

Explainable Artificial Intelligence: An Overview

Sargur N. Srihari
Department of Computer Science and Engineering
University at Buffalo, The State University of New York
srihari@buffalo.edu

Abstract

With a wide range of applications, Artificial Intelligence (AI) has spawned a spectrum of research activity on AI-related topics. One such area is that of explainable AI. It is a vital component of trustworthy AI systems. This paper provides an overview of explainable AI methods describing both post-hoc AI systems, which provide explanations with previously built conventional AI systems, and ante-hoc AI systems, which are configured from the start to provide explanations. The explanations take various forms: explanation based on features, explanation based on illustrative training samples, explanation based on embedded representations, and explanation based on heat-maps. There are also probabilistic explanations which combine neural network models with graphical models. Explainable AI is closely associated with many AI research topic frontiers such as neuro-symbolic AI and machine teaching.

Keywords: Trustworthy AI, Explainable AI, Post-hoc Explainability, Ante-hoc Explainability

Contents

1 Introduction

1.1 AI sub-disciplines

1.2 Trustworthy AI

1.3 AI methods

2 Explainable Artificial Intelligence

2.1 Need for XAI

2.2 Measures of explanation effectiveness

2.3 Taxonomy of XAI methods

3 Post-hoc XAI

3.1 Measures of Explanation Quality

3.1.1 Sensitivity Analysis

3.1.2 Layerwise Relevance Propagation

3.1.3 Evaluation of SA and LRP

3.2 Input Features as Explanation

3.3 Examples as Explanation

3.3.1 Machine Teaching

3.3.2 Bayesian Teaching

4 Ante-Hoc XAI

4.1 Explanation from Representation

4.1.1 Reverse time Attention Model

4.1.2 Explanations from Embeddings

4.2 Probabilistic Explanations

4.2.1 Bayesian Deep Learning

4.2.2 Graphical Model Inference

5 Concluding Remarks

6 References

1 Introduction

Artificial Intelligence (AI) is everywhere. There are billions of searches on handheld devices everyday. Smart phones use facial recognition. Alexa cutely answers our questions. The key element in Tic Toc is its recommender system. More generally, AI enables performing tasks requiring human cognition as well as decision-making.

While AI systems already incorporate intelligent behavior, they continue to be improved with the need to exhibit flexibility, resourcefulness, creativity, real-time responsiveness, and long-term reflection to demonstrate competence in complex environments and social contexts. This paper is an overview of one of the several sub-areas of active AI research known as Explainable AI (XAI). First we set the stage for where XAI fits into the spectrum of AI research topics and methods.

1.1 AI sub-disciplines

While AI has already been incorporated into a wide range of applications, it is also a topic of a great deal of current research. AI research areas can be divided into five areas as follows according to National Science Foundation (2019):

1. *Core AI*: Theory and methods for: (i) learning, abstraction, and inference (ii) architectures for intelligence and multi-agent systems. ML has made great advances through algorithms, computing power, and growing data. Other technologies include knowledge representation, logical and probabilistic reasoning, planning, search, constraint satisfaction, and optimization.
2. *Biologically-inspired AI*: Models may be inspired by living systems: connectionism, behavior, emergence. Computational neuroscience which deals with the theory of computation in the nervous system Behavioral and cognitive science Typical of human perceptual, motor, and cognitive processes and their interactions
3. *Perception and Communication*: Computer vision methods to sense and reason about visual world. Human language technologies (also called NLP, NLU) To analyze, produce, translate, and respond to human text and speech.

4. *Embodied AI*: Intelligent systems may be able to act upon the world through embodiment. Robotics is closely aligned with but not identical to embodied AI. An embodied AI may be a robot.
5. *Trustworthy AI*: AI amplifies human capabilities to accomplish individual and collective goals, Need to assess benefits, effects, and risks. How human, technical, and contextual aspects of systems interact to shape those effects. Relevant aspects of trustworthy AI are: Explainable AI (XAI), Validation of AI-enabled systems, AI safety, security, and privacy (including, for example, role of emotion and affect in the design and perception of AI).

1.2 Trustworthy AI

Trust is key in adoption of AI for economic growth and innovations to benefit society. Today, ability to understand AI decisions and measure their trustworthiness is limited. For it to be trustworthy, AI has to be: trusted to function reliably, trusted to be able to explain conclusions, trusted not to violate privacy, and trusted not to exhibit socially harmful bias.

It is the explainability aspect of trustworthy AI that we explore further here. Explanations are vital in decision making. Establishing human trust in the outcome requires the exchange of reasons for that outcome. Explanations must be in terms as appropriate to the task and as needed by users. Explanations help pinpoint errors or data. Research challenges include finding ways to make “black box” AI systems explainable. Models and frameworks for learning and reasoning that are both inherently explainable and powerful. Integrating psychology, cognitive science, to better understand and acceptability of an explanation.

1.3 AI methods

It has long been understood that designing AI needs knowledge. On a historic time-scale, efforts at doing this may be characterized as consisting of several overlapping waves. The first wave of AI began with the knowledge-based approach, where an input is transformed by a hand-designed program to the desired output, e.g., a rule-based expert system. On a parallel track the machine learning approach was developed, where the input is first transformed by a hand-designed program into features but the features were mapped to the output by a program that learns from examples.

The second wave of AI began by replacing feature engineering with representation learning, where the features are learnt automatically. Deep learning involves several representation layers. (Goodfellow, Bengio, and Courville, 2016). First simple features are learnt. Additional layers extract more abstract features. The final layers map the abstract features to the output. Deep learning allows AI systems to rapidly adapt to new tasks, since designing features can take great human effort – often decades for a community of researchers. It does not need

programmer to have deep knowledge of the problem domain. The two waves of AI are shown in Fig. 1.

Figure 1: Two waves of AI: (a) The first wave consisted of knowledge-based and machine-learning approaches, and (b) The second wave consists of representation learning and deep learning approaches.

