

# Prime implicants as a versatile tool to explain robust classification

**Hénoïk Willot, Sébastien Destercke & Khaled Belahcene**

*13th International Symposium on Imprecise Probabilities:  
Theories and Applications*

# Our team



# Our team



**CID Team**

## Our team



# Prime implicants as a versatile tool to explain robust classification

**Hénoïk Willot, Sébastien Destercke & Khaled Belahcene**

*13th International Symposium on Imprecise Probabilities:  
Theories and Applications*

# Introduction

## Classification problem

Recommend : class  $\mathbf{y} \in \mathcal{Y} = \{y_1, \dots, y_m\}$

Features :  $\mathcal{X}^N = \prod_{i=1}^n \mathcal{X}_i$

Discrete domains :  $\mathcal{X}_i = \{x_i^1, \dots, x_i^{k_i}\}$

Observation :  $\mathbf{x} \in \mathcal{X}^N$

# Introduction

## Classification problem

Recommend : class  $\mathbf{y} \in \mathcal{Y} = \{y_1, \dots, y_m\}$

Features :  $\mathcal{X}^N = \prod_{i=1}^n \mathcal{X}_i$

Discrete domains :  $\mathcal{X}_i = \{x_i^1, \dots, x_i^{k_i}\}$

Observation :  $\mathbf{x} \in \mathcal{X}^N$

### Crisp case :

One probability distribution  $p$

$$\mathbf{y} \succeq_{p,(\mathbf{x})} \mathbf{y}' \text{ if } p(\mathbf{y}|\mathbf{x}) \geq p(\mathbf{y}'|\mathbf{x})$$

⇒ Explanations by prime implicants are known

# Introduction

## Classification problem

### Credal case :

Probability distribution  $p$  replaced by convex sets of probabilities  $\mathcal{P}$



# Introduction

## Classification problem

### Credal case :

Probability distribution  $p$  replaced by convex sets of probabilities  $\mathcal{P}$

### Robust classification :

Necessary recommendation  $\mathbf{y} \succ_{\mathcal{P},(\mathbf{x})} \mathbf{y}'$ ,

$$\mathbf{y} \succ_{\mathcal{P},(\mathbf{x})} \mathbf{y}' \Leftrightarrow \forall p \in \mathcal{P}, p(\mathbf{y}|\mathbf{x}) \geq p(\mathbf{y}'|\mathbf{x}) \Leftrightarrow \inf_{p \in \mathcal{P}} \frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}'|\mathbf{x})} \geq 1$$

# Introduction

## Classification problem

### Credal case :

Probability distribution  $p$  replaced by convex sets of probabilities  $\mathcal{P}$

### Robust classification :

Incomparability  $\mathbf{y} \succ_{\mathcal{P},(\mathbf{x})} \mathbf{y}'$ ,

$$\exists p \in \mathcal{P} \frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}'|\mathbf{x})} < 1 \text{ and } \exists p' \in \mathcal{P} \frac{p'(\mathbf{y}'|\mathbf{x})}{p'(\mathbf{y}|\mathbf{x})} < 1$$

# Introduction

## Classification problem

### Credal case :

Probability distribution  $p$  replaced by convex sets of probabilities  $\mathcal{P}$

### Robust classification :

Incomparability  $\mathbf{y} \succ_{\mathcal{P},(\mathbf{x})} \mathbf{y}'$ ,

$$\begin{aligned} \exists p \in \mathcal{P} \frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}'|\mathbf{x})} < 1 \text{ and } \exists p' \in \mathcal{P} \frac{p'(\mathbf{y}'|\mathbf{x})}{p'(\mathbf{y}|\mathbf{x})} < 1 \\ \Leftrightarrow \inf_{p \in \mathcal{P}} \frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}'|\mathbf{x})} < 1 \text{ and } \inf_{p' \in \mathcal{P}} \frac{p'(\mathbf{y}'|\mathbf{x})}{p'(\mathbf{y}|\mathbf{x})} < 1 \end{aligned}$$

# Introduction

## *Running example*

**Data** : ZOO dataset from UCI repository

**Objective** : predict  $\mathbf{y}$  in  $\mathcal{Y} = \{ \text{Mammal } (M), \text{ Bird } (B), \text{ Reptile } (R), \text{ Fish } (F) \text{ or Invertebrate } (I) \}$

# Introduction

## Running example

**Data** : ZOO dataset from UCI repository

**Objective** : predict  $\mathbf{y}$  in  $\mathcal{Y} = \{ \text{Mammal } (M), \text{ Bird } (B), \text{ Reptile } (R), \text{ Fish } (F) \text{ or Invertebrate } (I) \}$

