

# A SURVEY ON VISUALIZATION FOR EXPLAINABLE CLASSIFIERS

by

**YAO MING**

Department of Computer Science and Engineering  
The Hong Kong University of Science and Technology  
Supervised by Prof. Huamin Qu

October 2017, Hong Kong

# TABLE OF CONTENTS

<b>Title Page</b>	i
<b>Table of Contents</b>	ii
<b>Abstract</b>	iv
<b>Chapter 1 Introduction</b>	1
1.1 Motivation	1
1.2 Challenges	2
1.3 Overview	3
<b>Chapter 2 Concepts and Definitions</b>	4
2.1 Classification	4
2.2 Classifiers	4
2.3 Explainability	5
<b>Chapter 3 Explainable Classifiers</b>	7
3.1 Interpretable Classifiers	7
3.1.1 Interpretable Architecture	8
3.1.2 Learning Sparse Models	9
3.2 Explanations of Classifiers	10
3.2.1 Local Explanations	11
3.2.2 Global Explanations	14
<b>Chapter 4 Visualization for Explainable Classifiers</b>	16
4.1 Life Cycle of an intelligent system	16
4.2 Visualization for Data Understanding	18
4.3 Visualization for Model Development	18
4.3.1 Understanding	18
4.3.2 Diagnosing	18
4.3.3 Assessment and Comparison	19
4.4 Visualization for Model Operation	19
4.4.1 Trust Establishment	19
4.4.2 Monitoring	19
4.5 Other Applications	19

4.5.1	Teaching and Communicating Models	19
4.5.2	Learn from the Model	19
4.6	Evaluation	19
<b>Chapter 5</b>	<b>Conclusion</b>	<b>20</b>
<b>Bibliography</b>		<b>21</b>

# A SURVEY ON VISUALIZATION FOR EXPLAINABLE CLASSIFIERS

by

**YAO MING**

Department of Computer Science and Engineering

The Hong Kong University of Science and Technology

## ABSTRACT

Classification is a fundamental problem in machine learning, data mining and computer vision. In practice, interpretability is a desirable property of classification models (classifiers) in critical areas, such as security, medicine and finance. For instance, a quantitative trader may prefer a more interpretable model with less expected return due to its predictability and low risk. Unfortunately, the best-performing classifiers in many applications (e.g., deep neural networks) are complex machines whose predictions are difficult to explain. Thus, there is a growing interest in using visualization to understand, diagnose and explain intelligent systems in both academia and in industry. Many challenges need to be addressed in the formalization of explainability, and the design principles and evaluation of explainable intelligent systems.

The survey starts with an introduction to the concept and background of explainable classifiers. Efforts towards more explainable classifiers are categorized into two: designing classifiers with simpler structures that can be easily understood; developing methods that generate explanations for already complicated classifiers. Based on the life circle of a classifier, we discuss the pioneering work of using visualization to improve its explainability at different stages in the life circle. The survey ends with a discussion about the challenges and future research opportunities of explainable classifiers.

# CHAPTER 1

## INTRODUCTION

### 1.1 Motivation

Classification is the problem of identifying if an observation or object belongs to a set or not, or which of several sets. It is a fundamental problem in machine learning, data mining and computer vision. With the support of the increasing capacity of computation resources and the growing volume of available data, the last few decades have witnessed an explosion of breakthroughs in these fields. Nowadays, classification models (classifiers) are widely adopted to solve real world tasks, including face recognition [], handwriting recognition [], sentiment analysis [] and spam filtering []. Take image classification for instance, a well-designed convolutional neural network can achieve human-level performance in a number of benchmark datasets [,].

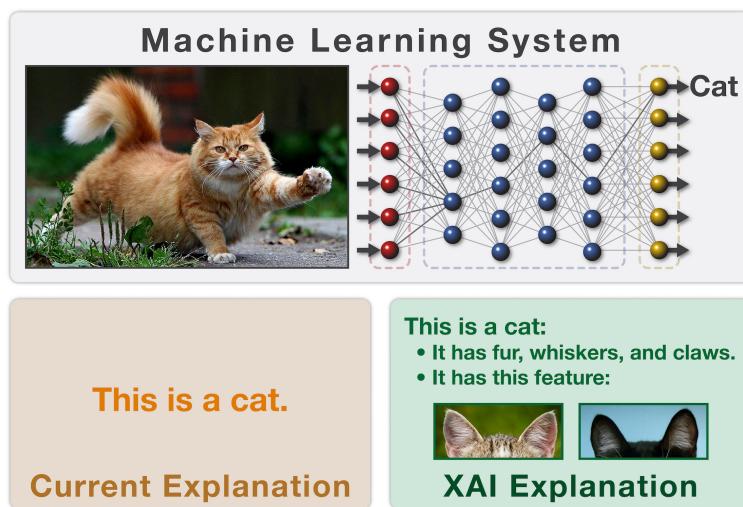


Figure 1.1. An illustration of an explainable image classifier [14].

Despite their promising capabilities, an often-overlooked aspect is the important role of humans [30]. When humans are to understand and collaborate with these autonomous systems, it is desirable if we have explanations of their outputs. For instance, a doctor using a machine classifier to assist in identifying early signs of lung cancer would need to know why the classifier “thinks” there might be a cancer so that he/she can make a more confident diagnosis. A natural way is to provide explanations (Figure 1.1). In machine learning, the term *explainability* does not have a standard and generally accepted definition. In some literature,

*interpretability* is used instead. Generally speaking, the explainability or interpretability of an intelligent system refers to the ability to explain its reasoning [8] to humans. For the sake of consistency, we use explainability as the ability to explain in this survey. Interpretability is used to refer to the property of how easily a model can be understood by humans.

The research for explainable intelligent systems can be traced back to the 1980s, when expert systems were created and proliferated [7, 27, 36]. These early works focused on reducing the difficulty of maintaining the complicated if-then rules by designing more explainable representations. A huge gap exists between today's state-of-the-art intelligent systems and the techniques that can make them explainable. The new challenges brought about by the new generation of intelligent systems have attracted growing research interests. DARPA have launched the Explainable Artificial Intelligence (XAI) project [14], which aims to develop new techniques to make these systems explainable. Google Inc. have initiated the People + AI Research Initiative (PAIR) [13] to advance humanistic considerations in AI.

