

Lab report - Interpretability evaluation  
framework:  
Aspects in Machine learning and hydrology  
domain applications

School of Computer Science and Engineering - Institute of  
Earth Sciences  
Hebrew University

Ronen Rojas

September 2023

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Literature review</b>	<b>5</b>
2.1	What is an attribute? . . . . .	5
2.2	Interpretability methods . . . . .	6
2.2.1	Grad . . . . .	6
2.2.2	Input $\times$ Gradient . . . . .	6
2.2.3	Integrated Gradients . . . . .	6
2.2.4	Local Integrated Gradients . . . . .	6
2.2.5	Smooth Grad . . . . .	7
2.2.6	VarGrad . . . . .	7
2.2.7	Smooth Grad Square SG-SQ . . . . .	7
2.2.8	Shapely Value Sampling . . . . .	8
2.2.9	LRP . . . . .	8
2.2.10	DeepLift . . . . .	9
2.2.11	Feature Ablation . . . . .	9
2.2.12	Deconvolution . . . . .	9
2.2.13	Guided backpropagation . . . . .	9
2.2.14	CAM . . . . .	9
2.2.15	Guided Grad CAM . . . . .	10
2.2.16	Lime . . . . .	10
2.3	Interpretability methods frameworks and metrics . . . . .	11
2.3.1	Feature importance removal framework . . . . .	11
2.3.2	Interpretability as a debug tool framework . . . . .	14
2.3.3	Human subject study . . . . .	16
2.3.4	Axiomatic approach framework . . . . .	17
2.3.4.1	Sensitivity, implementation invariance, linear- ity and Symmetry preserving . . . . .	17

2.3.4.2	Continuity . . . . .	17
2.3.4.3	Input invariance . . . . .	18
2.3.4.4	Monotonicity . . . . .	18
2.3.4.5	Completeness and sensitivity-n . . . . .	19
2.3.5	Adversarial framework . . . . .	20
2.4	Recognized gap . . . . .	21
<b>3</b>	<b>Research objective</b>	<b>22</b>
<b>4</b>	<b>Data and Models</b>	<b>24</b>
4.1	Synthetic Data . . . . .	24
4.1.1	Weather Generator . . . . .	24
4.1.2	DREAM . . . . .	25
4.2	Models . . . . .	26
4.2.1	Feed Forward Neural Network . . . . .	26
4.2.2	The Support Vector Classifier (SVC) . . . . .	28
4.2.3	RNN . . . . .	29
4.2.3.1	RNN - Introduction . . . . .	29
4.2.4	LSTM . . . . .	30
4.2.4.1	LSTM - Principles . . . . .	30
4.2.4.2	LSTM- Multi-cell . . . . .	32
4.2.4.3	LSTM - in context of rainfall . . . . .	32
<b>5</b>	<b>Results</b>	<b>34</b>
5.1	Sanity Check . . . . .	34
5.2	Dream synthetic Framework . . . . .	35
5.2.1	Dream synthetic model . . . . .	35
5.2.2	Interpretability method . . . . .	35
5.2.3	Framework evaluation results . . . . .	35
<b>6</b>	<b>Conclusions</b>	<b>36</b>

# Chapter 1

## Introduction

In recent years, machine learning models have become the standard tool for almost every task of estimation, prediction, and classification. Machine learning models learn for themselves what the important features in the data are by an optimization process called training. This training process has replaced the tedious task of manually extracting insights from the data .

The immense progress in tools, software packages, and hardware enhancements has made machine learning widely accessible. All that is needed is a strong computational framework and good, clean data. The performance of these kinds of machine learning efforts is unprecedented and they outperform humans on various tasks ([17], [22], [56], [32], [46]).

However, with great power comes great responsibility (as they say). How do we explain a model prediction in a sensible manner? Most machine learning models are considered to be "*black boxes*" and humanly incomprehensible. When considering a trained machine learning model we are now interested more in "why" it predicted something and less in "what" it predicted. This shift has given rise to significant developments in the field of interpretability methods ([33], [62], [60], [59], [40], [30], [39], [57]).

Interpretability tackles the issue of trying to gain more insight of the prediction mechanism of the model. It also tries to answer the following questions:

- *Why did the model made the prediction?*

- *Do we have a human-friendly explanation for this prediction?*
- *Can we visualize this explanation?*

Some interpretability methods try to tackle the idea of "why?" instead of "what?" from an axiomatic approach ([53], [19], [34], [21]). Some do so from a visualization perspective ([59], [61], [58], [41]). Even though visualization results can be insightful to humans, such evaluations are not objective and might contain biases.

There are evaluation frameworks that try to quantify how good an interpretability method is ([54], [38] [20] [34],[27], [6]). But since evaluation of interpretability methods can get convoluted and untrustworthy, it is difficult to define a good explanation. As there is no ground truth for that, there is no right or wrong answer to this conundrum.

Evaluation of interpretability methods is the main motivation for this research. In this paper I will detail methods for refinement and exploration of better understanding our interpretability methods and their limitations. I will evaluate the interpretability methods in a concrete manner and apply them to the scientific domain of hydrology.

# Chapter 2

## Literature review

In this chapter I will review interpretability methods that have gained popularity in recent years. I will also review frameworks of evaluating interpretability methods and what evaluation metrics are used. I will constrain the scope of this work only to model attribution-based explanations [62], these kind of explanations are local and post-hoc i.e., given a trained model and an input for the prediction, attribution measures or ranks each input feature to explain the model’s prediction.

### 2.1 What is an attribute?

Formally, if we have a model that was trained to perform a specific task:

$$\begin{aligned} y &= f_{\theta}(x) \\ x &\in \mathbb{R}^n, y \in \mathbb{R} \end{aligned}$$

An attribution method provides an explanation  $E$  for the input  $x$ :

$$\begin{aligned} E &= \text{attrib}(f_{\theta}, x) \\ x &\in \mathbb{R}^n, E \in \mathbb{R}^n \end{aligned}$$

$E$  gives an explanation power rank for each input coordinate for the prediction  $y = f_{\theta}(x)$ . Each coordinate of the input for a model is usually called a feature.

## 2.2 Interpretability methods

For this section we'll assume we have a model  $y = f(x)$  and an attribution-based interpretability method that outputs an explanation  $E$  for the prediction  $y$  with the same dimensions as the input  $x$ .

### 2.2.1 Grad

The most straight forward interpretability method [8] [47], also known as saliency map:

$$E_{\text{GRAD}}(f, x) = \nabla_x f(x)$$

### 2.2.2 Input $\times$ Gradient

Saliency map can be substantially improved [45] by simply multiplying the gradient with the input, this is basically first-order Taylor approximation in a sense:

$$E_{\text{INPUT} \times \text{GRAD}}(f, x) = x \odot \nabla_x f(x)$$

### 2.2.3 Integrated Gradients

Sundararajan et al. [53] proposed this interpretability method, abbreviated IG, denote the coordinate  $i$  as of input  $x$  as  $x_i$ :

$$[E_{\text{IG}}(f, x)]_i = (x_i - \bar{x}_i) \times \int_{\alpha=0}^{\alpha=1} \frac{\partial f(x + \alpha(\bar{x} - x))}{\partial x_i} d\alpha$$

where  $\bar{x}$  is the baseline input for the method, usually something neutral in the data-set e.g. the black image in images data-sets.

### 2.2.4 Local Integrated Gradients

Ancona et al. [5] suggested a variant of any method that has an multiplication term by the input. If the aforementioned term is omitted then it's call local method whereas the original is called global, in case of the IG method we'll get:

$$[E_{\text{local-IG}}(f, x)]_i = \int_{\alpha=0}^{\alpha=1} \frac{\partial f(x + \alpha(\bar{x} - x))}{\partial x_i} d\alpha$$

According to [5] *global attribution methods* describe marginal effect of a feature on the output with respect to a baseline, whereas *local attribution methods* describes changes for infinitesimally small perturbations around the original input.

