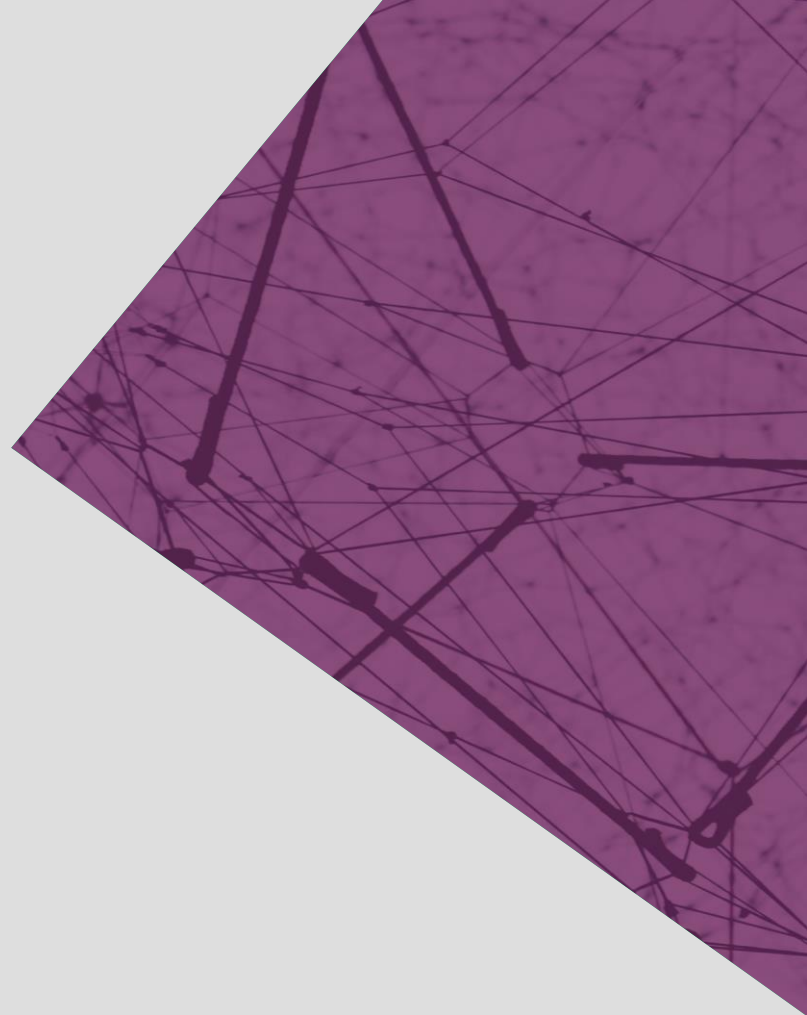# The Alan Turing Institute

# Explainability

**Milestone 5: Trade-offs and Interactions with other verticals in Trustworthy AI**

Roseline Polle

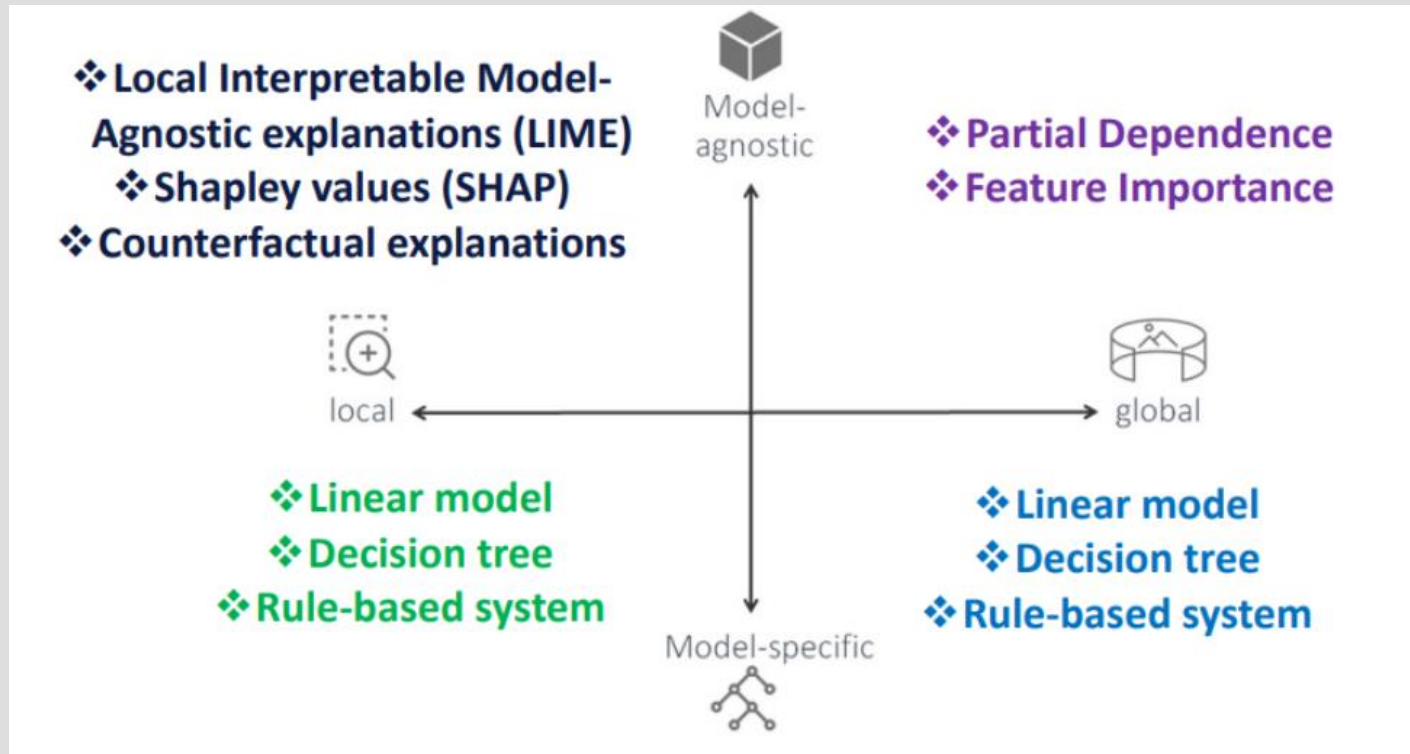Postgraduate, UCL

roseline.polle.19@ucl.ac.uk

# 1. Intro on Explainability

# Interpretability/Explainability

- **Interpretability**: "the degree to which a human can understand the cause of a decision" [Milller, 2018]