An emerging third wave of AI may be defined as neurosymbolic AI. It is essentially the combination of deep learning with symbolic reasoning. A symbolic reasoning process is used to bridge the learning of visual concepts, words and semantic parsing of sentences without explicit annotations for any of them (Mao, Gan, Kohli, Tenenbaum, and Wu, 2019). Symbolic approaches are usually constructed using graphical models (Koller and Friedman, 2009). Probabilities have been very much a part of the earlier waves, both in discriminative machine learning approaches as well as in generative approaches such as adversarial networks. The neurosymbolic approach relies on generative models at the symbolic level, where the symbols are computed using deep learning. Generative models, such as generative adversarial networks can be used to construct distributions. The symbolic approach calls for algorithms/architectures for: representation (such as Bayesian nets, Markov Random Fields) and inference (Exact, Approximate, Monte Carlo, Variational).

2 Explainable Artificial Intelligence

Explanations are vital in decision making. Establishing human trust in the outcome requires the exchange of reasons for that outcome. Explanations must be in terms as appropriate to the task and as needed by users. Explanations help pinpoint errors or data. Explanations help to detect bias in the system thereby leading to ethical AI.

Research challenges of XAI include finding ways to make “black box” AI systems explainable. Models and frameworks for learning and reasoning that are both inherently

explainable and powerful. Integrating psychology, cognitive science, to better understand and acceptability of an explanation.

One of the main criticisms of deep learning is opaqueness, i.e., having the characteristics of a blackbox. Formally, a blackbox is a function that is too complicated for any human to comprehend, a function that is proprietary, or model that is difficult to trouble-shoot. Deep Learning Models are blackbox models because they are recursive, non-intuitive and difficult for people to understand. See Fig. 2.

Figure 2: AI as a blackbox. Source image: Wikipedia

The role of explanation in a human-machine interactive scenario is given in Fig. 3. The explanation interface is one capable of answering the following types of queries Turek (2018):

1. Why did you do that?
2. Why not something else?
3. When do you succeed?
4. When do you fail?
5. When can I trust you?
6. How do I correct an error?

Figure 3: Explanation Framework. Source: (Turek, 2018).

2.1 Need for XAI

There are numerous reasons for certain AI deployments to be explainable. Some important ones are: justification, control, discovery, and improvement. We briefly describe each need here.

Explain to justify: The ability to explain one's decision to other people is an important aspect of human intelligence. Understanding the rationale behind the model's predictions would help users decide when to trust or not to trust their predictions.

Explain to control: This refers to compliance to legislation. For instance, in making credit decisions, how did the model decide to provide or deny credit to an individual? Was there bias: ethnicity, race religion? In the US the lender must provide reasons for adverse decision, such as take-home insufficient, insufficient collateral, poor credit rating. In the European Union, GDPR (General Data Protection Regulation): right to explanation for high-stakes automated decisions. Another example is healthcare, which is highly regulated due to HIPAA. How did AI predict grade 3 or grade 4 tumor?

Explain to improve: The first step in improving a system is to understand its weaknesses, such as detecting bias in the system. For instance, in the medical domain, an anecdote is that medical AI decisions were worse with AI, e.g., patient discharge to a nursing home did not take into account personal circumstances. An explanation would help the human decision-maker over-ride the system decision.

Explain to discover: AI systems are trained with millions of examples. They may observe previously unseen patterns in data that people may find to be useful.

An example of explanation generation in a deep network in the computer vision domain is given in Fig. 4. Here the goal of the system is to classify the type of bird. After it has generated the class to be a downy woodpecker it also utilizes the definition of the bird as follows: "This bird has a white breast, black wings and a red spot on its head" to generate the explanation "This is a Downy Woodpecker because it is a black and white bird with a red spot in its crown."

*Figure 4: An example of explanation generated by a deep network:
"This is a Downy Woodpecker because it is a black and white bird
with a red spot in its crown." Source: (Hendricks et al, 2016).*

After training a deep network we have large networks that work very well, but hard to tell how. They can fail unintuitively. Adversarial examples show this (See Fig. 5).

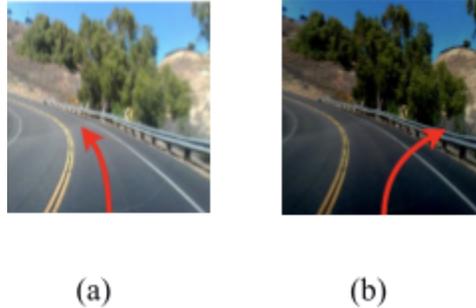


Figure 5: Unexpected failure: (a) correct steering in daytime lighting and (b) wrong steering in fading light.
Source image: (Pei, Cao, Yang, and Jana, 2017).

2.2 Measures of explanation effectiveness

Several measures of explanation effectiveness have been proposed Turek (2018):

User satisfaction: clarity of explanation (user rating), utility of explanation (user rating)

Mental model: understanding individual decisions, understanding the overall model, strength/weakness assessment, “what will it do” prediction, “How do I intervene” prediction

Task performance: Does it improve the user’s decision, task performance? Artificial decision tasks introduced to diagnose the user’s understanding

Trust assessment: appropriate future use and trust

Correctability: identifying errors, correcting errors, continuous training

2.3 Taxonomy of XAI methods

With wide adoption of AI in industry and government, the need for XAI has also grown commensurately. Existing XAI methods can be divided into two broad categories (See Fig. 6):

1. Post-hoc (Explain the Blackbox): Explainability based on test cases and results
2. Ante-hoc (Build a new learning model): Seeding explainability into model from the start.

Figure 6: Ante- and Post-Hoc Explainable AI. Source: (Marselis, 2019).

Some examples of post-hoc XAI systems are Sensitivity Analysis (SA), Layer-wise Relevance Propagation (LRP), and Local Interpretable Model-Agnostic Explanations (LIME). Examples of Ante-hoc XAI systems are: Reversed Time Attention Model (RETAIN), and Bayesian Deep Learning (BDL). We discuss each of these types of XAI systems next.

