# From Contrastive to Abductive Explanations and Back Again

**Alexey Ignatiev[1]**, Nina Narodytska[2], Nicholas Asher[3], and Joao Marques-Silva[3]

November 9, 2021 | **KR**

[1]Monash University, Melbourne, Australia
[2]VMware Research, CA, USA
[3]IRIT, CNRS, Toulouse, France
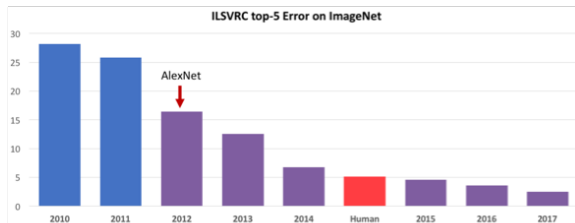
# Motivation

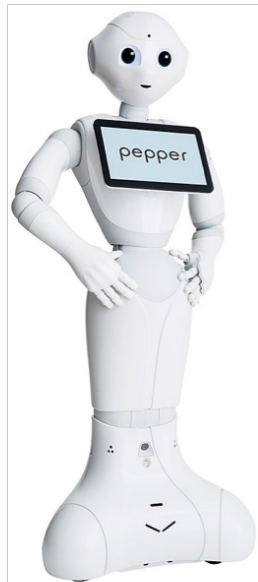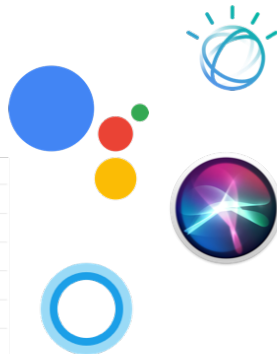# Ongoing ML Revolution



EHang184 Passenger Drone

https://en.wikipedia.org/wiki/Waymo

## Image & Speech Recognition



ILSVRC top-5 Error on ImageNet

AlexNet

2010 2011 2012 2013 2014 Human 2015 2016 2017

http://gradientscience.org/intro_adversarial/



DeepMind

AlphaGo

AlphaGo Zero & **Alpha Zero**



pepper

https://fr.wikipedia.org/wiki/Pepper_(robot)

# And yet...

| | A parrot | Machine learning algorithm |
|---|---|---|
| |  |  |
| **Learns random phrases** | ☑ | ☑ |
| **Doesn't understand s\*\*t about what it learns** | ☑ | ☑ |
| **Occasionally speaks nonsense** | ☑ | ☑ |

©Internet memes

eXplainable AI

Machine Learning System

Cat

This is a cat.

Current Explanation

This is a cat:
- It has fur, whiskers, and claws.
- It has this feature:

XAI Explanation

eXplainable AI

**Machine Learning System**

Cat

Which features?

Why cat?

This is a cat.

Explain?!

**Current Explanation**

This is a cat:
- It has fur, whiskers, and claws.
- It has this feature:

**XAI Explanation**

©DARPA

3/9

# Formal explanations

**classifier** $\tau : \mathbb{F} \rightarrow \mathcal{K}$, **instance** $\mathbf{v}$ **s.t.** $\tau(\mathbf{v}) = c$

**classifier $\tau : \mathbb{F} \to \mathcal{K}$, instance $\mathbf{v}$ s.t. $\tau(\mathbf{v}) = c$**

---

**abductive explanation $\mathcal{X}$**

$$\forall (\mathbf{x} \in \mathbb{F}) . \bigwedge\nolimits_{j \in \mathcal{X}} (x_j = v_j) \to (\tau(\mathbf{x}) = c)$$

**classifier  $\tau : \mathbb{F} \to \mathcal{K}$,  instance  $\mathbf{v}$  s.t.  $\tau(\mathbf{v}) = c$**

---

**abductive explanation $\mathcal{X}$**  *"why?"*

$$\forall (\mathbf{x} \in \mathbb{F}). \bigwedge\nolimits_{j \in \mathcal{X}} (x_j = v_j) \to (\tau(\mathbf{x}) = c)$$

