

Xplique A Neural Networks Explainability Toolbox

pip install xplique

Thomas FEL*, Lucas HERVIER*, David VIGOUROUX, Antonin POCHE, Thibault BOISSIN, Philippe DEJEAN, Benjamin DEPORTE, David PETITEAU, Justin PLAKOO and all the DEEL team

🛑 🔴 🦲







\land C R I A Q



Thomas FEL

PhD Student, **ANITI & DEEL** Team. SNCF Research & Innovation. **Working on Explainable AI** *under the supervision of the Prof. Thomas Serre (Brown University,* **ANITI**)













Summary

1.A short introduction to XAI2.Attribution Methods + Metrics module3.Feature Visualization module4.Concepts module



1. Introduction **The Black-box problem**



1. Introduction **A Conceptual challenge**

"An Explanation is a set of statements [...] which clarifies the cause, the context and consequences of those facts. [...] The component of an explanation can be implicit and interwoven with one another"

Jess DRAKE, LOGIC

- An explanation provides information
- An explanation depends on domain knowledge
- An explanation helps to complete the knowledge of the domain

Lombrozo (2006), Hempel & Oppenheim (1948), Bechtel and Abrahamsen (2005), Chater and Oaksford (2006), Glennan (2002), Keil (2006)



1. Introduction Why do we need Explanations ?

Build trust in the model prediction

Elucidate important aspects of learned models

Help satisfy regulatory requirements and Certification process

Reveal bias or other unintended effects learned by a model

Detect and prevent failure cases

Debug & Train better model

1. Introduction **A Technical challenge**

Model

Feature Viz, Concept Activation Vector Explanation 'by design'

> **Predictions** Feature Attribution Feature Inversion



Nearest Neighbourhood Influence Function

Data

. . .

SNCF



1. Introduction **Taxonomy**



Consider the following general supervised learning setting: input space $\mathcal{X} \subseteq \mathbb{R}^d$, an output space $\mathcal{Y} \subseteq$, and a black-box predictor \mathbf{f} , which for some test input \mathbf{x} predicts the output $\mathbf{f}(\mathbf{x})$. We then define a feature attribution explanation as a function $\Phi : \mathcal{F} \times \mathbb{R}^d \to \mathbb{R}^d$ to

We then define a feature attribution explanation as a function Φ : $\mathcal{F} \times \mathbb{R}^d \to \mathbb{R}^d$, that given a black-box predictor **f**, and a test point **x**, provides importance scores $\Phi(\mathbf{f}, \mathbf{x})$.



Saliency Maps Symonyan & al (2013)[1]

$$\Phi = \nabla f(x) \implies \phi_i = \frac{\partial f(x)}{\partial x_i}$$

SmoothGrad Smilkov & al (2017)[2]

$$\Phi = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I\sigma)} [\nabla f(x + \epsilon)]$$
$$\Phi = \frac{1}{N} \sum_{i=0}^{N} \nabla f(x + \epsilon)$$

In an infinitesimal neighborhood (often not feasible), what are the features that most impact the output score ? Fast, really noisy, not really meaningfull (except for Lipschitz networks?) Free Parameters

As the name suggests, smoothes out the gradient by averaging the effect of small perturbations within a neighborhood around each pixels. <u>N~80 to have good results on 224x224 images. Slower than</u> <u>Saliency. Lot of variants (VarGrad, Squaregrad...)</u> Parameters: N, epsilon

\land C R I A O

Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps
 SmoothGrad: removing noise by adding noise



Integrated Gradients Sundarajan & al (2017)[1]

$$\Phi = (x - x_0) \int_0^1 \frac{\partial f(x_0 + \alpha (x - x_0))}{\partial x} d\alpha$$
$$\Phi = (x - x_0) \frac{1}{N} \sum_{i=0}^N \frac{\partial f(x_0 + \frac{i}{N} (x - x_0))}{\partial x}$$

Occlusion Ancona & al (2017)[2] $\phi_i = f(x) - f(x_{[x_i = x_0]})$ Averages the gradient values along the path from a baseline state to the current value. The baseline state is often set to zero.