3 Post-hoc XAI

Among methods for visualizing, interpreting and explaining deep learning models, two popular techniques for explaining predictions are Sensitivity Analysis and Layerwise Relevance Propagation. We discuss each of these followed by objectively comparing the quality of the explanations provided.

3.1 Measures of Explanation Quality

Here we describe two measures of explanation quality for evaluating the performance of a deep network: Sensitivity Analysis and Layerwise Relevance Propagation. We then compare their efficacy on different tasks.

3.1.1 Sensitivity Analysis

Assumes that most relevant features are those to which output is most sensitive. Consider the input image x in Fig. 7. The system correctly classifies the input image as “rooster”. Then, an explanation method is applied to explain the prediction in terms of input variables. The result of this explanation process is a heatmap visualizing the importance of each input variable I (pixel) for the prediction $f(x)$. In this example the rooster’s red comb and wattle are the basis for the AI system’s decision. Sensitivity analysis (SA) explains a prediction based on the model’s locally evaluated gradient (partial derivative). This amounts to which pixels need to be changed to make image look more/less like the predicted class, e.g., changing yellow occluding pixels improves score, but does not explain rooster.

How changes in each pixel affect the score is given by the partial derivatives

$$R_I = \left\| \frac{\partial}{\partial x_I} f(x) \right\|$$

SA explains the prediction based on a locally evaluated gradient. However, it does not explain $f(x)$ but a variation of the input.

Figure 7: XAI using Sensitivity Analysis: Changing yellow (occluding pixels) improves score, but does not explain rooster. Source: (Samek, Wiegand, and Muller, 2018).

3.1.2 Layerwise Relevance Propagation

Layerwise Relevance Propagation (LRP) explains the classifier's prediction using decomposition. See Fig. 8. It redistributes the prediction $f(x)$ backwards using local redistribution rules until it assigns a relevance score R_i to each input variable (e.g., image pixel). The key property of this redistribution process is referred to as relevance conservation. It can be summarized by a sequence of sums as follows:

$$\sum R_i = \dots = \sum R_j = \sum R_k = f(x)$$

At every step of the redistribution process (e.g., at every layer of a deep neural network), the total amount of relevance (i.e., the prediction $f(x)$) is conserved. No relevance is artificially added or removed during redistribution. The relevance scores R_i of each input variable determines how much this variable has contributed to the prediction. Thus, in contrast to sensitivity analysis, LRP truly decomposes the function value $f(x)$.

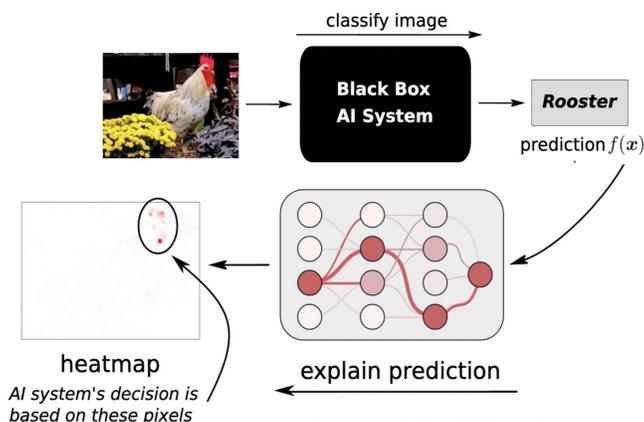


Figure 8: Layerwise Relevance Prediction. In this example the rooster's red comb and wattle are the basis for the AI system's decision. With the heatmap one can verify that the AI system works as intended. Source: (Samek et al, 2018).

The LRP redistribution process for feed-forward neural networks is as follows. Let x_j be the neuron activations at layer l , R_k be the relevance scores associated to the neurons at layer

$l + 1$ and w_{jk} be the weight connecting neuron j to neuron k . The simple LRP rule redistributes relevance from layer $l + 1$ to layer l in the following way:

$$R_j = \sum_k \frac{x_j w_{jk}}{\sum_j x_j w_{jk} + \epsilon} R_k$$

where the small stabilization term ϵ prevents division by zero. Intuitively, this rule redistributes relevance proportionally from layer $l + 1$ to each neuron in layer l based on two criteria, namely (i) the neuron activation x_j , i.e., more activated neurons receive a larger share of relevance, and (ii) the strength of the connection w_{jk} , i.e., more relevance flows through more prominent connections. Note that relevance conservation holds for $\epsilon = 0$.

3.1.3 Evaluation of SA and LRP

Heatmaps produced by different explanation methods can be used to measure the quality of explanation using *perturbation analysis*. It is based on the idea that perturbing input variables important for prediction leads to a steeper prediction score decline. Input variables are sorted by relevance score, and iteratively perturbed (starting from the most relevant ones). The prediction score is tracked after every perturbation step. The average decline of the prediction score is a measure of explanation quality; a large decline indicates successful explanation.

Evaluation of SA and LRP explanations on three different problems/classifiers are described next: (i) image classification using GoogleNet, (ii) document text classification using a convolutional neural network and (iii) recognition of human actions in videos using a Fisher Linear Discriminant/SVM.

1. Image Classification

A deep neural network, GoogleNet, was used to classify general objects from the ILSVRC2012 dataset. Fig. 9(a) shows two images correctly classified as “volcano” and “coffee cup”. The accompanying heatmaps visualize the explanations obtained with SA and LRP. The LRP heatmap of the coffee cup image shows that the model has identified the ellipsoidal shape of the cup to be a relevant feature for the category. In the volcano example, the shape of the mountain is regarded as evidence for a volcano. The SA heatmaps are much noisier than the ones computed with LRP and large values R_i are assigned to regions consisting of pure background, e.g., the sky, although these pixels are not really indicative for image category “volcano”. In contrast to LRP, SA does not indicate how much every pixel contributes to the prediction, but it rather measures the sensitivity of the classifier to changes in the input. Therefore, LRP produces subjectively better explanations of the model’s predictions than SA. Perturbation analysis (lower part of Fig. 9(a)) shows that LRP provides better explanations than SA— due to faster prediction score decrease using LRP heatmaps than using SA heatmaps.