**classifier $\tau : \mathbb{F} \to \mathcal{K}$, instance $\mathbf{v}$ s.t. $\tau(\mathbf{v}) = c$**

**abductive explanation $\mathcal{X}$**   *"why?"*

$$\forall (\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \to (\tau(\mathbf{x}) = c)$$

**contrastive explanation $\mathcal{Y}$**

$$\exists (\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\tau(\mathbf{x}) \neq c)$$

**classifier** $\tau : \mathbb{F} \rightarrow \mathcal{K}$, **instance** $\mathbf{v}$ **s.t.** $\tau(\mathbf{v}) = c$

abductive explanation $\mathcal{X}$ — *"why?"*

$$\forall (\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}} (x_j = v_j) \rightarrow (\tau(\mathbf{x}) = c)$$

contrastive explanation $\mathcal{Y}$ — *"why not?"*

$$\exists (\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\tau(\mathbf{x}) \neq c)$$

this work!

$$\mathbb{F} = \{0, 1, 2\}^5 \quad \mathcal{K} = \{\ominus, \oplus\}$$

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| | | | |
|---|---|---|---|
| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| | | | |
|---|---|---|---|
| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

**observe** $\tau(1, 1, 1, 1, 1) = \ominus$

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| | | | |
|---|---|---|---|
| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

**observe $\tau(1, 1, 1, 1, 1) = \ominus$**

**AXps $\mathbb{X} = \{\{1, 2\}, \{3\}\}$**

$$\mathbb{F} = \{0, 1, 2\}^5 \quad \mathcal{K} = \{\ominus, \oplus\}$$

| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
|---|---|---|---|
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

observe $\tau(1, 1, 1, 1, 1) = \ominus$

AXps $\mathbb{X} = \{\{1, 2\}, \{3\}\}$

CXps $\mathbb{Y} = \{\{1, 3\}, \{2, 3\}\}$

**Abductive XPs**

# Minimal hitting set duality

Abductive XPs

Contrastive XPs

AXps are minimal hitting sets of CXps, and vice versa

# CXp computation

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

# CXp computation – example

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| | | | |
|---|---|---|---|
| $R_0$: | IF | $x_1 = 1 \land x_2 = 1$ | THEN $\ominus$ |
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
|---|---|---|---|
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

**observe** $\tau(1, 1, 1, 1, 1) = \ominus$

# CXp computation – example

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| | | | |
|---|---|---|---|
| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

observe  $\tau(1, 1, 1, 1, 1) = \ominus$  **– why not $\oplus$?**

## CXp computation – example

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| | | | |
|---|---|---|---|
| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

**observe** $\tau(1, 1, 1, 1, 1) = \ominus$ **– why not $\oplus$?**

1. *can we drop feature 1?* $\quad \tau(1, *, *, *, *) \not\equiv \ominus$?

$$\exists(\mathbf{x} \in \mathbb{F}). \bigwedge\nolimits_{j \not\in \mathcal{Y}} (x_j = v_j) \wedge (\tau(\mathbf{x}) \neq c)$$

## CXp computation – example

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| | | | |
|---|---|---|---|
| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

**observe** $\tau(1, 1, 1, 1, 1) = \ominus$ **– why not $\oplus$?**

1. *can we drop feature 1?* $\quad \tau(1, *, *, *, *) \not\equiv \ominus$? **Yes**

$$\exists(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\tau(\mathbf{x}) \neq c)$$

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
|---|---|---|---|
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

**observe** $\tau(1, 1, 1, 1, 1) = \ominus$ **– why not $\oplus$?**

1. *can we drop feature 1?* $\quad \tau(1, *, *, *, *) \not\equiv \ominus$? **Yes**

2. *can we drop feature 2?* $\quad \tau(1, 1, *, *, *) \not\equiv \ominus$?

$$\exists(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\tau(\mathbf{x}) \neq c)$$

