

Explainable AI

Adam Noack

1 Introduction

AI is popping up everywhere. Examples are ubiquitous: Facial recognition software on our phones' cameras, natural language processing capabilities in personal assistants such as Siri, Alexa, and Google Assistant, self driving cars almost necessarily have a large number of different machine learning algorithms to help in the navigation task, Netflix, Spotify, and Amazon all use AI to deliver the content to its users that is most engaging, the list goes on and on.

AI's boom over the past decade or so has in large part been because of deep learning. Deep learning is a method of learning that relies on very large and deep neural networks. These networks can have millions of parameters and connections are extremely complex. Because of this, they are increasingly impossible for humans to understand without some sort of simplifying apparatus. The problem is that there aren't many of these apparatuses that are functional.

As more and more tasks are offloaded onto machines controlled by these extremely complex black-box algorithms, the more we will come to rely on predictions and approximations that we can in no way verify. That is, unless we devise some method by which we can interpret the output of any given model.

The approaches defined in the following section explain some of the more recent and successful attempts at making these black boxes a bit whiter.

2 Approaches

2.1 Anchors: High-Precision Model-Agnostic Explanations

Ribeiro, the creator of the LIME approach, introduces anchors in this paper, a new method by which black box machine learning models can be explained.

What are anchors? Take a rule, a set of predicates for a given dataset, and for each data sample determine which of those predicates are present in the sample. The rule is an anchor for a particular sample if each predicate is present in the sample and the presence of all the predicates in the rule is highly correlated with a particular output class. For example, if the set of words $A = \{\text{'not'}, \text{'bad'}\}$ is rule for a set of text data, $A(x_i)$ will return 1 and be an anchor for x_i if all of the feature predicates are present in a sample x_i . Assuming that $A(x_i) = 1$ predicts the sentiment “positive” for this dataset with sufficiently high probability, A is an *anchor* for the sample $x_1 = \text{“This movie is not bad.”}$ because $A(x_1)$. The anchor can then be interpreted as an explanation for the class prediction.

How are anchors computed? One method by which Ribeiro computes anchors for a given feature space begins with an empty set of predicates $B = \{\}$ and iteratively adds new features to B . In other words, each time through the search loop, the most explanatory, unused feature predicate is selected and added to B . The search loop terminates when the class for a given sample x_i can be predicted with high probability given only the fact that $B(x_i) = 1$.

Anchors are useful because they are extremely easy and intuitive for humans to interpret. Simple rules such as the presence of a few words, in the case of textual data, or the presence of a few superpixels, in the case of image data, are easily communicated to and understood by the user.

2.2 Extracting Rules from Artificial Neural Networks with Distributed Representations

In this paper, Sebastian Thrun uses validity interval analysis to generate easily understood rules for a neural network that predicts whether or not a robotic arm is in a safe position or not.

Thrun starts off by providing the motivation for this line of research. He explains that many common methods for generating rules for neural networks are not applicable in cases where the internal representation of the neural network is distributed in nature, ie, the network is not sparse, and the activations are continuous as opposed to binary.

Validity interval analysis works by attaching intervals to every neuron’s activation value. For example, for neuron x_i , $\sigma(x_i) \in [a_i, b_i]$. For the robotic arm example, there are two classes of output, “safe” and “unsafe”. A single output neuron can be used to cover the output space: the value for x_{output} will be predicted “safe” if $\sigma(x_{output}) \geq .5$ and “unsafe” if $\sigma(x_{output}) < .5$. It’s true then that any combination of input features that satisfies the constraint that the output neuron’s activation value must lie in the interval $[0, .5)$ to will always be classified as “unsafe”.

In practice, constructing rules for neural networks using validity interval analysis is usually done in one of two ways, specific to general or general to specific. Thrun took a specific to general approach. He would start with a single data sample, initially assigning each neuron x_i ’s interval to be $[\sigma(x_i), \sigma(x_i)]$. From there, the intervals for each neuron were gradually widened until the output prediction class changed. The range of input values that “worked” in that range became a rule for that output class. For example, if the input to the net is a five element array, a , and the i th feature of a can be accessed by a_i , then a learned rule for the network might be written as follows: *if $a_2 < 50$ and $a_4 > 75$ then “safe”*.

2.3 Learning with Interpretable Structure from RNN

RNNs are very difficult for humans to understand – arguably even more difficult than standard feed-forward neural networks. This is in large part due to the fact that their hidden unit values are represented numerically rather than semantically. Interestingly though, the task that RNNs are so good at performing, interpreting sequences of data, can also be performed by finite state automata (FSA). This is particularly interesting because states in FSA are represented semantically and are therefore much easier for humans to understand than RNN states.

It was previously known that the values of the hidden states of classical RNNs tend to clump together. This fact led Bo-Jian Hou to believe that the hidden states of classical RNNs could be represented as semantically separate states in an FSA. I.e. each clumping of hidden unit values could be interpreted as a separate state in an FSA.

However, classical RNN units suffer from the vanishing or exploding gradient issue, so they have in large part been deprecated in favor of gated recurrent units. It was not known if the values of gated RNN units such as the long short-term memory (LSTM) unit, gated recurrent units (GRUs), or the minimal gated unit (MGU) clumped together. After training several separate gated unit RNNs Hou determined that the values of gated units also clump together. With this information at hand, Hou set to work converting trained gated RNNs into equivalent FSA.

Hou did this in the following manner: He first trained gated RNNs on data. He then used k-means to find clusters of semantically similar hidden state values. These clusters were then made into the states of the FSA. The FSA’s vocabulary was then created by simply enumerating all of the possible inputs to the RNN. A transition matrix was then learned that defined how one should move given the current state and the next input.

Hou obtained interesting results. He found that of the three different types of gated unit

RNNs tested, MGU RNNs could be represented by an FSA with the fewest number of states. This implies that MGUs are much easier to interpret via FSA than LSTMs and GRUs. Given the fact that MGUs tend to be only slightly outperformed by LSTMs and GRUs, MGUs might be the better option if one desires both interpretability and performance.

2.4 Learning to Explain: An Information-Theoretic Perspective on Model Interpretation

This approach to model explanation, Learning to Explain or L2X as it's called, functions at a high level in the following manner. For a given input sample x existing in \mathbb{R}^d , L2X seeks to find a vector x_s existing in \mathbb{R}^k where $k < d$ and the exact value for k is defined by the user of the model. In other words, like many of the other approaches to model explanation, L2X produces explanations by rank-ordering the features for a particular sample that are most influential in that sample's prediction.

L2X is a bit different, though, in the way it decides on the most important features for a sample. It finds the subset of features X_s that maximizes the mutual information between X_s and the correct output value Y . Mutual information can be thought of as the amount of information gained about one event given the observation of a different event. If the mutual information between X_s and Y is high, conceptually, that means that observing X_s gives us a lot of insight into Y 's value.

A direct method for finding the subset of features that maximizes mutual information is not possible, so in this paper the solution is approximated using a variational approach.