Visualization is an effective and efficient technique for communicating information and understanding complex datasets for humans. The visual system is a proxy with a very large bandwidth to human brains [25]. Thus, visualization can be an ideal weapon to help explain complicated classifiers to humans. Early related research can be traced back to the software and algorithm visualization for computer science education in the 1980s and 1990s [6, 35, 28]. Visualization, especially interactive visualization, was proved to be very effective in facilitating people's understanding of complex softwares and algorithms. Little research has been done to visualize the increasingly complicated classifiers, which are actually algorithms learned from the data. It has not been until recently that visualization was popularized as a media for understanding classification models, especially for image classifiers [33, 43, 2, 44]. However, these methods have limited applications to neural networks for image data. There is also a lack of a unified and convenient evaluation method for the generated visualizations.

## 1.2 Challenges

The need for visually explaining classifiers is actually a result of the successes and advances of AI. The major challenges of visually explaining classifiers results from the complexity of the model and data, and the limits of humans.

First, it is challenging to explain complex classification models both concisely and precisely. The best-performing classifiers (e.g., neural networks) are becoming increasingly complex, in terms of the number of parameters and operations they employed, which makes

them difficult to explain. A convolutional network typically employs thousands of neurons and millions of parameters. A random forest used for classification may employ hundreds of decision trees, each containing hundreds of nodes. Sampling a small number of parameters/neurons/nodes to explain might be easier to understand for humans, but it brings with it the risk of misunderstanding as well. The variety of model architectures also increases the difficulty of an explanation framework for classifiers to be effective and general at the same time.

Another challenge is the volume and variety of the data used for training the classifiers. To explain a classifier, a most common strategy is to trace back to the input data. Which part of the input data contributes to the prediction? How does a model behave on this subset of data? Some explanation methods require computations over the entire training set, which may become impractical if the data set is very large. Different data types may require different forms of visual explanation. Image data are readily interpretable, but how to effectively explain classifiers on categorical, text and speech data is still a problem.

These challenges are, to some extent, due to compromises owing to the limits of humans' cognitive ability. If humans can make sense of the meaning of thousands of parameters and complex model structures by merely looking at the raw data or code, there is no need to struggle with how to better visualize them. There are already some studies discussing the structure, function and effectiveness of explanations in cognitive science. However, it is still unclear how we can effectively evaluate the quality of an explanation, and the load that its visual representations exert over humans.

### 1.3 Overview

This survey mainly focuses on how visualization techniques can be used to support explainable classifiers. In Chapter 3, we first introduce the definition of classification and classifiers, and the concept of explainable classifiers. Two major research directions towards more explainable classifiers are identified: designing classifiers that are readily interpretable, and methods that generate explanations for a classifier without modifying the model. In Chapter 4, we first articulate the life cycle of a classifier into different stages, that is, the recursive procedures of data collection and processing, model development and testing, and operations and maintenance. Then, we illustrate how visualization can be applied at different stages to provide explainability for classifiers. Based on the specified life cycle, we categorize the surveyed literature and discuss the challenges and opportunities for future research in visualization for explainable classifiers.

# CHAPTER 2

## CONCEPTS AND DEFINITIONS

In this section, we first briefly introduce the problem of classification, as an instance of supervised learning, and a few popular classifications models (classifiers). To clarify the scope of this survey, we discuss the concepts of explainability of classifiers and illustrate in which circumstances explainability are desirable or needed.

### 2.1 Classification

Given an input space  $\mathcal{X}$  and an output space  $\mathcal{Y} = \{1, 2, \dots, K\}$  with  $K$  classes, **classification** is the problem of identifying any **observation**  $\mathbf{x} \in \mathcal{X}$  to a class  $y \in \mathcal{Y}$ . For multi-label classification, where class labels are not exclusive, we can view it as multiple related binary classification. For simplicity, we only consider the basic formulation in this survey.

A **classifier** is an algorithm  $f$  that implements classification, *i.e.*,  $y = f(\mathbf{x})$ . To handle ambiguity, a classifier is often used in a probabilistic setting. That is, the output of  $f$  a probabilistic distribution  $p(y | \mathbf{x}, \mathcal{D})$  over all possible classes in  $\mathcal{Y}$ .  $\mathcal{D}$  is the training set, which is a subset of  $\mathcal{X} \times \mathcal{Y}$ , that have already been observed. Thus, in practice, a classifier will often take the form of  $\mathbf{y} = f(\mathbf{x})$ , where  $\mathbf{y} = (y_i) \in \mathbb{R}^K$  is a vector denoting the probabilistic distribution. Then the final classification will be the class  $i$  with largest probability  $\arg \max_i y_i$ .

### 2.2 Classifiers

Classification is now widely applied in solving many real world applications. A few examples are: face recognition [], handwritten recognition [], sentiment analysis [] and spam filtering. Here we briefly present a few popular models for classifications, including k-nearest neighbors, support vector machines, decision trees and neural networks.

**K-nearest neighbor.**

**Support vector machine.**

**Decision trees.**

**Neural networks.** CNN, RNN.

## 2.3 Explainability

What is explainability? What is the explainability of a classifier? There is no commonly agreed definitions so far. Doshi-Velez and Kim defines interpretability (or explainability) as the ability to explain or to present in understandable terms to a human [8], which is already a good general definition. To clarify the scope of this survey, we define the **explainability** of a classifier as the ability to explain the reasoning of its predictions so that a human can understand. Simple models such as a linear classifier already has good explainability since humans can easily understand the model’s reasoning by simply looking at the coefficients of each feature. For a complicated classifier like a deep neural network, a human may find it difficult to understand due to layer-wise structure and the nonlinearity of the computation. Thus, the key of explainability is the humans.

An immediate question is: why do we need explainability? The need of explainability for a full automated classifier mainly comes from three aspects: humans’ curiosity about knowledge, limitations of current intelligent algorithms, and moral and legal issues:

- **Curiosity of human.** Humans are curious about new knowledge. Often, a classifier is not developed solely for performing the classification tasks, but also for knowledge discovery. For todays popular neural networks, humans are curious of how the impressive human-level classification accuracy is achieved. There are also examples of how insights learned from the behavior of a model lead to improvement on the design of a classifier [43, 1]. Besides, given that AlphaGo Zero [32] can learn to master the game of Go much better than human players, it is desirable that the machine can explain its learned strategy (knowledge) to us.
- **Limitations of machines.** The current state-of-the-art intelligent systems are usually not fully testable. Human knowledge are required as a complement in case the machines fail. In the seeable future, machines are expected to assist rather than replace humans in many domains, such as security, medical services, education and design. By providing explainability, users’ trust can be more easily established. Besides, explainability can provide an interface for humans to monitor machine.
- **Moral and legal issues.** The “right to explanation”, which is a regulation included in the GDPR <sup>1</sup> of the European Union, has recently raised a debate on to which extent we should require automatic decision-making systems to provide explanations to the

---

<sup>1</sup><https://www.privacy-regulation.eu/en/r71.htm>

subjects of the decision. If one's application for a loan is denied by a automatic classifier, he/she has the right to ask why. A doctor may need to know why a patient is classified to have a lung cancer to give the final diagnosis. Another issue is the fairness or the discrimination problem of a classifier, which may be easily neglected during the development phase.

