

An attempt to understand our models' predictions

Parul Pandey, Machine Learning Engineer



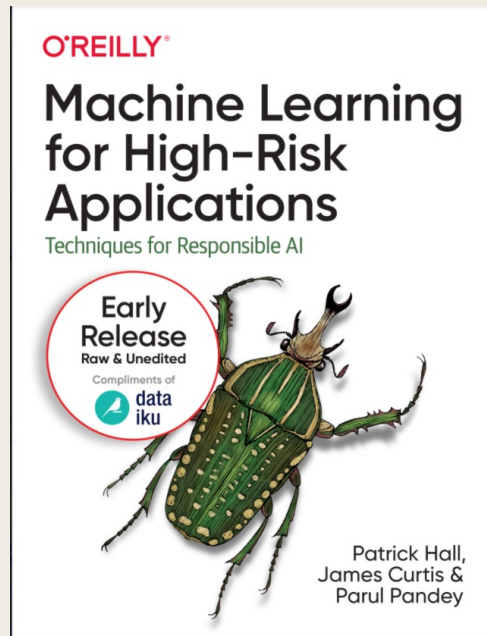
Who am I?

- ML Engineer @ Weights & Biases
- Electrical Engineer
- Kaggle Grandmaster(Notebooks)
- Author



Who am I?

- ML Engineer @ Weights & Biases
- Electrical Engineer
- Kaggle Grandmaster(Notebooks)
- Author



Motivation

Complex systems tend to drift toward unsafe conditions unless constant vigilance is maintained.

— Closing the AI Accountability Gap, Google Research

Motivation

Machine Bias

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

by Ju

Amit Datta*, Michael Carl Tschantz, and Anupam

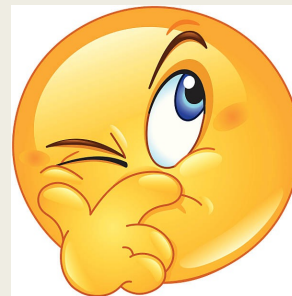
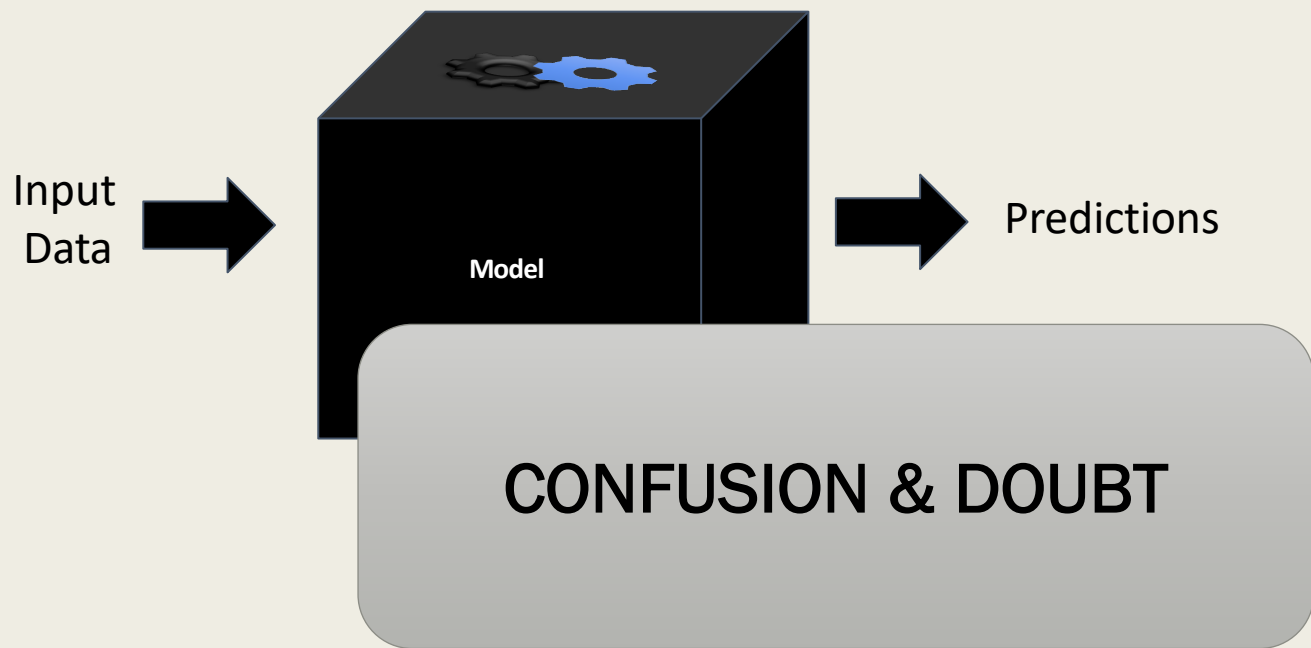
Automated Experiments

