Verification and Explanation of Deep Neural Networks

Nina Narodytska, VMware Research

Joint work with Alexey Ignatiev, Joao Marques-Silva, Aditya Shrotri, Kuldeep Meel









Motivation

source: wikipedia





AlphaGo Zero & Alpha Zero

Image & Speech Recognition



My photos

OK, but first you'll need to set up Photos to do that.

Open Photos

YOU CAN ALSO TRY

"Launch the Photos app"





A **Level 4** autonomous car is one defined as a car that can completely drive itself, from start to finish, within a specifically-designated area.







Does ML-controller comply with a specification?





Does ML-controller comply with a specification?

Why does ML controller drive a car into a wall?







Verification/robustness

Why does ML controller drive a car into a wall?





Verification/robustness

Explainability/interpretability





Robustness of ML models

Robustness of ML models

Original image



88% tabby cat

Original image + Perturbation =





88% tabby cat

Original image + Perturbation = Perturbed image







88% tabby cat

Original image + Perturbation = Perturbed image







88% tabby cat

99% guacamole

Robustness of ML: STOP sign attack



Eykholt et al'18





Aung et al'17

Robustness of ML: Autonomous car attack



Adversarial machine learning, Y. Vorobeychik, B. Li

Standardization: ISO

Standardization: ISO



ISO/IEC NP TR 24029-1

Artificial Intelligence (AI) — Assessment of the robustness of neural networks — Part 1: Overview



Statue - A Under development

Robustness of ML models









The Challenge of Crafting Intelligible Intelligence

By Daniel S. Weld, Gagan Bansal Communications of the ACM, June 2019, Vol. 62 No. 6, Pages 70-79 10.1145/3282486 Comments (1)

VIEW AS: 🚊 🛄 🏟 🔂 🔂 SHARE: 🖂 🥶 외 <table-cell> 🗈 💽



Artificial Intelligence (ai) systems have reached or exceeded human performance for many circumscribed tasks. As a result, they are increasingly deployed in mission-critical roles, such as credit scoring, predicting if a bail candidate will commit another crime, selecting the news we read on social networks, and selfdriving cars. Unlike other mission-critical software, extraordinarily complex AI systems are difficult to test: AI decisions are context specific and often based on thousands or millions of factors. Typically, AI behaviors are generated by searching vast action spaces or learned by the opaque optimization of mammoth neural networks operating over prodigious amounts of training data. Almost by definition, no clear-cut method can accomplish these AI tasks.

Explainable Artificial Intelligence (XAI)

Standardization: GDPR

Standardization: GDPR

European Union regulations on algorithmic decision-making and a "right to explanation"

Bryce Goodman,1* Seth Flaxman,2

We summarize the potential impact that the European Union's new General Data Protection Regulation will have on the routine use of machine-learning algorithms. Slated to take effect as law across the European Union in 2018, it will place restrictions on automated individual decision making (that is, algorithms that make decisions based on user-level predictors) that "significantly affect" users. When put into practice, the law may also effectively create a right to explanation, whereby a user can ask for an explanation of an algorithmic decision that significantly affects them. We argue that while this law may pose large challenges for industry, it highlights opportunities for computer scientists to take the lead in designing algorithms and evaluation frameworks that avoid discrimination and enable explanation.

Robustness of ML models

Robustness of ML models



Robustness of ML models



Robustness of ML models







Robustness of ML models






Verification and explanation of ML models are important topics!

Robustness of ML models





I.Goodfellow at AAAI 2019, "Adversarial Machine Learning"

Interpretation of ML models







Neural Networks

Neural Networks



- features
- images

Linear transformation



- features
- images

Non-Linear transformation



max(0, Wz + b)

- features
- images

Non-Linear transformation



• images

Block structure



• images

Repeated block structure



- features
- images

Repeated block structure



- features
- images

 $\frac{e^{a_i}}{\sum_{j=1}^k e^{a_j}}, i \in \{1, \dots, k\}$







Interpretable Explanations (Discussion)







How would you explain this decision?



How would you explain this decision?

From: salem@pangea.Stanford.EDU (Bruce Salem) Subject: Re: Science and theories

How is it possible for us to believe in God (or god, I gues when science has shown his existence to be impossible? I think atheism is the way to go forward.



Atheism

From: salem@pangea.Stanford.EDU (Bruce Salem) Subject: Re: Science and theories

How is it possible for us to believe in God (or god, I gues when science has shown his existence to be impossible? I think atheism is the way to go forward.



How would you explain this decision?



How would you explain this decision?







How would you explain this decision?



How would you explain this decision?