Though explainability is a desirable property, it should be noted that it is not always necessary. Explainability is not required if 1) the application domain has high resistance to errors, and thus, unexpected errors is acceptable; 2) the application domain has been well studied and the classifier has been well tested in production, and thus, it is unlikely to have unexpected results.

# CHAPTER 3

## EXPLAINABLE CLASSIFIERS

In this section, we discuss methods that provide explainability to classifiers. We categorized related work into two depending on how the explainability of a classifier is achieved.

The first type of work develops more interpretable classifiers that are easy to understand for humans. The second generates explanations for a classifier without modifying the model, either by explaining the classifier locally on specific instances, or by explaining the behavior of the classifier globally. A summary is shown in Table 3.1.

Categories		Related Papers	Remarks
Interpretable Classifiers	Interpretable Architecture	Decision Trees [5], Rule Lists [21, 40], Rule Sets [41]	rule-based
		Linear Models [4]	linear
		kNNs [10, 18]	instance-based
	Learning Sparse Models	Decision Trees [29], Sparse SVMs [9], Sparse CNNs [23]	simplification
		Sparsity by Bayesian [38], Integer Models [37, 39]	direct-sparsity
Explanations of Classifiers	Local	Model-unaware	Sensitivity Analysis [33, 22, 34]
			LIME [30]
			Generate Visual Explanations [16]
	Model-aware	Global	De-convolution [43],
			Layer-wise Propagation [2],
			Prediction Difference [44],
			Output Decomposition [26],
			Direct Mapping [17]
			CNN
			CNN
			Image
			LSTM
			RNN
		Model-unaware	cell3
			cell5
			cell8
			cell8
			cell9
			cell9

Table 3.1. Towards explainable classifiers.

### 3.1 Interpretable Classifiers

**Interpretable classifiers** are the classifiers that are commonly recognized to be more understandable than others, and hence, do not need extra explicit explanations. Summarizing

```

if hemiplegia and age > 60 then stroke risk 58.9% (53.8%–63.8%)
else if cerebrovascular disorder then stroke risk 47.8% (44.8%–50.7%)
else if transient ischaemic attack then stroke risk 23.8% (19.5%–28.4%)
else if occlusion and stenosis of carotid artery without infarction then
stroke risk 15.8% (12.2%–19.6%)
else if altered state of consciousness and age > 60 then stroke risk
16.0% (12.2%–20.2%)
else if age ≤ 70 then stroke risk 4.6% (3.9%–5.4%)
else stroke risk 8.7% (7.9%–9.6%)

```

Figure 3.1. A decision list learned by the BRL algorithm [21].

existing work, we find two major strategies for creating interpretable classifiers: developing interpretable models with easy-to-understand structures, and learning simpler or sparser models.

### 3.1.1 Interpretable Architecture

To create more interpretable classifiers, a natural way is to use simple computation structures (*e.g.*, if-then rules). Most classifiers that fall into this category are rule-based.

**Rule-based.** A widely adopted type of models are the decision trees [5]. A decision tree classifier uses internal nodes and branches to represent its classification reasoning as conjunctions of rules. A human can trace back a specific classification from a leaf to the root to understand the prediction of the classifier. However, the difficulty of constructing a high-accuracy and interpretable decision tree has long been criticized.

Focused on balancing among performance, explainability and computation, a few recent studies introduce the Bayesian framework in rule-based classifiers. Letham *et al.* [21] develop the Bayesian Rule List (BRL) which employs a prior structure that encourages sparsity in the generated decision lists with a good accuracy. The rule lists have the form of if-then-else structures, as shown in Figure 3.1. Wang and Rudin [40] design the Falling Rule Lists that use an ordered if-then rule list so that the most at-risk occasion will be handled first. Wang *et al.* [41] construct rule sets based on AND and OR operations and highlight its low computation cost and on-par accuracy compared with SVM and random forest.

The most serious problem of these interpretable models with easy-to-understand structures is the scalability. The performance of the rule-based models increases as the number of rules increases or the non-linearity increases. Although the rule-based models are easy to learn and understand at the first glance, it is intractable to understand the classifier as

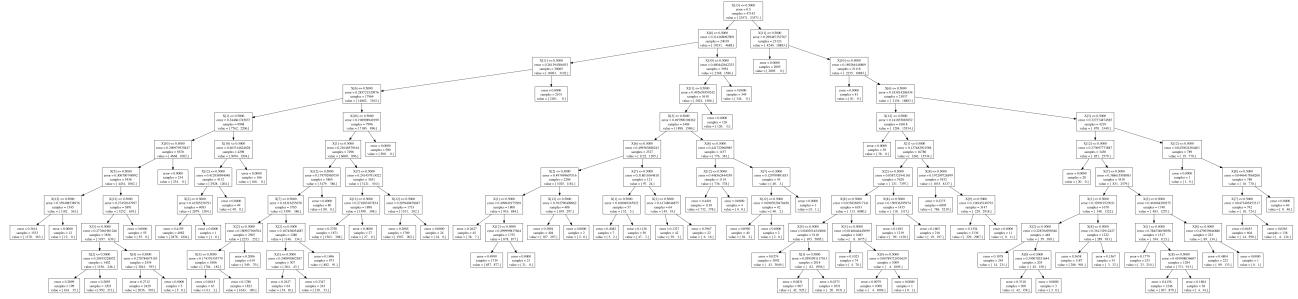


Figure 3.2. A decision tree with over one hundred nodes, which is hard to explain its reasoning.

a whole when the number of nodes of rules grows up to a few hundreds. An example is shown in Figure 3.2.