A Tale of Opacity, Choice, and Discrimination



- Source: <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>
- <https://www.andrew.cmu.edu/user/danupam/dtd-pets15.pdf>
- Washington Post

Complex & Opaque Models

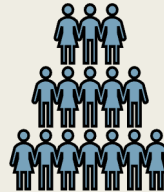


Why Explainability

- To inculcate Trust in models
- To debug predictions
- To detect bias
- Ensure suitability of models for deployment
- Global AI Regulation



Data Scientists and ML Engineers



End Users



Regulators

Explanations

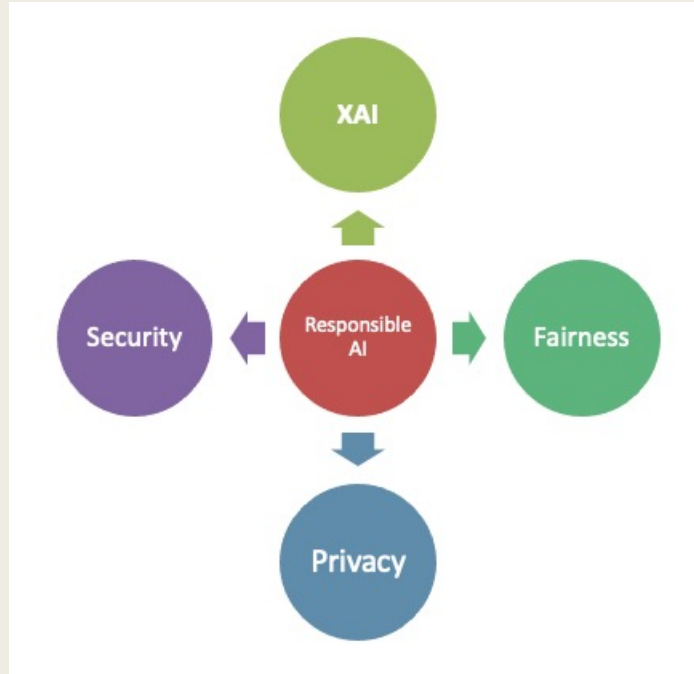
The degree to which a humans can understand and trust a ML model's predictions.

- Miller, Tim. "Explanation in artificial intelligence: Insights from the social sciences"
- Kim, Been, Rajiv Khanna, and Oluwasanmi O. Koyejo. "Examples are not enough, learn to criticize! Criticism for interpretability." Advances in Neural Information Processing Systems (2016).

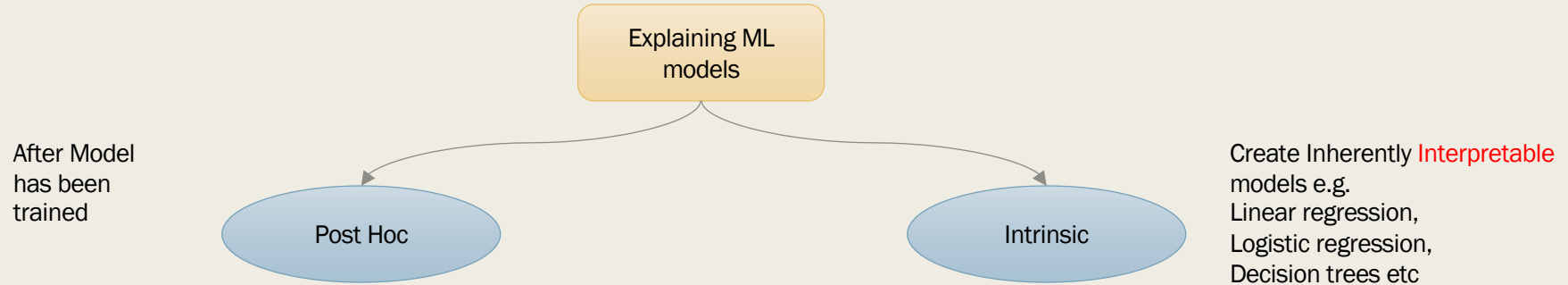
Interpretable vs Explainable AI (XAI)

