Hybrid techniques based on black-box classifiers and interpretable control modules

Luis Galárraga Seminaire DKM Lannion, 12/07/2018

- Interpretability in classifiers: What and Why?
- Black-box vs. interpretable classifiers
- Explaining the black-box
- Conclusion & open research questions

- Interpretability in classifiers: What and Why?
- Black-box vs. interpretable classifiers
- Explaining the black-box
- Conclusion & open research questions

- Interpretability in **classifiers**: What and Why?
- Black-box vs. interpretable classifiers
- Explaining the black-box
- Conclusion & open research questions

Machine Learning Classifiers

- Models that assign classes to instances
 - The model is learned and trained from labeled data
 - Labels are predefined: supervised learning



(Hopefully) a lot of labeled data

Machine Learning Classifiers

- Models that assign classes to instances
 - The model is learned and trained from labeled data
 - Labels are predefined: supervised learning



- Interpretability in classifiers: What and Why?
- Black-box vs. interpretable classifiers
- Explaining the black-box
- Conclusion & open research questions

Some ML classifiers can be really complex



Some ML classifiers can be really complex



A classifier is *interpretable* if the rationale behind its answers can be easily *explained*



A classifier is *interpretable* if the rationale behind its answers can be easily *explained*



interpretability \cong explainability \cong comprenhensibility



• Classifiers are used to make critical decisions



How a Self-Driving Uber Killed a Pedestrian in Arizona

By TROY GRIGGS and DAISUKE WAKABAYASHI UPDATED MARCH 21, 2018

A woman was <u>struck and killed</u> on Sunday night by an autonomous car operated by Uber in Tempe, Ariz. It was believed to be the first pedestrian death associated with selfdriving technology.

What We Know About the Accident



- Classifiers are used to make critical decisions
- Need to know the rationale behind an answer
 - For debugging purposes
 - To tune the classifier
 - To spot biases in the data
 - For legal and ethical reasons
 - General Data Protection Regulation (GDPR)
 - To understand the source of the classifier's decision bias
 - To generate trust



Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublicc May 23, 2016

O N A SPRING AFTERNOON IN 2014, Brisha Borden was running late to pick up her god-sister from school when she spotted an unlocked kid's blue Huffy bicycle and a silver Razor scooter. Borden and a friend grabbed the bike and scooter and tried to ride them down the street in the Fort Lauderdale suburb of Coral Springs.

Just as the 18-year-old girls were realizing they were too big for the tiny conveyances which belonged to a 6-year-old boy — a woman came running after them saying, "That's my kid's stuff." Borden and her friend immediately dropped the bike and scooter and walked away.

But it was too late — a neighbor who witnessed the heist had already called the police. Borden and her friend were arrested and charged with burglary and petty theft for the items, which were valued at a total of \$80.

Amazon just showed us that 'unbiased' algorithms can be inadvertently racist

Rafi Letzter Apr. 21, 2016, 4:50 PM

A Bloomberg report Thursday revealed that Amazon's sameday delivery service offered to Prime users around major US cities seems to routinely, if unintentionally, exclude black neighborhoods.

The maps, which you should check out on Bloomberg's site, show that in cities like Chicago, New York, and Atlanta, sameday delivery covers just about every zip code at this point except the majority black ones.



Chicago was one of the cities highlighted in Bloomberg's report. Klichiro Sato/AP



BUSINESS INSIDER INTELLIGENCE EXCLUSIVE ON ARTIFICIAL INTELLIGENCE

DISCOVER THE FUTURE OF FINTECH WITH THIS EXCLUSIVE SLIDE DECK

3

Tay: Microsoft issues apology over racist chatbot fiasco



Microsoft has apologised for creating an artificially intelligent chatbot that quickly turned into a holocaust-denying racist.

But in doing so made it clear Tay's views were a result of nurture, not nature. Tay confirmed what we already knew: people on the internet can be cruel.

Tay, aimed at 18-24-year-olds on social media, was targeted by a "coordinated attack by a subset of people" after being launched earlier this week.

Within 24 hours Tay had been deactivated so the team could make "adjustments".



