# Bringing Order to Chaos: Probing the Disagreement Problem in XAI

#### Hima Lakkaraju





Model understanding is absolutely critical in several domains -- particularly those involving *high stakes decisions*!





# Motivation: Why Model Understanding?



# Motivation: Why Model Understanding?



# Achieving Model Understanding

Take 1: Build *inherently interpretable* predictive models



# Achieving Model Understanding

Take 2: *Explain* pre-built models *in a post-hoc manner* 



[Cireşan et. al. 2012, Caruana et. al. 2006, Frosst et. al. 2017, Stewart 2020]

# Inherently Interpretable Models vs. Post hoc Explanations



In *certain* settings, *accuracy-interpretability trade offs* may exist.

In *certain* settings, you may just have access to a (proprietary) black box.

# Feature Attribution Based Local Explanations

- Local explanations
  - explain individual predictions of any classifier
- Output feature attributions for individual instances, which capture the effect/contribution of each feature on the black box prediction
- Examples: LIME, SHAP, Gradient, Gradient times Input, SmoothGrad, Integrated Gradients

# **Disagreement Problem in XAI: Overview**

- Study to understand:
  - if and how often feature attribution based explanation methods disagree with each other in practice
  - What constitutes disagreement between these explanations, and how to formalize the notion of explanation disagreement based on practitioner inputs?
  - How do practitioners resolve explanation disagreement?

### Practitioner Inputs on Explanation Disagreement

- 30 minute semi-structured interviews with 25 data scientists
- 84% of participants said they often encountered disagreement between explanation methods
- Characterizing disagreement:
  - Top features are different
  - Ordering among top features is different
  - Direction of top feature contributions is different
  - Relative ordering of features of interest is different

### Practitioner Inputs on Explanation Disagreement

- Participants typically characterize explanation disagreement based on factors such as:
  - mismatch in top features,
  - feature ordering, and
  - directions of feature contributions,
  - But NOT on the feature importance values output by different explanation methods
- 24 out of 25 participants (96%) in our study opine that feature importance values output by different explanation methods are not directly comparable

#### Practitioner Inputs on Explanation Disagreement

• Quote: "The values generated by different explanation methods are clearly different. So, I would not characterize disagreement based on that. But, I would at least want the explanations they output to give me consistent insights. The explanations should agree on what are the most important features, the ordering among them and so on for me to derive consistent insights. But, they don't!"

# Formalizing the Notion of Explanation Disagreement (Top K)

 $FeatureAgreement(E_{a}, E_{b}, k) = \frac{|top\_features(E_{a}, k) \cap top\_features(E_{b}, k)|}{k}$ 

 $RankAgreement(E_a, E_b, k)$ 

 $\left|\bigcup_{s\in S} \{s \mid s \in top\_features(E_a, k) \land s \in top\_features(E_b, k) \land rank(E_a, s) = rank(E_b, s)\}\right|$ 

#### k

 $SignAgreement(E_a, E_b, k)$ 

 $\frac{|\bigcup_{s\in S} \{s \mid s \in top\_features(E_a, k) \land s \in top\_features(E_b, k) \land sign(E_a, s) = sign(E_b, s)\}|}{k}$ 

 $SignedRankAgreement(E_a, E_b, k)$ 

$$|\bigcup_{s \in S} \{s \mid s \in top\_features(E_a, k) \land s \in top\_features(E_b, k) \land sign(E_a, s) = sign(E_b, s) \land rank(E_a, s) = rank(E_b, s)\}|$$

# Formalizing the Notion of Explanation Disagreement (Features of Interest)

Spearman rank correlation coefficient computed over features of interest

 $RankCorrelation(E_a, E_b, F) = r_s(Ranking(E_a, F), Ranking(E_b, F))$ 

 $PairwiseRankAgreement(E_a, E_b, F) = \frac{\sum_{i,j \text{ for } i < j} \mathbb{1} [RelativeRanking(E_a, f_i, f_j) = RelativeRanking(E_b, f_i, f_j)]}{\binom{|F|}{2}}$ 

