

# Interpretable Machine Learning

## Global explainability techniques

Vasilis Gkolemis<sup>12</sup>

<sup>1</sup>ATHENA Research and Innovation Center

<sup>2</sup>Harokopio University of Athens

February 2022

# Motivation

- The credit assessment system leads to the rejection of an application for a loan - the client suspects racial bias<sup>1</sup>

---

<sup>1</sup><https://www.technologyreview.com/2021/06/17/1026519/racial-bias-noisy-data-credit-scores-mortgage-loans-fairness-machine-1>

<sup>2</sup><https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

# Motivation

- The credit assessment system leads to the rejection of an application for a loan - the client suspects racial bias<sup>1</sup>
- A model that assesses the risk of future criminal offenses (and used for decisions on parole sentences) is biased against black prisoners<sup>2</sup>

---

<sup>1</sup><https://www.technologyreview.com/2021/06/17/1026519/>

[racial-bias-noisy-data-credit-scores-mortgage-loans-fairness-machine-1](https://www.propublica.org/article/racial-bias-noisy-data-credit-scores-mortgage-loans-fairness-machine-1)

<sup>2</sup>[https://www.propublica.org/article/](https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing)

[machine-bias-risk-assessments-in-criminal-sentencing](https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing)

# Interpretability

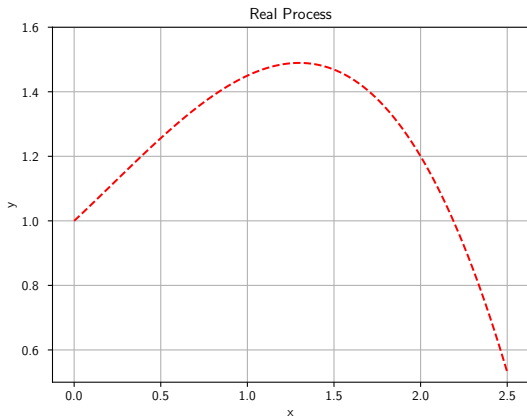
- Questions:
  - Why did a model make a specific decision?
  - Can we summarize the model's behavior?
- Answer:
  - Interpretable Machine Learning (IML)
  - “Extraction of relevant knowledge from a machine-learning model concerning relationships either contained in data or learned by the model”<sup>3</sup>

---

<sup>3</sup>Murdoch, W. J., Singh, C., Kumbier, K., Abbasi-Asl, R. and Yu, B. “Definitions, methods, and applications in interpretable machine learning.” Proceedings of the National Academy of Sciences, 116(44), 22071-22080. (2019)

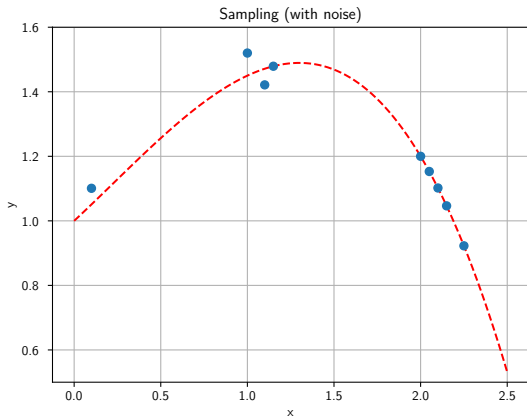
# Example

Consider the following mapping  $x \rightarrow y$



# Example

Process unknown  $\rightarrow$  we only have samples



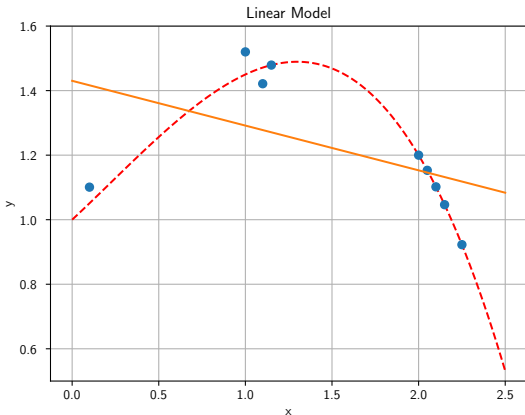
# Example

Our goal is to model the process using the available samples  
(regression)

# Example

Linear model → Underfitting!

$$y = w_1 \cdot x + w_0$$

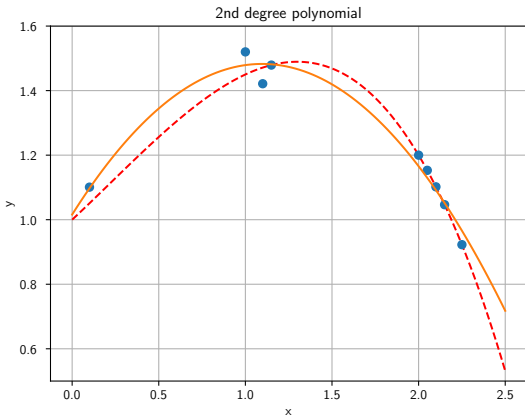




# Example

2<sup>nd</sup> degree polynomial → Decent Fit!

