

Zoo data example

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

Feature	feathers	eggs	aquatic	toothed	backbone	breathes with lungs	venomous	fins	legs	tail
Domain	{羽毛, ✕}	{卵, ✕}	{脊椎, ✕}	{歯, ✕}	{喙, ✕}	{肺呼吸, ✕}	{毒, ✕}	{鰐類, ✕}	{2*, 4*, 5*, 6*, 8*}, ✕	{尾, ✕}

If the value is left unspecified, we use the question mark symbol as for instance 🦴 for feathers

Observation : $\mathbf{x} = (\text{羽毛}, \text{卵}, \text{脊椎}, \text{歯}, \text{喙}, \text{肺呼吸}, \text{毒}, \text{鰐類}, \text{2*}, \text{4*}, \text{5*}, \text{6*}, \text{8*}, \text{尾})$

Crisp Case

A class y dominates a class y' given a probability distribution p

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

A Validatory Prime Implicant (PI) of $\mathbf{y} \succeq_{p,(\mathbf{x})} \mathbf{y}'$ is a subset E of features whose value in the observation is sufficient to obtain dominance no matter the other values

Example : $M \succeq_{p,(\mathbf{x})} B$, $E = \{\text{feathers, eggs}\}$ is a Validatory Prime Implicant because $M \succeq_{p,(\mathbf{x}_{-E}, ?_{-E})} B$ with $(\mathbf{x}_{-E}, ?_{-E}) = (\text{脊椎}, \text{歯}, \text{喙}, \text{肺呼吸}, \text{毒}, \text{鰐類}, \text{2*}, \text{4*}, \text{5*}, \text{6*}, \text{8*})$

A Contrastive Prime Implicant (PI) of $\mathbf{y} \succeq_{p,(\mathbf{x})} \mathbf{y}'$ is a subset E of features such that changing the observation on these values is sufficient to obtain a reversed dominance, while keeping other values untouched

Example : $M \succeq_{p,(\mathbf{x})} B$, $E = \{\text{feathers, eggs, legs}\}$ is a Contrastive Prime Implicant because $B \succeq_{p,(\mathbf{x}_{-E}, (\text{羽毛}, \text{卵}, \text{2*}))} M$ i.e. $(\text{羽毛}, \text{卵}, \text{脊椎}, \text{歯}, \text{喙}, \text{肺呼吸}, \text{毒}, \text{鰐類}, \text{2*}, \text{4*}, \text{5*}, \text{6*}, \text{8*})$ is classified as a Bird

Robust prediction with credal sets

Credal set: Probability distribution p replaced by convex sets of probabilities \mathcal{P}

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

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

Incomparability : When neither classes is necessarily better, incomparability $\mathbf{y} \succ_{\mathcal{P},(\mathbf{x})} \mathbf{y}'$ arises :

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

Validator PI

$$\begin{aligned} \forall x_{-E} \in \mathcal{X}^{-E} \inf_{p \in \mathcal{P}} \frac{p(\mathbf{y}|(\mathbf{x}_E, x_{-E}))}{p(\mathbf{y}'|(\mathbf{x}_E, x_{-E}))} &\geq 1 \\ \Leftrightarrow \phi_{\mathbf{y}, \mathbf{y}'}^v(E) = \inf_{x_{-E} \in \mathcal{X}^{-E}} \frac{p(\mathbf{y}|(\mathbf{x}_E, x_{-E}))}{p(\mathbf{y}'|(\mathbf{x}_E, x_{-E}))} &\geq 1 \end{aligned}$$

Contrastive PI

$$\begin{aligned} \exists x_{-E} \in \mathcal{X}^{-E} \inf_{p \in \mathcal{P}} \frac{p(\mathbf{y}|(x_E, x_{-E}))}{p(\mathbf{y}'|(x_E, x_{-E}))} &< 1 \\ \Leftrightarrow \phi_{\mathbf{y}, \mathbf{y}'}^c(E) = \inf_{x_E \in \mathcal{X}^E} \frac{p(\mathbf{y}|(x_E, x_{-E}))}{p(\mathbf{y}'|(x_E, x_{-E}))} &< 1 \end{aligned}$$