- Describes the internals of a system in a way which is understandable to humans
- Summarize the reasons for their behaviour.

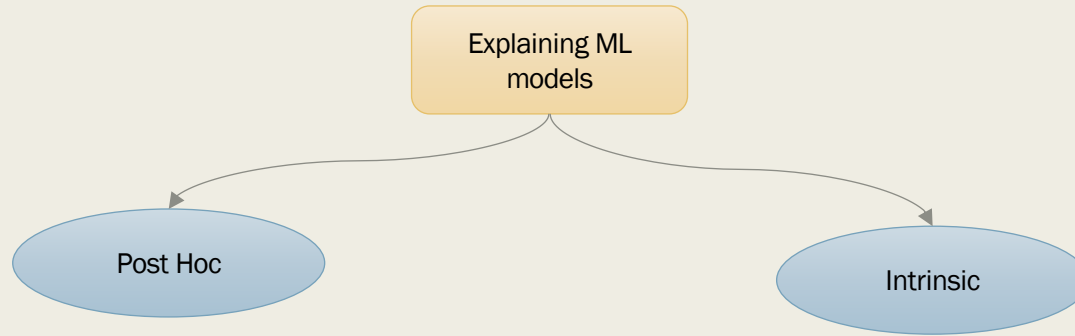
Responsible AI



Taxonomy

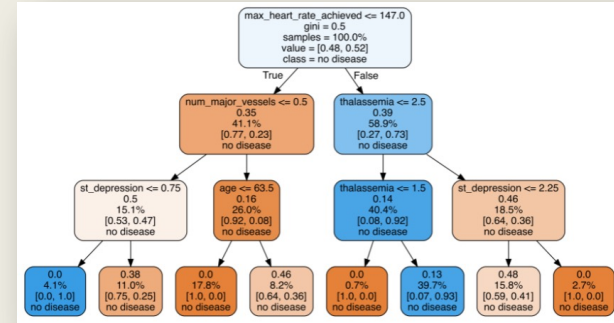


Taxonomy



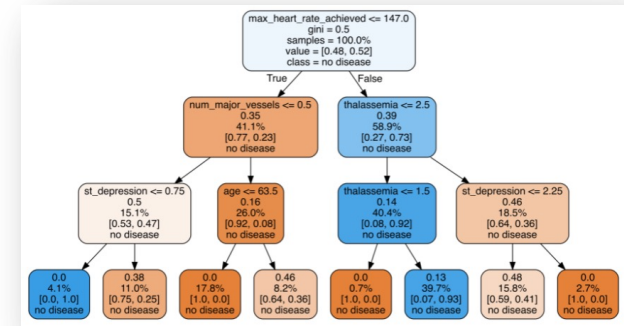
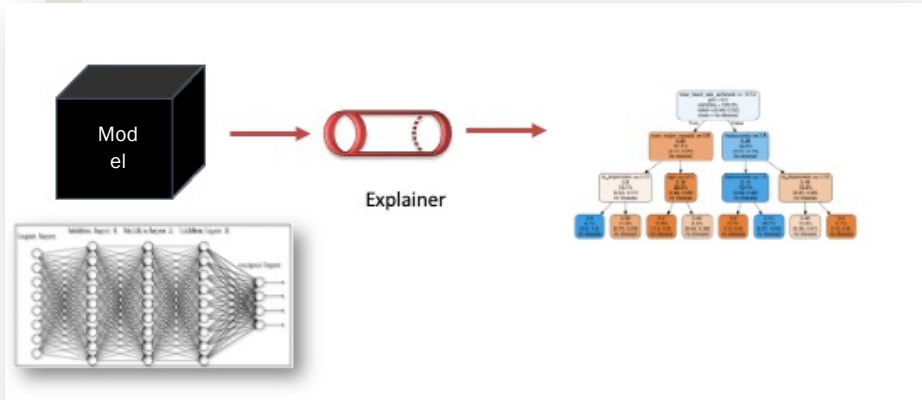