Figure 9: Explanation Quality Measures: (a) Image classification. Two images correctly classified as volcano, coffee cup by a deep learning network. Heat maps and perturbation analysis show better explanatory power of LRP than SA. (b) Text Classification: SA and LRP heatmaps identify words such as “discomfort”, “body” and “sickness” as relevant ones for explaining the prediction of medicine. In contrast to SA, LRP distinguishes between positive (red) and negative (blue) relevances. (c) Human Action Recognition. Explaining prediction sit-up. The LRP heatmaps of a video which was classified as “sit-up” show increased relevance on frames in which the person is performing an upwards and downwards movement. Source: (Samek et al, 2018).

2. Text Classification

In this experiment a word-embedding based convolutional neural network was trained to classify text documents from the 20Newsgroup dataset.

Fig. 9(b) shows SA and LRP heatmaps (e.g., a relevance score R_i is assigned to every word) overlaid on top of a document, which was classified as “sci.med”, i.e., medical topic. Both SA and LRP indicate that words such as “sickness”, “body” or “discomfort” are the basis for this classification decision. In contrast to SA, LRP distinguishes between positive (red) and negative (blue) words, i.e., words which support “sci.med” and words which speak for another category (e.g., “sci.space”). Words such as “ride”, “astronaut” and “shuttle” strongly speak for space, but not necessarily for medicine. With the LRP heatmap we can see that although the classifier decides for the correct “sci.med” class, there is evidence in the text which contradicts this decision. The SA method does not distinguish between positive and negative evidence. As before, perturbation analysis shows that LRP provides more informative heatmaps than SA, because these heatmaps lead to a larger decrease in classification accuracy compared to SA heatmaps.

3. Human Action Recognition

In this experiment a Fisher Vector / SVM classifier was trained for predicting human actions from compressed videos. To reduce computation, the classifier was trained on block-wise motion vectors (not individual pixels). The evaluation was performed on the HMDB51 dataset. Fig. 9(c) shows LRP heatmaps overlaid onto five exemplar frames of a video sample. The video was correctly classified as showing the action “sit-up”. The model focuses on blocks surrounding the upper body of the person as this part of the frame shows motion indicative of

“sit-up”, i.e., upward and downward body movements. The curve at the bottom of Fig. 9(c) displays the distribution of relevance over (four consecutive) frames. The relevance scores are larger for frames in which the person is performing an upwards and downwards movement. Thus, LRP heatmaps not only visualizes the relevant locations of the action within a video frame (i.e., where relevant action happens), but it also identifies the most relevant time points within a video sequence (i.e., when relevant action happens).

3.2 Input Features as Explanation

One approach is to attempt to explain the predictions of any machine learning classifier by having access to its input features. The goal is to provide explanations of the form “A is something because of B,C and D.” For example, “This is a bird because it has feathers, wings and a beak.” Such an explanation is concise—there are not a hundred reasons. It relies on B,C,D which are also high level concepts.

Local Interpretable Model-Agnostic Explanations (LIME) is such a system (de Sousa, Vellasco, Sun, and da Silva, 2019). The use of LIME in a medical diagnostic application is shown in Fig. 10. Here the model predicts that a certain patient has the flu. The prediction is then explained by an “explainer” that highlights the symptoms that are most important to the model. The physician is thereby empowered whether to trust the model or not.

To

Figure 10: Input features in explaining medical diagnosis. Source: (de Sousa et al, 2019).

explain a classifier that predicts whether an image contains a tree frog (Fig. 11(a)) by means of the super-pixels (a group of pixels with the same value) in the input image. The explainer generate a set of perturbed instances by turning some interpretable components “off” (making them gray) as illustrated in Fig. 11 (b). For each instance, we get the probability that a tree frog is in the image according to the model. We then learn a simple (linear) model on this data set, which is locally weighted—that is, we care more about making mistakes in perturbed instances that are more similar to the original image. In the end, we present the superpixels with highest positive weights as an explanation, graying out everything else.

Figure 11: Input features in explaining in prediction with



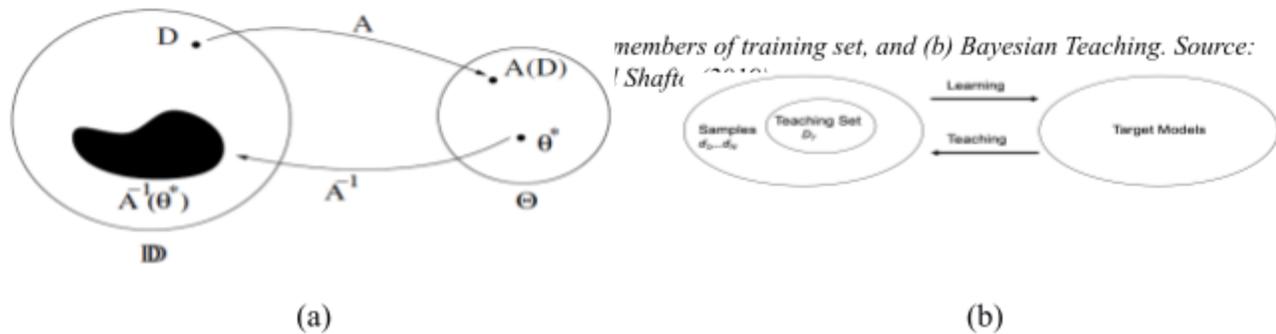
3.3 Examples as Explanation

The aim to select subset of the dataset that leads to similar conclusions as the entire dataset. The intuition is that subsets of training data that lead a model to the same (or approximately similar) inference as the model trained on all the data should be useful to understand the fitted model.

Explanation is viewed as the inverse of modeling. It is based on two fundamental observations: (i) all machine learning models are trained on data, and (ii) data is the common language of the user and the model. Examples may not capture what philosophers and cognitive scientists call as explanation, but they harness people's proclivity to inductive inference.

It is related to machine teaching in which the teacher designs the optimal training data to drive the learning algorithm to a target model.