### 2.2.5 Smooth Grad

Smilkov et al. [49] proposed adding some noise to the attribution process that will help smooth the saliency maps :

$$E_{\text{SGRAD}}(f, x) = \frac{1}{N} \sum_{i=1}^N \nabla_x f(x + g_i)$$

Where  $g_i$  are noise vectors  $g_i \sim \mathcal{N}(0, \sigma^2)$  are drawn i.i.d. from a normal distribution. One can generalize this method, given an attribution method  $E(f, x)$  its smooth counterpart is:

$$E_{\text{Smooth}}(f, x) = \frac{1}{N} \sum_{i=1}^N E(f, x + g_i)$$

### 2.2.6 VarGrad

Adebayo et al. [3] suggests using the variance operator similarly to the SGRAD. obviously this method can be generalized to any attribution method, denote the variance operator  $\mathcal{V}$ :

$$E_{\text{Variance}}(f, x) = \mathcal{V}[E(f, x + g_i)]$$

### 2.2.7 Smooth Grad Square SG-SQ

Hooker et al. [20] proposed the method element-wise square of SGrad

$$E_{\text{SG-SQ}}(f, x) = E_{\text{SGRAD}}(f, x) \odot E_{\text{SGRAD}}(f, x)$$

Where  $\odot$  is the element-wise product.



### 2.2.8 Shapely Value Sampling

Shapley values provide is a concept in game theory to calculate "fair" distribution of winnings for each player members of the winning coalition in a cooperative game by measuring marginal contribution to the final outcome ([10], [51] [42], [31], [1]).

Denote  $\mathcal{M}$ , the set of players i.e. the grand coalition,  $\mathcal{S}$  the set of players in the partial coalition and  $v$  is the contribution map between subsets of players, the Shapley value  $\phi_i$  for coalition member  $i$  is defined as follows:

$$\phi_i(v) = \sum_{\mathcal{S} \subseteq \mathcal{M} \setminus \{i\}} \binom{|\mathcal{M}| - |\mathcal{S}| - 1}{|\mathcal{S}|} \left( v(\mathcal{S} \cup \{i\}) - v(\mathcal{S}) \right)$$

Based on this definition an efficient perturbation process for an attribution-based explanation can be devised [52]. By taking a random permutation of the input features  $x_i, i \in [N]$ , adding them one-by-one to the given baseline  $x'$  and defining  $v = f(x) - f(x')$ . Repeating this process  $n$  times, each time choosing a new random permutation of the input features will yield to an empiric estimation of the real Shapley value.

### 2.2.9 LRP

Layer-wise relevance propagation abbreviated LRP [27] [34], is a method that tries to estimate every neurons' layer relevancy. The model output score represents the initial relevance which then is being propagated backwards in a set of pre-defined rules the are applied sequentially to all layers of the model. When the propagation reaches the first layer we will get the wanted attribution of the input. This backward propagation mechanism resembles somewhat of a numerical differentiation [60]:

$$[E_{\text{LRP}}(f, x)]_i = x_i \cdot \frac{\partial^g f(x)}{\partial x_i}, \quad g(z) = \frac{\sigma(z)}{z}$$

where  $\sigma(z)$  is the non-linearity in the network and  $g(z)$  is a replacement the derivative of  $\sigma(z)$ .

### 2.2.10 DeepLift

Shrikumar et al. [44] tried to devise a more intricate backward propagation mechanism called Deep Learning Important FeaTures, DeepLift in short.

The mechanism used pre-defined rules for each type of non-linearity and layer type. Their novelty was introducing a reference point in their propagation mechanism.

$$[E_{\text{DeepLift}}(f, x)]_i = (x_i - x_i^{\text{ref}}) \cdot \frac{\partial g f(x)}{\partial x_i}, \quad g(z) = \frac{\sigma(z) - \sigma(z^{\text{ref}})}{z - z^{\text{ref}}}$$

Both LRP and DeepLIFT have gained the nickname *discrete gradient methods* [60].

### 2.2.11 Feature Ablation

A perturbation based approach to computing attribution, replacing each input feature  $x_i$  with a given baseline  $x'_i$ , and computing the difference in output.

### 2.2.12 Deconvolution

Deconvolution [58] computes the gradient of the target output with respect to the input, but instead of using the gradients of ReLU functions it is computed taking ReLU of the output gradient, essentially only propagating non-negative gradients [29].

### 2.2.13 Guided backpropagation

Guided backpropagation is a slight variant of the Deconvolution method [50]. The difference is mainly how the gradients of non-linearities are being calculated. The method adds an additional guidance signal from the higher layers to the usual backpropagation, this in turn prevents backward flow of negative gradients.

### 2.2.14 CAM

Class Activation Mapping [61], CAM in short, is an attribution method which tries to estimate the importance of each unit in the feature maps

i.e. the rectified result of the convolutional by using global average pooling (GAP) in CNNs.

Since the feature maps have a lower resolution the activation maps are up-sampled, usually the last convolutional is the one that is considered for the attribution.

### 2.2.15 Guided Grad CAM

Gradient-weighted Class Activation Mapping [41], in short Guided Grad CAM, combines two methods. Guided Grad CAM Computes element-wise product of guided backpropagation attributions with upsampled (non-negative) GradCAM attributions:

$$E_{\text{Guided Grad CAM}}(f, x) = E_{\text{Guided Backprop}}(f, x) \odot E_{\text{CAM}}(f, x)$$

### 2.2.16 Lime

Local Interpretable Model-agnostic Explanations [37] (LIME). This method trains an interpretable surrogate model around the input example and using model evaluations at these points to train a simpler interpretable ‘surrogate’ model, for example a linear model.

## 2.3 Interpretability methods frameworks and metrics

### 2.3.1 Feature importance removal framework

There has been a line of study to treat attribution-based explanations as feature importance ([38] [20] [34],[27], [6], [13], [28], [7], [35]). This line of works asks what will happen if we "remove" these features from the input strategically? what will happen to model predictions?, if we remove from most to least important, ranked from the explanations) we will probably see degradation. This can be evaluated by a certain metric. On the other we can ask the opposite question what will happen if we remove the least important features ? can the model sustain it's merits ? until what percentage of the input is removed ?

By comparing the degradation or robustness of a certain model we can compare different interpretability method attributes, if explanations from one method "broke" the model quicker in a sense this method is better and vice versa if we were able to sustain accuracy or other performance metrics of a model when removing least important features from explanations ranks. Motavon et al. [34] coined this desired property as *Explanation Selectivity*.

Samec et al. [38] have devised an input perturbation process called *most relevant first*, abbreviated MoRF for the trained model  $f(x)$ . In this process we remove information from the input by setting the coordinates (or a region of surrounding pixels e.g.  $9 \times 9$  neighborhood) to a uniform distribution. MoRF process is done by the individual attribution ranking  $e$  of a given input  $x$ , this is done for  $L$  iterations:

$$\begin{aligned} x_{\text{MoRF}}^{(0)} &= x \\ \forall 1 \leq k \leq L \quad x_{\text{MoRF}}^{(k)} &= \text{Remove}(x_{\text{MoRF}}^{(k-1)}, e) \end{aligned}$$

The metric used is the area over the MoRF perturbation curve (AOPC):

$$\text{AOPC} = \left\langle \sum_{k=0}^L f(x) - f(x_{\text{MoRF}}^{(k)}) \right\rangle_{p(x)}$$

Where  $\langle \cdot \rangle_{p(x)}$  means we average over all inputs in the data set, a larger AOPC means a sharper decrease of the MoRF process ( $f(x_{\text{MoRF}}^{(k)})$ ) thus a better feature importance for the attribution-based explanation.