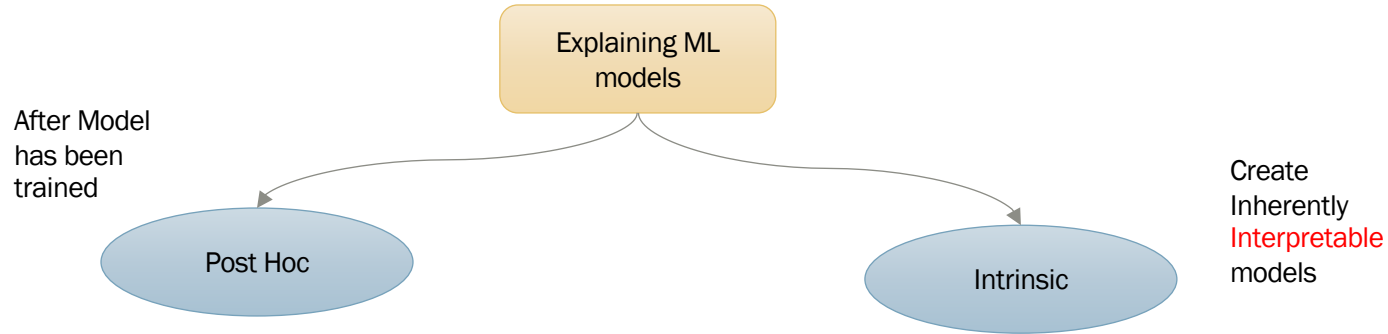
After Model
has been
trained

Create Inherently **Interpretable**
models e.g.
Linear regression,
Logistic regression,
Decision trees etc



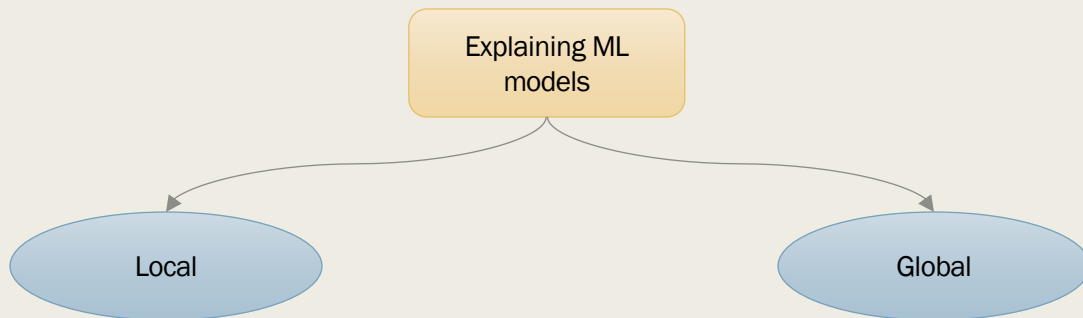
Does a person have a heart disease ?

Taxonomy



Does a person have a heart disease ?