**Others.** Except for rule-based models, there are a few other models with more complicated models are recognized to be interpretable. One family of interpretable models worth noticing are the generalized linear models [4], which are pervasive in statistics and finance. Although these models can have highly nonlinear computations, the additive relation between nonlinear functions of features are believed to be easy-to-understand. However, the generalized linear classifiers can also be hard to understand when their non-linearity increased to a certain extent. The other non-probabilistic family of classifiers are the k-nearest neighbors (kNN) classifiers, whose prediction can be easily understood by presenting the observation’s k-nearest neighbors. Numerous work has been done to boost the performance the kNN classifiers, including weighted kNN with different kernels [10] and fuzzy kNN [18]. The explainability of kNN classifiers may easily fail when there lacks good neighbors for certain observations.

### 3.1.2 Learning Sparse Models

As discussed above, the explainability often decreases as the complexity (*i.e.*, number of parameters or nodes) of the model increases. Thus, we can improve the explainability by learning a sparser model with the same architecture. Two common strategies are used to learn a sparse model: simplifying a “dense” model through pruning or approximation; introducing sparsity as a prior to learn a sparse model from scratch.

**Simplification.** The methods for simplifying classifiers are usually developed in a model-specific manner. Quinlan summarizes four techniques for simplifying decision trees [29], *i.e.*, cost-complexity pruning, reduced error pruning, pessimistic pruning and simplification to rule sets. Downs *et al.* [9] recognize and eliminate dependent support vectors while leaving the outputs of SVM unchanged. Liu *et al.* [23] use a sparse decomposition method to zero

out redundant parameters in a CNN, achieving about 10-times speedup while only losing about 1% accuracy. Though these methods can practically simplify classifiers and speed up computations, they do not directly provide explainability. A simplified SVM with fewer support vectors but utilizing a complicated kernel is still difficult to explain. A simplified decision tree with 200 nodes instead of 1,000 nodes is still hard to interpret.

**Learning from scratch.** To directly learn sparser models from scratch, Tipping [38] proposed a general Bayesian framework that treats sparsity as a prior and specialized this method on SVM. Instead of restricting the complexity of the parameters, Tan *et al.* [37] uses a 0-1 control variable to each input feature, and convert the learning to a mixed integer programming problem. Similar idea can be found in the sparse linear integer models proposed by Ustun *et al.* [39]. Although these methods can learn sparse classifiers without losing much performance, they mainly focus on reducing the computation costs instead of providing explainability. They do not guarantee explainability if the classification problem is difficult.

In most cases, the efforts of developing more interpretable classifiers are tradeoffs between performance and explainability. For performance-critical applications, it is always difficult to train a interpretable classifier that do not need extra explanations.

## 3.2 Explanations of Classifiers

Generating explanations of a classifier without modifying the classifier itself is preferred when the underlying model is already too complicated, *e.g.*, neural networks and SVMs, and we don't want to sacrifice performance. There is also no common-recognized definition for what an **explanation** of a classifier is. Most existing work uses a subset or a weighted subset of input features to explain a single prediction of a classifier, *e.g.*, a mask over the input image, a heatmap with the same size as the input image, a bag of words or categorical fields. Some work [30] proposed to induct a simpler classifier (*e.g.*, linear classifier) as the explanation of a prediction. Here we only discuss explanations for complex classifiers. Thus, illustrative diagrams for simple classifiers are not included here.

In cognitive science, explanations are characterized as arguments that demonstrating all or a subset of the causes of the **explanandum** (the subject being explained), usually following deductions from natural laws or empirical conditions [15, 24]. Here we give a general definition:

**Explanations of a classifier** is the human-understandable representations that identify the causes of the classifier's prediction(s). A typical form of human-understandable rep-



Figure 3.3. Images (first row) and their saliency maps (second row) for the top-1 predicted class in ILSVRC-2013 test images [33].

resentations is the visualization. As introduced in Section 2.1, a classifier can be regarded as a function  $f$ , which is in general learned from a training dataset  $\mathcal{D}$ , specified by learned parameters  $\theta$ . Thus, the **causes** can traced to 1) parts of the training data  $\mathcal{D}$ , 2) parts of the parameters  $\theta$  as some components of the classifier, or 3) parts of the input,  $x$ , of a prediction or predictions.

If an explanation is provided to explain  $f$ 's prediction in a small region around a given input point  $x$ , we call it as a **local explanation**. If it is generated to explain  $f$  on the whole input space  $\mathcal{X}$  in general or as a summary, we call it as a **global explanation**. Sometime, we also want to have an intermediate **subset explanation** that is performed on a subset  $S$  of  $\mathcal{X}$ , in which the inputs share some common features. Next, we discuss the related work on explanations of classifiers based on the above taxonomy.

### 3.2.1 Local Explanations

As we have discussed, the causes used in an explanation of a classifier can be the training data, model parameters, and inputs. For local explanation, the input  $x$  is always specified and used in the explanation. Depending on whether an explanation is generated directly using model's parameters or structures, we categorize local explanation methods into model-aware and model-unaware methods.

**Model-unaware explanations** only require the input  $x$  and a computable  $f$ . Most work uses sensitivity or saliency-based techniques to derive explanations. Simonyan *et al.* [33]

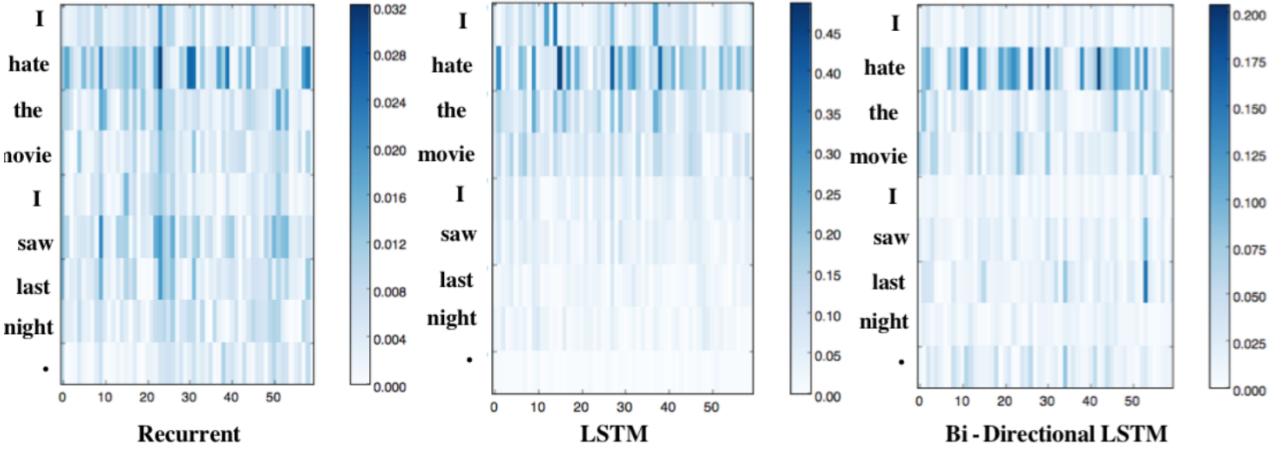


Figure 3.4. Saliency matrix maps for “I hate the movie I saw last night.” of three RNN sentiment classifiers. *Left:* a vanilla RNN; *Middle:* an LSTM; *Right:* a bi-directional LSTM [33].