3.3.1 Machine Teaching



Given a training dataset $D \in \mathbf{D}$, the process of machine learning returns a model $A(D) \in \theta$. A is in general many-to-one. Conversely, given a target model $\theta^* \in \Theta$, the inverse function A^{-1} returns a set of training examples that will result in θ^* . Machine Teaching aims to identify optimal member(s) among $A^{-1}(\theta^*)$. See Fig. 12(a).

3.3.2 Bayesian Teaching

The teaching problem is to select a small subset of data that with high probability leads the learner to model to the correct inference. See Fig. 12(b).

This requires two kinds of inference: (i) Teacher's inference (T) which is done in the space of possible teaching sets, and (ii) Learner's inference (L) which is done in the space of possible target models. For any subset of the training data \mathbf{x} the probability assigned by the model can be written as

$$P_T(\mathbf{x}|\theta) = \frac{P_L(\theta|\mathbf{x})P(\mathbf{x})}{\int_{\mathbf{x}} P_L(\theta|\mathbf{x})P_L(\mathbf{x})d\mathbf{x}}$$

where θ denotes the target model, which can be an entire model or a particular substructure, such as latent features, relations, grammars, programs, or combinations of these; $P_T(\mathbf{x}|\theta)$ is the probability of choosing \mathbf{x} as the teaching examples for explaining target model θ , $P_L(\theta|\mathbf{x})$

is learner's posterior inference after receiving \mathbf{x} , $P(\mathbf{x})$ describes bias for certain kind of examples (e.g., favoring smaller subsets), and the integral is over all partitions of the training data (i.e., if the size of \mathbf{x} is m and the size of the entire training corpus is N , there are ${}^N C_m$ partitions).

A Bayesian teacher:

With training data $D=\{d_1, d_2, \dots, d_N\}$ and teaching set size $n < N$ teaches a target model θ^* by sampling a teaching set $D_T \subseteq D$ from $\mathcal{D} = \{D | D \in \mathcal{P}(D) \wedge |D| \leq n\}$ according to

$$p(D_T|\theta^*) = \frac{p(D_T)p_L(\theta^*|D_T)}{p(\theta^*)} = \frac{p(D_T)p_L(\theta^*|D_T)}{\sum_{D \in \mathcal{D}} p(D)p_L(\theta^*|D)}$$

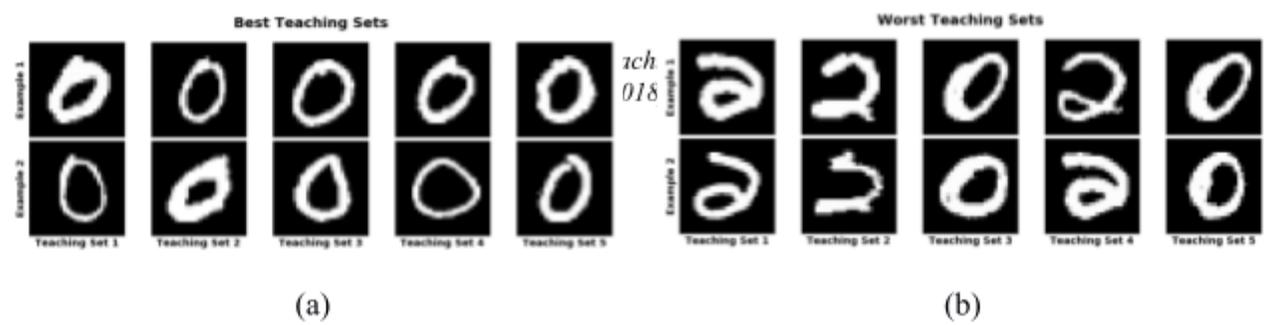
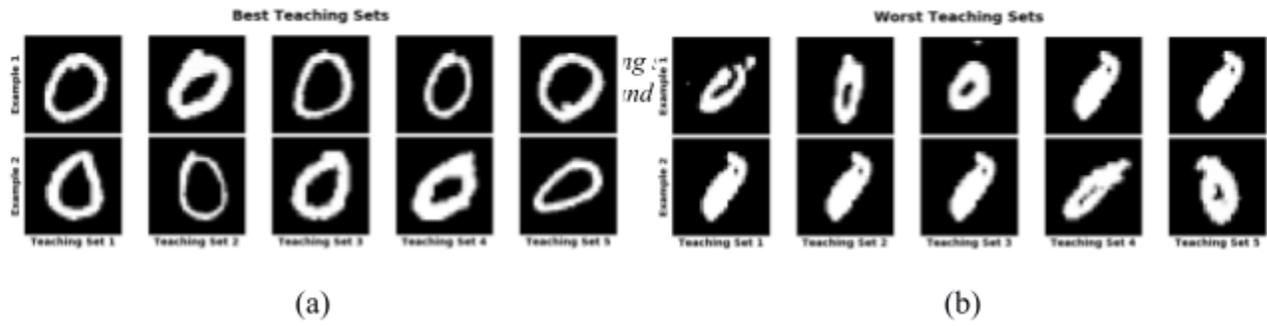
where $\mathcal{P}(D)$ is the power set of D and D is the space of teaching sets. $p_L(\theta^*|D)$ is the probability the learner will infer the target model θ^* given a particular teaching set D (i.e. the learner's posterior probability given that teaching set), and $p(D)$ is the teacher's prior probability on the same teaching set D . Priors assign higher probabilities to smaller teaching sets.

Teaching Inference for model explanation is as follows. Pick a teaching set size (e.g., 2) to constrain the search space. Perform teaching inference for the category to be understood. Rank the teaching sets based on the teaching probabilities. See Fig. 13.

Figure 13: Teaching Inference for model explanation. Source: (Yang and Shafto, 2019).

Consider the explanation of a new model at a conference. Authors show examples classified correctly and incorrectly. They mention that humans would find some to be hard. Rather than cherry-pick, we care about those classified correctly with high confidence, those classified correctly with low confidence, those classified incorrectly with low confidence, those classified incorrectly with high confidence. But some methods don't offer certainty estimates. Bayesian teaching offers finding/ranking such examples as seen below.