N~80, Axiom Grounded, lots of tricks to leverage Integral approximation. What is a good baseline (x0)? Parameters: N, baseline

Sweeps a patch that occludes pixels over the image, and uses the variations of the model prediction to identify critical areas. hard to tune, black-box, need a baseline. Always good to have. Parameters: patch size, patch stride, baseline

[1] Axiomatic Attribution for Deep Networks
[2] Towards better understanding of gradient-based attribution methods for Deep Neural Networks





Rise Petsiuk & al (2018)[1]

$$\phi_i = \mathbb{E}[f(x \odot m) | m = 1]$$

Sobol Attribution Method

Thomas Fel*, Rémi Cadène*, Mathieu Chalvidal, Matthieu Cord, David Vigouroux, Thomas Serre, NeurIPS (2021)[2]

$$\phi_i = \frac{\mathbb{E}_{M_{\sim i}}(\operatorname{Var}_{M_i}(\boldsymbol{f} \circ \zeta(\boldsymbol{x}, \boldsymbol{M}) | \boldsymbol{M}_{\sim i}))}{\operatorname{Var}(\boldsymbol{f} \circ \zeta(\boldsymbol{x}, \boldsymbol{M}))}$$

[1] RISE: Randomized Input Sampling for Explanation of Black-box Models[2] Look at the Variance! Efficient Black-box Explanations with Sobol-based Sensitivity Analysis

Probes the model with randomly masked versions of the input image.



\land C R I A O

N~8000, Hard to compute. Baseline sensitive. Parameter sensitive. Random binary masks (i.i.d) are upsampled from low resolution Parameters : N, Low Res Grid, M=1 Probability

Estimates the variance of the model output (Sobol' indices) w.r.t a perturbation function using Quasi-Monte Carlo sampling. <u>N~4000. No Baseline.</u>

Parameters : Grid size, Perturbation function (e.g inpainting, blurring...)



CAM Zhou & al (2016)[1] & Grad-CAM Selvaraju & al (2017)[2]

$$\Phi = ReLU(\sum_{k=0}^{K} w^{(k)} A^{(k)})$$

W^(k) weight for each feature map

BROWN

Learning Deep Features for Discriminative Localization
 Visual Explanations from Deep Networks via Gradient-based Localization

For CAM (Conv + Global Average Pooling, one unit per class), the weight is 1 only for the feature map of the class else 0.

For Grad-CAM (any ConvNet), the weight is the avg of the gradients of each feature maps.

\land C R I A Q

 $w^{(k)} = \frac{1}{Z} \sum_{k=1}^{\infty} \sum_{k=1}^{\infty} w^{(k)} = \frac{1}{Z} \sum_{k$

Grad-CAM Selvaraju & al (2017)

$$w^{(k)} = \frac{1}{Z} \sum_{i} \sum_{j} \frac{\partial f(x)}{\partial A_{ij}^{(k)}}$$
$$\Phi = ReLU(\sum_{k=0}^{K} w^{(k)} A^{(k)})$$

Uses the gradients flowing back into the last convolution layer to generate a weight. This weight is used for the feature map. Aggregates all the weighted feature maps and removes the negative values. The results are then extrapolated (gives smooth results)

Quick to compute (1 forward, ~1/2 backward) and give good results. Lots of variants (Ablation-CAM, Grad-CAM++, Score-CAM, Shapley-CAM...). Parameters : a conv layer ~white-box state of the art.

\land C R I A Q

2. Feature Attributions 'Cool, but how to use that with Xplique ?'

import xplique from xplique.attributions import (Saliency, SmoothGrad,

GradCAM)

import xplique from xplique.attributions import (Saliency,

SmoothGrad, GradCAM)

2. Feature Attributions

•••

import numpy as np
from xplique.attributions import Saliency, GradCAM, SmoothGrad

model = tf.keras.Model(...)

initialize your explainers

explainers = [
 Saliency(model, batch_size=50),
 GradCAM(model, conv_layer='mixed3d'),
 SmoothGrad(model, nb_samples=50)

inputs = np.random.rand(10, 224, 224, 3)
labels = np.random.rand(10, 1000)

explainer share the same api once initialized
explanations = [
 explainer(inputs, labels) for explainer in explainers

ANITI DE

2. Feature Attributions Other data format!

Bounding Box (notebook coming soon!)

🐝 ΙΨΑΦΟ 🗛 C ΒΙΑ Ο

Confirmation bias.