Taxonomy



explaining an individual prediction or a small part of the model's prediction space

explaining the entire model behaviour



Interpretable Models

- Penalized Regression
- Additive Models – GAMs
- Decision Trees
- Constrained XGBoost Models



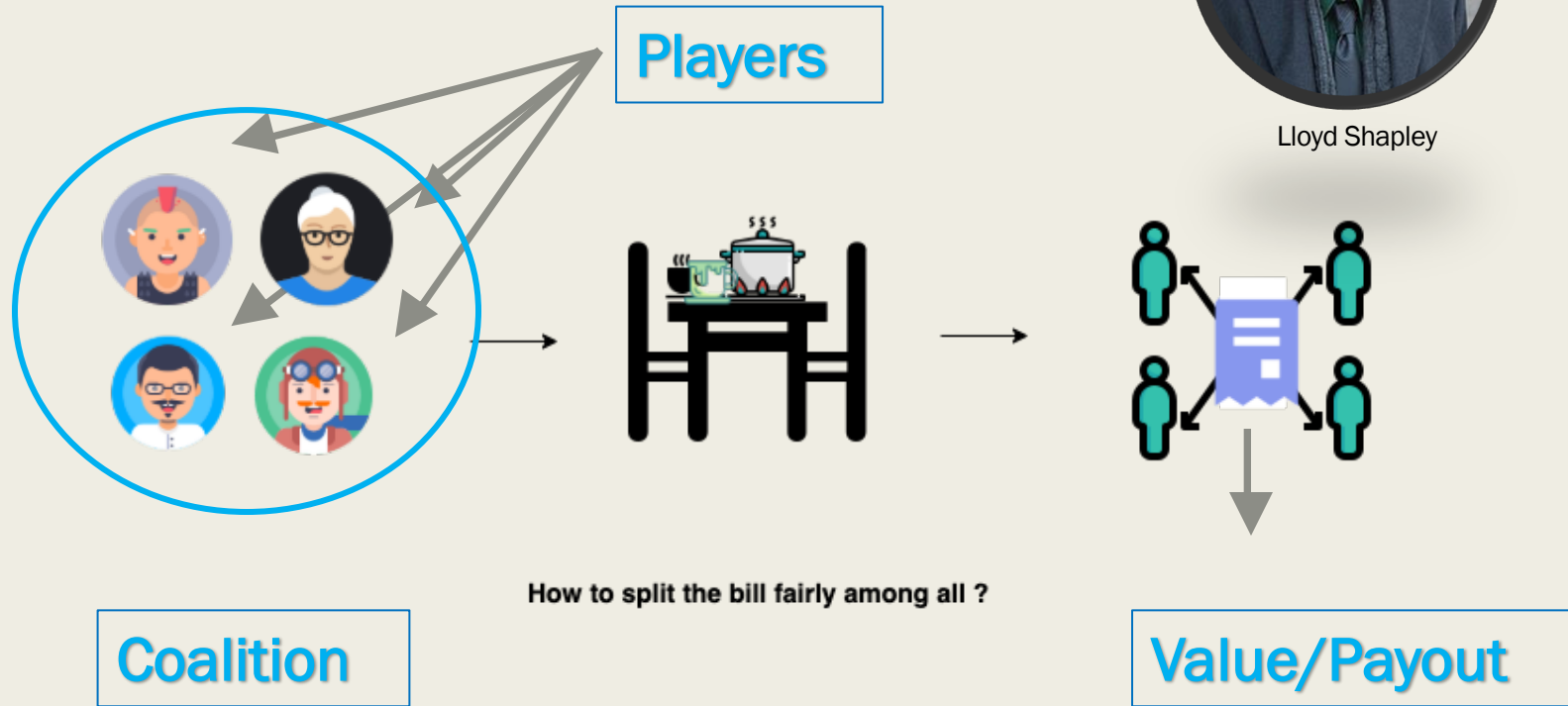
Post Hoc Techniques

- Partial Dependence Plots
- SHAP
- LIME
- Counterfactuals
- Saliency Maps
- Occlusion