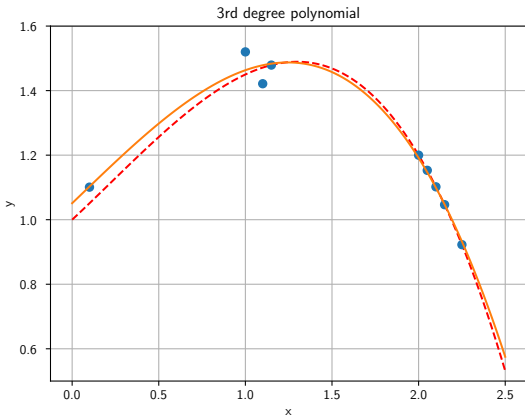
$$y = w_2 \cdot x^2 + w_1 \cdot x + w_0$$



# Example

3<sup>rd</sup> degree polynomial → Good Fit!

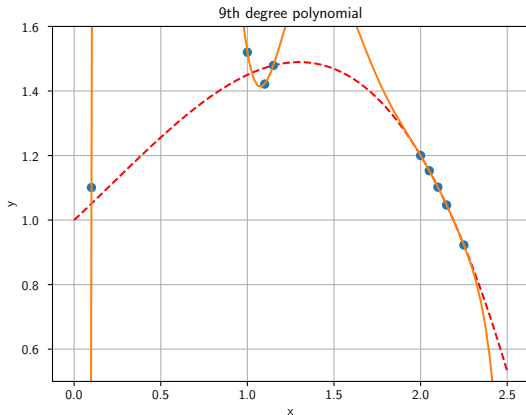
$$y = w_3 \cdot x^3 + w_2 \cdot x^2 + w_1 \cdot x + w_0$$



# Example

9<sup>th</sup> degree polynomial → Overfitting!

$$y = \sum_{i=0}^9 w_i \cdot x^i$$



# Problem diagnosis

- Model behavior is **explained** by the shape of the function
- Overfitting, Underfitting are easily diagnosed
- If the input has multiple dimensions  $D$ ?
  - We often have tens or hundreds of features
  - Images and signals: Several thousands of input dimensions

# Bike Sharing Problem

- Predict Bike rentals per hour in California
- We have 11 features
  - e.g., month, hour, temperature, humidity, windspeed
- We fit a Neural Network  $y = \hat{f}(\mathbf{x})$
- How to make a plot like before?
  - Feature Effect methods

# Feature effect methods

- High-dimensional input space  $\mathbf{x} \in \mathbb{R}^D$ 
  - $x_s \rightarrow$  feature of interest
  - $\mathbf{x}_c \rightarrow$  other features
- How do we isolate the effect of  $x_s$ ?

# Partial Dependence Plots (PDP)

- Proposed by J. Friedman on 2001<sup>4</sup> and is the marginal **effect** of a feature to the model output:

$$f_s(x_s) = E_{X_c} \left[ \hat{f}(x_s, X_c) \right]$$

- Computation:

$$\hat{f}_s(x_s) = \frac{1}{n} \sum_{i=1}^n \hat{f}(x_s, \mathbf{x}_c^{(i)})$$

---

<sup>4</sup>J. Friedman. "Greedy function approximation: A gradient boosting machine." Annals of statistics (2001): 1189-1232

# Partial Dependence Plots (PDP)

## Bike sharing Dataset:

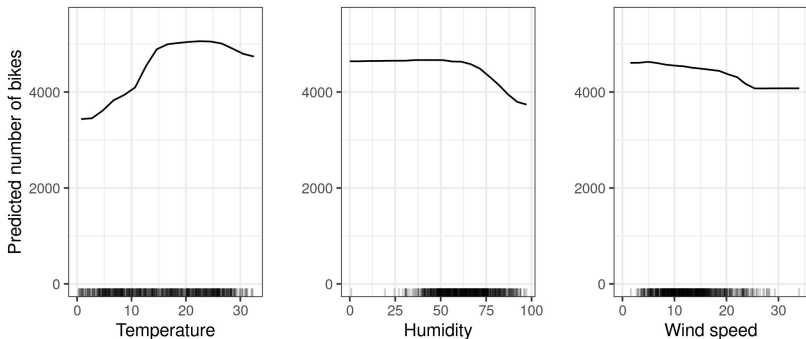


Figure: C. Molnar, IML book

<sup>4</sup>J. Friedman. "Greedy function approximation: A gradient boosting machine." *Annals of statistics* (2001): 1189-1232