Just because it makes sense to humans doesn't mean it reflects the evidence for prediction.

2. Feature Attributions Attribution methods can be manipulated

Fairwashing Explanations with Off-Manifold Detergent

Input Grad x o Grad IntGrad LRP Image: Strate strate

Figure 2. Example explanations from the original model g (left) and the manipulated model \tilde{g} (right). Images from the test sets of FashionMNIST (top) and CIFAR10 (bottom).

Interpretation of Neural Networks is Fragile

Feature-Importance Map

Interpretable Deep Learning under Fire

Figure 1: Sample (a) benign, (b) regular adversarial, and (c) dual adversarial inputs and interpretations on ResNet [22] (classifier) and CAM [64] (interpreter).

ͺ 🕍 ΙΫΑΟΟ

2. Feature Attributions Formal method for robust & efficient explainability

A C R I A Q

Don't Lie to Me! Robust and Efficient Explainability with Verified Perturbation Analysis

Thomas FEL*, Mélanie DUCOFFE*, David VIGOUROUX, Rémi CADENE, Mikael CAPELLE, Claire NICODÈME, Thomas SERRE, Pre-print under review

2. Feature Attributions **Fidelity metric**

ANITI

-1 99 0 45 2 6

SNCF

DEEL BROWN

-1.68 0.67 2.49

A C R I A Q

"pixel-flipping" procedure explanation techniques examples heatmaps 16 - sensitivity - simple Taylor score LRP 12 classification 10 (1) ? i C (1) average ີ 0 (1) compute current heatmap 10 0 5 15 20 (2) remove most relevant features # features removed

comparing

Figure 8: Illustration of the "pixel-flipping" procedure. At each step, the heatmap is used to determine which region to remove (by setting it to black), and the classification score is recorded.

Evaluating the visualization of what a Deep Neural Network has learned, Samek & al, 2015.

2. Feature Attributions Consistency & Representativity metrics

1-Lipschitz model gives 'better' explanations!

How Good is your Explanation? Algorithmic Stability Measures to Assess the Quality of Explanations for Deep Neural Networks Thomas FEL, David VIGOUROUX, Rémi CADENE, Thomas SERRE, WACV (2022)

2. Feature Attributions Metrics

•••

import xplique
from xplique.attributions import GradCAM
from xplique.metrics import Deletion

model = tf.keras.Model(...)

inputs = np.random.rand(10, 224, 224, 3) # images
labels = np.random.rand(10, 1000) # one-hot

explainer = GradCAM(model, batch_size=32)
explanations = explainer(inputs, labels)

2. Feature Attributions **Fidelity metric**

2. Feature Attributions

Attribution Method	Type of Model	Source	Tabular	Images	Time-	Tutorial	Attribution Metrics	Type of Model	Property	Source
Deconvolution	TF	Paper			WIP	CO Open in Colab	MuFidelity	TF	Fidelity	Paper
Grad-CAM	TF	Paper		✓ ✓	WIP	CO Open in Colab	Deletion	TF	Fidelity	Paper
Grad-CAM++	TF	Paper		~	WIP	CO Open in Colab	Insertion	TF	Fidelity	Paper
Gradient Input	TF	Paper	~	✓	WIP	CO Open in Colab	Average Stability	тг	Ctobility	Demor
Guided Backprop	TF	Paper	~	 ✓ 	WIP	CO Open in Colab	Average Stability		Stability	Paper
Integrated	TF	Paper	J		WIP	CO Open in Colab	MeGe	TF	Representativity	Paper
Gradients							ReCo	TF	Consistency	Paper
Kernel SHAP	Callable*	Paper	✓	✓	WIP	CO Open in Colab			, ,	
Lime	Callable*	Paper	✓	✓	WIP	CO Open in Colab	(WIP) e-robustness			
Occlusion	Callable*	Paper	✓	✓	WIP	CO Open in Colab				
Rise	Callable*	Paper	WIP	~	WIP	CO Open in Colab				
Saliency	TF	Paper	~	✓	WIP	CO Open in Colab				
SmoothGrad	TF	Paper	~	✓	WIP	CO Open in Colab				
SquareGrad	TF	Paper	~	✓	WIP	CO Open in Colab				
VarGrad	TF	Paper	~	~	WIP	CO Open in Colab				