SHAP values

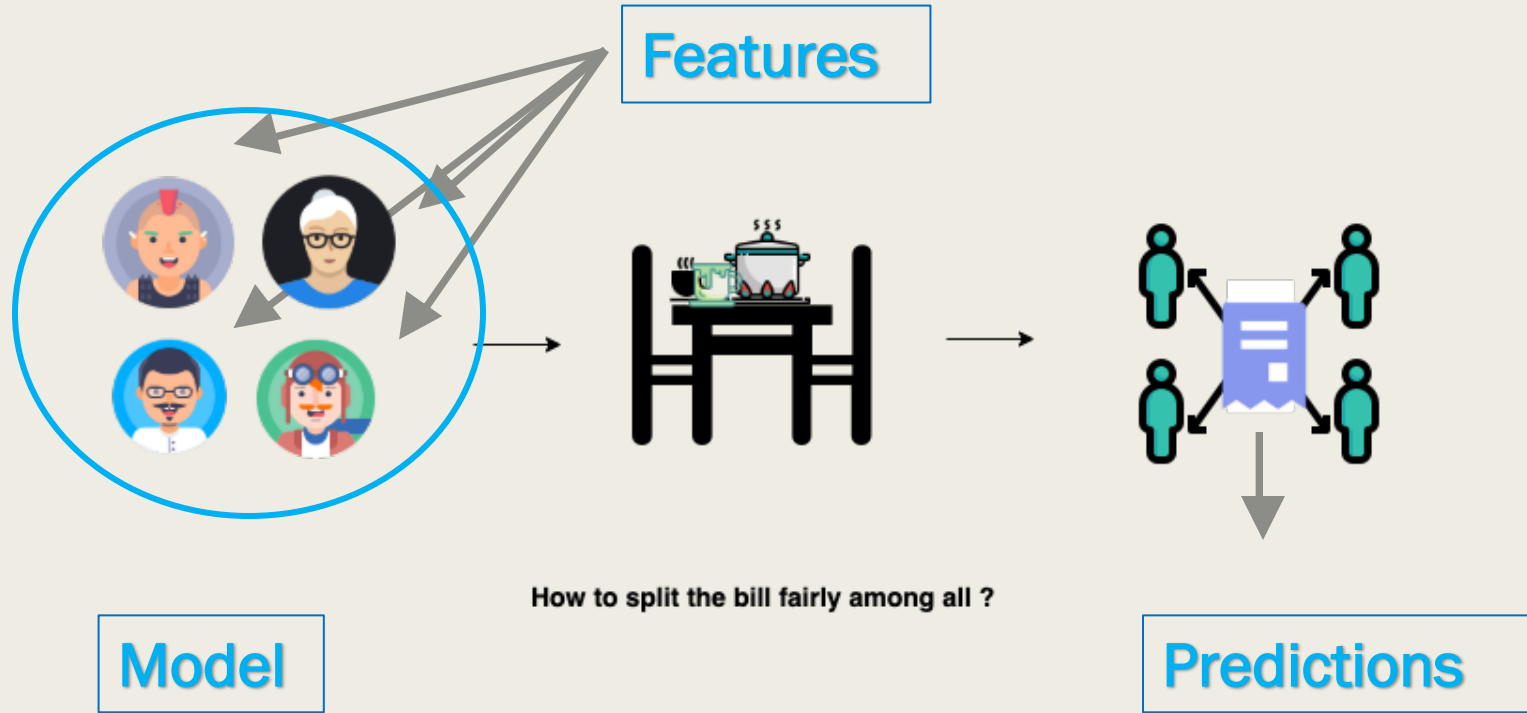
- SHAP stands for **SH**apley **A**dditive ex**P**lanations.
- It leverages the concept of **Cooperative Game theory** to break down a prediction to measure the impact of each feature.

