VVSOR Symposium on Transparent Machine Learning, October 5, 2022

The Limits of Explainable Machine Learning: Some Things Are Simply Impossible

Tim van Erven



Joint work with:



Hidde Fokkema



Rianne de Heide

Explainable Machine Learning

The Need for Explanations:

Why did the machine learning system

- Classify my company as high risk for money laundering?
- Reject my bank loan?
- Give a certain medical diagnosis?
- Make a certain mistake?
- Reject the profile picture I uploaded to get a new OV chipcard?¹

Explainable Machine Learning

The Need for Explanations:

Why did the machine learning system

- Classify my company as high risk for money laundering?
- Reject my bank loan?
- Give a certain medical diagnosis?
- Make a certain mistake?
- Reject the profile picture I uploaded to get a new OV chipcard?¹

▶ ...

A Communication Limit:

- Cannot communicate millions of parameters!
- Can communicate only some relevant aspects and/or need high-level concepts in common with user

¹Personal experience

Binary Classification



Binary Classification



Local Post-hoc Explanations



Local: only explain the part of f that is (most) relevant for x.
Post-hoc: ignore explainability concerns when estimating f.

Local Explanations via Attributions



 $\phi_f(x) \in \mathbb{R}^d$ attributes a weight to each feature, which explains how important the feature is for the classification of x by f.

Examples of Local Attribution Methods

Example Attribution Method: LIME

LIME: Do local linear approximation of f near x (optionally in dimensionality reduced space), and report coefficients

LIME for tabular data:²



(classifying edibility of mushrooms)

²Image source: https://github.com/marcotcr/lime

Example Attribution Method: LIME

LIME: Do local linear approximation of f near x (optionally in dimensionality reduced space), and report coefficients

LIME for images:²



(a) Original Image

(b) Explaining *Electric guitar* (c) Explaining *Acoustic guitar*

(d) Explaining Labrador

Example: Gradient-based Explanations





³Image source: [Smilkov et al., 2017]

Example: Counterfactual Explanations

"If you would have had an income of €40 000 instead of €35 000, your loan request would have been approved."



Counterfactual explanation: $\tilde{x} = \underset{x': \operatorname{sign}(f(x')) \neq \operatorname{sign}(f(x))}{\operatorname{arg min}} \operatorname{dist}(x', x)$

Example: Counterfactual Explanations

"If you would have had an income of €40 000 instead of €35 000, your loan request would have been approved."



Counterfactual explanation: $\tilde{x} = \underset{x': \text{sign}(f(x')) \neq \text{sign}(f(x))}{\arg \min} \operatorname{dist}(x', x)$ Viewed as attribution method: $\phi_f(x) = \tilde{x} - x$

How Do We Evaluate Explanations?

- When are they good? Are some better than others?
- What is even the goal they are trying to achieve?

Explanations with Recourse as their Goal

"If you change your current income of €35 000 to €40 000, then your loan request will be approved."



Attribution methods provide recourse if they tell the user how to change their features such that f takes their desired value. Impossibility:

No Single Method Can Be Both Recourse Sensitive and Robust

Theorem

For any $\delta > 0$ there exists a continuous function f such that no attribution method ϕ_f can be both recourse sensitive and continuous.

Recourse Sensitivity

• Our definition: weakest possible requirement for providing recourse.



Recourse Sensitivity

Our definition: weakest possible requirement for providing recourse.



1. Assume user can change their features by at most some $\delta > 0$

Recourse Sensitivity

• Our definition: weakest possible requirement for providing recourse.



- 1. Assume user can change their features by at most some $\delta > 0$
- φ_f(x) can point in any direction that provides recourse within distance δ, and length does not matter as long as it is > 0.
- 3. If no direction provides recourse, then $\phi_f(x)$ can be arbitrary.

Recourse Sensitivity: Example

Profile picture is accepted if contrast between profile and background is large enough:



(a) Accepted profile picture



(b) Rejected profile picture

Recourse Sensitivity: Example

Profile picture is accepted if contrast between profile and background is large enough:



(a) Accepted profile picture



(b) Rejected profile picture



Recourse Sensitivity: Example

Profile picture is accepted if contrast between profile and background is large enough:



(a) Accepted profile picture







Robustness of Explanations

Compare:

- 1. "If you change your current income of €35 000 to €40 000, then your loan request will be approved."
- "If you change your current income of €35 001 to €45 000, then your loan request will be approved."

Minor changes in x should not cause big changes in explanations!

Robustness of Explanations

Compare:

- 1. "If you change your current income of €35 000 to €40 000, then your loan request will be approved."
- "If you change your current income of €35 001 to €45 000, then your loan request will be approved."

Minor changes in x should not cause big changes in explanations!

Robustness: If f is continuous, then ϕ_f should also be continuous. (e.g. survey of recourse by [Karimi et al., 2021])

Conclusion

Summary:

- In binary classification: exist f for which recourse sensitivity + robustness is impossible
- Further extensions in the paper:
 - Generalization to multiclass and regression using utility functions
 - Include constraints on user actions
 - Exact characterization of impossible f when user can only change a single feature

Conclusion

Summary:

- In binary classification: exist f for which recourse sensitivity + robustness is impossible
- Further extensions in the paper:
 - Generalization to multiclass and regression using utility functions
 - Include constraints on user actions
 - Exact characterization of impossible f when user can only change a single feature

Discussion:

Is the field of explainable machine learning in trouble? Not, but need to **refine goals** of explainability for recourse. E.g.:

- Accept that robustness sometimes fails
- Set-valued explanations
- Randomized explanations



References

H. Fokkema, R. de Heide and T. van Erven. Attribution-based Explanations that Provide Recourse Cannot be Robust, ArXiv:2205.15834 preprint, 2022.

Other references:

- A.-H. Karimi, G. Barthe, B. Schölkopf, and I. Valera. A survey of algorithmic recourse: definitions, formulations, solutions, and prospects. arXiv preprint arXiv:2010.04050, 2021.
- 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, 2016.
- D. Smilkov, N. Thorat, B. Kim, F. Viégas, and M. Wattenberg. Smoothgrad: removing noise by adding noise. *ArXiv:1706.03825*, 2017.



 $L = \{x : \text{recourse possible by moving at most } \delta \text{ left} \}$ $R = \{x : \text{recourse possible by moving at most } \delta \text{ right} \}$



 $L = \{x : \text{recourse possible by moving at most } \delta \text{ left} \}$ $R = \{x : \text{recourse possible by moving at most } \delta \text{ right} \}$

Recourse sensitivity implies:

$$\phi_f(x) \begin{cases} < 0 & \text{for } x \in L \setminus R \\ > 0 & \text{for } x \in R \setminus L \\ \neq 0 & \text{for } x \in L \cap R \end{cases}$$



 $L = \{x : \text{recourse possible by moving at most } \delta \text{ left} \}$ $R = \{x : \text{recourse possible by moving at most } \delta \text{ right} \}$

Recourse sensitivity implies:

$$\phi_f(x) \begin{cases} < 0 & \text{for } x \in L \setminus R \\ > 0 & \text{for } x \in R \setminus L \\ \neq 0 & \text{for } x \in L \cap R \end{cases}$$

But this **contradicts continuity**! (by the mean-value theorem)

Can embed 1D example in higher dimensions as well.