State-of-the-art methods













State-of-the-art

Heuristic approaches (e.g. LIME)

Why Should I Trust You?" Explaining the Predictions of Any Classifier, Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin, KDD'16
























































LIME: build a local neighborhood





LIME: build a local neighborhood



Heuristic approaches

1. local explanations

Heuristic approaches

- 1. local explanations
- 2. no guarantees about quality

Heuristic approaches

- 1. local explanations
- 2. no guarantees about quality
- 3. robustness of explanations



Logic-based approach to explanations

Abduction-Based Explanations for Machine Learning Models, AAAI'19 Alexey Ignatiev, Nina Narodytska and Joao Marques-Silva

Andy Shih and Arthur Choi and Adnan Darwiche A Symbolic Approach to Explaining Bayesian Network Classifiers, IJCAI'18 Compiling Bayesian Network Classifiers into Decision Graphs, AAAI'19

The main idea comes from work on model diagnosis to explain failures of the systems from a given set of hypotheses

R. Reiter. A theory of diagnosis from first principles. Artif. Intell., 32

Age is (37, 48], White, Male, Married







Given a classifier \mathbf{M} , represented by some logic encoding, a cube Z and a prediction p, compute a subset-minimal $\mathcal{Z} \subseteq Z$ s.t.

$1.(\mathcal{Z} \land \mathbf{M} \not\models \bot)$ $2.(\mathcal{Z} \land \mathbf{M} \models p)$

Given a classifier \mathbf{M} , represented by some logic encoding, a cube Z and a prediction p, compute a subset-minimal $\mathcal{Z} \subseteq Z$ s.t.

$1.(\mathcal{Z} \land \mathbf{M} \not\models \bot)$ $2.(\mathcal{Z} \land \mathbf{M} \models p)$



















$$\begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} -1 & 1 \\ -1 & -1 \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$
$$a_1 = -z_1 + z_2$$

$$a_2 = -z_1 - z_2$$

NN

model \mathbf{M}





$$\begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} -1 & 1 \\ -1 & -1 \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

$$a_1 = -z_1 + z_2 \qquad r_1 = max(0, a_1)$$

$$a_2 = -z_1 - z_2 \qquad r_2 = max(0, a_2)$$





NN







Given a classifier \mathbf{M} , represented by some logic encoding, a cube Z and a prediction p, compute a subset-minimal $\mathcal{Z} \subseteq Z$ s.t.

$1.(\mathcal{Z} \land \mathbf{M} \not\models \bot)$ $2.(\mathcal{Z} \land \mathbf{M} \models p)$

 $1.(\mathcal{Z} \land \mathbf{M} \not\models \bot)$ $2.(\mathcal{Z} \land \mathbf{M} \models p)$





 $\mathcal{Z} \land \mathbf{M} \not\models \bot$



$\mathcal{Z} \land \mathbf{M} \not\models \bot$

By definition: $Z \wedge \mathbf{M} \models p$



 $Z \wedge M \not\models_{autology}$

By definition:
$$Z \wedge \mathbf{M} \models p$$

 $1.(\mathcal{Z} \land \mathbf{M} \not\models \bot)$ $2.(\mathcal{Z} \land \mathbf{M} \models p)$









$\mathcal{Z} \land \mathbf{M} \models p$



 $\mathcal{Z} \land \mathbf{M} \models p$ $\Leftrightarrow (\mathcal{Z} \models (\mathbf{M} \rightarrow p))$



$\mathcal{Z} \models (\mathbf{M} \to p)$



$\mathcal{Z} \models (\mathbf{M} \to p)$
Propositional abduction problem



An **explanation** is a subset of input features so that changes to the rest of inputs do not affect the prediction.

Subset-minimal explanation

- **Input:** M, initial cube Z, prediction p
- 1 begin
- 2 foreach $l \in Z$ do
- 3 4 if $Z \setminus \{l\} \vDash \mathbf{M} \to p$ then $\ \ \left\lfloor Z \leftarrow Z \setminus \{l\} \right\rfloor$
- 5 return Z
- 6 end

Algorithm 1: Computing a subset-minimal explanation

Subset-minimal explanation

- **Input:** M, initial cube Z, prediction p
- 1 begin
- 2 foreach $l \in Z$ do 3 $| if Z \setminus \{l\} \vDash M \rightarrow p$ then 4 $| Z \leftarrow Z \setminus \{l\}$
- 5 return Z
- 6 end

Algorithm 1: Computing a subset-minimal explanation

- 1. local explanations
- 2. no guarantees about quality
- 3. robustness of explanations

- 1. global explanations
- 2. no guarantees about quality
- 3. robustness of explanations

- 1. global explanations
- 2. provide guarantees about quality
- 3. robustness of explanations

On Validating, Repairing and Refining Heuristic ML Explanations, CoRR report'19, Alexey Ignatiev, Nina Narodytska, and Joao Marques-Silva

- 1. global explanations
- 2. provide guarantees about quality
- 3. evaluate robustness of explanations

Assessing Heuristic Machine Learning Explanations with Model Counting, SAT'19 Nina Narodytska, Aditya A. Shrotri, Kuldeep S. Meel, Alexey Ignatiev, João Marques-Silva (Used approximate model counting to evaluate quality of the explanations)



Figure 1: Possible minimal explanations for digit one.