@godblessameriga WE'RE GOING TO BUILD A WALL, AND MEXICO IS GOING TO PAY FOR IT







Following

(a) Husky classified as wolf

(b) Explanation

https://www.bbc.com/news/technology-35902104

- Interpretability in classifiers: What and Why?
- Black-box vs. interpretable classifiers
- Explaining the black-box
- Conclusion & open research questions

Black-box



- Neural Networks
 (DNN, RNN, CNN)
- Ensemble methods
 - Random Forests
- Support Vector
 Machines

• Interpretable



- Decision Trees
- Classification Rules
 - If-then rules
 - *m*-of-*n* rules
 - Lists of rules
 - Falling rule lists
 - Decision sets
- Prototype-based methods

Decision Trees



- If-then Rules
 - Past-Depression Λ Melancholy \Rightarrow Depressed
- m-of-n rules
 - Predict a class if at least *m* out of *n* relevant attributes are present
 - If 2-of-{Past-Depression, ¬Melancholy, ¬Insomnia} ⇒ Healthy
- Lists of rules
 - Select the top-k rules for each class, and predict the class with the rule set of highest expected accuracy

• Falling rule lists

Falling Rule Lists

	Conditions		Probability	Support
IF	IrregularShape AND Age ≥ 60	THEN malignancy risk is	85.22%	230
ELSE IF	Spiculated Margin AND Age ≥ 45	THEN malignancy risk is	78.13%	64
ELSE IF	IllDefinedMargin AND Age ≥ 60	THEN malignancy risk is	69.23%	39
ELSE IF	IrregularShape	THEN malignancy risk is	63.40%	153
ELSE IF	LobularShape AND Density ≥ 2	THEN malignancy risk is	39.68%	63
ELSE IF	RoundShape AND Age ≥ 60	THEN malignancy risk is	26.09%	46
ELSE		THEN malignancy risk is	10.38%	366

Decision sets

If Respiratory-Illness=Yes and Smoker=Yes and Age≥ 50 then Lung Cancer

If Risk-LungCancer=Yes and Blood-Pressure≥ 0.3 then Lung Cancer

If Risk-Depression=Yes and Past-Depression=Yes then Depression

If BMI ≥ 0.3 and Insurance=None and Blood-Pressure ≥ 0.2 then Depression

If Smoker=Yes and BMI ≥ 0.2 and Age ≥ 60 then Diabetes

If Risk-Diabetes=Yes and BMI > 0.4 and Prob-Infections > 0.2 then Diabetes

If Doctor-Visits ≥ 0.4 and Childhood-Obesity=Yes then Diabetes

- Prototype-based methods
 - Predict a class and provide a prototypical instance labeled with the same class
 - Challenge: pick a set of prototypes per class such that
 - The set is of minimal size
 - It provides full coverage, i.e., every instance should have a close prototype
 - They are far from instances of other classes
 - (a) formulates prototype selection as an optimization problem and uses it to classify images of handwritten digits

(a) J. Bien and R. Tibshirani. Prototype selection for interpretable classication. The Annals of Applied Statistics, pages 2403-2424, 2011.

• Prototype-based methods



FIG. 7. The first 88 prototypes (out of 3,372) of the greedy solution. We perform MDS (R function sammon) on the tangent distances to visualize the prototypes in two dimensions. The size of each prototype is proportional to the log of the number of correct-class training images covered by this prototype.

17



- Random Forests
 - Bagging: select *n* random samples (with replacement) and fit n decision trees.
 - Prediction: aggregate the decisions of the different trees to make a prediction



Support Vector Machines



- Interpretability in classifiers: What and Why?
- Black-box vs. interpretable classifiers
- Explaining the black-box
- Conclusion & open research questions

Design an interpretation layer between the classifier and the human user



Design an interpretation layer between the classifier and the human user



Design an interpretation layer between the classifier and the human user

Evaluation

accuracy(classifier,truth)= $\frac{\# \text{ examples such that } classifier = truth}{\# \text{ all examples}}$ fidelity = accuracy(interpretable classifier, classifier)

complexity = f(interpretable classifier)



- Methods can be classified into three categories:
 - Methods for global explainability
 - Methods for local (outcome) explainability
 - Methods for classifier inspection
- Methods can be black-box dependent or black-box agnostic

The interpretable approximation is a classifier that provides explanations for all possible outcomes



- Global explanations for NNs date back to the 90s
 - Trepan⁽¹⁾ is a black-box agnostic method that induces decision trees by querying the black box



(1) M. Craven and J. W. Shavlik. Extracting tree-structured representations of trained networks. In Advances in neural information processing systems, pages 24-30, 1996.