use the derivative of a image classifier  $f_i$  w.r.t. the input image  $\mathbf{x}$  as the saliency score of the class  $i$ , and map the score of each pixel to a saliency map as the explanation of  $\mathbf{x}$ . Li *et al.* [22] also calculate the derivative of a text sentiment RNN classifier w.r.t. the embeddings of an input sentence (which is a matrix), as shown in Figure 3.4. The heatmap matrices are used as explanations for users to identify salient dimensions of the embedding vector and salient words that contribute the most to the prediction. Although these sensitivity-based methods are intuitive and can be efficiently approximated, their generated explanations are often noisy (as shown in Figure 3.3), due to the high nonlinearity of the complicated classifier. Recently, Smilkov *et al.* [34] propose a random sampling technique to smooth the gradient, which achieve more meaningful visual explanations. However, this smoothing technique is computationally expensive and non-deterministic.

Instead of sensitivity analysis, some other work trains another model to explain the explanandum. Ribeiro *et al.* [30] approximate a complicated classifier locally using a simple linear classifier, and proceed to generate super-pixel patches as explanations. This method is actually similar to the gradient smoothing, since the training of the linear classifier is also done by sampling around the current input. Forming the problem as image captioning, Hendricks *et al.* [16] use extra labeled explanation texts of images to train a explainer that generate explanatory texts of a image classifier. This method highly depends on the quality of text explanations labels, which require extra expensive labeling. Besides, it introduces another model, which is another potential explanandum that need to explain.

**Model-aware explanations** utilize another cause – the parameters of the model,  $\theta$ , to amplify the information in the explanation. Zeiler and Fergus [43] developed a de-convolution method that map the outputs of a CNN classifier back to the input space, utilizing the in-

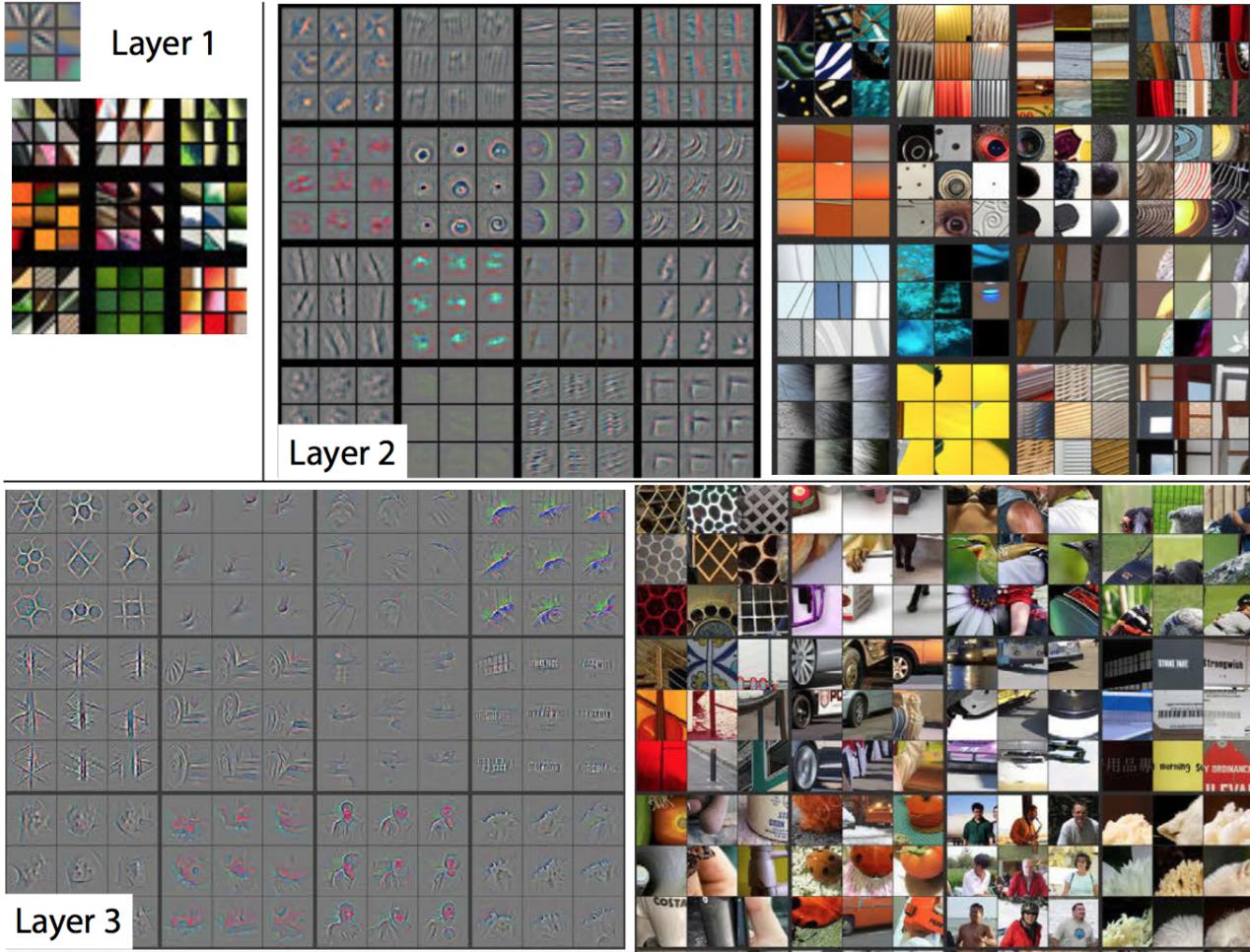


Figure 3.5. Visualization of a CNN generated using the de-convolution method [43]. Gray images in the left are the top nine activations in a random subset of neurons across validation data, projected back to image space. In the right are corresponding image patches.

verse operations of different layers. The projected images Figure 3.5 can be used as an explanation of what different neurons are used for. Bach *et al.* [2] uses a layer-wise relevance propagation (LRP) that improves the sparsity of the image heatmap. Zintgraf *et al.* [44] developed a visualization that highlights evidence for and against a prediction separately through prediction difference analysis. Murdoch and Szlam [26] decomposed the output of an LSTM classifier into multiplicative contribution scores of input words and uses the scores to explain how important the words are for the prediction. Though these methods typically result in much better (sparse, meaningful) explanations than model-unaware methods, they are developed in a per-model manner, which are hard to generalize for other classifiers. There is a lack of general explanatory frameworks to guide and evaluate the development of model-aware methods. Flooded by the interest in deep learning, we can hardly find methods developed for classification models other than neural networks.

