To cite this presentation: Pendyala, V.S. (2024) "Explainability of Al/ML models for a socially sustainable Al growth". University of Bolton's November to Remember 2024 Inaugural Lecture

## Explainability of AI/ML models for a socially sustainable AI growth

Vishnu S. Pendyala, Ph.D. San Jose State University

Video Recording: https://www.youtube.com/watch?v=R8HR9Y3vG-8

©Vishnu S. Pendyala This work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 International License

PROFILES

Source: The New Yorker, November 13, 2023

## WHY THE GODFATHER OF A.I. FEARS WHAT HE'S BUILT

Geoffrey Hinton has spent a lifetime teaching computers to learn. Now he worries that artificial brains are better than

ours.



## AI makes mistakes sometimes

Decides unfairly

Overfits causing privacy concerns

Hallucinates giving incorrect answers

And often, no one can accurately explain its behavior!

Pendyala, V., & Kim, H. (2024). Assessing the Reliability of Machine Learning Models Applied to the Mental Health Domain Using Explainable Al. *Electronics*, *13*(6), 1025.

"This work proves that merely achieving superlative evaluation metrics can be dangerously misleading and may infringe upon ethical horizons. A future direction is to investigate methods to quantify the effectiveness of machine learning models in terms of insights from their explainability."

©Vishnu S. Pendyala This work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 Inte

Y. Liu, C. Tantithamthavorn, L. Li and Y. Liu, "Explainable AI for Android Malware Detection: Towards Understanding Why the Models Perform So Well?," 2022 IEEE 33rd International Symposium on Software Reliability Engineering (ISSRE), Charlotte, NC, USA, 2022, pp. 169-180, "our results indicate that ML models classify malware based on temporal differences between malware and benign, rather than the actual malicious behaviors."

"We discover that temporal sample inconsistency in the training dataset brings over-optimistic classification performance (up to 99%F1 score and accuracy)."

Tian, Y., Ma, S., Wen, M. et al. To what extent do DNN-based image classification models make unreliable inferences?. Empir Software Eng 26, 84 (2021).

"we applied our approach to 18 pre-trained single-label image classification models and 3 multi-label classification models, and then examined their inferences on the ImageNet and COCO datasets. We found that unreliable inferences are pervasive. Specifically, for each model, *more than thousands of correct classifications are actually made using irrelevant features.*"

©Vishnu S. Pendyala This work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 International License

The "Right to explanation" in the GDPR

#### Recital 71

"...to obtain an explanation of the decision reached after such assessment and to challenge the decision..."

Articles 13, 14, and 15

"...of automated decision-making, including profiling, referred to in Article 22(1) and (4) and, at least in those cases, *meaningful information about the logic involved*"





# Explaining human behavior: Credit Denied!

**Poor Credit Score**: Applicant had a history of late payments and defaults

**Insufficient Income**: Applicant's income is too low relative to their existing debt and living expenses

**Too Many Open Credit Accounts**: Applicant had multiple credit cards and loans - ability to manage additional credit is questionable

#### ...and more

## Machine Learning Model for Credit Approval: Features



#### Logistic Regression model to approve or deny credit applications



We can generate similar explanations as humans generate, based on the features and their weights



## K Nearest Neighbors for Credit Approval



## Approval using a Decision Tree Machine Learning Model



#### **Prior Credit Approval Applications**



# Can we create a model-agnostic method for explaining predictions of ML models?

Explanations are not specific to the ML models – sometimes nearest neighbors or decision trees are not the best way to explain

Easily replace or upgrade models without changing the explanation approach

Explainability is now a layer on top of the black box ML model

Easier benchmarking and comparisons of model interpretability

©Vishnu S. Pendyala This work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 International License



#### What if the features are not human-interpretable?

#### What if the features are not human-interpretable?

embedding = [

0.234, -0.115, 0.587, 0.020, -0.367, 0.200, 0.150, 0.405, 0.100, -0.075, 0.289, 0.451, -0.201, 0.120, 0.340, 0.460, -0.069, 0.370, -0.199, 0.023, 0.135, -0.020, 0.250, 0.015, -0.233, 0.415, 0.056, 0.089, -0.305, 0.670, 0.222, 0.300, -0.045, 0.090, 0.159, -0.028, 0.016, 0.100, -0.220, 0.321, 0.067, -0.112, 0.245, 0.091, -0.030, 0.375, 0.190, -0.177, 0.022, 0.060, 0.289, -0.098, 0.155, -0.074, 0.118, 0.034, 0.205, -0.005, 0.019, 0.350, 0.088, -0.093, 0.417, 0.204, -0.154, 0.136, -0.105, 0.209, 0.095, -0.050, 0.072, -0.189, # ... (remaining values) 0.237, 0.152, 0.043, 0.092, -0.203, 0.155, 0.046, 0.360, -0.112, 0.244, -0.087, 0.271, 0.048, 0.112, 0.190

## Human interpretable knowledge representation



<sup>] ©</sup>Vishnu S. Pendyala This work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 International License

### Local Interpretable Model-agnostic Explanations (LIME)

Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "Why should i trust you?" Explaining the predictions of any classifier." Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining. 2016. Provides local explanations for individual predictions (not global)

Uses a surrogate interpretable model like the ones discussed earlier

Local: Surrogate model is trained on data in a local context, typically around a specific data item

Interpretable: Human-understandable even if the features are not

Outputs feature importance scores for the prediction

©Vishnu S. Pendyala Thi work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 International License