The opposite approach then is the perturbation process *least relevant first*, abbreviated LeRF. The proposed metric is the area between the perturbations curves:

$$\text{ABPC} = \left\langle \sum_{k=0}^L f(x_{\text{LeRF}}^{(k)}) - f(x_{\text{MoRF}}^{(k)}) \right\rangle_{p(x)}$$

In this case the LeRF process needs to sustain the information as much as possible thus **larger** area is wanted.

In [40] and [27] to evaluate and compare different interpretability methods "pixel-flipping" was used. Arras et al. [6] also used LRP evaluation and simply "deleted" word for their model into by setting the corresponding word embedding to zero in order, Lundberg et al. using different strategies on how to remove the input features [28] for evaluations explanations for Tree-based machine learning models.

Hooker et al. [20] opine that any removal of information from the input is disruptive since the metrics will be used on modified input that is not from the data set original distribution, thus the generalization errors of "out of sample" inputs for the model becomes dominant.

To mitigate this case [20] presented a training strategy, RemOve and Re-train (ROAR). This ensures that the metric used, which in the case was a simple test-set accuracy, is used in sample from the distribution the model was trained on. The opposite approach was to Keep And Retrain (KAR). In this approach, **minimizing** degradation to test-set accuracy is desirable

DeYoung et al. [13] used a similar approach in NLP explanations, called rationales in the NLP jargon, evaluation benchmark. [13] proposed ERASER, which stands for **E**valuation **R**ationales **A**nd **S**imple **E**nglish **R**easoning, a framework that utilized the metric of AOPC (among others metrics) with the only exception that re-training is not necessary since the input size for the model (BERT [12]) can vary so removing tokens (input features) can be

done seamlessly.

Chefer et al. [11] used the ERASER framework for the NLP domain and the "removing" pixels for the vision domain in their evaluation of interpretability methods of self-attentions transformers. In both cases, the metric used was the area-under-the-curve (AUC) perturbation process.

### 2.3.2 Interpretability as a debug tool framework

Another line of study is to let the interpretability method to be a debug tool for the model or training process itself ([4], [2]). If an interpretability method is sensitive to a model that trained or tested in an ill manner compared to a model that was trained in a nominal way it means that it can better explain the model and in a way could be a sanity check for the interpretability method itself.

Adebayo et al. [4] proposed 2 test frameworks:

1. *Model parameter randomization test*: Apply attributions-based method on a trained model and on the same architecture but on a randomly initialized untrained model.
2. *Data randomization test*: Apply attributions-based method on a trained model and on a copy of the mode trained on a copy of the data set in which all labels were randomly permuted.

*Model parameter randomization* tests whether attributions method outputs differ substantially between the 2 models, if the output are similar it means the method is insensitive to model parameters which in a sense does not bode well to the explanation goals which are to understand why did the model made its prediction. It's important to note that the randomization of the model parameters was done in a cascading fashion from top to bottom layers.

*Data randomization* tests again as before the difference between the 2 models, An insensitivity to the permuted labels reveals that the method does not depend on the relationship between input and labels which is not a desired property for an attribution-based method.

[4] used visualization to show the the differences in the attributes, but also used more rigorous ways namely *similarity metrics* as follows:

1. Spearman rank correlation with / without absolute value.
2. Structural similarity index measure (SSIM [55]).
3. Pearson correlation of the histogram of gradients (HOGs)

In his later work Abedeyo et al. [2] used a more elaborate scheme for tampering with the data-set, the training or evaluation process itself, simply named "bugs":

1. *Data contamination bugs*
  - (a) *Labeling errors* - Incorrectly labeled data , Similar to previous work.
  - (b) *Spurious Correlation* - Make the model associate uncorrelated reason to the task, e.g. blue sky backgrounds with the bird class.
2. *Model contamination bugs* - re-initialization of model weights, similar to previous work.
3. *Test-time contamination bugs* - Out of distribution (OOD) samples, domain shift for the data-set.

*Spurious bug implementation* was implemented by placing all birds onto one of the sky backgrounds and all dogs onto a bamboo forest background. Logically, explanations on a model that was trained on a data-set like this would identify this correlation i.e. attributing most of the background pixels to the class prediction. In a sense this is a ground truth for the explanation output itself and a "debug" tool for the attribution method.

*Test-time contamination bugs* assess the ability of attributions to diagnose domain shift, e.g. the attribution for an MNIST input from a model trained on MNIST, to an attribution for the same input but derived from an output of a model trained on a different data-set.

The metric used for the attributes is SSIM to compare the similarity between the explanation in the aforementioned tests.



### 2.3.3 Human subject study

Incorporating human study for evaluation of interpretability methods is a very natural ([2], [16], [37]). The explanations coming out of these methods are meant to make black-box machine learning in a sense more understandable and human comprehensible.

Adebayo et al. [2] conduct a 54-person study to assess whether end-users can recognize the bugs in the tests according to the attributions. Evidently this approach showed the people are more biased towards the model predictions even in the presence of its attributions.

Hase et al. [16] tried to test *simulatability* on interpretability methods for machine learning models. A model is simulatable when a person can predict its behavior on new inputs. For this task [16] conducted a two-fold 39-person human-subject study:

1. *Forward simulation* - given an input and an “explanation,” users must predict what a model would output for the given input
2. *Counterfactual simulation* - are given an input, a model’s output for that input, and an “explanation” of that output, and then they must predict what the model will output when given a perturbation of the original input.

### 2.3.4 Axiomatic approach framework

In this section we'll present the approach of desired properties which attribution-based interpretability ought to have. These properties are mostly common sense and stem from a more theoretical thought process when devising a scheme for an attribution method hence we call them *axioms* in this context.

Although these *axioms* are not comparable between methods per se they're still worth to mentioned as it can be used as a rigorous framework to evaluate these methods and some can become more quantifiable in a sense as we'll explain later in this section.

#### 2.3.4.1 Sensitivity, implementation invariance, linearity and Symmetry preserving

Sundararajan et al. [53] introduced these axioms in his novel attributed-based method *integrated gradients*:

- *Sensitivity* - If the function implemented by the deep network does not depend (mathematically) on some variable, then the attribution to that variable is always zero.
- *Implementation invariance* - Attributions should be identical for two functionally equivalent networks.
- *Linearity* - Attributions method should preserve any linearity within the network.
- *Symmetry preserving* - If 2 inputs to the network are symmetrical  $F(x, y) = F(y, x)$ , so should be their corresponding attributions.

#### 2.3.4.2 Continuity

Montavon et al [34] introduced the desired axiomatic property of *continuity* as follows:

- *Continuity* - If two data points are nearly equivalent, then the explanations of their predictions should also be nearly equivalent.

This axiom can be also quantified as follows, denote an attribution method  $E$ , two inputs  $x, x'$  and their explanations respectively  $E, E'$  :

$$\max_{x \neq x'} \frac{\|E - E'\|_1}{\|x - x'\|_2}$$

Bhatt et al. [9] used a similar definition but used the term *low/average sensitivity*

#### 2.3.4.3 Input invariance

Kindermans et al. [21] introduced the *Input invariance* axiom:

- *Input invariance* - The attribution-based method should mirror the sensitivity of the model with respect to transformations of the input

For example a constant shift in the input with two model that were trained on the original data and the shifted data and have the **same predictions** should have the **same attributions**.