### Features :

- feathers : { , ~~~~ }
- eggs : { , ~~~~ }
- aquatic : { , ~~~~ }
- toothed : { , ~~~~ }
- backbone : { , ~~~~ }
- breathes with lungs : { , ~~~~ }
- venomous : { , ~~~~ }
- fins : { , ~~~~ }
- legs :  
{ ~~~~, 2 \* , 4 \* , 5 \* , 6 \* , 8 \*  }
- tail : { , ~~~~ }

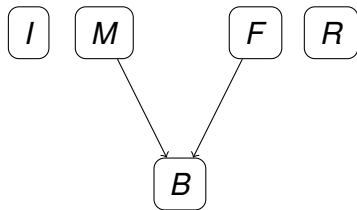
# Introduction

## Running example

**Observation** : 🦘 = (🚫🔪, 🚫🩸, 🌊, 🦷, 🦴, 🦷, 🦷, 🚫🦷, 🦷, 🚫🦷, 🚫🦷)

**Model** :

- $\log p(M|\text{🦘}) \in [-1.6, -0.012]$
- $\log p(B|\text{🦘}) \in [-8.36, -3]$
- $\log p(R|\text{🦘}) \in [-3.08, -0.016]$
- $\log p(F|\text{🦘}) \in [-2.59, -0.029]$
- $\log p(I|\text{🦘}) \in [-4.82, -0.45]$



# Introduction

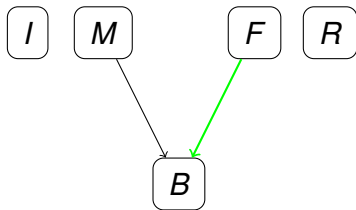
## Running example

Observation :  $\mathcal{M} = (\text{✂}, \text{✂}, \text{👉}, \text{🦷}, \text{👉}, \text{👉}, \text{☠}, \text{👉}, \text{✂}, \text{✂})$

Model :

- $\log p(M|\mathcal{M}) \in [-1.6, -0.012]$
- $\log p(B|\mathcal{M}) \in [-8.36, -3]$
- $\log p(R|\mathcal{M}) \in [-3.08, -0.016]$
- $\log p(F|\mathcal{M}) \in [-2.59, -0.029]$
- $\log p(I|\mathcal{M}) \in [-4.82, -0.45]$

$F \succ_{\mathcal{P}, (\mathcal{M})} B$  because  $-2.59 + 3 \geq 0$



# Introduction

## Running example

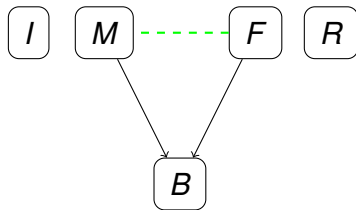
Observation :  $\mathcal{M} = (\text{✂}, \text{✂}, \text{👉}, \text{🦷}, \text{👉}, \text{👉}, \text{☠}, \text{👉}, \text{✂}, \text{✂})$

Model :

- $\log p(M|\mathcal{M}) \in [-1.6, -0.012]$
- $\log p(B|\mathcal{M}) \in [-8.36, -3]$
- $\log p(R|\mathcal{M}) \in [-3.08, -0.016]$
- $\log p(F|\mathcal{M}) \in [-2.59, -0.029]$
- $\log p(I|\mathcal{M}) \in [-4.82, -0.45]$

$F >_{\mathcal{P}, (\mathcal{M})} B$  because  $-2.59 + 3 \geq 0$

$M >_{\mathcal{P}, (\mathcal{M})} F$  because  $-1.6 + 0.029 < 0$





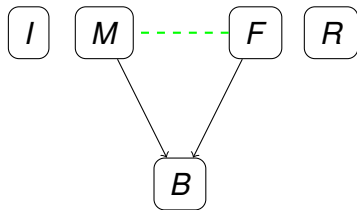
# Introduction

## Running example

Observation :  $\mathcal{M} = (\text{✂}, \text{✂}, \text{👉}, \text{🦷}, \text{👉}, \text{👉}, \text{☠}, \text{👉}, \text{✂}, \text{✂})$

Model :

- $\log p(M|\mathcal{M}) \in [-1.6, -0.012]$
- $\log p(B|\mathcal{M}) \in [-8.36, -3]$
- $\log p(R|\mathcal{M}) \in [-3.08, -0.016]$
- $\log p(F|\mathcal{M}) \in [-2.59, -0.029]$
- $\log p(I|\mathcal{M}) \in [-4.82, -0.45]$