Figure 1: Possible minimal explanations for digit one.





Figure 1: Possible minimal explanations for digit one.





Figure 1: Possible minimal explanations for digit one.



Explanations are not unique



Figure 1: Possible minimal explanations for digit one.



Some look more sensible than others!





Verification of NNs

 Certification of Neural Networks (train a network that satisfies a property)

 Verification of Neural Networks (complete, incomplete verification)

Verification of NNs

 Certification of Neural Networks (train a network that satisfies a property)

 Verification of Neural Networks (complete, incomplete verification)

> Algorithms for Verifying Deep Neural Networks Changliu Liu, Tomer Arnon, Christopher Lazarus, Clark Barrett, Mykel J. Kochenderfer

Given a classifier \mathbf{M} , a *counterexample* to a prediction p is a subset-minimal set of feature literals \mathcal{T} , such that

$$\mathcal{T} \models \lor_{t,t \neq p} (\mathbf{M} \to t)$$





 $\mathcal{T} \models \lor_{t,t \neq p} (\mathbf{M} \to t)$





\mathcal{T} breaks \mathcal{Z}





\mathcal{T} breaks \mathcal{Z}







\mathcal{T} breaks \mathcal{Z}





Given a model \mathbf{M} , represented by some logic encoding, and a prediction p,

- every explanation \mathcal{Z} of p breaks every counterexample of p, and
- every counterexample \mathcal{T} of p breaks every explanation of p.

On Relating Explanations and Adversarial Examples, NeurIPS'2019 Alexey Ignatiev, Nina Narodytska, Joao Marques-Silva

Input: formula \mathbf{M} and prediction p**Output:** set \mathbb{E} of all absolute explanations of prediction p1 $(\mathbb{C}, \mathbb{E}, \mathcal{Z}) \leftarrow (\emptyset, \emptyset, \emptyset)$ 2 do: if $\mathcal{Z} \vDash (\mathbf{M}
ightarrow p)$: 3 4 $\mathbb{E} \leftarrow \mathbb{E} \cup \{\mathcal{Z}\}$ # \mathcal{Z} is an explanation; save it 5 else: $(\mathcal{T}, t) \leftarrow \texttt{ExtractInstance}()$ # get an instance \mathcal{T} with a 6 prediction t, $t \neq p$ for $l\in\mathcal{T}$: 7 if $(\mathcal{T} \setminus \{l\}) \vDash (\mathbf{M} \to t)$: 8 $\mathcal{T} \leftarrow \mathcal{T} \setminus \{l\}$ 9 $\mathbb{C} \leftarrow \mathbb{C} \cup \{\mathcal{T}\}$ # update \mathbb{T} with a new counterexample \mathcal{T} 10 $\mathcal{E} \leftarrow \texttt{MinimumHS}(\mathbb{C})$ # get a new hitting set of $\mathbb C$ 11 12 while $\mathcal{Z} \neq \emptyset$ 13 return \mathbb{E}