#### 2.3.4.4 Monotonicity

Although *monotonicity* defined as metric in [36] it stems from an axiomatic approach on the mechanism of attribution methods that assign an importance value to each feature. From a more theoretical approach a more concrete metric was proposed namely *monotonicity*.

Nguyen et al. [36] posit that feature importance rank attributions should follow a desired property, denote the explanation  $E_i$  for feature  $i \in [N]$ ,  $y^* = f(\mathbf{x}^*)$ :

$$|E_i| \propto \mathbb{E}(l(y^*, f_i) | \mathbf{x}_{-i}^*) = \int_{\mathcal{X}_i} l(y^*, f_i(x_i)) p(x_i) dx_i$$

Where  $l$  is the loss function, density probability function for feature  $x_i$  is  $p(x_i)$  and  $f_i$  is the restriction of the function  $f$  to the feature  $i$  obtained by fixing the other features at the values  $\mathbf{x}_{-i}^* = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_N)$ , so the

monotonicity metric for attribution  $E$  is defined as the Spearman's correlation coefficient  $\rho_S$ :

$$\begin{aligned}\mathbf{e} &= [|E_1|, \dots, |E_i|, \dots, |E_N|] \\ \mathbf{f} &= [\mathbb{E}(l(y^*, f_1)|\mathbf{x}_{-1}^*), \dots, \mathbb{E}(l(y^*, f_i)|\mathbf{x}_{-i}^*), \dots, \mathbb{E}(l(y^*, f_N)|\mathbf{x}_{-N}^*)] \\ \text{monotonicity}(E, f, \mathbf{x}^*, y^*) &= \rho_S(\mathbf{e}, \mathbf{f})\end{aligned}$$

#### 2.3.4.5 Completeness and sensitivity-n

Sundararajan et al. [53] proved that the attribution-based method *iterated gradients* holds the axiom of *completeness*:

$$\sum_{i=1}^N IG(\mathbf{x}) = f(\mathbf{x}) - f(\mathbf{x}')$$

Where  $\mathbf{x}'$  is a baseline input that is used for the attribution calculation and we'll be elaborated in later sections.

Ancona et al. [5] generalized the *completeness* axioms namely, *sensitivity-n*:

- *sensitivity-n*: For any subset of features of cardinality  $n$ ,  $\mathbf{x}_s = [x_1, \dots, x_n] \subseteq \mathbf{x}$  it holds  $\sum_{i=1}^n E_i = f(\mathbf{x}) - f(\mathbf{x}_{[\mathbf{x}_s=0]})$

It easy to see that *completeness* axiom is a private case of the *sensitivity-n* axiom when  $n = N$ . Bhatt et al.

[9] introduced a similar notion of *faithfulness* where the sum of the attributions and the difference in output when setting those features to a reference baseline should be proportionality correlated

### 2.3.5 Adversarial framework

Adversarial framework is a well researched subject for finding vulnerabilities in neural networks performance. Adversarial framework for interpretability methods asks whether the method at hand can be fooled via some sort of manipulation. Although some adversarial techniques did not quantify numerically how an interpretability method is susceptible to a specific attack it is rather important to mention this framework as they give some notion of robustness to the methods.

Heo et al. [18] showed that LRP, Grad-CAM, and simple saliency map , can be easily fooled with model manipulation. They also suggested a quantitative metric , *Fooling Success Rate* (FSR). FSR measures how much an interpretability method is prone to aforementioned adversarial manipulation.

Slack et al. [48] proposed a novel scaffolding technique that hides classifier biases by allowing an adversarial entity to craft an arbitrary desired explanation for the LIME and methods. Dombrowski et al. [14] exploited certain geometrical properties of neural networks to manipulate an explanation arbitrarily with hardly perceptible differences in the input.

## 2.4 Recognized gap

Many frameworks tackle tasks like images classification and NLP tasks. They also use popular widely researched existing data sets. I recognized a twofold gap:

1. There is no attempt to extract more than the perceived feature importance of the data and model prediction in a method
2. There is no attempt to utilize and evaluate interpretability methods in hydrological setup.

In a sense I want to know whether the interpretability method was able to gain some more insights from the data, for instance some hidden modalities in the data and labels that are more easily perceived with the explanations rather than the data or can be learned with fewer samples.

The incentive of using hydrological domain, is that sometimes the explanation of a trained model can reveal a pattern that even expert judgment could not. Also incorporating machine learning models in general is sparse so it is rather interesting to explore and evaluate the use of these models for regression task in hydrology.

# Chapter 3

## Research objective

In this report I will present a novel framework for evaluating an interpretability method. In this framework I will try to see if inside the explanations there is more than meets eyes.

I will ingrain the data set with different sets of scenarios that adhere to a general mechanism, this is the way to generate the data-label distribution. But also the process will insure that the different scenarios will create different characteristics that are not so noticeable in the data itself but rather apparent in their counterpart explanations. I will quantify the contrast between these different point of views.

These scenarios characteristics can be perceived differently in different method of attributes. it can be easily compared by devising a classification tasks for the scenarios and compare attribution methods.

Taking advantage of the fact that the attributes of model and the inputs for the model have the same dimensions, one can devise two type classifiers that have almost the same architecture. One type that uses the real data and one that uses to explanations. It can be stated that by comparing the two types of the classifiers is comparing apples to apples so to speak.

The aforementioned evaluation approach can be summarized in the following scheme 3.1:

1. Generate a 2-modal data-set

2. Train a machine learning model on the data-set
3. ‘Train two classifiers for the modalities of the data-set:
  - (a) A clean view – training only with the data-set
  - (b) An interpreted view – training only with attributes of the trained model
4. Evaluate interpretability method by comparison of the clean classifier to interpreted classifier with simple classification metrics

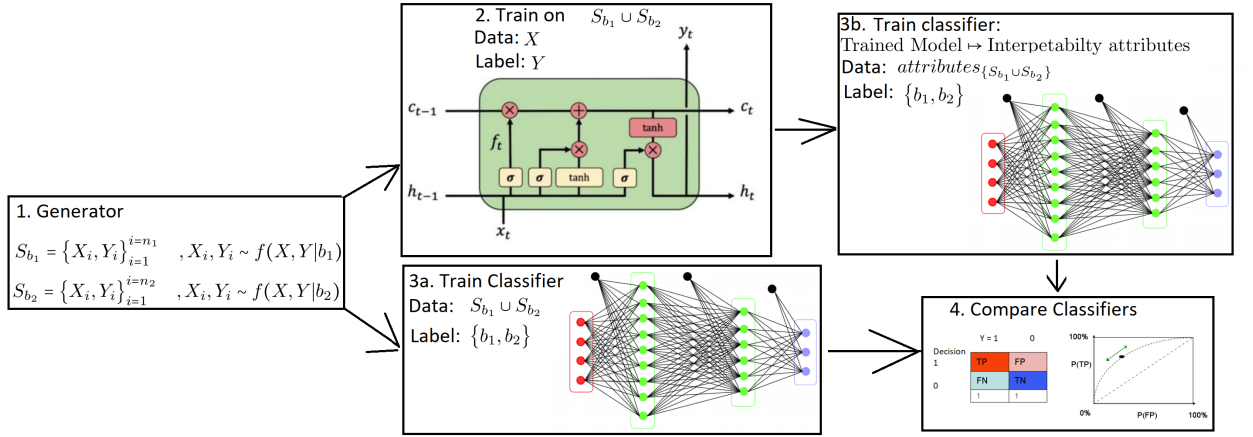


Figure 3.1: Scheme



# Chapter 4

## Data and Models

### 4.1 Synthetic Data

Our machine learning model task is to predict the overground runoff from synthetic weather data for precipitation and potential evapotranspiration (PET) inputs. The creation of the data is done in two phases:

1. Weather Generator (Inputs) - daily time series of the precipitation  $\left[\frac{mm}{day}\right]$  and PET  $\left[\frac{mm}{day}\right]$ .
2. DREAM (labels) - daily time series of runoff  $\left[\frac{mm}{day}\right]$ .