	Local Interpretability	Global Interpretability
Focus	Explains individual predictions	Explains the overall model behavior
Scope	Specific instances or data points	Entire model or dataset
Methods	LIME, SHAP	Decision Trees, Rule-based models
Goal	Understand why a model made a specific prediction for a data item	Understand how the model works in general across all data
Application	Debugging, identifying biases, building trust	Model design, feature engineering, risk assessment
Use cases	Clinicians, loan officers needing explanations for specific decisions	Regulatory, where trust / compliance are crucial; scientific research
Level of Detail	High level of detail for specific predictions	Some level of detail sacrificed for comprehensiveness





The value of the RBF Kernel is maximum at the center where the distance = 0  $e^0 = 1$ 

This Photo by Unknown Author is licensed under CC BY-SA

$$K(\mathbf{x},\mathbf{x}') = \exp\!\left(-rac{\|\mathbf{x}-\mathbf{x}'\|^2}{2\sigma^2}
ight)$$

LIME – the math  

$$L(f,g,\pi_x) = \sum_{z \in \mathbb{Z}} \pi_x(z) (f(z) - g(z'))^2, L = loss function$$
f is the blackbox model, g is the surrogate interpretable model  

$$\pi_x(z) = \exp\left(\frac{-D(x,z)^2}{\sigma^2}\right); z \in \mathbb{R}^d \text{ is a sample in the perturbed set, } Z$$
f(z) is the label generated by the blackbox model  
z' is the interpretable version of z; g(z') is the label from g  

$$\xi(x) = \arg\min_{g \in G} L(f,g,\pi_x) + \Omega(g) \text{ is the generated explanation}$$
G is the set of interpretable models,  $\Omega$  is the regularization  
evision of z = Generated by the blackbox model

#### Pendyala, Vishnu, and Hyungkyun Kim. "Assessing the Reliability of Machine Learning Models Applied to the Mental Health Domain Using Explainable AI." *Electronics* 13.6 (2024): 1025.



©Vishnu S. Pendyala This work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 International License

### **Explaining Object Detection**



## SHapley Additive exPlanations (SHAP)

Scott, M., & Su-In, L. (2017). A	Based on cooperative game theory, specifically using Shapley values; treats each feature as a player in a game
approach to interpreting model	Can be used for explanations at both global and local levels
predictions. Advances in neural information	The model prediction is the sum of the SHAP values for each feature
processing systems, 30, 4765-4774.	Calculates the average contribution of each feature to the prediction by considering all possible subsets of feature

©Vishnu S. Pendyala This work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 International License



Shapley, Lloyd S. "A value for n-person games." Contribution to the Theory of Games 2 (1953)

Goal: Distribute total gains among players in a coalition based on their contributions

Each player gets their fair share of the total gains based on their marginal contributions

Considers all possible coalitions of players and their interactions

Order independent, proportional weighting - fairness is key



## Feature Subsets

Full set of features, F: Credit Score (CS), Annual Income (AI), Loan Amount (LA), Debt-to-Income Ratio (DIR), and Employment Length (EL).

Feature Subset	Features Included		
{LA}	Loan Amount		
{EL}	Employment Length		
{CS, AI}	Credit Score, Annual Income		
{AI, EL}	Annual Income, Employment Length		
{LA, DIR}	Loan Amount, Debt-to-Income Ratio		
{CS, AI, EL}	Credit Score, Annual Income, Employment Length		
{CS, LA, EL}	Credit Score, Loan Amount, Employment Length		
{CS, DIR, EL}	Credit Score, Debt-to-Income Ratio, Employment Length		
{CS, AI, DIR, EL}	Credit Score, Annual Income, Debt-to-Income Ratio, Employment Length		
{AI, LA, DIR, EL}	Annual Income, Loan Amount, Debt-to-Income Ratio, Employment Length		
{CS, AI, LA, DIR, EL}	Credit Score, Annual Income, Loan Amount, Debt-to-Income Ratio, Employment Length		

©Vishnu S. Pendyala This work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 International License

## SHAP – the math



#### SHAP for mental health prediction model explainability

Pendyala, Vishnu, and Hyungkyun Kim. "Assessing the Reliability of Machine Learning Models Applied to the Mental Health Domain Using Explainable AI." Electronics 13.6 (2024): 1025.







## Explainable Misinformation Classification using Large Language Models

Legend: Negative 🗆 Neutral 🗖 Positive				
Predicted Label	Experiment Type	Word Importance		
		#s Please select the option that most closely describes the following claim by Donald Trump : H ill ary Cl inton has spoken such lies		
Barely True	IG	about my foreign policy . They said I want Japan to get nuclear weapons . Give me a break . A ) True B ) Most ly True C ) Half True		
(0.03)		D) B are ly True E) False F) P ants on Fire ( abs urd lie ) Choice : (		
		#s Please select the option that most closely describes the following claim by Donald Trump : H ill ary Cl inton has spoken s		
	LIME	about my foreign policy , They said I want Japan to get nuclear weapons . Give me a break . A ) True B ) Most ly True C ) Half True		
		D ) B are ly True E ) False F ) P ants on Fire ( abs urd lie ) Choice : (		
		#s Please select the option that most closely describes the following claim by Donald Trump : H ill ary Cl inton has spoken		
	SHAP	about my foreign policy. They said I want Japan to get nuclear weapons . Give me a break . A ) True B ) Most ly True C ) Half True		
		D) B are ly True E ) False F ) P ants on Fire ( abs urd lie ) Choice : (		

©Vishnu S. Pendyala This work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 International License

## **Future Directions**