The five best teaching sets using *ground truth labels* are shown in Fig. 14 (a). Each pair of images represents a teaching set; pairs are sorted by teaching probabilities in descending order, from left to right (leftmost is best). The five worst teaching sets using ground truth labels are in Fig. 14 (b). Each pair of images represents a teaching set; pairs sorted by teaching probabilities in ascending order, from left to right, (leftmost set is worst). The five best teaching sets using *model predictions as labels* are shown in Fig. 15(a). The five worst teaching sets using model predictions as labels are in are shown in Fig. 15(b).



The model for Category 0 is explained by: (i) those classified correctly with high confidence (ii) those classified correctly with low confidence (iii) those classified incorrectly with low confidence (iv) those classified incorrectly with high confidence.

In summary, Bayesian teaching leverages the common understanding of model behavior—the data—to explain opaque models through the examples from the original data that are most representative of the inference. In doing so, it integrates learning and explanation by taking the learning model as input into the explanation process, and outputs an explanation in terms of the examples from the original data.

4 Ante-Hoc XAI

While post-hoc explanation techniques allow models to be trained normally, with explainability being incorporated as an afterthought, ante-hoc techniques entail baking explainability into a model from the beginning. The goal of transitioning from a conventional AI system to an ante-hoc AI system, in the context of image recognition, is illustrated in Fig. 16.



The goal of ante-hoc AI is to produce more explainable models, while maintaining a high level of performance, say prediction accuracy. The goal is to enable human users to understand,

trust and manage emerging AI partners. Do we want a complex black box model such as an RNN or a less accurate traditional model with better interpretation, say logistic regression., i.e., a 90% accurate model we understand versus a 99% accurate model we don't. The role of performance in ante-hoc AI is illustrated in Fig. 17.

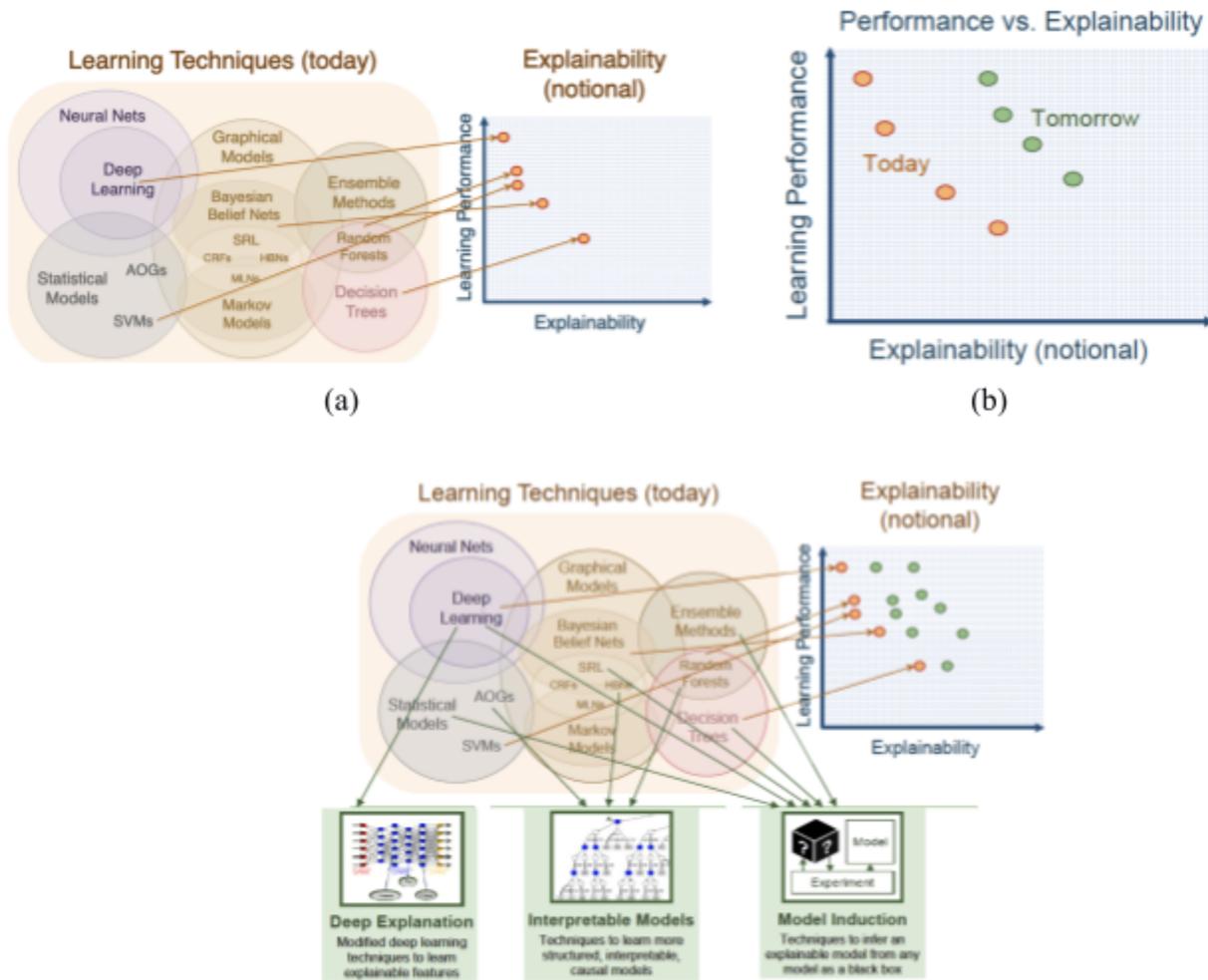


Figure 17: Performance of Explainable AI: (a) today, (b) tomorrow and (c) tomorrow's methods. Source: (Marselis, 2019).

Some active research efforts in ante-hoc AI are:

Explanation from Representation. Techniques to identify the most salient input features used in a decision, e.g., it is a cat because it has whiskers and fur. They include techniques to select the training examples most influential in a decision. The explanation can also be based on embedded or computed features, e.g., network dissection techniques to identify meaningful features inside the layers of a deep net. Intertwined are deep learning techniques to generate explanations.