#### 4.1.1 Weather Generator

The weather generator is essentially a Markov chain  $\{X_i\}_{i=0}^n, X_i \in \{0, 1\}$  where:

$$X_{i|i-1} \sim \begin{cases} \text{Bernoulli}(\delta_1) & X_{i-1} = 0 \\ \text{Bernoulli}(\delta_2) & X_{i-1} = 1 \end{cases}$$

The precipitation daily values are:

$$Y_i = X_i \times Z_i$$

Where  $Z_i \stackrel{i.i.d}{\sim} \text{Weibull}(\lambda, k)$  and its density function is defined as follows:

$$f(x; \lambda, k) = \begin{cases} \frac{k}{\lambda} \left(\frac{x}{\lambda}\right)^{k-1} e^{-(x/\lambda)^k}, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

The PET values are drawn  $\overset{i.i.d}{\sim} \mathcal{N}(\mu, \sigma^2)$ .

For the two modalities of the data set we use different distribution parameters mode 1 (dry) and mode 2 (wet):

Table 4.1: Weather generation parameters

<b>Parameter</b>	<b>Mode 1</b>	<b>Mode 2</b>
$\delta_1$	0.2	0.8
$\delta_2$	0.5	0.8
$k$	0.5	1.5
$\mu$	3.0	3.0
$\sigma$	1.0	1.0

$\lambda$  was used as a normalization factor to make the annual mean precipitation equal between the two modes.

#### 4.1.2 DREAM

For the runoff model we’ve used the well established Hydrometeorological daily recharge assessment model (DREAM) [43]. In short this model tries to delineate between the runoff discharge and the ground water recharge, this model has displayed good empirical results.

## 4.2 Models

### 4.2.1 Feed Forward Neural Network

We've used Deep Neural Networks (DNNs) for classification tasks. The network consist of multiple layers of interconnected nodes. In the context of classification, DNNs can automatically learn to extract relevant patterns and representations from complex data.

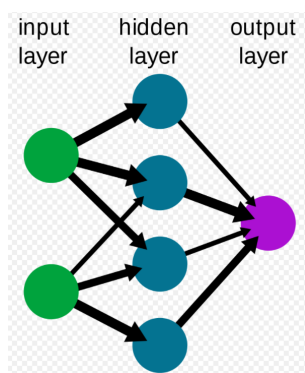


Figure 4.1: Neural Network Architecture

Every layer is an affine layer that ends with a non-linear point wise activation function e.g.: sigmoid, Relu.

An Affine layer is defined as follows:

$$y = Ax + b, x \in \mathbb{R}^n, b, y \in \mathbb{R}^m, A \in \mathbb{R}^{m \times n}$$

Activation functions:

- $\sigma_g(z) = \frac{1}{1+e^{-z}}$  - Sigmoid
- $ReLU(z) = \max(0, z)$  - Rectified linear unit
- $\sigma_g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$  - hyperbolic tangent function

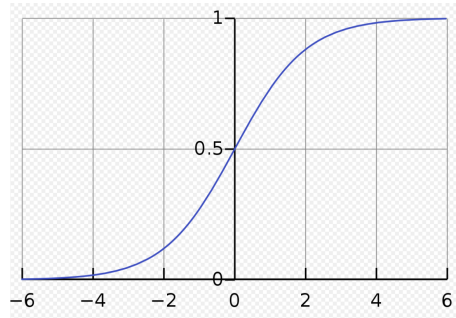


Figure 4.2: Sigmoid function

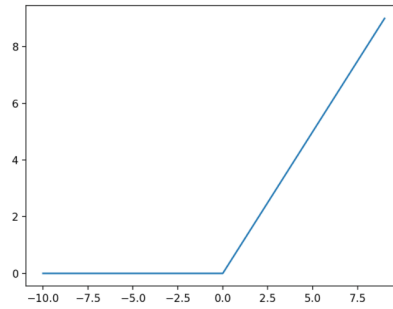


Figure 4.3: ReLU function

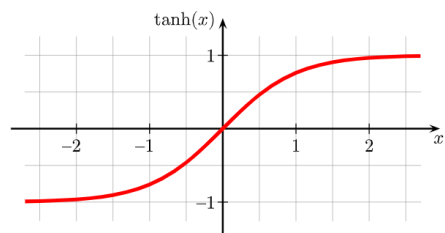


Figure 4.4: Hyperbolic tangent function

The process of training involves adjusting the network's parameters through backpropagation, where errors are propagated backward to fine-tune the connections between nodes.

There are different loss function for different tasks, MSE, L1, Cross entropy. But for the classification task we've used the binary cross-entropy loss:

$$l(x, y) = -x \log(y) - (1 - x) \log(1 - y)$$

### 4.2.2 The Support Vector Classifier (SVC)

The Support Vector Classifier (SVC) aims to find the optimal hyperplane that separates classes in a binary classification problem. One can use a linear kernel to find the hyperplane separator but some non-linear kernels  $\phi(x)$  can be utilized to get a non-linear separation.

Given a set of training samples  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$  where  $x_i$  is the feature vector and  $y_i \in \{-1, 1\}$  is the class label, the optimization problem for SVM with a kernel function:

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(w \cdot \phi(x_i) + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \\ & \xi_i \geq 0, \quad i = 1, \dots, n \end{aligned}$$

where:

$w$  is the weight vector perpendicular to the hyperplane,

$b$  is the bias term,

$\xi_i$  are slack variables,

$C$  is the regularization parameter,

$\phi(x_i)$  is the feature mapping of  $x_i$  into a higher-dimensional space using the kernel function .

## 4.2.3 RNN

### 4.2.3.1 RNN - Introduction

Recurrent Neural Networks (RNNs) have connections that loop back, this enables the network to have memory of past inputs so to speak. The input to RNNs are signals with timestamp (time series) -  $\{x_t, y_t\}_{t=1}^k$ .

RNNs utilise the hidden state to predict the labels  $h_t = f_W(x_t, h_{t-1})$ .

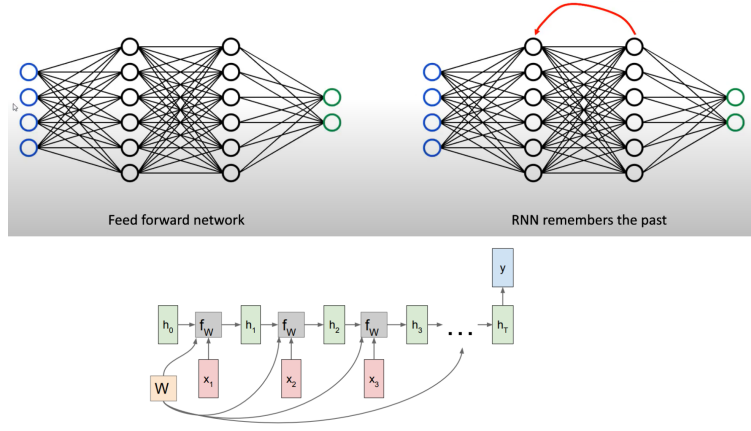


Figure 4.5: RNN

However, standard RNNs can suffer from vanishing or exploding gradients, limiting their ability to capture long-range dependencies. To overcome these issues, variants like Long Short-Term Memory (LSTM) were developed which will be presented in the next section as it will be our main workhorse for the hydrology modeling.