2. Feature Attributions Fidelity doesn't mean Usefulness

KNITI

What I Cannot Predict, I Do Not Understand: A Human-Centered Evaluation Framework for Explainability Methods

 Thomas Fel^{1,3,4 *}
 Julien Colin^{1,3 *}
 Rémi Cadène^{1,2 †}
 Thomas Serre^{1,3}

 ¹Carney Institute for Brain Science, Brown University, USA
 ²Sorbonne Université, CNRS, France

 ³Artificial and Natural Intelligence Toulouse Institute, Université de Toulouse, France

 ⁴ Innovation & Research Division, SNCF

 {thomas.fel, julien.colin, remi.cadene}@brown.edu

Which explanation is the most useful to humans?

Prediction: **Red Fox** Saliency(0.92)

92) Occlusion (0.89)

 $\mathsf{Grad} ext{-}\mathsf{CAM}(0.89)$

4.2. Faithfulness metric as a proxy for Usefulness?

Figure 6. *Utility* vs Faithfulness correlation. The utility scores on the two datasets Husky vs. Wolf (point marker) and Leaves (square marker) are plotted showing a poor or anti-correlation between the two measures. Concerning the ImageNet dataset (triangle marker), the *Utility* scores are insignificant since none of the methods improves the baseline.

D 🛦 CRIA Q

2. Feature Attributions Tips & Tricks

- Always use multiple methods
- Sobol, Rise, Grad-Cam and Smoothgrad work in more cases
- Clip percentile before using `imshow(.)`
- Beware of confirmation bias
- The *quality* of the explanation does not depend only on the method! Robust models (e.g. 1-Lipschitz model) seem to give better explanations

3. Feature Visualization

3. Feature Visualization **Overview**

3. Feature Visualization **Parameterization**

FIGURE 1: As long as an image parameterization is differentiable, we can backpropagate (<--) through it.

3. Feature Visualization

7.19 - Betulaceae

5.56 - Betulaceae

5.31 - Betulaceae

5.20 - Betulaceae

🐝 IVADO 🗛 CRIAQ

3. Feature Visualization

5.20 - Betulaceae

Where

Paleo Al Project, Serre LAB Ivan Felipe Rodriguez*, Thomas FEL*, Thomas Serre, Peter Wilf

BROWN ANITI BROWN CONCEPTION

What

3. Feature Visualization **Explaining logits**

4. Testing with Concept Activation Vectors **Beyond feature Attribution**

4. Testing with Concept Activation Vectors **TCAV**

5. What's Next

Research

• Guiding model explanations during training – for models with human-like explanations

A C R I A Q

- Explaining models using data (Recursive Flow Explainable AI project)
- Influenciae Library (Influence Function library)

Implementation

- Adapt Xplique for NLP
- More Attribution methods
- Automatic Concept Extraction

🖉 🖬 BROWN

5. See also

Libraries

- <u>Lucid</u> the wonderful library specialized in feature visualization from OpenAI.
- <u>Captum</u> the Pytorch library for Interpretability research
- <u>Tf-explain</u> that implement multiples attribution methods and propose callbacks API for tensorflow.
- <u>Alibi</u> Explain for model inspection and interpretation
- <u>SHAP</u> a very popular library to compute local explanations using the classic Shapley values from game theory and their related extensions implementation

Tutorials

• Interpretable Machine Learning by Christophe Molnar.

BROWN

- Interpretability Beyond Feature Attribution by Been Kim.
- Explaining ML Predictions: State-of-the-art, Challenges, and Opportunities by *Himabindu Lakkaraju, Julius Adebayo and Sameer Singh*.
- A Roadmap for the Rigorous Science of Interpretability by *Finale Doshi-Velez*.
- **DEEL** White paper a summary of the **DEEL** team on the challenges of certifiable AI and the role of explainability for this purpose

• (AAAI 2022) On Explainable AI - XAI Coding And Engineering Practices – Showcase Xplique ! 💐

Thank you for your attention! Any questions?

Thomas FEL*, Lucas HERVIER*, David VIGOUROUX, Antonin POCHE, Thibault BOISSIN, Philippe DEJEAN, Benjamin DEPORTE, David PETITEAU, Justin PLAKOO and all the DEEL team

ACRIAO

Sivado