- Global explanations for NNs date back to the 90s
 - Trepan⁽¹⁾ is a black-box agnostic method that induces decision trees by querying the black box
 - Trepan's split criterion depends on entropy and fidelity



(1) M. Craven and J. W. Shavlik. Extracting tree-structured representations of trained networks. In Advances in neural information processing systems, pages 24-30, 1996.

- (2) uses genetic programming to enhance explanations
 - By modifying the tree
 - By mixing the original and the BB labeled data



(2) U. Johansson and L. Niklasson. Evolving decision trees using oracle guides. In Computational Intelligence and Data Mining, 2009. CIDM'09., pages 238-244. IEEE, 2009.

- ⁽³⁾ uses notions from ensemble methods to improve accuracy
 - Dataset is augmented via bagging



(3) P. Domingos. Knowledge discovery via multiple models. Intelligent Data Analysis, 2(1-4):187-202, 1998.

- Other methods generate sets of rules as explanations
 - (4) learns m-of-n rules from the original data plus BB labeled data (NNs)

(4) M. Craven and J. W. Shavlik. Using sampling and queries to extract rules from trained neural networks. In ICML, pages 37-45, 1994.

39

- Other methods generate sets of rules as explanations
 - (4) learns m-of-n rules from the original data plus BB labeled data (NNs)
 - BETA⁽⁵⁾ applies reverse engineering on the BB plus itemset mining to extract if-then rules.
 - Rules are restricted to two levels
 - If two contradictory rules apply to an example, the one with higher fidelity wins

```
If Age > 50 and Gender = Male Then

If Past-Depression = Yes and Insomnia = No and Melancholy = No ⇒ Healthy

If Past-Depression = Yes and Insomnia = No and Melancholy = Yes ⇒ Depressed
```

(4) M. Craven and J. W. Shavlik. Using sampling and queries to extract rules from trained neural networks. In ICML, pages 37-45, 1994.
 (5) H. Lakkerain and E. Kamar and P. Carnana and Jure Laskeras. Interpretable & Evalurable Approximations of

- BETA⁽⁵⁾ applies reverse engineering on the BB plus itemset mining to extract if-then rules.
 - Conditions (Gender = Male) obtained via pattern mining
 - Rule selection is formulated as an optimization problem

$$\underset{\mathcal{R}\subseteq\mathcal{ND}\times\mathcal{DL}\times\mathcal{C}}{\arg\max}\sum_{i=1}^{5}\lambda_{i}f_{i}(\mathcal{R})$$
(1)

s.t.
$$size(\mathcal{R}) \le \epsilon_1, maxwidth(\mathcal{R}) \le \epsilon_2, numdsets(\mathcal{R}) \le \epsilon_3$$
 (2)

 $f_{1}(\mathcal{R}) = \mathcal{P}_{max} - numpreds(\mathcal{R}), \text{ where } \mathcal{P}_{max} = \mathcal{P}_{max} = 2 * \mathcal{W}_{max} * |\mathcal{ND}| * |\mathcal{DL}|$ $f_{2}(\mathcal{R}) = \mathcal{O}_{max} - featureoverlap(\mathcal{R}), \text{ where } \mathcal{O}_{max} = \mathcal{W}_{max} * |\mathcal{ND}| * |\mathcal{DL}|$ $f_{3}(\mathcal{R}) = \mathcal{O}'_{max} - ruleoverlap(\mathcal{R}), \text{ where } \mathcal{O}'_{max} = N \times (|\mathcal{ND}| * |\mathcal{DL}|)^{2}$ $f_{4}(\mathcal{R}) = cover(\mathcal{R})$ $f_{5}(\mathcal{R}) = \mathcal{F}_{max} - disagreement(\mathcal{R}), \text{ where } \mathcal{F}_{max} = N \times |\mathcal{ND}| * |\mathcal{DL}|$

(4) M. Craven and J. W. Shavlik. Using sampling and queries to extract rules from trained neural networks. In ICML, pages 37-45, 1994.