## 4.2.4 LSTM

LSTM is a special kind of RNN, designed to overcome the short-comes of RNN such as complicated training procedures and difficulty to processes long sequences. The architecture of LSTM networks facilitates (as its eponymous name) the learning of both short-term and long-term patterns, remembering information for long periods of time is intrinsic to LSTM.

### 4.2.4.1 LSTM - Principles

LSTMs has three gating mechanisms: the input gate, forget gate, and output gate. These gates enable selective retention or discard of information over time.

Each gate is essentially an Affine transformation followed by a non-linear point wise function:

$$\begin{aligned}f_t &= \sigma_g(W_f x_t + U_f h_{t-1} + b_f) \\i_t &= \sigma_g(W_i x_t + U_i h_{t-1} + b_i) \\o_t &= \sigma_g(W_o x_t + U_o h_{t-1} + b_o) \\\tilde{c}_t &= \sigma_h(W_c x_t + U_c h_{t-1} + b_c) \\c_t &= f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \\h_t &= o_t \circ \sigma_h(c_t)\end{aligned}$$

where:

$x_t \in \mathbb{R}^d$  : input vector to the LSTM unit

$f_t \in \mathbb{R}^h$  : forget gate's activation vector

$i_t \in \mathbb{R}^h$  : input/store gate's activation vector

$o_t \in \mathbb{R}^h$  : output gate's activation vector

$h_t \in \mathbb{R}^h$  : hidden state vector also known as output vector of the LSTM unit

$\tilde{c}_t \in \mathbb{R}^h$  : cell input activation vector

$c_t \in \mathbb{R}^h$  : cell state vector

$\sigma_g$  : sigmoid function

$\sigma_h$  : hyperbolic tangent function

The initial values are  $c_0 = 0$  and  $h_0 = 0$  and the operator  $\circ$  denotes the Hadamard product (element-wise product).

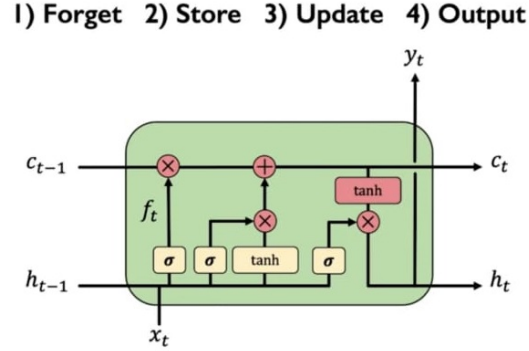


Figure 4.6: LSTM - Scheme

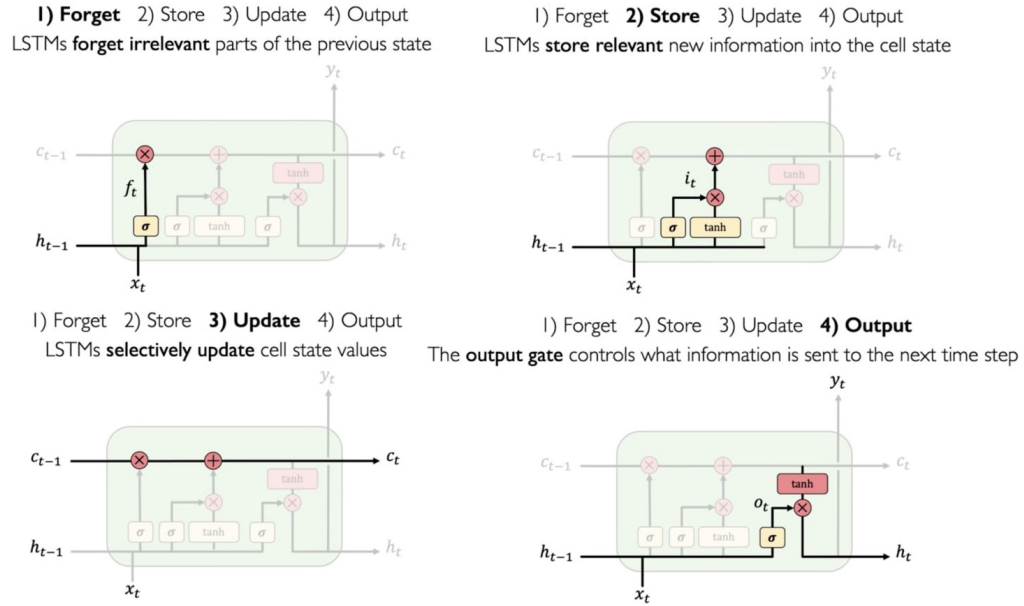


Figure 4.7: LSTM - Scheme



#### 4.2.4.2 LSTM- Multi-cell

LSTM Multi-cell architecture stack memory cells in parallel. This enables having a more expressive networks making them more easy to deal with intricate temporal short-term and long-term dependencies in complex sequences:

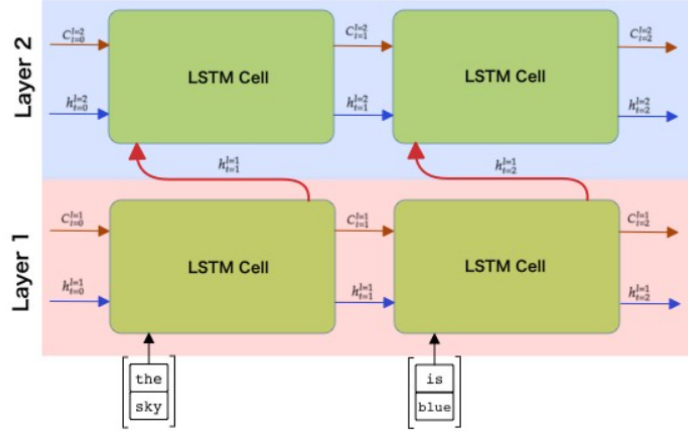


Figure 4.8: LSTM - Multi layer cells

#### 4.2.4.3 LSTM - in context of rainfall

LSTM in hydrology have gained popularity in recent years, specifically the success of incorporating LSTM in rainfall-runoff modeling is well established [24],[25],[26] and [15]. There was even an attempt to try interpretability (IG) on these models [23].

It's worth mentioning a widely used statistical metric for assessing the performance of hydrological models namely the Nash-Sutcliffe model efficiency coefficient, often denoted as NSE.

Given the observed discharge  $Q_o(t)$  and the predicted discharge  $Q_p(t)$  NSE is defined as follows:

$$NSE = 1 - \frac{\sum_{t=1}^T (Q_o(t) - Q_p(t))^2}{\sum_{t=1}^T (Q_o(t) - \bar{Q}_o)^2}$$

Where  $\overline{Q}_o = \sum_{t=1}^N \frac{Q_o(t)}{T}$  i.e. the mean of the observed discharge.

NSE values range from negative infinity to 1, where an NSE of 1 signifies a perfect match, while values below zero indicate that the observed mean outperforms the model i.e. better than just predicting the mean of the sequence.

# Chapter 5

## Results

### 5.1 Sanity Check

paragraph on the rationale of the example

Formulation and figures of the results

paragraph on explaining the results

## **5.2 Dream synthetic Framework**

### **5.2.1 Dream synthetic model**

Figures and explanation of the Rainfall runoff data

Paragraph on how the 2 scenarios are distinguished from each other

Paragraph on how the LSTM was trained, event-driven, MSE and NSE

Figures and explanation of several events prediction Vs. label

### **5.2.2 Interpretability method**

Paragraph on how what methods were used and what were the considerations or constraints

### **5.2.3 Framework evaluation results**

Figures and explanation of the metric used, the process in which they were evaluated and the actual results on each method, lots of comparisons!