Doubt PI

$$\begin{aligned} \forall x_{-E} \in \mathcal{X}^{-E} \inf_{p \in \mathcal{P}} \frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}'|\mathbf{x})} &< 1 \text{ and } \inf_{p' \in \mathcal{P}} \frac{p'(\mathbf{y}'|\mathbf{x})}{p'(\mathbf{y}|\mathbf{x})} < 1 \\ \Leftrightarrow \phi_{\mathbf{y}, \mathbf{y}'}^d(E) = \sup_{x_{-E} \in \mathcal{X}^{-E}} \inf_{p \in \mathcal{P}} \frac{p(\mathbf{y}|(x_E, x_{-E}))}{p(\mathbf{y}'|(x_E, x_{-E}))} &< 1 \\ \text{and } \phi_{\mathbf{y}', \mathbf{y}}^d(E) = \sup_{x_{-E} \in \mathcal{X}^{-E}} \inf_{p' \in \mathcal{P}} \frac{p'(\mathbf{y}'|(x_E, x_{-E}))}{p'(\mathbf{y}|(x_E, x_{-E}))} &< 1 \end{aligned}$$

Application with NCC [1, 2]

Naive Bayes leads to $\frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}'|\mathbf{x})} = \frac{p_y(\mathbf{y})}{p_y(\mathbf{y}') \prod_{i \in N} p_i(\mathbf{x}_i|\mathbf{y}')} \prod_{i \in N} p_i(\mathbf{x}_i|\mathbf{y})$ and with **credal sets** for all $p_i(\cdot|\mathbf{y}')$ distributions we obtain $\inf_{p \in \mathcal{P}} \frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}'|\mathbf{x})} = \frac{p_y(\mathbf{y})}{\bar{p}_y(\mathbf{y}') \prod_{i \in N} \frac{p_i(\mathbf{x}_i|\mathbf{y})}{\bar{p}_i(\mathbf{x}_i|\mathbf{y}')}}$ with \underline{p} and \bar{p} lower/upper bounds of $p \in \mathcal{P}$

Validator + Contrastive PI

$$\begin{aligned} \mathbf{x}^{v, \mathbf{y}, \mathbf{y}'} = \mathbf{x}^{c, \mathbf{y}, \mathbf{y}'} &= \times_{i=1}^n \arg \inf_{x_i^a \in \mathcal{X}_i} \frac{p_i(x_i^a|\mathbf{y})}{\bar{p}_i(x_i^a|\mathbf{y}')} \\ \forall E \subset N \log \phi_{\mathbf{y}, \mathbf{y}'}^v(E \cup \{i\}) - \log \phi_{\mathbf{y}, \mathbf{y}'}^v(E) &= G_{\mathbf{y}, \mathbf{y}'}^v(i) \\ \Rightarrow \log \phi_{\mathbf{y}, \mathbf{y}'}^v(E) &= \log \phi_{\mathbf{y}, \mathbf{y}'}^v(\emptyset) + \sum_{i \in E} G_{\mathbf{y}, \mathbf{y}'}^v(i) \geq 0 \end{aligned}$$

We have exactly the same properties for $\phi_{\mathbf{y}, \mathbf{y}'}^c$. It is *item selection* problem [3]

$$G_{\mathbf{y}, \mathbf{y}'}^v(i) = \left(\log \frac{p_i(\mathbf{x}_i|\mathbf{y})}{\bar{p}_i(\mathbf{x}_i|\mathbf{y}')} \right) - \left(\log \frac{p_i(\mathbf{x}_i^{v, \mathbf{y}, \mathbf{y}'}|\mathbf{y})}{\bar{p}_i(\mathbf{x}_i^{v, \mathbf{y}, \mathbf{y}'}|\mathbf{y}')} \right)$$

$$G_{\mathbf{y}, \mathbf{y}'}^c(i) = \left(\log \frac{p_i(\mathbf{x}_i^{c, \mathbf{y}, \mathbf{y}'}|\mathbf{y})}{\bar{p}_i(\mathbf{x}_i^{c, \mathbf{y}, \mathbf{y}'}|\mathbf{y}')} \right) - \left(\log \frac{p_i(\mathbf{x}_i|\mathbf{y})}{\bar{p}_i(\mathbf{x}_i|\mathbf{y}')} \right)$$

