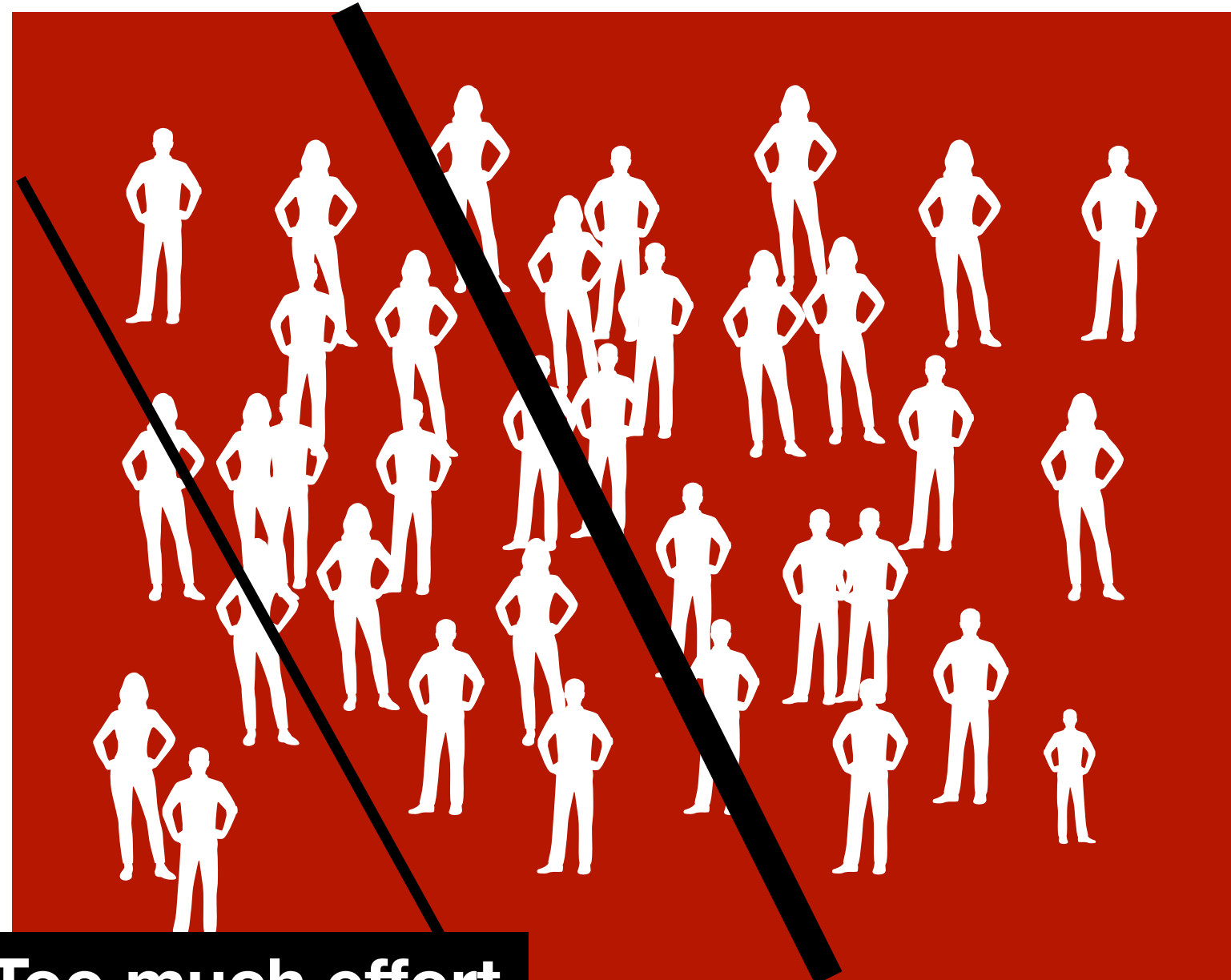
***Yes and No***

**Too much effort**



# Strategising

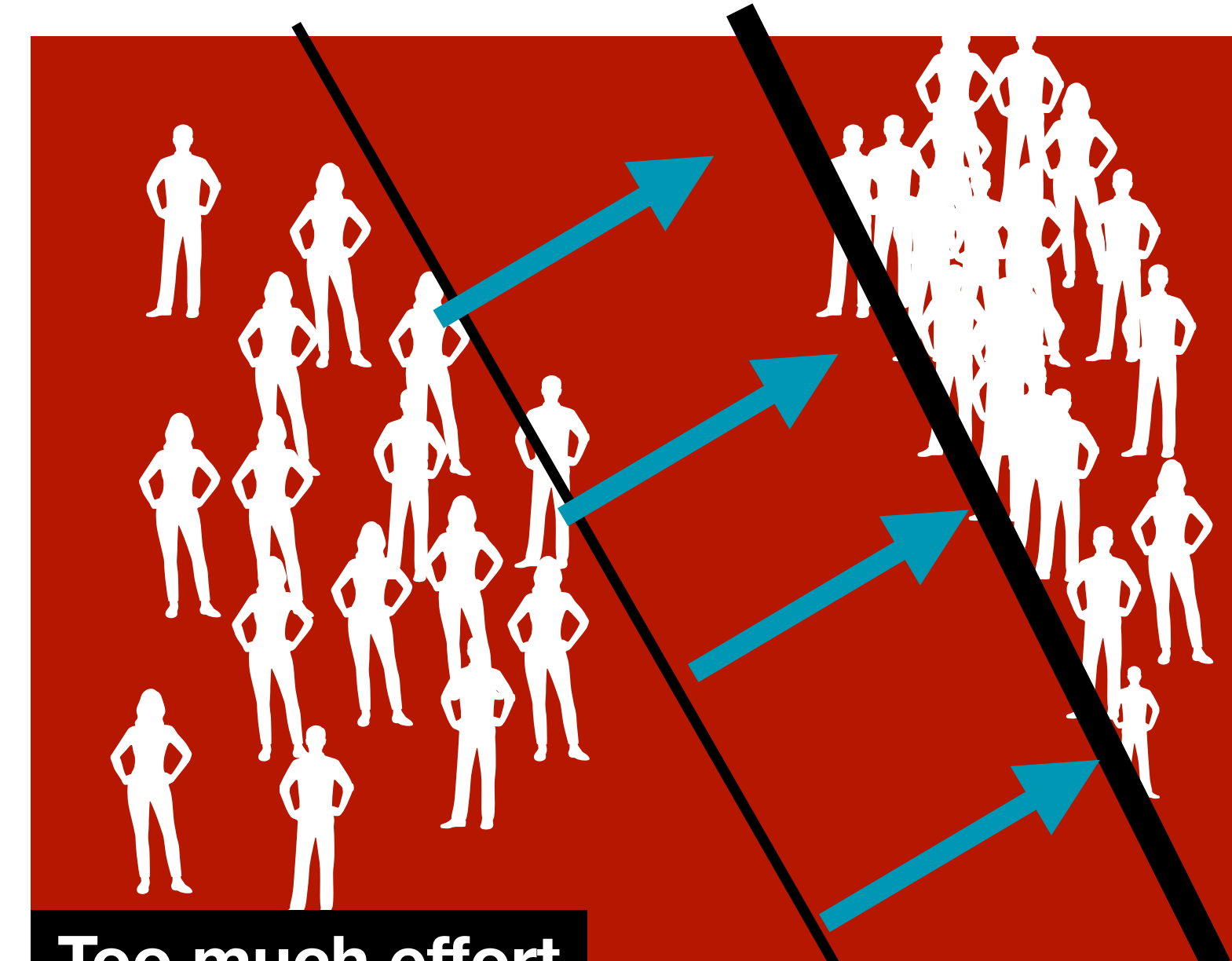
## Example



**Too much effort**



**Too much effort**



**Too much effort**

*Exactly the same people get a loan*

*More effort*

# Strategising

- ▶ Can generalise this result
  - ▶ Defiant case: Can only be as good as before
  - ▶ Compliant case: In principle, could have arbitrary improvement, but would increase the cost for every applicant substantially



# Strategising

## Formal result, a new definitions

### *Definition*

Let  $\mathcal{F}$  be some model class,  $r(x_0) \in \{0,1\}$  and  $\varphi^r(x_0) = r(x_0)\varphi(x_0) + (1 - r(x_0))\varphi(x_0)$ , define

$$\mathcal{F}_\varphi^r := \{f' : x_0 \mapsto f(\varphi^r(x_0)) \mid f \in \mathcal{F}\}.$$

The set of functions induced by the recourse map.

If  $\mathcal{F}_\varphi^r = \mathcal{F}$ , we call  $\mathcal{F}$  *invariant under recourse*

For example,

$$\mathcal{F} = \{f(x) = \text{sign}(a^\top x + b) \mid a, b \in \mathbb{R}^d\} \text{ and } r(x_0) = 1_{\{\|\varphi(x_0) - x_0\| < D\}}.$$

Then  $\mathcal{F} = \mathcal{F}_\varphi$

# Strategising

## Formal result (Defiant)

### *Theorem*

If  $\mathcal{F}$  is a recourse invariant model class, then

$$\min_{f \in \mathcal{F}} R_P(f) = \min_{f \in \mathcal{F}} R_Q(f).$$

***You can only do as well as before***

# Strategising

## Formal result (Compliant)

### *Theorem*

If  $\mathcal{F}$  is a recourse invariant model class, then

$$\min_{f \in \mathcal{F}} R_Q(f) \leq \min_{f \in \mathcal{F}} R_P(f) - \gamma.$$

Where  $\gamma \in \mathbb{R}$  depends on  $P$  and  $f$ .

***Improvement is possible when  $\gamma > 0$ !***

For example, in the Gaussian example

***However, there will be a cost for every individual***

Some perspectives

# Some perspectives

## When can Recourse still be beneficial?

### Some possible answers:

- ▶ In situations where Accuracy is not the most important metric
- ▶ When counterfactuals have other explanatory properties
- ▶ Should Recourse be replaced by Contestability?

**Thank you for your attention!**