Algorithm 1: Duality-based computation of all absolute explanations

Input: formula M and prediction p	
Output: set \mathbb{E} of all absolute explanations of prediction p	
1 $(\mathbb{C},\mathbb{E},\mathcal{Z}) \leftarrow (\emptyset,\emptyset,\emptyset)$	
2 do:	
$\mathbf{if} \ \mathcal{Z} \vDash (\mathbf{M} ightarrow p)$:	
4 $\mathbb{E} \leftarrow \mathbb{E} \cup \{\overline{\mathcal{Z}}\}$ # \mathcal{Z} is an explanation; save it	
5 else:	
6 $(\mathcal{T}, t) \leftarrow \texttt{ExtractInstance}()$ # get an instance \mathcal{T} with a	
prediction t , $t \neq p$	
7 for $l \in \mathcal{T}$:	
8 if $(\mathcal{T} \setminus \{l\}) \vDash (\mathbf{M} \to t)$:	
9 $\mathcal{T} \leftarrow \mathcal{T} \setminus \{l\}$	
10 $\mathbb{C} \leftarrow \mathbb{C} \cup \{\mathcal{T}\}$ # update \mathbb{T} with a new counterexample \mathcal{T}	
11 $\mathcal{E} \leftarrow \texttt{MinimumHS}(\mathbb{C})$ # get a new hitting set of \mathbb{C}	
12 while $\mathcal{Z} \neq \emptyset$	
13 return $\mathbb E$	
Algorithm 1: Duality-based computation of all absolute explanations	

```
Input: formula \mathbf{M} and prediction p
    Output: set \mathbb{E} of all absolute explanations of prediction p
 1 (\mathbb{C}, \mathbb{E}, \mathcal{Z}) \leftarrow (\emptyset, \emptyset, \emptyset)
 2 do:
     if \mathcal{Z} \vDash (\mathbf{M} 
ightarrow p) :
 3
     \mathbb{E} \leftarrow \mathbb{E} \cup \{\mathcal{Z}\}
                                                           # \mathcal{Z} is an explanation; save it
 4
       else:
 5
               (\mathcal{T}, t) \leftarrow \texttt{ExtractInstance}() # get an instance \mathcal{T} with a
 6
                prediction t, t \neq p
              for l \in \mathcal{T}:
 7
                    if (\mathcal{T} \setminus \{l\}) \vDash (\mathbf{M} \to t):
 8
                         \mathcal{T} \leftarrow \mathcal{T} \setminus \{l\}
 9
              \mathbb{C} \leftarrow \mathbb{C} \cup \{\mathcal{T}\} # update \mathbb{T} with a new counterexample \mathcal{T}
10
         \mathcal{E} \leftarrow \texttt{MinimumHS}(\mathbb{C})
                                                               # get a new hitting set of \mathbb C
11
12 while \mathcal{Z} \neq \emptyset
13 return \mathbb{E}
   Algorithm 1: Duality-based computation of all absolute explanations
```



(a) digit "6"



(a) digit "6"

(b) patch area



(a) digit "6"

(b) patch area

(c) an XP



(a) digit "6"

(b) patch area

(c) an XP

(d) all XP's



(a) digit "6"

(b) patch area

(c) an XP

(d) all XP's

(e) a CE

140







Scalability



Heuristic methods for NN explainability do not suffer from scalability issues


Heuristic methods for NN explainability do not suffer from scalability issues

no guarantees about the quality of explanations

Logic-based methods on explanability of NNs do suffer from scalability issues

Input: formula \mathbf{M} and prediction p**Output:** set \mathbb{E} of all absolute explanations of prediction p1 $(\mathbb{C}, \mathbb{E}, \mathcal{Z}) \leftarrow (\emptyset, \emptyset, \emptyset)$ 2 do: if $\mathcal{Z} \models (\mathbf{M} \rightarrow p)$: 3 $\mathbb{E} \leftarrow \mathbb{E} \cup \{\mathcal{Z}\}$ # \mathcal{Z} is an explanation; save it $\mathbf{4}$ else: 5 $(\mathcal{T}, t) \leftarrow \texttt{ExtractInstance}()$ # get an instance \mathcal{T} with a 6 prediction t, $t \neq p$ for $l \in \mathcal{T}$: 7 if $(\mathcal{T} \setminus \{l\}) \vDash (\mathbf{M} \to t)$: 8 $\mathcal{T} \leftarrow \mathcal{T} \setminus \{l\}$ 9 $\mathbb{C} \leftarrow \mathbb{C} \cup \{\mathcal{T}\}$ # update \mathbb{T} with a new counterexample \mathcal{T} 10 $\mathcal{E} \leftarrow \texttt{MinimumHS}(\mathbb{C})$ # get a new hitting set of $\mathbb C$ 11 12 while $\mathcal{Z} \neq \emptyset$ 13 return \mathbb{E}

Algorithm 1: Duality-based computation of all absolute explanations



It is an issue!

Many papers on verification of NNs is battling scalability issue, e.g.

- add constraints to the training procedure
- use approximate reasoning
- simplify networks during the training

Many papers on verification of NNs is battling scalability issue, e.g.

- add constraints to the training procedure
- use approximate reasoning
- simplify networks during the training

In Search for a SAT-friendly Binarized Neural Network Architecture

Hongce Zhang (summer @ VMware), Aarti Gupta, Toby Walsh

Training for Faster Adversarial Robustness Verification via Inducing ReLU Stability Kai Y. Xiao, Vincent Tjeng, Nur Muhammad (Mahi) Shafiullah, Aleksander Madry













Can we train a NN so that it is easier to analyze?

Conclusion



We showed a tight connection between explanations and counterexamples



There is hope to battle scalability issues!

Conclusion

Thanks!