- BETA⁽⁵⁾ applies reverse engineering on the BB plus itemset mining to extract if-then rules.
 - Conditions (Gender = Male) obtained via pattern mining
 - Rule selection is formulated as an optimization problem

$$\underset{\mathcal{R}\subseteq\mathcal{ND}\times\mathcal{DL}\times\mathcal{C}}{\arg\max}\sum_{i=1}^{5}\lambda_{i}f_{i}(\mathcal{R})$$
(1)

Sorcery!!!

s.t.
$$size(\mathcal{R}) \le \epsilon_1, maxwidth(\mathcal{R}) \le \epsilon_2, numdsets(\mathcal{R}) \le \epsilon_3$$
 (2)

$$\begin{split} f_{1}(\mathcal{R}) &= \mathcal{P}_{max} - numpreds(\mathcal{R}), \text{ where } \mathcal{P}_{max} = \mathcal{P}_{max} = 2 * \mathcal{W}_{max} * |\mathcal{ND}| * |\mathcal{DL}| \\ f_{2}(\mathcal{R}) &= \mathcal{O}_{max} - featureoverlap(\mathcal{R}), \text{ where } \mathcal{O}_{max} = \mathcal{W}_{max} * |\mathcal{ND}| * |\mathcal{DL}| \\ f_{3}(\mathcal{R}) &= \mathcal{O}'_{max} - ruleoverlap(\mathcal{R}), \text{ where } \mathcal{O}'_{max} = N \times (|\mathcal{ND}| * |\mathcal{DL}|)^{2} \\ f_{4}(\mathcal{R}) &= cover(\mathcal{R}) \\ f_{5}(\mathcal{R}) &= \mathcal{F}_{max} - disagreement(\mathcal{R}), \text{ where } \mathcal{F}_{max} = N \times |\mathcal{ND}| * |\mathcal{DL}| \end{split}$$

(4) M. Craven and J. W. Shavlik. Using sampling and queries to extract rules from trained neural networks. In ICML, pages 37-45, 1994.

- BETA⁽⁵⁾ applies reverse engineering on the BB plus itemset mining to extract if-then rules.
 - Conditions (Gender = Male) obtained via pattern mining
 - Rule selection is formulated as an optimization problem
 - Maximize fidelity and coverage
 - Minimize rule, feature overlap, and complexity
 - Constrained by number of rules, maximum width, and number of first level conditions

(4) M. Craven and J. W. Shavlik. Using sampling and queries to extract rules from trained neural networks. In ICML, pages 37-45, 1994.

- RxREN⁽⁶⁾ learns rule-based explanations for NNs
 - First, it iteratively prunes the *insignificant* input neurons while the accuracy loss is less than 1%



(6) M. G. Augasta and T. Kathirvalavakumar. Reverse engineering the neural networks for rule extraction in classification problems. Neural processing letters, 35(2):131-150, 2012.

- RxREN⁽⁶⁾ learns rule-based explanations for NNs
 - Second, build a matrix with the [min, max] of the values of the significant neurons for # misses > θ

	Cat	Dog	Fish	Koala
N_2	[3, 5]	0	[1, 4]	0
N ₃	[6, 7]	[8, 9]	0	[3, 6]

(6) M. G. Augasta and T. Kathirvalavakumar. Reverse engineering the neural networks for rule extraction in classification problems. Neural processing letters, 35(2):131-150, 2012.

- RxREN⁽⁶⁾ learns rule-based explanations for NNs
 - Third, learn rules from the matrix
 - Sort classes by number of non-zero entries

Cat Dog Fish Koala N_2 [3, 5] 0 [1, 4] 0 N_3 [6, 7] [8, 9] 0 [3, 6] If $N_2 \in [6, 7] \land N_3 \in [3, 5] \Rightarrow Cat$ Else If $N_2 \in [8, 9] \Rightarrow Dog$ Else

(6) M. G. Augasta and T. Kathirvalavakumar. Reverse engineering the neural networks for rule extraction in classification problems. Neural processing letters, 35(2):131-150, 2012.