# Empirical Analysis of Explanation Disagreement





• We carried out empirical analysis with 6 post hoc explanation methods, 4 real world datasets (tabular, NLP, images), 8 model classes, and found several disagreements between explanation methods

### How do Practitioners Resolve Disagreements?



#### Below, you see a data point, as well as its explanation using methods LIME and KernelSHAP.

As a reminder, the 7 features of the COMPAS dataset are age, two\_year\_recid (whether the defendant recidivated after 2 years of the original crime, priors\_count (number of prior crimes committed), length\_of\_stay (length the defendant stayed in jail), c\_charge\_degree (whether the previous charge was a Misdemeanor or Felony), sex, and race

To what extent do you think the two explanations shown above agree or disagree with each other? Completely agree 
Mostly agree Completely disagree

Please explain why you chose the above answer.

Since you believe that the above explanations disagree (to some extent), which explanation would you rely on?

Please explain why you chose the above answer.

# How do Practitioners Resolve Disagreements?

- Online user study where 25 users were shown explanations that disagree and asked to make a choice, and explain why
- Practitioners are choosing methods due to:
  - Associated theory or publication time (33%)
  - Explanations matching human intuition better (32%)
  - Type of data (23%)
    - E.g., LIME or SHAP are better for tabular data

# How do Practitioners Resolve Disagreements?

Algorithm	Reasons that algorithm was chosen in disagreement		
	• [36%] SHAP is better for tabular data ("SHAP is more commonly used		
KernelSHAP	[than Gradient] for tabular data")		
	• [25%] SHAP is more familiar ("More information present + more		
	familiarity")		
	• [14%] SHAP is a better algorithm overall ("SHAP seems more method- ical than LIME", "SHAP is a more rigorous approach [than LIME] in		
	theory")		
SmoothGrad	• [33%] SmoothGrad paper is newer or better ("SmoothGrad is apparently		
	more robust", "SmoothGrad is often considered improved verison of grad")		
	• [58%] Reasons based on the explainability map shown ("directionality		
	of the attributions [agree] with intuition", "gradient has unstability		
	problems [, so] smoothgrad")		
LIME	• [54%] LIME is better for tabular data ("I use LIME for structured		
	data.")		
	• [15%] LIME is more familiar/easier to interpret ("I am more familiar		
	with LIME", "LIME is easy to interpret")		
Integrated	• [86%] Integrated Gradients paper is better ("IG came after gradi-		
Gradients	ents and paper shows improvements", "integrated gradients paper showed		
	improvements [over Gradient $\times$ Input]"		

# Insights and Moving Forward

- Feature attribution methods often disagree in practice w.r.t. basic insights, and practitioners adopt ad hoc heuristics to resolve those disagreements!
- Why do feature attribution methods disagree?
- Given that feature attribution methods disagree, which explanation method should we choose for different kinds of data and applications?

# Why do Feature Attribution Methods Disagree?

 Various feature attribution methods (e.g., LIME, C-LIME, KernelSHAP, Occlusion, Vanilla Gradients, Gradient times Input, SmoothGrad, Integrated Gradients) are essentially local function approximations.

$$g^* = \operatorname*{arg\,min}_{g \in \mathcal{G}} \mathop{\mathbb{E}}_{\xi \sim \mathcal{Z}} \ell(f, g, \mathbf{x}_0, \xi)$$

• But...