# Chapter 6

## Conclusions

TBD

# Bibliography

- [1] Kjersti Aas, Martin Jullum, and Anders Løland. *Explaining individual predictions when features are dependent: More accurate approximations to Shapley values*. 2019. DOI: 10.48550/ARXIV.1903.10464. URL: <https://arxiv.org/abs/1903.10464>.
- [2] Julius Adebayo et al. “Debugging Tests for Model Explanations”. In: *CoRR* abs/2011.05429 (2020). arXiv: 2011.05429. URL: <https://arxiv.org/abs/2011.05429>.
- [3] Julius Adebayo et al. “Local Explanation Methods for Deep Neural Networks Lack Sensitivity to Parameter Values”. In: *CoRR* abs/1810.03307 (2018). arXiv: 1810.03307. URL: <http://arxiv.org/abs/1810.03307>.
- [4] Julius Adebayo et al. *Sanity Checks for Saliency Maps*. 2018. DOI: 10.48550/ARXIV.1810.03292. URL: <https://arxiv.org/abs/1810.03292>.
- [5] Marco Ancona et al. *Towards better understanding of gradient-based attribution methods for Deep Neural Networks*. 2017. DOI: 10.48550/ARXIV.1711.06104. URL: <https://arxiv.org/abs/1711.06104>.
- [6] Leila Arras et al. “‘What is Relevant in a Text Document?’: An Interpretable Machine Learning Approach”. In: *CoRR* abs/1612.07843 (2016). arXiv: 1612.07843. URL: <http://arxiv.org/abs/1612.07843>.
- [7] Leila Arras et al. “Evaluating Recurrent Neural Network Explanations”. In: *CoRR* abs/1904.11829 (2019). arXiv: 1904.11829. URL: <http://arxiv.org/abs/1904.11829>.

- [8] David Baehrens et al. “How to explain individual classification decisions”. In: *The Journal of Machine Learning Research* 11 (2010), pp. 1803–1831.
- [9] Umang Bhatt, Adrian Weller, and Jose M. F. Mour. “Evaluating and Aggregating Feature-based Model Explanations”. In: *CoRR* abs/2005.00631 (2020). arXiv: 2005.00631. URL: <https://arxiv.org/abs/2005.00631>.
- [10] Javier Castro, Daniel Gómez, and Juan Tejada. “Polynomial calculation of the Shapley value based on sampling”. In: *Computers and Operations Research* 36.5 (2009). Selected papers presented at the Tenth International Symposium on Locational Decisions (ISOLDE X), pp. 1726–1730. ISSN: 0305-0548. DOI: <https://doi.org/10.1016/j.cor.2008.04.004>. URL: <https://www.sciencedirect.com/science/article/pii/S0305054808000804>.
- [11] Hila Chefer, Shir Gur, and Lior Wolf. “Transformer Interpretability Beyond Attention Visualization”. In: *CoRR* abs/2012.09838 (2020). arXiv: 2012.09838. URL: <https://arxiv.org/abs/2012.09838>.
- [12] Jacob Devlin et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, June 2019, pp. 4171–4186. DOI: 10.18653/v1/N19-1423. URL: <https://aclanthology.org/N19-1423>.
- [13] Jay DeYoung et al. “ERASER: A Benchmark to Evaluate Rationalized NLP Models”. In: *CoRR* abs/1911.03429 (2019). arXiv: 1911.03429. URL: <http://arxiv.org/abs/1911.03429>.
- [14] Ann-Kathrin Dombrowski et al. *Explanations can be manipulated and geometry is to blame*. 2019. DOI: 10.48550/ARXIV.1906.07983. URL: <https://arxiv.org/abs/1906.07983>.
- [15] Martin Gauch et al. “Rainfall–runoff prediction at multiple timescales with a single Long Short-Term Memory network”. In: *Hydrology and Earth System Sciences* 25.4 (2021), pp. 2045–2062.

- [16] Peter Hase and Mohit Bansal. “Evaluating Explainable AI: Which Algorithmic Explanations Help Users Predict Model Behavior?” In: *CoRR* abs/2005.01831 (2020). arXiv: 2005.01831. URL: <https://arxiv.org/abs/2005.01831>.
- [17] Kaiming He et al. *Deep Residual Learning for Image Recognition*. 2015. arXiv: 1512.03385 [cs.CV].
- [18] Juyeon Heo, Sunghwan Joo, and Taesup Moon. “Fooling Neural Network Interpretations via Adversarial Model Manipulation”. In: *CoRR* abs/1902.02041 (2019). arXiv: 1902.02041. URL: <http://arxiv.org/abs/1902.02041>.
- [19] Robin Hesse, Simone Schaub-Meyer, and Stefan Roth. “Fast Axiomatic Attribution for Neural Networks”. In: *CoRR* abs/2111.07668 (2021). arXiv: 2111.07668. URL: <https://arxiv.org/abs/2111.07668>.
- [20] Sara Hooker et al. *A Benchmark for Interpretability Methods in Deep Neural Networks*. 2018. DOI: 10.48550/ARXIV.1806.10758. URL: <https://arxiv.org/abs/1806.10758>.
- [21] Pieter-Jan Kindermans et al. *The (Un)reliability of saliency methods*. 2017. DOI: 10.48550/ARXIV.1711.00867. URL: <https://arxiv.org/abs/1711.00867>.
- [22] Kamran Kowsari et al. “RMDL: Random Multimodel Deep Learning for Classification”. In: *CoRR* abs/1805.01890 (2018). arXiv: 1805.01890. URL: <http://arxiv.org/abs/1805.01890>.
- [23] Frederik Kratzert et al. “NeuralHydrology – Interpreting LSTMs in Hydrology”. In: *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*. Springer International Publishing, 2019, pp. 347–362. DOI: 10.1007/978-3-030-28954-6\_19. URL: [https://doi.org/10.1007/978-3-030-28954-6\\_19](https://doi.org/10.1007/978-3-030-28954-6_19).
- [24] Frederik Kratzert et al. “Rainfall–runoff modelling using long short-term memory (LSTM) networks”. In: *Hydrology and Earth System Sciences* 22.11 (2018), pp. 6005–6022.
- [25] Frederik Kratzert et al. “Toward improved predictions in ungauged basins: Exploiting the power of machine learning”. In: *Water Resources Research* 55.12 (2019), pp. 11344–11354.



- [26] Frederik Kratzert et al. “Towards learning universal, regional, and local hydrological behaviors via machine learning applied to large-sample datasets”. In: *Hydrology and Earth System Sciences* 23.12 (2019), pp. 5089–5110.
- [27] Sebastian Lapuschkin et al. “On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation”. In: *PLoS ONE* 10 (July 2015), e0130140. DOI: 10.1371/journal.pone.0130140.
- [28] Scott M. Lundberg et al. “Explainable AI for Trees: From Local Explanations to Global Understanding”. In: *CoRR* abs/1905.04610 (2019). arXiv: 1905.04610. URL: <http://arxiv.org/abs/1905.04610>.
- [29] Aravindh Mahendran and Andrea Vedaldi. “Salient Deconvolutional Networks”. In: *Computer Vision – ECCV 2016*. Ed. by Bastian Leibe et al. Cham: Springer International Publishing, 2016, pp. 120–135. ISBN: 978-3-319-46466-4.
- [30] Aniek F. Markus, Jan A. Kors, and Peter R. Rijnbeek. “The role of explainability in creating trustworthy artificial intelligence for health care: a comprehensive survey of the terminology, design choices, and evaluation strategies”. In: *CoRR* abs/2007.15911 (2020). arXiv: 2007.15911. URL: <https://arxiv.org/abs/2007.15911>.
- [31] Rory Mitchell et al. *Sampling Permutations for Shapley Value Estimation*. 2021. DOI: 10.48550/ARXIV.2104.12199. URL: <https://arxiv.org/abs/2104.12199>.
- [32] Volodymyr Mnih et al. “Playing Atari with Deep Reinforcement Learning”. In: *CoRR* abs/1312.5602 (2013). arXiv: 1312.5602. URL: <http://arxiv.org/abs/1312.5602>.
- [33] Christoph Molnar. *Interpretable Machine Learning. A Guide for Making Black Box Models Explainable*. <https://christophm.github.io/interpretable-ml-book/>. 2019.
- [34] Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. “Methods for Interpreting and Understanding Deep Neural Networks”. In: *CoRR* abs/1706.07979 (2017). arXiv: 1706.07979. URL: <http://arxiv.org/abs/1706.07979>.