## CXp computation – example

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| | | | |
|---|---|---|---|
| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{\text{DEF}}$: | ELSE | | THEN $\ominus$ |

**observe** $\tau(1, 1, 1, 1, 1) = \ominus$  **– why not $\oplus$?**

1. *can we drop feature 1?*  $\tau(1, *, *, *, *) \not\equiv \ominus$? **Yes**
2. *can we drop feature 2?*  $\tau(1, 1, *, *, *) \not\equiv \ominus$? **No**

$$\exists(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{y}}(x_j = v_j) \wedge (\tau(\mathbf{x}) \neq c)$$

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| | | | |
|---|---|---|---|
| $R_0$: | **IF** | $x_1 = 1 \wedge x_2 = 1$ | **THEN** $\ominus$ |
| $R_1$: | **ELSE IF** | $x_3 \neq 1$ | **THEN** $\oplus$ |
| $R_{DEF}$: | **ELSE** | | **THEN** $\ominus$ |

**observe** $\tau(1, 1, 1, 1, 1) = \ominus$ **– why not** $\oplus$**?**

1. *can we drop feature 1?* $\quad \tau(1, *, *, *, *) \not\equiv \ominus$? **Yes**
2. *can we drop feature 2?* $\quad \tau(1, 1, *, *, *) \not\equiv \ominus$? **No**
3. *can we drop feature 3?* $\quad \tau(1, *, 1, *, *) \not\equiv \ominus$?

$$\exists (\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\tau(\mathbf{x}) \neq c)$$

## CXp computation – example

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
|---|---|---|---|
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

**observe** $\tau(1, 1, 1, 1, 1) = \ominus$ **– why not $\oplus$?**

1. *can we drop feature 1?* $\quad \tau(1, *, *, *, *) \not\equiv \ominus$? **Yes**
2. *can we drop feature 2?* $\quad \tau(1, 1, *, *, *) \not\equiv \ominus$? **No**
3. *can we drop feature 3?* $\quad \tau(1, *, 1, *, *) \not\equiv \ominus$? **No**
4. *can we drop feature 4?* $\quad \tau(1, *, *, 1, *) \not\equiv \ominus$?

$$\exists(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{Y}}(x_j = v_j) \wedge (\tau(\mathbf{x}) \neq c)$$

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| | | | |
|---|---|---|---|
| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

**observe** $\tau(1, 1, 1, 1, 1) = \ominus$ **– why not $\oplus$?**

1. *can we drop feature 1?*    $\tau(1, *, *, *, *) \not\equiv \ominus$? **Yes**
2. *can we drop feature 2?*    $\tau(1, 1, *, *, *) \not\equiv \ominus$? **No**
3. *can we drop feature 3?*    $\tau(1, *, 1, *, *) \not\equiv \ominus$? **No**
4. *can we drop feature 4?*    $\tau(1, *, *, 1, *) \not\equiv \ominus$? **Yes**
5. *can we drop feature 5?*    $\tau(1, *, *, 1, 1) \not\equiv \ominus$?

$$\exists (\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\tau(\mathbf{x}) \neq c)$$

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
|---|---|---|---|
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{DEF}$: | ELSE | | THEN $\ominus$ |

**observe** $\tau(1, 1, 1, 1, 1) = \ominus$ **– why not $\oplus$?**

1. *can we drop feature 1?* $\quad \tau(1, *, *, *, *) \not\equiv \ominus$? **Yes**
2. *can we drop feature 2?* $\quad \tau(1, 1, *, *, *) \not\equiv \ominus$? **No**
3. *can we drop feature 3?* $\quad \tau(1, *, 1, *, *) \not\equiv \ominus$? **No**
4. *can we drop feature 4?* $\quad \tau(1, *, *, 1, *) \not\equiv \ominus$? **Yes**
5. *can we drop feature 5?* $\quad \tau(1, *, *, 1, 1) \not\equiv \ominus$? **Yes**

$$\exists (\mathbf{x} \in \mathbb{F}). \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\tau(\mathbf{x}) \neq c)$$