Using the cause of training data for explanations has not attracted interest until recently. Koh and Liang proposed a fast approximation of the influence function, which is a well-

studied method in statistics. The influence functions can help identifying training points that are most responsible for a given prediction [19], and thus can explain the prediction from the aspect of training data.

### 3.2.2 Global Explanations

The global explanations are not dependent on any specific inputs. A global explanation is actually a summary of the reasoning of how a classifier generally behaves. Unlike local explanations which are defined around a certain point, the global explanations are considered to be ill-posed and are much harder to achieve. Similarly to local explanations, we divide existing work into model-aware and model-unaware methods.

**Model-unaware explanations.** To our knowledge, very few methods has been proposed to generate global explanations for general classifiers. Ribeiro *et al.* [30] proposed to select a collection of representative local explanations and present to the users one local explanation at a time to give a global understanding. This method will easily fail when the dataset (training data) is too large. The users will not be able to remember a lot representative local explanations to form a global understanding.

**Model-aware explanations.** The first attempt to understand a complex classifier globally is done by Féraud and Clérot [12]. They partitioned the hidden representation space through clustering the representation of the whole training set, where each cluster represent a semantic concept learned by the classifier. To qualitatively understand a CNN, Erhan *et al.* [11] proposed the activation maximization method. Each neuron in the CNN can be explained using an image patch that will maximize its activation. To provide conceptual meanings of the explanation, Bau *et al.* [3] align hidden units with human understandable concepts (objects) through a dissection process. However, these methods often require explorations over multiple nodes or neurons. The big picture is often neglected. Additionally, a common issue is that it is hard to compare these methods due to the lack of an evaluation framework.

In summary, there are two common strategies to make a classifier explainable. First, we develop simpler or sparser classifiers that can meet the performance requirements. Second, we build another human-understandable interface for explanation on top of a classifier. Both strategies are useful in different scenarios.

However, an important, but neglected aspect of existing methods is the human. Few have paid attention to model human. Most of them study the classifiers and develop techniques for explainability and then argue that their methods help humans to understand the

classifier, without studying how humans exactly response to the results generated by these techniques.

Another related problem raised by Doshi-velez and Kim [8] is the lack of the evaluation methods for explainability. Without a rigorous evaluation, it is hard to compare which method is better in a certain setting. It will also be infeasible to clarify the gap of current research and direction for future research.

## CHAPTER 4

# VISUALIZATION FOR EXPLAINABLE CLASSIFIERS

As discussed at the end of Chapter 3, current research only focuses on one subject of the problem of explainable classifier and neglects the other subject – human. Thus, we study explainable classifier from the aspects of visualization and human computer interaction in this chapter. Broadly speaking, the visualization for explainable classifiers can be viewed as a special case of algorithm visualization or software visualization. The former aims to provide better understanding of algorithms for education purposes in computer science. The latter focus on assisting developers and operation engineers for the development and maintenance of complex software. Here, the subject of visualization is the classifiers, which can be treated as algorithms learned from the data, or complex systems that need assistance in understanding.

We view the development and the operation of an intelligent system as a system engineering problem, and divide the life cycle of an intelligent system into different stages. The classification system can be treat as a specification of the general intelligent system. Then, we identify the current issues in explaining classifiers and discuss the research opportunities of visualization regarding different stages.

### 4.1 Life Cycle of an intelligent system

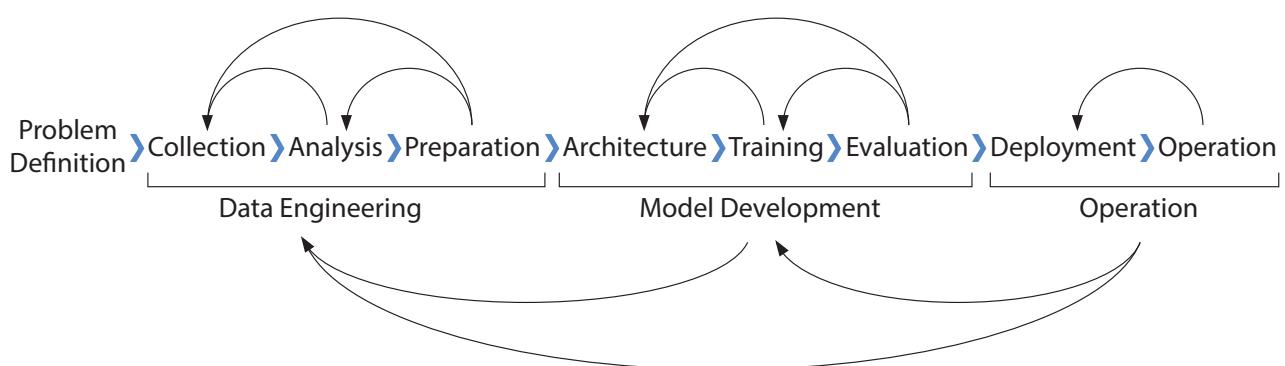


Figure 4.1. The life cycle of a classifier.

A classification system can be thought as a specialized case of an intelligence system. An intelligent system is developed to perform certain tasks with artificial intelligence (here we only consider data-driven systems). The development of a data-driven intelligent system is an iterative process. In this survey, the entire life cycle of an intelligent system is divided into three major stages (data engineering, model development, and operation) and eight sub-procedures (see Figure 4.1). This definition is formed based on the cross-industry standard process for data mining [42], a professional advice from Gartner, Inc. [31], the life cycle for expert system [20], and the machine learning workflow of Google Cloud<sup>1</sup>.