# Issues with PDPs

- The marginal distribution ignores correlated features!
- To compute the effect of temperature = 33 degrees it will (also) use an instance with month = January

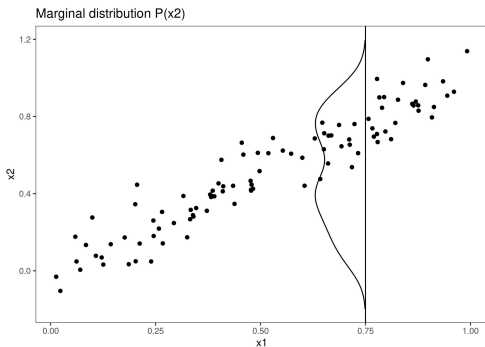



Figure: C. Molnar, IML book

# Accumulated Local Effects (ALE)<sup>5</sup>

- Resolves problems that result from the feature correlation by computing differences over a (small) window
- Definition:  $f(x_s) = \int_{x_{min}}^{x_s} \mathbb{E}_{\mathbf{x}_c|z} \left[ \frac{\partial f}{\partial x_s}(z, \mathbf{x}_c) \right] dz$

---

<sup>5</sup>D. Apley and J. Zhu. “Visualizing the effects of predictor variables in black box supervised learning models.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 82.4 (2020): 1059-1086. 

# ALE approximation

$$\text{Approximation: } f(x_s) = \underbrace{\sum_{k=1}^{k_x} \frac{1}{|\mathcal{S}_k|}}_{\text{bin effect}} \sum_{i: x^i \in \mathcal{S}_k} \underbrace{[f(z_k, \mathbf{x}_c^i) - f(z_{k-1}, \mathbf{x}_c^i)]}_{\text{point effect}}$$

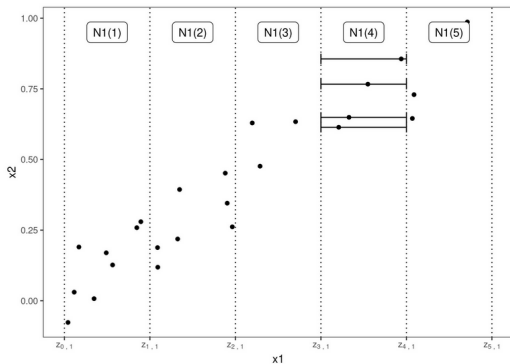


Figure: C. Molnar, IML book

# ALE plots - examples

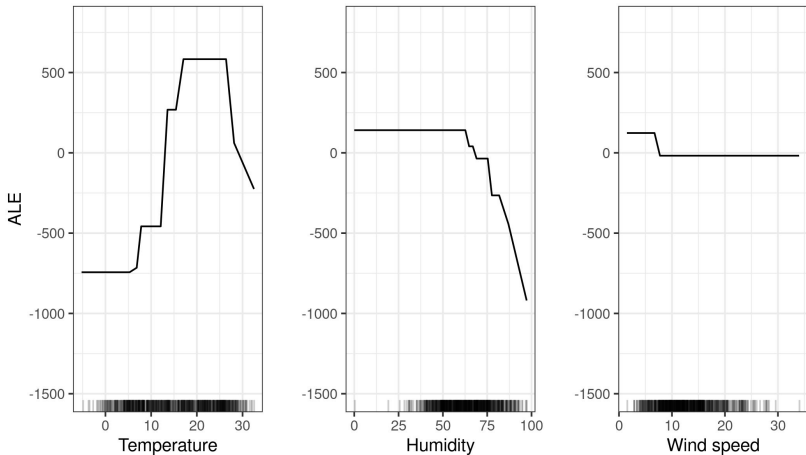


Figure: C. Molnar, IML book

# Our work

- Differential Accumulated Local Effects (DALE)
  - Asian Conference in Machine Learning (ACML 2022)
  - work done with my supervisors: Christos Diou, Theodore Dalamagas
- More efficient and accurate extension of ALE
- Works only with differential models (like Neural Networks)
- <https://arxiv.org/abs/2210.04542>