# Introduction to interpretable AI

SETI Master 2024

---

Julien Girard-Satabin (CEA LIST): julien.girard2@cea.fr

February 8, 2024

# Preliminaries

## You and I

Myself:

1. researcher at CEA on formal methods for software safety and security applied to machine learning;

2. working on case-based reasoning and out-of-distribution detection in industrial use cases;

3. informed citizen;

You:

1. M2 SETI master students;

2. future practitionners of AI systems: designer, developers, debuggers;

3. informed citizens;

Preliminaries
○○●○○○○

Post-hoc explanations
○○○○○○○○

Explanable by design programs
○○○○○○○

# Hands on TP

1. `https://git.frama-c.com/pub/seti_master/-/`
   `archive/xai_tp/seti_master-xai_tp.zip`
2. `bash setup.sh`
3. wait some time
4. `bash launch.sh`

This will download the required python environment and several other dependencies

Preliminaries
○○○●○○

Post-hoc explanations
○○○○○○○○

Explanable by design programs
○○○○○○○

# Definitions

### Explanation

"An explanation is a presentation of (aspects of) the reasoning, functioning and/or behavior of a machine learning model in human-understandable terms" [Nau+23]

"The **belief** (by the trustor) in the ability (of the trustee) to achieve **something**"

Preliminaries
○○○○●○

Post-hoc explanations
○○○○○○○○

Explanable by design programs
○○○○○○○

# Explanation is a spectrum

Social science have quite a big corpus on what constitutes a good explanation ([Mil19])?

1. *contrastive*: why P instead of Q?
2. *a social process*: A explains P to B
3. *more generic* (cover more facts), *simpler* (quote less causes), and *coherent* (related to previous knowledge) are more easily understood

Preliminaries
○○○○○●

Post-hoc explanations
○○○○○○○○

Explanable by design programs
○○○○○○○

# Why it matters

1. debugging and audit

2. refutability

3. compliance with regulation (GDPR article 13.f [SP17])

# Post-hoc explanations

Preliminaries
oooooo

Post-hoc explanations
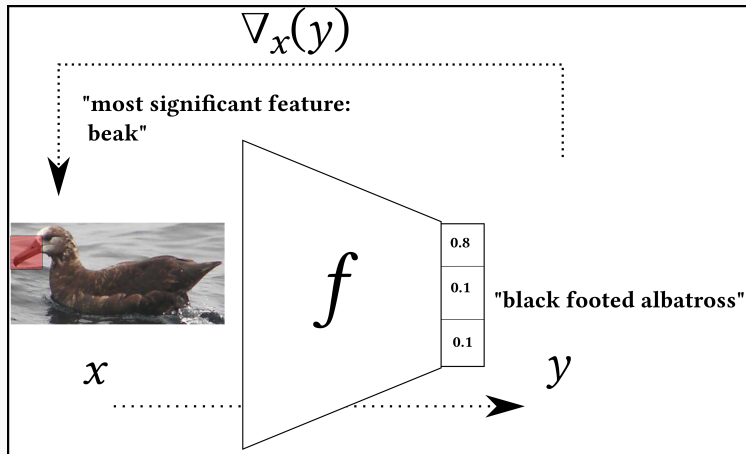o●oooooo

Explanable by design programs
oooooooo

## Notations

1. samples $x \in \mathcal{X} \subseteq \mathbb{R}^d$ an input space, $i^{th}$ feature $x_i$

2. an output $y \in \mathcal{Y} \subseteq \mathbb{R}^p$, the $i^{th}$ feature $y_i$

3. a program $f : \mathcal{X} \mapsto \mathcal{Y}$ trained on a $\mathcal{X}$
   - we can usually decompose $f = h \circ g$
   - in the following, $h(x)$ is the output of an intermediate layer for neural network

4. $\nabla_x(y)$ is the gradient of $y$ at $x$

Preliminaries
○○○○○○

Post-hoc explanations
○○●○○○○○

Explanable by design programs
○○○○○○○

# Framework of feature attribution

Preliminaries
oooooo

Post-hoc explanations
ooooooooo

Explanable by design programs
ooooooo

## Some caveats

1. gradient based approaches may not capture variations
   - given $f(x) = 1 - ReLu(1 - x)$, $\nabla_0 f$ and $\nabla_2 f$ have the same value
2. strong, local variations without any regularization scheme

Preliminaries
oooooo

Post-hoc explanations
oooo●ooo

Explanable by design programs
ooooooo

## Smoothgrad

SMOOTHGRAD [Smi+17] $\nabla_{x^*}(y)$ where $x^*$ is a gaussian neighborhood of $x$

$$\nabla_{x^*}(y) \approx \frac{1}{n} \sum_0^n \nabla_x f(x + \mathcal{N}(0, \sigma))$$

Preliminaries
oooooo
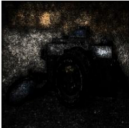
Post-hoc explanations
ooooo●oo