The first stage, data engineering, is defined to include any procedures that are data-related, namely, data collection, exploratory analysis, data preparation. The second stage, model development, includes procedures such as designing the architecture of the classifier (*e.g.*, what type of model to use and parameters), training the model using data prepared, and evaluating whether the model meets certain requirements. After developing a classifier, the model is deployed, and in certain cases, is operated by some people. As we can see from Figure 4.1, there are back-links from each stage/step to its previous stages/steps. This is the nature development. For example, we have a model with unsatisfactory performance after the training. This might due to model used is not suitable for certain tasks (go back to architecture setting), or it is because the volume of the data is too small (go back to collect more data). Similar problems might occur in other stages or procedures, which force us to go back and improve.

Visualization and visual analytics systems are semi-automatic solutions at different stages to make a classifier more explainable. At the stage of data engineering, data visualization can help humans explore the data and get a qualitative sense of the nature of the data. Since the training of a classifier is actually extracting information from the training data, with more knowledge of the data in mind, humans (*i.e.* developers or data scientists) can better understand if a failure results from the quality or volume issues of the data. At the second stage, visual analytics systems serve as development tools, which make the development more transparent and understandable. When designing or selecting model architectures, visual analytics systems can help humans better understand the characteristics of different classifiers, and even inspire improvements in the architecture. Visual diagnosing tools can help identify the problems in the training process and improve the debugging efficiency. For evaluations, visual analytics systems can help compare different classifiers and qualitatively evaluate the robustness and fairness of a classifier. At the last stage, when a classification system is deployed, visualization can help explain the inner workings of the system to end-

---

<sup>1</sup><https://cloud.google.com/ml-engine/docs/ml-solutions-overview>

users. For routine operations, visualization can better explain the predictions of the system, which make the monitoring and management easier. Also note that, visualization can be used in the life-cycle for other purposes instead of explainability. For example, monitoring the training process by plotting loss curves, or visualization for crowd sourcing to collect data with better quality.

---

A table summarizing all related papers

---

## 4.2 Visualization for Data Understanding

At the stage of data engineering, visualization can be used mainly in the procedure of data analysis to assist humans' understanding of the data. Data plays an important role in the success of machine learning advances. A trained classifier can be viewed as a machine that has extracted the information in the training data and abstracted the information as its parameters. Thus, understand the data is the first step to understand a classifier.

There is a long history of research in visualization for exploratory data analysis. When it comes to the

## 4.3 Visualization for Model Development

There is a surge of interest to use visualization for explainable classifiers, focusing on the stage of model development. We summarized three tasks for visualization, namely, model understanding, model diagnosing, and assessment and comparison, which are corresponding to the three stages: architecture design, training, and evaluation.

### 4.3.1 Understanding

Scientific understanding. Investigate the characteristic of the model.

Existing work:

### 4.3.2 Diagnosing

Diagnose model and data.

Existing work:

### **4.3.3 Assessment and Comparison**

Unquantifiable assessments, Fairness (e.g., discrimination), Vulnerability

Existing work:

## **4.4 Visualization for Model Operation**

### **4.4.1 Trust Establishment**

### **4.4.2 Monitoring**

## **4.5 Other Applications**

### **4.5.1 Teaching and Communicating Models**

Narrative, Interactive, etc. to explain your model to others.

### **4.5.2 Learn from the Model**

Knowledge Discovery; Learn lessons from what the model learned (Alpha Go)

## **4.6 Evaluation**

Review methods and standards of evaluating visualization.

Address the problem of the lack of evaluation standards for visualization for explainable classifiers.

Proposed?

1. Fidelity. How visualization reflects the real model. (The relativity and faithfulness of explanation)
2. Understandability. How easy the visualization is to be understood.

# **CHAPTER 5**

# **CONCLUSION**

Placeholder for Conclusion.

## BIBLIOGRAPHY

- [1] B. Alsallakh, A. Jourabloo, M. Ye, X. Liu, and L. Ren, "Do convolutional neural networks learn class hierarchy?" *IEEE Transactions on Visualization and Computer Graphics*, vol. PP, no. 99, pp. 1–1, 2017.
- [2] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek, "On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation," *PLOS ONE*, vol. 10, no. 7, p. e0130140, 2015. [Online]. Available: <http://dx.doi.org/10.1371/journal.pone.0130140>
- [3] D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba, "Network dissection: Quantifying interpretability of deep visual representations," in *Computer Vision and Pattern Recognition*, 2017.
- [4] K. D. Bock, K. Coussement, and D. V. den Poel, "Ensemble classification based on generalized additive models," *Computational Statistics & Data Analysis*, vol. 54, no. 6, pp. 1535 – 1546, 2010.
- [5] L. Breiman, J. Friedman, C. J. Stone, and R. A. Olshen, *Classification and regression trees*. CRC press, 1984.
- [6] M. H. Brown, "Algorithm animation," Ph.D. dissertation, Providence, RI, USA, 1987, uMI Order No. GAX87-15461.
- [7] W. Clancey, "The epistemology of a rule-based expert system: A framework for explanation," Stanford, CA, USA, Tech. Rep., 1981.
- [8] F. Doshi-Velez and B. Kim, "Towards a rigorous science of interpretable machine learning," *arXiv*, 2017. [Online]. Available: <https://arxiv.org/abs/1702.08608>
- [9] T. Downs, K. E. Gates, and A. Masters, "Exact simplification of support vector solutions," *Journal of Machine Learning Research*, vol. 2, no. Dec, pp. 293–297, 2001.
- [10] S. A. Dudani, "The distance-weighted k-nearest-neighbor rule," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. SMC-6, no. 4, pp. 325–327, April 1976.
- [11] D. Erhan, Y. Bengio, A. Courville, and P. Vincent, "Visualizing higher-layer features of a deep network," University of Montreal, Tech. Rep. 1341, Jun. 2009, also presented at the ICML 2009 Workshop on Learning Feature Hierarchies, Montréal, Canada.