- [35] Dong Nguyen. “Comparing Automatic and Human Evaluation of Local Explanations for Text Classification”. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. New Orleans, Louisiana: Association for Computational Linguistics, June 2018, pp. 1069–1078. DOI: 10.18653/v1/N18-1097. URL: <https://aclanthology.org/N18-1097>.
- [36] An-phi Nguyen and María Rodríguez Martínez. “On quantitative aspects of model interpretability”. In: *CoRR* abs/2007.07584 (2020). arXiv: 2007.07584. URL: <https://arxiv.org/abs/2007.07584>.
- [37] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ““Why Should I Trust You?”: Explaining the Predictions of Any Classifier”. In: *CoRR* abs/1602.04938 (2016). arXiv: 1602.04938. URL: <http://arxiv.org/abs/1602.04938>.
- [38] Wojciech Samek et al. *Evaluating the visualization of what a Deep Neural Network has learned*. 2015. DOI: 10.48550/ARXIV.1509.06321. URL: <https://arxiv.org/abs/1509.06321>.
- [39] Wojciech Samek et al. *Explainable AI: interpreting, explaining and visualizing deep learning*. Vol. 11700. Springer Nature, 2019.
- [40] Wojciech Samek et al. “Toward Interpretable Machine Learning: Transparent Deep Neural Networks and Beyond”. In: *CoRR* abs/2003.07631 (2020). arXiv: 2003.07631. URL: <https://arxiv.org/abs/2003.07631>.
- [41] Ramprasaath R. Selvaraju et al. “Grad-CAM: Why did you say that? Visual Explanations from Deep Networks via Gradient-based Localization”. In: *CoRR* abs/1610.02391 (2016). arXiv: 1610.02391. URL: <http://arxiv.org/abs/1610.02391>.
- [42] L. S. Shapley. “17. A Value for n-Person Games”. In: *Contributions to the Theory of Games (AM-28), Volume II*. Ed. by Harold William Kuhn and Albert William Tucker. Princeton: Princeton University Press, 2016, pp. 307–318. ISBN: 9781400881970. DOI: doi:10.1515/9781400881970-018. URL: <https://doi.org/10.1515/9781400881970-018>.

- [43] Nathan Sheffer et al. “Hydrometeorological daily recharge assessment model (DREAM) for the Western Mountain Aquifer, Israel: Model application and effects of temporal patterns”. In: *Water Resources Research* 46 (May 2010). DOI: 10.1029/2008WR007607.
- [44] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. “Learning Important Features Through Propagating Activation Differences”. In: *CoRR* abs/1704.02685 (2017). arXiv: 1704.02685. URL: <http://arxiv.org/abs/1704.02685>.
- [45] Avanti Shrikumar et al. “Not Just a Black Box: Learning Important Features Through Propagating Activation Differences”. In: *CoRR* abs/1605.01713 (2016). arXiv: 1605.01713. URL: <http://arxiv.org/abs/1605.01713>.
- [46] David Silver et al. “Mastering the Game of Go with Deep Neural Networks and Tree Search”. In: *Nature* 529.7587 (Jan. 2016), pp. 484–489. DOI: 10.1038/nature16961.
- [47] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. *Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps*. 2013. DOI: 10.48550/ARXIV.1312.6034. URL: <https://arxiv.org/abs/1312.6034>.
- [48] Dylan Slack et al. “How can we fool LIME and SHAP? Adversarial Attacks on Post hoc Explanation Methods”. In: *CoRR* abs/1911.02508 (2019). arXiv: 1911.02508. URL: <http://arxiv.org/abs/1911.02508>.
- [49] Daniel Smilkov et al. “SmoothGrad: removing noise by adding noise”. In: *CoRR* abs/1706.03825 (2017). arXiv: 1706.03825. URL: <http://arxiv.org/abs/1706.03825>.
- [50] Jost Tobias Springenberg et al. *Striving for Simplicity: The All Convolutional Net*. 2014. DOI: 10.48550/ARXIV.1412.6806. URL: <https://arxiv.org/abs/1412.6806>.
- [51] Erik Strumbelj and Igor Kononenko. “An Efficient Explanation of Individual Classifications using Game Theory”. In: *J. Mach. Learn. Res.* 11 (2010), pp. 1–18.
- [52] Erik Strumbelj and Igor Kononenko. “An efficient explanation of individual classifications using game theory”. In: *The Journal of Machine Learning Research* 11 (2010), pp. 1–18.

- [53] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. “Axiomatic Attribution for Deep Networks”. In: *CoRR* abs/1703.01365 (2017). arXiv: 1703.01365. URL: <http://arxiv.org/abs/1703.01365>.
- [54] Richard Tomsett et al. “Sanity Checks for Saliency Metrics”. In: *CoRR* abs/1912.01451 (2019). arXiv: 1912.01451. URL: <http://arxiv.org/abs/1912.01451>.
- [55] Zhou Wang et al. “Image quality assessment: from error visibility to structural similarity”. In: *IEEE Transactions on Image Processing* 13.4 (2004), pp. 600–612. DOI: 10.1109/TIP.2003.819861.
- [56] Martin Wistuba, Ambrish Rawat, and Tejaswini Pedapati. “A Survey on Neural Architecture Search”. In: *CoRR* abs/1905.01392 (2019). arXiv: 1905.01392. URL: <http://arxiv.org/abs/1905.01392>.
- [57] Hao Yuan et al. “Explainability in Graph Neural Networks: A Taxonomic Survey”. In: *CoRR* abs/2012.15445 (2020). arXiv: 2012.15445. URL: <https://arxiv.org/abs/2012.15445>.
- [58] Matthew D. Zeiler and Rob Fergus. “Visualizing and Understanding Convolutional Networks”. In: *CoRR* abs/1311.2901 (2013). arXiv: 1311.2901. URL: <http://arxiv.org/abs/1311.2901>.
- [59] Quanshi Zhang and Song-Chun Zhu. “Visual Interpretability for Deep Learning: a Survey”. In: *CoRR* abs/1802.00614 (2018). arXiv: 1802.00614. URL: <http://arxiv.org/abs/1802.00614>.
- [60] Yu Zhang et al. “A Survey on Neural Network Interpretability”. In: *CoRR* abs/2012.14261 (2020). arXiv: 2012.14261. URL: <https://arxiv.org/abs/2012.14261>.
- [61] Bolei Zhou et al. “Learning Deep Features for Discriminative Localization”. In: *CoRR* abs/1512.04150 (2015). arXiv: 1512.04150. URL: <http://arxiv.org/abs/1512.04150>.
- [62] Jianlong Zhou et al. “Evaluating the Quality of Machine Learning Explanations: A Survey on Methods and Metrics”. In: *Electronics* 10.5 (2021). ISSN: 2079-9292. DOI: 10.3390/electronics10050593. URL: <https://www.mdpi.com/2079-9292/10/5/593>.