Cooperative Game Theory

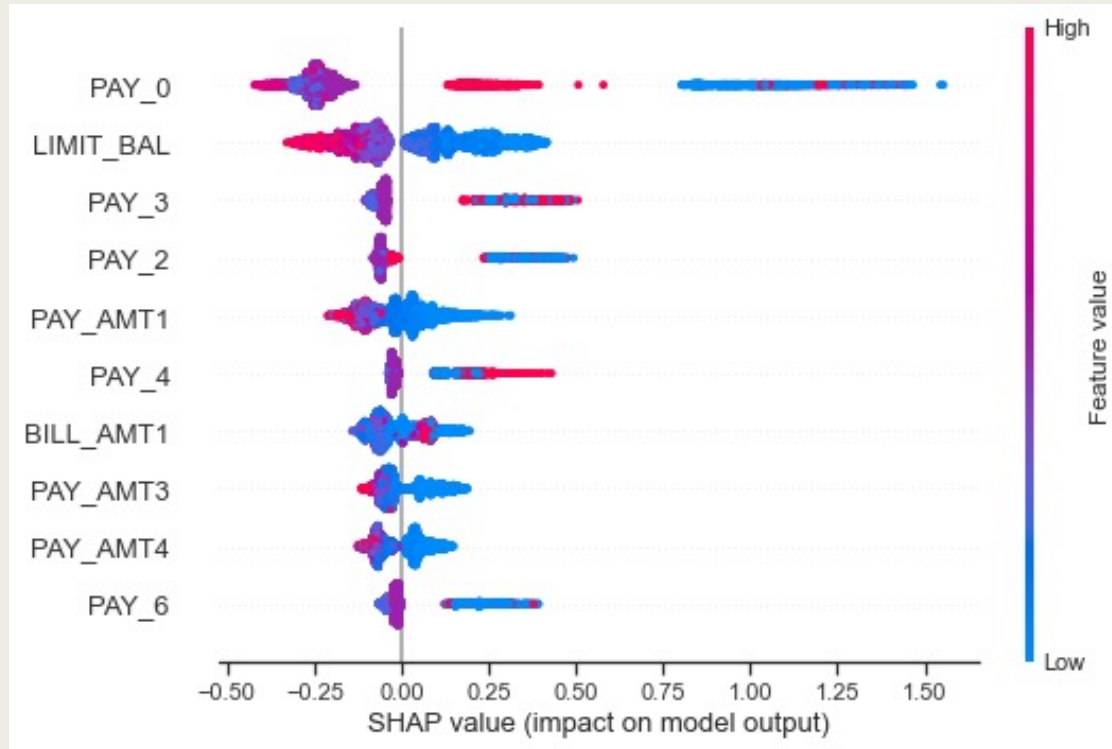


Lloyd Shapley

SHAP values : Machine Learning context



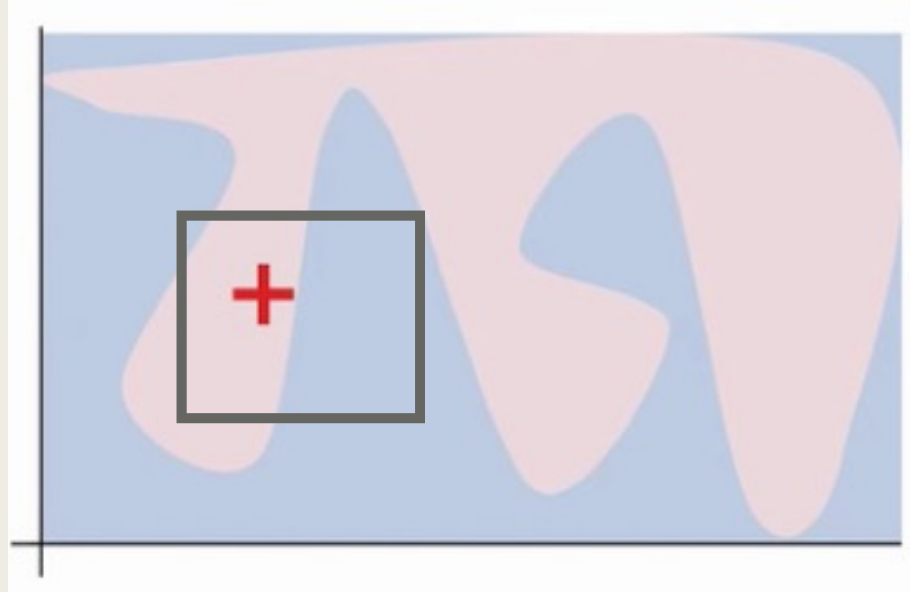
Who will default on a Credit Card Payment?



