20S2 RDICHI **Characterizing Data Scientists' Mental Models** of Local Feature Importance

Dennis Collaris, Hilde Weerts, Daphne Miedema, Jarke J. van Wijk, Mykola Pechenizkiy

Eindhoven University of Technology

Aarhus, Denmark -October 10, 2022



Motivation Why do we need explanations?



Undertrusting the model → Overtrusting the model → Unfair treatment of groups

Motivation How do we explain models?

eXplainable Artificial Intelligence

Motivation **Low do we explain models?**

-1





Country = Denmark

Attempts = 29

Feature • importance definition

-0,5

A quantitative assignment of *importance* to the *feature* of a prediction made by a machine learning model.

Model prediction: FRAUD!



Gradient-based (e.g., LIME) Sensitivity to feature value changes

 $\hat{y} = \alpha + \sum \beta_i X_i$

 $\hat{y} = \epsilon + \sum_{i=1}^{n} \phi_{i}$

Gradient-based (e.g., LIME) Sensitivity to feature value changes



Gradient-based (e.g., LIME) Sensitivity to feature value changes



Gradient-based (e.g., LIME) Sensitivity to feature value changes



Gradient-based (e.g., LIME) Sensitivity to feature value changes



Motivation Which technique do l use?

Motivation Which technique do l use?



Method and results What Exploratory mixed-methods survey



Method and results What Exploratory mixed-methods survey





$\bigotimes \bigotimes$

Qualitative – Quantitative

Method and results Who Demographics

Туре	Answers (count)	
Gender	Male (20), Female (12), Prefer not	
Age	5 - 0 -	
Location	Netherlands (24) , United States (2 Singapore (1) , Spain (1) , Switzerla	
Role	Data scientist (18) , (Data science) (Data science) consultancy (2) , AV	
Data science experience (years)	5 0	
Familiarity	Linear regression coefficients (30) LIME (23), Permutation importanc Anchors (2), DeepLift (1)	





2), Prefer not to disclose (2), Colombia (1), India (1), and (1), United Kingdom (1)

researcher (6), Software/data/AI engineer (3), PhD candidate (3), /P (Assistant Vice President) (1), Prefer not to disclose (1)



, Random Forest feature importance (29), Shapley values (23), ce (18), Saliency maps (e.g., GradCAM) (11), treeinterpreter (9),



Open questions:

- Explain in your own words what feature importance is, in the context of machine learning.
- What does it mean to you when a feature is important in a (trained) machine learning model? (global)
- What does it mean to you when a feature is important for an individual prediction? (local)



Three main themes were identified:



Explanandum



3



Underlying mechanism

Three main themes were identified:



Explanandum



3



Underlying mechanism

Three main themes were identified:



Explanandum



3



Underlying mechanism

Three main themes were identified:



Explanandum



3



Underlying mechanism

Method and results RQ2 What are the data scientists' expectations of properties of local feature importance?

Strongly disagree

Disagree

I expect that if the model is relatively **certain** about a risk prediction, the sum of feature importances is relatively **high**.

O Strongly disagree	 Disagree 	O Neutral	O Agree	Strongly agree

I expect that if the model is relatively **uncertain** about a risk prediction, the sum of feature importances is relatively **high**.

Method and results RQ2 What are the data scientists' expectations of properties of local feature importance?

Expected properties:

- Robustness and stability
- Causality and sampling distribution
- Additivity

Varying expectations for:

- Selectivity
- Actionability and proportionality

Method and results RQ2 What are the data scientists' expectations of properties of local feature importance?

Expected properties:

- Robustness and stability
- Causality and sampling distribution
- Additivity

Varying expectations for:

- Selectivity
- Actionability and proportionality

Method and results RQ2 O Robustness & 2 Stability







Person **B**

Method and results RQ2 O Robustness & O Stability





4 Robustness

Method and results RQ2 O Robustness & 2 Stability





Method and results RQ2 O Robustness & 2 Stability



Person A







4 Robustness

Method and results RQ2 1 Robustness & 2 Stability







4 Robustness

Method and results RQ2 **B** Selectivity



4 Selectivity

Method and results RQ2 **B** Selectivity





✓ Selectivity



	Actionable	Proportional
Gradient- Based (LIME)		
Ablation- Based (SHAP)		



	Actionable	Proportional
Gradient- Based (LIME)		
Ablation- Based (SHAP)		



Conclusions and future work What do we conclude

- The definition and computation of "feature importance" can vary between methods resulting in different properties of feature importance values
- Data scientists expect feature importance values to be robust and stable
- Data scientists have varying expectations for selectivity, actionability, and proportionality.
- Data scientists have expectations of properties that are incompatible with existing techniques

Conclusions and future work What technique should you pick?

None of the existing techniques matched all participants' expectations!

Conclusions and future work What technique should you pick?

None of the existing techniques matched all participants' expectations!

Feature importance



Feature sensitivity for gradient-based techniques.

Feature attribution for ablation-based techniques.

20S2 RDICHI **Characterizing Data Scientists' Mental Models** of Local Feature Importance

Dennis Collaris, Hilde Weerts, Daphne Miedema, Jarke J. van Wijk, Mykola Pechenizkiy

Eindhoven University of Technology

Aarhus, Denmark -October 10, 2022