# Why do Feature Attribution Methods Disagree?

But, they adopt different loss functions, and local neighborhoods

<b>Explanation Method</b>	Local Neighborhood $\mathcal{Z}$ around $\mathbf{x}_0$	Loss Function $\ell$
C-LIME SmoothGrad Vanilla Gradients	$ \begin{array}{l} \mathbf{x}_0 + \xi; \ \xi(\in \mathbb{R}^d) \sim \operatorname{Normal}(0, \sigma^2) \\ \mathbf{x}_0 + \xi; \ \xi(\in \mathbb{R}^d) \sim \operatorname{Normal}(0, \sigma^2) \\ \mathbf{x}_0 + \xi; \ \xi(\in \mathbb{R}^d) \sim \operatorname{Normal}(0, \sigma^2), \sigma \to 0 \end{array} $	Squared Error Gradient Matching Gradient Matching
Integrated Gradients Gradients × Input	$ \begin{aligned} & \xi \mathbf{x}_0; \ \xi (\in \mathbb{R}) \sim \text{Uniform}(0,1) \\ & \xi \mathbf{x}_0; \ \xi (\in \mathbb{R}) \sim \text{Uniform}(a,1), a \to 1 \end{aligned} $	Gradient Matching Gradient Matching
LIME KernelSHAP Occlusion	$ \begin{array}{l} \mathbf{x}_0 \odot \xi; \ \xi (\in \{0,1\}^d) \sim \text{Exponential kernel} \\ \mathbf{x}_0 \odot \xi; \ \xi (\in \{0,1\}^d) \sim \text{Shapley kernel} \\ \mathbf{x}_0 \odot \xi; \ \xi (\in \{0,1\}^d) \sim \text{Random one-hot vectors} \end{array} $	Squared Error Squared Error Squared Error

# Why Do Feature Attribution Methods Disagree?

• *No Free Lunch Theorem for Explanation Methods*: No single method can perform optimally across all neighborhoods

**Theorem 3** (No Free Lunch for Explanation Methods). Consider the scenario where we explain a black-box model f around point  $\mathbf{x}_0$  using an interpretable model g from class  $\mathcal{G}$  and a valid loss function  $\ell$  where the distance between f and  $\mathcal{G}$  is given by  $d(f, \mathcal{G}) = \min_{g \in \mathcal{G}} \max_{\mathbf{x} \in \mathcal{X}} \ell(f, g, 0, \mathbf{x})$ . Then, for any explanation  $g^*$  on a neighborhood distribution  $\xi_1 \sim \mathcal{Z}_1$  such that  $\max_{\xi_1} \ell(f, g^*, \mathbf{x}_0, \xi_1) \leq \epsilon$ , we can always find another neighborhood  $\xi_2 \sim \mathcal{Z}_2$  such that  $\max_{\xi_2} \ell(f, g^*, \mathbf{x}_0, \xi_2) \geq d(f, \mathcal{G})$ .

# Which Method Should We Choose?: Take 1

- A guiding principle based on function approximation: choose a method which recovers the underlying model when the model is a member of the explanation function class
- For continuous data, use additive continuous noise methods (e.g. SmoothGrad, Vanilla Gradients, C-LIME) or multiplicative continuous noise methods (e.g. Integrated Gradients, Gradient x Input). For binary data, use binary noise methods (e.g. LIME, KernelSHAP, Occlusion).

# Which Method Should We Choose?: Take 2

- OpenXAI: open-source framework to readily evaluate and benchmark post hoc explanation methods
- Systematic, efficient, and reproducible evaluations of post hoc explanation methods on various datasets
- Assessing reliability of post hoc explanation methods from diverse perspectives (e.g., faithfulness, stability, fairness)
- (Customizable) dashboards to compare existing and new methods across various datasets easily

# **Conclusions and Summary**

- Several methods proposed to "explain" machine learning models in prior research
- Important to characterize these methods, and understand which methods can be useful under what circumstances
- Critical to bridge the gaps between researchers and practitioners

# Thank You!

- Email: <u>hlakkaraju@hbs.edu</u>; <u>hlakkaraju@seas.harvard.edu</u>;
- Webpage: <u>https://himalakkaraju.github.io</u>
- Course on interpretability and explainability: <u>https://interpretable-ml-class.github.io/</u>
- Multiple tutorials on explaining ML models (ranging from 1 hour to 3 hours): <u>explainml-tutorial.github.io</u>
- Trustworthy ML Initiative: <u>https://www.trustworthyml.org/</u>
  - Lots of resources and seminar series on topics related to explainability, fairness, adversarial robustness, differential privacy, causality etc.