# CXp computation – example

$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| | | | |
|---|---|---|---|
| $R_0$: | **IF** | $x_1 = 1 \wedge x_2 = 1$ | **THEN** $\ominus$ |
| $R_1$: | **ELSE IF** | $x_3 \neq 1$ | **THEN** $\oplus$ |
| $R_{DEF}$: | **ELSE** | | **THEN** $\ominus$ |

**observe** $\tau(1, 1, 1, 1, 1) = \ominus$ **– why not $\oplus$?**

1. *can we drop feature 1?* $\quad \tau(1, *, *, *, *) \not\equiv \ominus$? **Yes**

2. *can we drop feature 2?* $\quad \tau(1, 1, *, *, *) \not\equiv \ominus$? **No**

> CXp $\mathcal{Y} = \{2, 3\}$

3. *can we drop feature 3?* $\quad \tau(1, *, 1, *, *) \not\equiv \ominus$? **No**

4. *can we drop feature 4?* $\quad \tau(1, *, *, 1, *) \not\equiv \ominus$? **Yes**

5. *can we drop feature 5?* $\quad \tau(1, *, *, 1, 1) \not\equiv \ominus$? **Yes**

$$\exists (\mathbf{x} \in \mathbb{F}) . \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\tau(\mathbf{x}) \neq c)$$

## Explanation Enumeration

**Function** XPENUM($\tau$, **v**, c)
  **Input:** $\tau$: ML model, **v**: Input instance, $c = \tau(\mathbf{v})$: Prediction
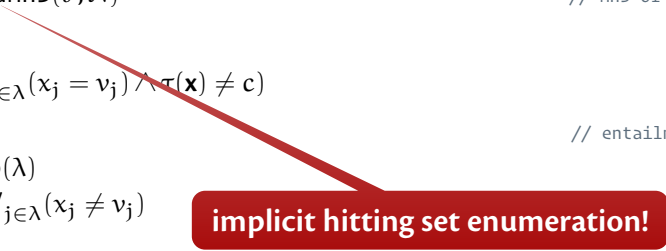
1     $\mathcal{K} = (\mathcal{N}, \mathcal{P}) \leftarrow (\emptyset, \emptyset)$            // Block AXps & CXps

2     **while** true:

3        $(st_\lambda, \lambda) \leftarrow \texttt{FindMHS}(\mathcal{P}, \mathcal{N})$          // MHS of $\mathcal{P}$ s.t. $\mathcal{N}$

4        **if** $\neg st_\lambda$**: break**

5        $st_{c'} \leftarrow \texttt{SAT}(\bigwedge_{j \in \lambda}(x_j = v_j) \wedge \tau(\mathbf{x}) \neq c)$

6        **if** $\neg st_{c'}$**:**                               // entailment holds

7           $\texttt{ReportAXp}(\lambda)$

8           $\mathcal{N} \leftarrow \mathcal{N} \cup \bigvee_{j \in \lambda}(x_j \neq v_j)$

9        **else:**

10          $\mu \leftarrow \texttt{ExtractCXp}(\tau, \mathbf{v}, c, \mathcal{P})$

11          $\texttt{ReportCXp}(\mu)$

12          $\mathcal{P} \leftarrow \mathcal{P} \cup \bigvee_{j \in \mu}(x_j = v_j)$

## Explanation Enumeration

**Function** XpEnum($\tau$, **v**, c)
    **Input:** $\tau$: ML model, **v**: Input instance, $c = \tau(\mathbf{v})$: Prediction

1    $\mathcal{K} = (\mathcal{N}, \mathcal{P}) \leftarrow (\emptyset, \emptyset)$                                           // Block AXps & CXps

2    **while** true:

3        $(st_\lambda, \lambda) \leftarrow$ FindMHS$(\mathcal{P}, \mathcal{N})$                          // MHS of $\mathcal{P}$ s.t. $\mathcal{N}$