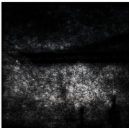Explanable by design programs
ooooooo

# Integrated gradients

Gradient on the line between $x$ and a baseline image $x^{'}$ [STY17]

$$\text{IG}_i = (x_i - x_i^{'}) \int_{\alpha=0}^{1} \nabla_{x_i} f(x^{'} + \alpha(x - x^{'})) d\alpha$$

usually computed using Riemann approaches

$$\text{IG}_i \approx (x_i - x_i^{'}) \sum_{k=0}^{m} \nabla_{x_i} f(x^{'} + \frac{m}{k}(x - x^{'})) * \frac{1}{m}$$

Preliminaries
oooooo

Post-hoc explanations
ooooooeo

Explanable by design programs
ooooooo

# Integrated gradients

Preliminaries
oooooo

Post-hoc explanations
ooooooo●

Explanable by design programs
ooooooo

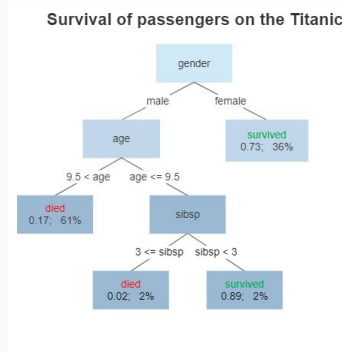# Wrapping up: empirical feature attribution approaches

1. usually only require gradient computation access;

2. provide attributions on the input space, but no direct exposition of the underlying decision process;

3. brittle, require sanity checks[Ade+18];

4. heavily rely on the program internal representation;

5. no validity domain;

# Explanable by design programs

Preliminaries
○○○○○○

Post-hoc explanations
○○○○○○○○

Explanable by design programs
○●○○○○○

# Decision trees



from Wikipedia https://en.wikipedia.org/wiki/Decision_tree_learning/

Issue: the deeper the tree, the less amenable it is to understand its decision

Integrating decision trees as the decision process for image classification

Preliminaries
○○○○○○

Post-hoc explanations
○○○○○○○○

Explanable by design programs
○○○●○○○

# Protoype based approaches - ProtoTrees



"White albatross"

Preliminaries
oooooo

Post-hoc explanations
oooooooo

Explanable by design programs
oooo●oo

# Protoype based approaches - ProtoTrees

1. learn "prototypes" $p$: part of the input set that are deemed representative for the prediction;

2. during inference, $M_i(x)$ are compared to $p_i$ using a similarity layer S;

Preliminaries
○○○○○○

Post-hoc explanations
○○○○○○○○

Explanable by design programs
○○○○○●○

# Protoype based approaches - Tackling tree complexity

ProtoTree have two hyperparameters that influence the decision tree:

1. the decision tree depth;
2. the pruning threshold;

Preliminaries
○○○○○○

Post-hoc explanations
○○○○○○○○

Explanable by design programs
○○○○○○●

# Prototype based approaches - Caveat

1. Still rely on the hypothesis that similarity in the feature space equals similarity in the input space;

2. Need retraining models;

## References

[Ade+18]   Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. "Sanity Checks for Saliency Maps". In: *Advances in Neural Information Processing Systems 32*. 2018, p. 11 (cit. on p. 15).

[Mil19]   Tim Miller. "Explanation in Artificial Intelligence: Insights from the Social Sciences". In: *Artificial Intelligence* 267 (2019), pp. 1–38 (cit. on p. 6).

## Bibliography ii

[Nau+23]   Meike Nauta, Jan Trienes, Shreyasi Pathak, Elisa Nguyen, Michelle Peters, Yasmin Schmitt, Jörg Schlötterer, Maurice van Keulen, and Christin Seifert. "From Anecdotal Evidence to Quantitative Evaluation Methods: A Systematic Review on Evaluating Explainable AI". In: *ACM Comput. Surv.* (Feb. 2023). Just Accepted. ISSN: 0360-0300. DOI: 10.1145/3583558. URL: https://doi.org/10.1145/3583558 (cit. on p. 5).

[Smi+17]   Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda B. Viégas, and Martin Wattenberg. "SmoothGrad: removing noise by adding noise". In: *ArXiv* abs/1706.03825 (2017) (cit. on p. 12).

[SP17]     Andrew D Selbst and Julia Powles. "Meaningful information and the right to explanation". In: *International Data Privacy Law* 7.4 (Dec. 2017), pp. 233–242. ISSN: 2044-3994. DOI: 10.1093/idpl/ipx022. eprint: https://academic.oup.com/idpl/article-pdf/7/4/233/22923065/ipx022.pdf. URL: https://doi.org/10.1093/idpl/ipx022 (cit. on p. 7).

## Bibliography iv

[STY17]     Mukund Sundararajan, Ankur Taly, and Qiqi Yan. "Axiomatic
            Attribution for Deep Networks". In: *Proceedings of the 34th
            International Conference on Machine Learning*. Ed. by Doina Precup
            and Yee Whye Teh. Vol. 70. Proceedings of Machine Learning
            Research. PMLR, June 2017, pp. 3319–3328. URL:
            `https://proceedings.mlr.press/v70/`
            `sundararajan17a.html` (cit. on p. 13).