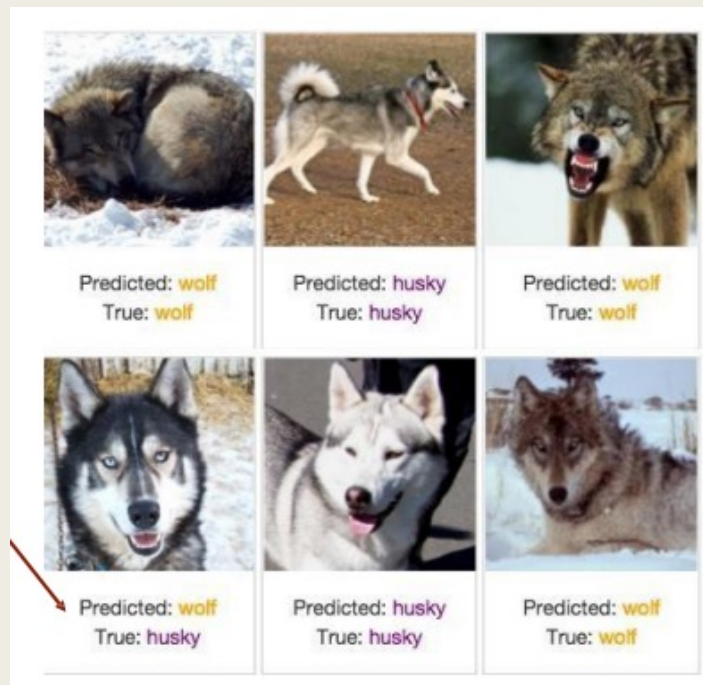
Lime

- LIME stands for **L**ocally **I**nterpretable **M**odel-Agnostic **E**xplanations.
- A technique to explain the predictions of any machine learning classifier – Model Agnostic.
- Approximates a given model by an interpretable one (such as a linear model)

Intuition

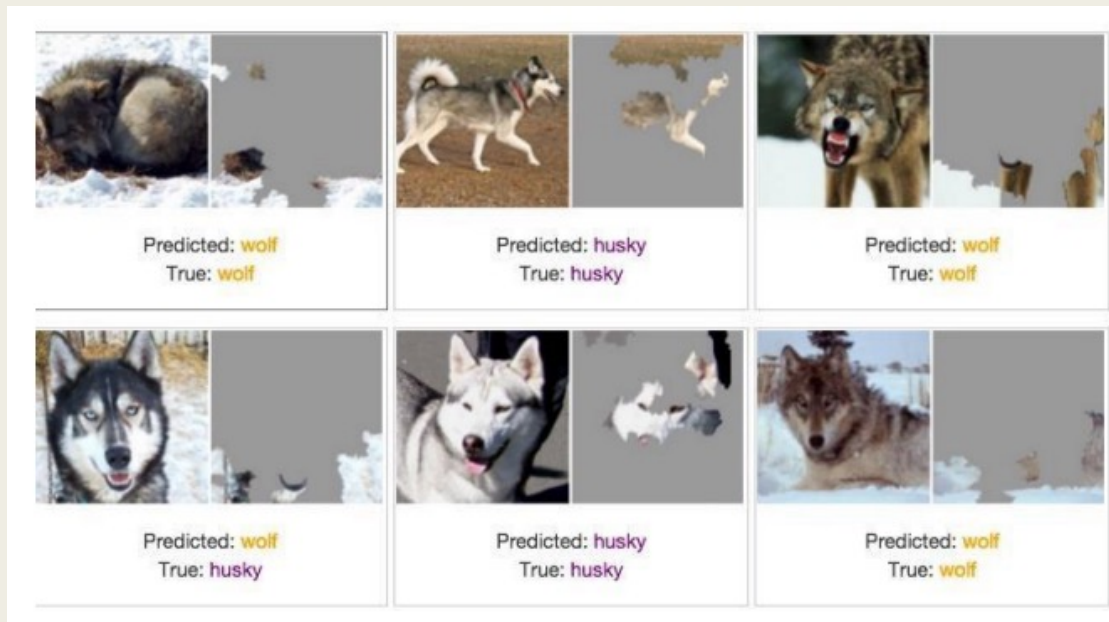


Lime for Image data



Husky or Wolf ?

Lime for Image data



A great SNOW detector



Edited by Murat Durmus (CEO AISOMA)

Thankyou