The interpretable approximation is a classifier that provides explanations for the answers of the black box in the vicinity of an individual instance.



- LIME⁽⁷⁾ is a black-box agnostic method that optimizes for local fidelity
 - First, write examples in an *interpretable* way



Original Image



Interpretable Components

(7) 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, pages 1135–1144. ACM, 2016.

- LIME⁽⁷⁾ is a black-box agnostic method that optimizes for local fidelity
 - First, write examples in an *interpretable* way

Délai de livraison parfait très bon état du livre en ce qui concerne le bouquin en lui-même c'est extraordinaire . un point de vue sur l'histoire de l'humanité qui fait voir les choses sous un nouvel angle



(7) 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, pages 1135–1144. ACM, 2016.

- LIME⁽⁷⁾ is a black-box agnostic method that optimizes for local fidelity
 - Then, learn an interpretable model (e.g., linear model) in the vicinity of the given instance.



(7) 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, pages 1135–1144. ACM, 2016.

- LIME⁽⁷⁾ is a black-box agnostic method that optimizes for local fidelity
 - Then, learn an interpretable model (e.g., linear model) in the vicinity of the given instance.



(7) 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, pages 1135–1144. ACM, 2016.

⁽⁸⁾ generates class activation maps (CAM) from NNs used for image classification



(8) B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning deep features for discriminative localization. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 2921–2929, 2016.

- SHAP⁽⁹⁾ is black-box agnostic and provides local explainability via additive feature attribution
 - Shap values quantify feature importance
 - They correspond to the average distance between the model's answer and **all possible reduced models** that omit a feature

 SHAP⁽⁹⁾ is black-box agnostic and provides local explainability via additive feature attribution



(9) Lundberg, Scott M. and Su-In Lee. A Unified Approach to Interpreting Model Predictions. NIPS 2017.

- SHAP⁽⁹⁾ is black-box agnostic and provides local explainability via additive feature attribution
 - It also offers some model-specific extensions such as DeepShap and TreeShap
 - It is written in Python and available at https://github.com/slundberg/shap

The goal is to *plot* the *correlations* between the input features and the output classes



• Sensitivity analysis⁽¹⁰⁾ explains the influence of the inputs on the classifier's output

- Sensitivity analysis⁽¹⁰⁾ explains the influence of the inputs on the classifier's output
 - Build a prototype vector with the average/median/mode of the input attributes
 - Vary each attribute value, apply the classifier



- Sensitivity analysis⁽¹⁰⁾ explains the influence of the inputs on the classifier's output
 - Express each output class as a binary variable
 - Compute metrics for each attribute: range, gradient, variance, importance



• Sensitivity analysis⁽¹⁰⁾ explains the influence of the inputs on the classifier's output

$$- \text{ range}_{\text{class}}$$
 (Age) $= 0$

- gradient_{class=} (PClass) = (1 + 0)/2 = 0.5



• Sensitivity analysis⁽¹⁰⁾ explains the influence of the inputs on the classifier's output

$$- \text{ range}_{\text{class}}$$
 (Age) $= 0$

- gradient_{class=} (PClass) = (1 + 0)/2 = 0.5



- Sensitivity analysis⁽¹⁰⁾ explains the influence of the inputs on the classifier's output
 - Importance of an input feature *a* according to metric *s*



- Interpretability in classifiers: What and Why?
- Black-box vs. interpretable classifiers
- Explaining the black-box
- Conclusion & open research questions

Conclusion

- Interpretability in ML classifier matters
 - For human, ethical and legal reasons
 - For technical reasons
- Interpretability has two dimensions: global & local
- Global interpretability for black boxes has been by far more studied
- The key of opening the black box is *reverse* engineering

Open research questions

- Is it possible to define interpretability and explanations in a rigorous mathematical way?
- Explanations are generated solely on the input features
 - The input features are always known
 - Some answers may depend on unknown latent factors.
 Can we generate explanations that consider such factors?