Probabilistic Explanation. These methods leverage inference methods from probabilistic graphical models. They are a natural approach to use with neuro-symbolic methods.

In the following two subsections we describe efforts based on each of these two approaches.

4.1 Explanation from Representation

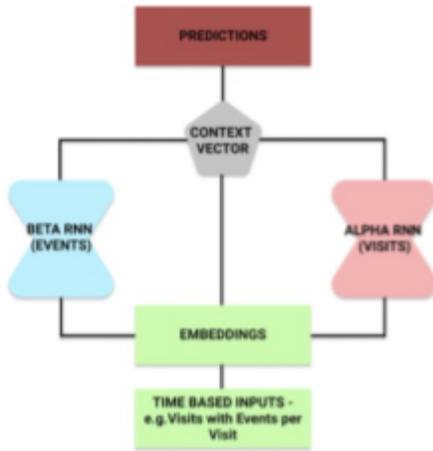
4.1.1 Reverse time Attention Model

The Reverse time Attention Model (RETAIN) mimics a physician in providing explanations (Choi et al, 2016). The goal is to help physicians understand the AI software's predictions. RETAIN uses an electronic health record (EHR) in reverse time order. It calculates the contribution of variables (medical codes) to diagnostic prediction using RNNs. See Fig. 18.

Patient hospital visit data is sent to two recurrent neural networks (RNNs) both of which have an attention mechanism. The concept of attention in RETAIN is analogous to attention in machine translation. Given a sentence of length S in the source, first generate h_1, \dots, h_S , to represent input words. To find the j^{th} target word, generate attention α_i for $i = 1, \dots, S$ for

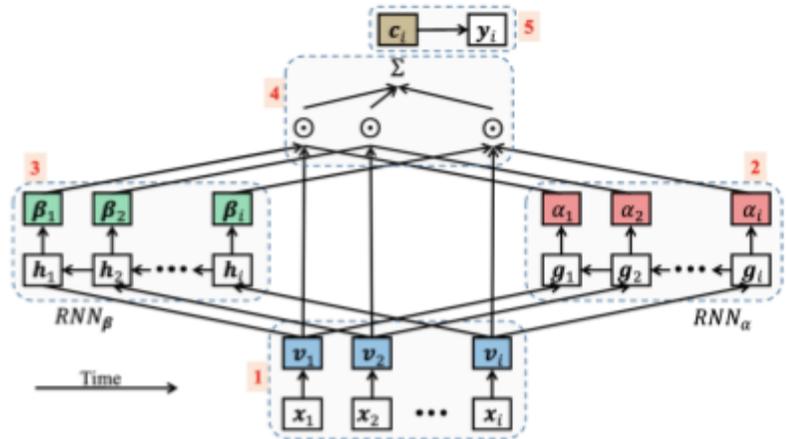
each word in the source sentence. Compute context $c_j = \sum_i \alpha_{ij} h_i$ and use it to predict the j^{th} target word i . Attention allows focus on specific words in the given sentence when generating each word in the target.

The attention mechanism in RETAIN helps explain which part the neural network was focusing on and which features helped influence its choice.



(a)

del: (a)



(b)

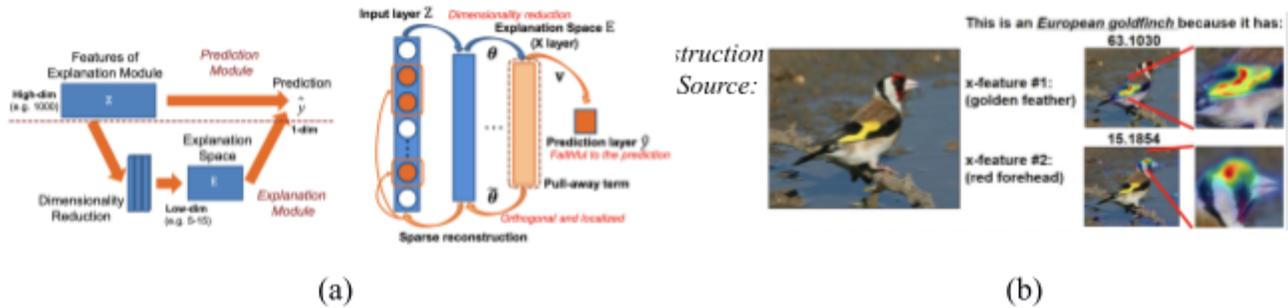
4.1.2 Explanations from Embeddings

The essence of deep learning is that of learning representations, or embeddings, that are useful to easily perform the final computation task, of, say classification or regression.

A proposed approach to harness what has been learnt by a deep network is to construct an explanation module by embedding a high-dimensional deep network layer nonlinearly into a low-dimensional explanation space. In this process a goal is to retain retain faithfulness, so that the original deep learning predictions can be constructed from the few concepts extracted by the explanation module.

The explanation module is a dimensionality reduction mechanism so that the original deep learning prediction \hat{y} can be reproduced from this low-dimensional space. It can be attached

to any layer in the prediction deep network. The network output can be faithfully recovered from this low-dimensional explanation space. A sparse Reconstruction Autoencoder is used as an explanation module (Qi, Khorram, and Li, 2018). See Fig. 19(a). An explanation generated by this model is shown in Fig. 19(b).



Caption-guided Image explanation

Deep image captioning systems learn to translate visual input into languages: potential map between visual concepts and words. Despite good captioning performance, they are hard to understand “black boxes.” A solution proposed is caption guided visual saliency: a top-down neural saliency map. See Fig. 20.

Figure 20: Caption-guided visual saliency. Source: (Ramanishka, Das, Zhang, and Saenko, 2017).

4.2 Probabilistic Explanations

4.2.1 Bayesian Deep Learning

Bayesian deep learning (BDL) enables one to gauge how uncertain a neural network is about its predictions. These deep architectures can model complex tasks by leveraging the hierarchical representation power of deep learning, while also being able to infer complex multi-modal posterior distributions. BDL models typically form uncertainty estimates by either placing distributions over model weights, or by learning a direct mapping to probabilistic outputs. By knowing the weight distributions of various predictions and classes, we can tell a lot about what feature led to what decisions and the relative importance of it.