- [12] R. Féraud and F. Clérot, "A methodology to explain neural network classification," *Neural Netw.*, vol. 15, no. 2, pp. 237–246, Mar. 2002.
- [13] Google Inc. (2017) PAIR | people + ai research initiative. [Online]. Available: <http://ai.google/pair>
- [14] D. Gunning, "Explainable artificial intelligence (XAI)," *Defense Advanced Research Projects Agency (DARPA)*, 2017.
- [15] C. Hempel and P. Oppenheim, "Studies in the logic of explanation," *Philosophy of Science*, vol. 15, no. 2, pp. 135–175, 1948.
- [16] L. Hendricks, Z. Akata, M. Rohrbach, J. Donahue, B. Schiele, and T. Darrell, "Generating visual explanations," *CoRR*, vol. abs/1603.08507, 2016. [Online]. Available: <http://arxiv.org/abs/1603.08507>
- [17] A. Karpathy, J. Johnson, and L. Fei-Fei, "Visualizing and understanding recurrent networks," in *International Conference on Learning Representations (ICLR) Workshop*, 2016.
- [18] J. M. Keller, M. R. Gray, and J. A. Givens, "A fuzzy k-nearest neighbor algorithm," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. SMC-15, no. 4, pp. 580–585, July 1985.
- [19] P. W. Koh and P. Liang, "Understanding black-box predictions via influence functions," in *Proceedings of the 34th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, D. Precup and Y. W. Teh, Eds., vol. 70. International Convention Centre, Sydney, Australia: PMLR, 06–11 Aug 2017, pp. 1885–1894. [Online]. Available: <http://proceedings.mlr.press/v70/koh17a.html>
- [20] A. J. La Salle and L. R. Medsker, "The expert system life cycle: What have we learned from software engineering?" in *Proceedings of the 1990 ACM SIGBDP Conference on Trends and Directions in Expert Systems*, ser. SIGBDP '90. New York, NY, USA: ACM, 1990, pp. 17–26.
- [21] B. Letham, C. Rudin, T. McCormick, and D. Madigan, "Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model," *Ann. Appl. Stat.*, vol. 9, no. 3, pp. 1350–1371, 09 2015. [Online]. Available: <https://doi.org/10.1214/15-AOAS848>
- [22] J. Li, X. Chen, E. Hovy, and D. Jurafsky, "Visualizing and understanding neural models in nlp," in *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. San

- Diego, California: Association for Computational Linguistics, June 2016, pp. 681–691. [Online]. Available: <http://www.aclweb.org/anthology/N16-1082>
- [23] B. Liu, M. Wang, H. Foroosh, M. Tappen, and M. Penksy, "Sparse convolutional neural networks," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015, pp. 806–814.
- [24] T. Lombrozo, "The structure and function of explanations," *Trends in Cognitive Sciences*, vol. 10, no. 10, pp. 464 – 470, 2006. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1364661306002117>
- [25] T. Munzner, *Visualization analysis and design*. CRC press, 2014.
- [26] W. J. Murdoch and A. Szlam, "Automatic rule extraction from long short term memory networks," in *International Conference on Learning Representations (ICLR)*, 2017.
- [27] R. Neches, W. Swartout, and J. Moore, "Enhanced maintenance and explanation of expert systems through explicit models of their development," *IEEE Transactions on Software Engineering*, vol. SE-11, no. 11, pp. 1337–1351, Nov 1985.
- [28] B. Price, I. Small, and R. Baecker, "A taxonomy of software visualization," vol. ii. IEEE Publishing, 1992, pp. 597–606.
- [29] J. R. Quinlan, "Simplifying decision trees," *International journal of man-machine studies*, vol. 27, no. 3, pp. 221–234, 1987.
- [30] M. T. Ribeiro, S. Singh, and C. Guestrin, ""why should i trust you?": Explaining the predictions of any classifier," in *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. New York, NY, USA: ACM, 2016, pp. 1135–1144.
- [31] C. Sapp. (2017) Preparing and architecting for machine learning. [Online]. Available: <https://www.gartner.com/doc/3573617/preparing-architecting-machine-learning>
- [32] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and D. Hassabis, "Mastering the game of go without human knowledge," *Nature*, vol. 550, no. 7676, pp. 354–359, Oct. 2017. [Online]. Available: <http://dx.doi.org/10.1038/nature24270>

- [33] K. Simonyan, A. Vedaldi, and A. Zisserman, "Deep inside convolutional networks: Visualising image classification models and saliency maps," in *International Conference on Learning Representations (ICLR) Workshop*, 2014.
- [34] D. Smilkov, N. Thorat, B. Kim, F. B. Viégas, and M. Wattenberg, "Smoothgrad: removing noise by adding noise," *CoRR*, vol. abs/1706.03825, 2017. [Online]. Available: <http://arxiv.org/abs/1706.03825>
- [35] J. T. Stasko, "Tango: a framework and system for algorithm animation," *Computer*, vol. 23, no. 9, pp. 27–39, Sept 1990.
- [36] W. Swartout, C. Paris, and J. Moore, "Explanations in knowledge systems: design for explainable expert systems," *IEEE Expert*, vol. 6, no. 3, pp. 58–64, June 1991.
- [37] M. Tan, L. Wang, and I. W. Tsang, "Learning sparse svm for feature selection on very high dimensional datasets," in *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, 2010, pp. 1047–1054.
- [38] M. E. Tipping, "Sparse bayesian learning and the relevance vector machine," *Journal of machine learning research*, vol. 1, no. Jun, pp. 211–244, 2001.
- [39] B. Ustun and C. Rudin, "Supersparse linear integer models for optimized medical scoring systems," *Machine Learning*, vol. 102, no. 3, pp. 349–391, Mar 2016. [Online]. Available: <https://doi.org/10.1007/s10994-015-5528-6>
- [40] F. Wang and C. Rudin, "Falling Rule Lists," in *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, G. Lebanon and S. V. N. Vishwanathan, Eds., vol. 38. San Diego, California, USA: PMLR, 09–12 May 2015, pp. 1013–1022. [Online]. Available: <http://proceedings.mlr.press/v38/wang15a.html>
- [41] T. Wang, C. Rudin, F. Doshi-Velez, Y. Liu, E. Klampfl, and P. MacNeille, "A bayesian framework for learning rule sets for interpretable classification," *Journal of Machine Learning Research*, vol. 18, no. 70, pp. 1–37, 2017. [Online]. Available: <http://jmlr.org/papers/v18/16-003.html>
- [42] R. Wirth and J. Hipp, "Crisp-dm: Towards a standard process model for data mining," in *Proceedings of the 4th international conference on the practical applications of knowledge discovery and data mining*, 2000, pp. 29–39.

- [43] M. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Computer Vision – ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part I*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds. Cham: Springer International Publishing, 2014, pp. 818–833.
- [44] L. Zintgraf, T. Cohen, T. Adel, and M. Welling, "Visualizing deep neural network decisions: Prediction difference analysis," in *International Conference on Learning Representations (ICLR)*, 2017.