Example

Worst opponent for $M \succ_{\mathcal{P},(\mathbf{x})} B$ and $M \succ_{\mathcal{P},(\mathbf{x})} R$						
$\mathbf{x}^{v, M, B} = (\text{羽毛}, \text{卵}, \text{脊椎}, \text{喙}, \text{肺呼吸}, \text{毒}, \text{2*}, \text{5*}, \text{尾})$						
$\mathbf{x}^{d, M, R} = (\text{羽毛}, \text{卵}, \text{脊椎}, \text{喙}, \text{肺呼吸}, \text{毒}, \text{4*}, \text{尾})$						
$\mathbf{x}^{d, R, M} = (\text{羽毛}, \text{卵}, \text{脊椎}, \text{歯}, \text{喙}, \text{肺呼吸}, \text{毒}, \text{尾})$						
$G_{M, B}^v$	2.64	2.29	0	2.29	.515	-6.77 3.27
$G_{M, R}^d$	0	0	-0.472	0	0	1.58
$G_{R, M}^d$	0	-1.49	0	-0.61	0	2.45
	?	?	?	?	?	
$G_{M, B}^v$.515	.515	.187	.91	.187	
$G_{M, R}^d$	0	0	-0.219	-1.31	-0.219	
$G_{R, M}^d$	-0.809	-1.25	-0.707	0	-0.707	

Doubt PI

$$(\alpha, \beta) \text{ is either } (\mathbf{y}, \mathbf{y}') \text{ either } (\mathbf{y}', \mathbf{y})$$

$$\mathbf{x}^{d, \alpha, \beta} = \times_{i=1}^n \arg \sup_{x_i^a \in \mathcal{X}_i} \frac{p_i(x_i^a|\alpha)}{\bar{p}_i(x_i^a|\beta)}$$

Like Validator or Contrastive PI, we obtain :

$$\log \phi_{\alpha, \beta}^d(E) = \log \phi_{\alpha, \beta}^d(\emptyset) + \sum_{i \in E} G_{\alpha, \beta}^d(i) < 0$$

$$G_{\alpha, \beta}^d(i) = \left(\log \frac{p_i(\mathbf{x}_i|\alpha)}{\bar{p}_i(\mathbf{x}_i|\beta)} \right) - \left(\log \frac{p_i(\mathbf{x}_i^{d, \alpha, \beta}|\alpha)}{\bar{p}_i(\mathbf{x}_i^{d, \alpha, \beta}|\beta)} \right)$$

Adding feature i in E changes both functions $\phi_{\mathbf{y}, \mathbf{y}'}^d$ and $\phi_{\mathbf{y}', \mathbf{y}}^d$ by adding in both a positive value, while both functions are bounded by 0.

It is a 2-dimentional knapsack problem !

Conclusion

Summary :

- Prime implicants for robust preferences
- Application to the NCC with convex domains
- Polynomial calculation of implicants

Perspectives :

- Pairwise or holistic explanations ?
- Implications on complexity to remove the independence hypothesis ?
- Refine explanations for incomparability (random or epistemic) ?

References

- [1] Marco Zaffalon. "The naive credal classifier". In: *Journal of Statistical Planning and Inference* 105.1 (2002), pp. 5–21.
- [2] Jean-Marc Bernard. "An introduction to the imprecise Dirichlet model for multinomial data". In: *International Journal of Approximate Reasoning* 39.2-3 (2005), pp. 123–150.
- [3] João Marques-Silva et al. "Explaining Naive Bayes and Other Linear Classifiers with Polynomial Time and Delay". In: *NeurIPS 2020*, December 6-12, 2020, virtual. 2020.