4.2.2 Graphical Model Inference

Neuro-symbolic models aim to use both neural mechanisms to infer symbolic entities, but also aim to incorporate symbolic reasoning mechanisms to answer queries of interest. For instance, in an image of a horse, the neural mechanism infers that we have a horse with high probability. Before performing classification it identifies its features such as its legs, tail, mane, etc. A simple linear classifier performs the final classification. Suppose one of the horses legs is occluded. Even though the neural model can infer the presence of a horse, it would require a symbolic mechanism to infer the presence and location of the invisible fourth leg, and provide an appropriate explanation (Fig. 21).



Figure 21: Image of two horses with one having its legs cropped out. While a neural network would recognize both horses, based on observed features, a symbolic reasoning mechanism would infer the presence and location of the cropped-out legs. Image source: Benoit photo.

Symbolic reasoning and inference is done efficiently using graphical models such as Bayesian networks and Markov networks. In fact the inference of the most likely probability of a set of variables in a graphical model given evidence is referred to as maximum a posteriori (MAP) probability inference equivalently referred to as the most probable explanation (MPE) of the evidence (Koller and Friedman, 2009).

Probabilistic graphical models represent joint probability distributions over a set of variables χ . They are used to answer queries of interest, given evidence variables whose known value is e , i.e., $E = e$. Assume that a query to the system yields a value y to the target variables Y . We assume that the explanatory variables are the rest of the variables $W \subseteq \chi - Y - E$.

Assume $Y = y$ is the response to the query variable. The conditional probability of explanation is $P(W|E = e, Y = y)$. The task of MPE can be expressed as follows.

$$\arg \max_{w \in W} P(W|E = e, Y = y).$$

which can be computed using

$$p(w|E = e, Y = y) = \frac{P(w, e, y)}{\sum_e \sum_y P(w, e, y)}.$$

Since the denominator has an exponential number of terms in the number of variables, the computation is NP-hard. The situation can be alleviated by using methods from probabilistic graphical model domain. They include exact inference algorithms for MPE such as *variable elimination* and *clique trees*. Also, approximate algorithms based on either optimization or sampling (particle-based).

5 Concluding Remarks

Among the large number of research areas spawned by the growth of applications of artificial intelligence, one that has attracted general public interest is that of trustworthiness. Among key attributes of trustworthiness is explainability. An attribute that can be useful to not only the user, e.g., to understand whether there is bias, but also to the designer, e.g., to improve design.

Explanations can be made in terms of identifying salient features in the input. They can also be in terms of computed features, known as embeddings in deep learning systems. Explanations of AI systems can be made in terms of identifying archetypes in the training data.

If explanations are generated from existing AI systems, they are called post-hoc systems. AI systems configured from the start to generate explanations are known as ante-hoc systems.

Post-hoc explanations can be generated by perturbing input variables in variable heat maps so as to determine as to which variables have the most effect in changing the generated

decision. They can also be generated from deep learning models using attention mechanisms that identify the context for each decision.

Ante-hoc explanations have been obtained by mapping outputs of layers to lower-dimensional spaces. They can also be generated using neuro-symbolic methods, where a neural network generates values of high level variables and a symbolic reasoning system generates an explanation in those terms.

Probabilistic graphical models offer computationally feasible algorithms to generate probabilistic explanations in terms of given evidence and the observed target value.

6 References

- [1] E. Choi, M. T. Bahadori, J. Sun, J. Kulas, A. Schuetz, and W. Stewart. RETAIN: An interpretable predictive model for healthcare using reverse time attention mechanism. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 3504–3512. Curran, 2016.
- [2] I.P. de Sousa, M. Vellasco, J. Sun, and E.C. da Silva. Local interpretable model-agnostic explanations for classification of lymph node metastases. *Sensors*, 19(3), 2019. <https://www.oreilly.com/content/introduction-to-local-interpretable-model-agnostic-explanation-s-lime/>.
- [3] National Science Foundation. National Artificial Intelligence (AI) Research Institutes, 2019.
- [4] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016.
- [5] L.A. Hendricks, Z. Akata, M. Rohrbach, J. Donahue, B. Schiele, and T. Darrell. Generating visual explanations. In *Proceedings European Conference on Computer Vision*, 2016.
- [6] D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, 2009.
- [7] J. Mao, C. Gan, P. Kohli, J.B. Tenenbaum, and J. Wu. The neurosymbolic concept learner. In *Proceedings of ICLR*, 2019.
- [8] R. Marselis. Make your Artificial Intelligence more trust- worthy with eXplainable AI.
- [9] K. Pei, Y. Cao, J. Yang, and S. Jana. Deepxplore automated whitebox testing of deep learning systems. In *Proceedings 26th Symposium on Operating Systems Principles*, 2017.
- [10] Z. Qi, S. Khorram, and F. Li. Embedding Deep Networks into Visual Explanations, 2018.
- [11] V. Ramanishka, A. Das, J. Zhang, and K. Saenko. Top-down Visual Saliency Guided by Captions, 2017. <https://arxiv.org/abs/1612.07360>.
- [12] W. Samek, T. Wiegand, and K.R. Muller. Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. *ITU Journal: ICT Discoveries*, pages 39–48, 2018.

- [13] M. Turek. Explainable Artificial Intelligence (XAI), 2018. <https://www.darpa.mil/program/explainable-artificial-intelligence>.
- [14] W.K. Vong, R.B. Sojitra, A. Reyes, S. Yang, and P. Shafto. Bayesian Teaching of Image Categories. scottchenghsinyang.com/paper/Vong-2018.pdf.
- [15] S.Y. Yang and P. Shafto. Explainable Artificial Intelligence via Bayesian Teaching. http://shaftolab.com/assets/papers/yangShafto_NIPS_2017_machine_teaching.pdf.