4        **if** $\neg st_\lambda$**: break**

5        $st_{c'} \leftarrow$ SAT$(\bigwedge_{j \in \lambda}(x_j = v_j) \wedge \tau(\mathbf{x}) \neq c)$

6        **if** $\neg st_{c'}$**:**                                              // entailment holds

7            ReportAXp$(\lambda)$

8            $\mathcal{N} \leftarrow \mathcal{N} \cup \bigvee_{j \in \lambda}(x_j \neq v_j)$

9        **else:**

10           $\mu \leftarrow$ ExtractCXp$(\tau, \mathbf{v}, c, \mathcal{P})$

11           ReportCXp$(\mu)$

12           $\mathcal{P} \leftarrow \mathcal{P} \cup \bigvee_{j \in \mu}(x_j = v_j)$

**implicit hitting set enumeration!**

**see paper for details**

# Conclusions

# Conclusions

- **formal definition of contrastive explanations**
  - **similar to abductive explanations**

# Conclusions

- **formal definition of contrastive explanations**
  - **similar to abductive explanations**

- **minimal hitting set duality between CXps and AXps**
  - **explanation enumeration algorithms**
  - **solving membership problems**

# Conclusions

- **formal definition of contrastive explanations**
  - **similar to abductive explanations**

- **minimal hitting set duality between CXps and AXps**
  - **explanation enumeration algorithms**
  - **solving membership problems**

**proved helpful in several papers!**

## Conclusions

- **formal definition of contrastive explanations**
  - **similar to abductive explanations**

- **minimal hitting set duality between CXps and AXps**
  - **explanation enumeration algorithms**
  - **solving membership problems**

  **proved helpful in several papers!**

- **experimental results**
  - **XP enumeration**
  - **CXp enumeration – *helps to debug SHAP***

**Questions?**

# References i

📄 Adnan Darwiche and Auguste Hirth.
   **On the reasons behind decisions.**
   In *ECAI*, pages 712–720, 2020.

📄 Xuanxiang Huang, Yacine Izza, Alexey Ignatiev, and Joao Marques-Silva.
   **On efficiently explaining graph-based classifiers.**
   In *KR*, 2021.

📄 Yacine Izza, Alexey Ignatiev, and João Marques-Silva.
   **On explaining decision trees.**
   *CoRR*, abs/2010.11034, 2020.

📄 Yacine Izza and João Marques-Silva.
   **On explaining random forests with SAT.**
   In *IJCAI*, pages 2584–2591, 2021.

# References ii

Alexey Ignatiev, Nina Narodytska, Nicholas Asher, and João Marques-Silva.
**From contrastive to abductive explanations and back again.**
In *AI\*IA*, pages 335–355, 2020.

Alexey Ignatiev, Nina Narodytska, and João Marques-Silva.
**On relating explanations and adversarial examples.**
In *NeurIPS*, pages 15857–15867, 2019.

Alexey Ignatiev, Nina Narodytska, and Joao Marques-Silva.
**Abduction-based explanations for machine learning models.**
In *AAAI*, pages 1511–1519, 2019.

Alexey Ignatiev and João P. Marques Silva.
**SAT-based rigorous explanations for decision lists.**
In *SAT*, pages 251–269, 2021.

📄 João Marques-Silva, Thomas Gerspacher, Martin C. Cooper, Alexey Ignatiev, and Nina Narodytska.
**Explaining naive bayes and other linear classifiers with polynomial time and delay.**
In *NeurIPS*, 2020.

📄 João Marques-Silva, Thomas Gerspacher, Martin C. Cooper, Alexey Ignatiev, and Nina Narodytska.
**Explanations for monotonic classifiers.**
In *ICML*, pages 7469–7479, 2021.

📄 Andy Shih, Arthur Choi, and Adnan Darwiche.
**A symbolic approach to explaining Bayesian network classifiers.**
In *IJCAI*, pages 5103–5111, 2018.