$F >_{\mathcal{P}, (\mathcal{M})} B$  because  $-2.59 + 3 \geq 0$

$M >_{\mathcal{P}, (\mathcal{M})} F$  because  $-1.6 + 0.029 < 0$  and  $-2.59 + 0.012 < 0$

## ... **validatory**

$E \subseteq N$ , as a subset of feature indices, is a **validatory implicant** of decision  $\mathbf{y} \succ_{\mathcal{P},(\mathbf{x})} \mathbf{y}'$  if :

$$\forall \mathbf{x}_{-E} \in \mathcal{X}^{-E} \quad \inf_{\rho \in \mathcal{P}} \frac{\rho(\mathbf{y} | (\mathbf{x}_E, \mathbf{x}_{-E}))}{\rho(\mathbf{y}' | (\mathbf{x}_E, \mathbf{x}_{-E}))} \geq 1$$

## ... **validatory**

$E \subseteq N$ , as a subset of feature indices, is a **validatory implicant** of decision  $\mathbf{y} \succ_{\mathcal{P},(\mathbf{x})} \mathbf{y}'$  if :

$$\forall \mathbf{x}_{-E} \in \mathcal{X}^{-E} \quad \inf_{\rho \in \mathcal{P}} \frac{\rho(\mathbf{y} | (\mathbf{x}_E, \mathbf{x}_{-E}))}{\rho(\mathbf{y}' | (\mathbf{x}_E, \mathbf{x}_{-E}))} \geq 1$$

$$\Leftrightarrow \inf_{\substack{\rho \in \mathcal{P} \\ \mathbf{x}_{-E} \in \mathcal{X}^{-E}}} \frac{\rho(\mathbf{y} | (\mathbf{x}_E, \mathbf{x}_{-E}))}{\rho(\mathbf{y}' | (\mathbf{x}_E, \mathbf{x}_{-E}))} \geq 1$$

## ... *validatory*

$E \subseteq N$ , as a subset of feature indices, is a *validatory implicant* of decision  $\mathbf{y} \succ_{\mathcal{P},(\mathbf{x})} \mathbf{y}'$  if :

$$\forall \mathbf{x}_{-E} \in \mathcal{X}^{-E} \inf_{\rho \in \mathcal{P}} \frac{\rho(\mathbf{y} | (\mathbf{x}_E, \mathbf{x}_{-E}))}{\rho(\mathbf{y}' | (\mathbf{x}_E, \mathbf{x}_{-E}))} \geq 1$$

$$\Leftrightarrow \inf_{\substack{\rho \in \mathcal{P} \\ \mathbf{x}_{-E} \in \mathcal{X}^{-E}}} \frac{\rho(\mathbf{y} | (\mathbf{x}_E, \mathbf{x}_{-E}))}{\rho(\mathbf{y}' | (\mathbf{x}_E, \mathbf{x}_{-E}))} \geq 1$$

$E \subseteq N$  is a *prime implicant* if  $\forall i \in E$ , the inequality does not hold, *i.e.*  $E$  is subset-minimal

For one decision, it might exist different prime implicants with different cardinals !

## ...validatory

*i.e.* observing  $\mathbf{x}_E$  is sufficient to conclude no matter the values on other attributes of  $\mathcal{X}^{-E}$

$E = \{\text{'feathers'}, \text{'eggs'}, \text{'toothed'}\}$  is a prime implicant of  $M \succ_{\mathcal{P}, (\text{🐾})} B$

**Observation :**

$(\text{🐾}) = (\text{✂️}, \text{👄}, \text{👁️}, \text{🦷}, \text{🦋}, \text{🦷}, \text{🦋}, \text{👉}, \text{✂️}, \text{✂️}) \Rightarrow M \succ_{\mathcal{P}, (\text{🐾})} B$

**Observation :**

$(\text{🐾}_E, \text{?}_{-E}) = ((\text{✂️}, \text{👄}, \text{🦷}), (\text{👁️}, \text{🦋}, \text{🦷}, \text{🦋}, \text{👉}, \text{?}, \text{?}))$   
 $\Rightarrow M \succ_{\mathcal{P}, (\text{🐾}_E, \text{?}_{-E})} B$

## ... contrastive and doubt ?

Prime implicant can also be used for constrastive explanation

## ... contrastive and doubt ?

Prime implicant can also be used for constrastive explanation

But also to explain incomparability !

## ... contrastive and doubt ?

Prime implicant can also be used for constrastive explanation

But also to explain incomparability !

To know more come visit my poster 😊

Thank you for your attention !