- **Explainability**: the degree to which the inner mechanics of an algorithm are understood by a human
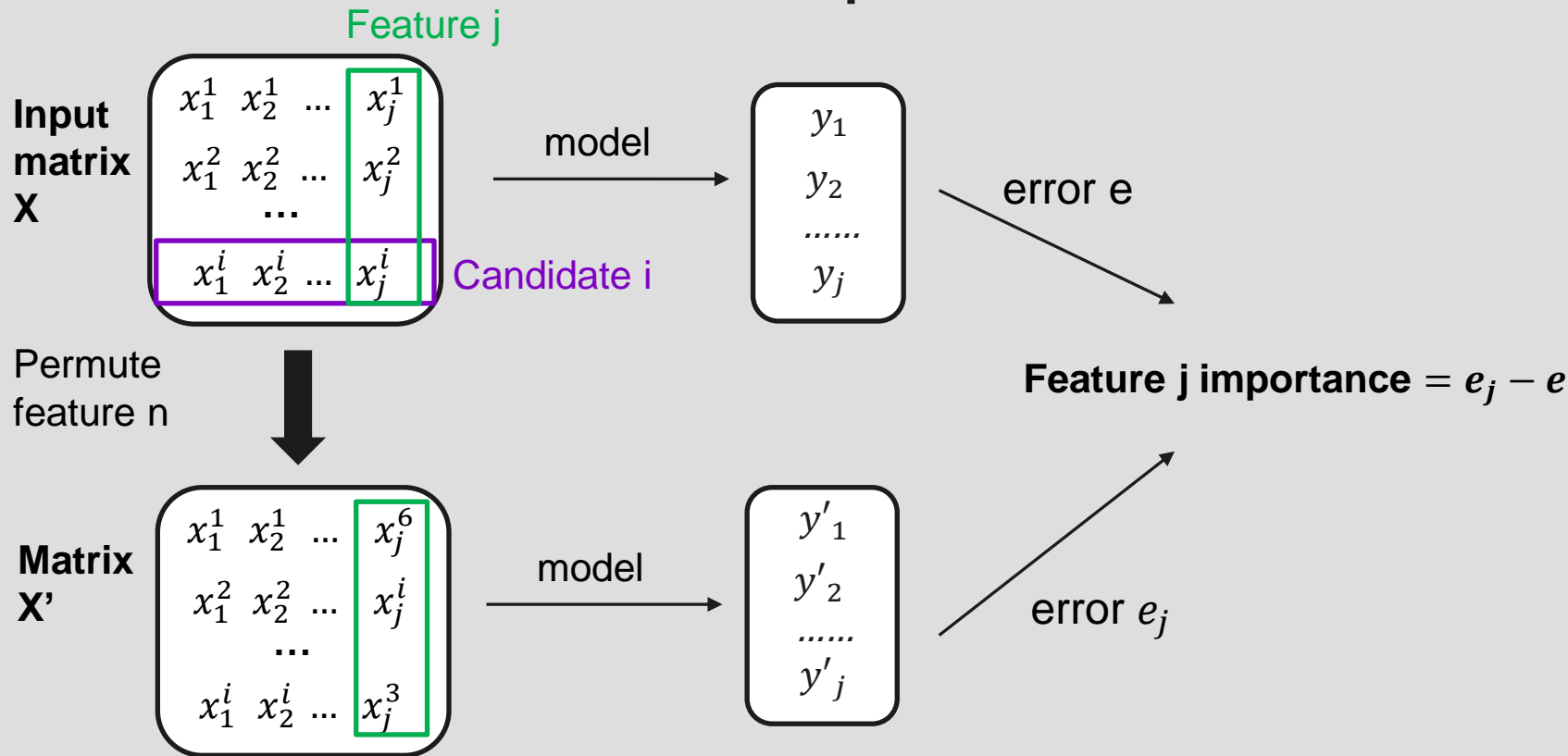
# Types of Explainability



[Koshiyama et al.,2021]

# Feature importance

- Looks to assign a score to each feature relative to its importance in the prediction.

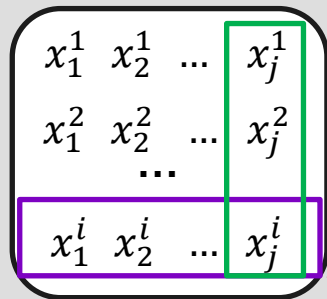- Simple example in linear regression:

$$y = w_1 * x_1 + w_2 * x_2 + \ldots + w_n * x_n$$

Importance of feature 1 in outcome y

# Permutation feature importance

# LIME

**Input matrix X**

$$\begin{matrix} x_1^1 & x_2^1 & ... & x_j^1 \\ x_1^2 & x_2^2 & ... & x_j^2 \\ & ... & & \\ x_1^i & x_2^i & ... & x_j^i \end{matrix}$$

Candidate i

Feature j

$N(.., 1)$

sample 1
sample 2
......
sample 5000

model

$y_1$
$y_2$
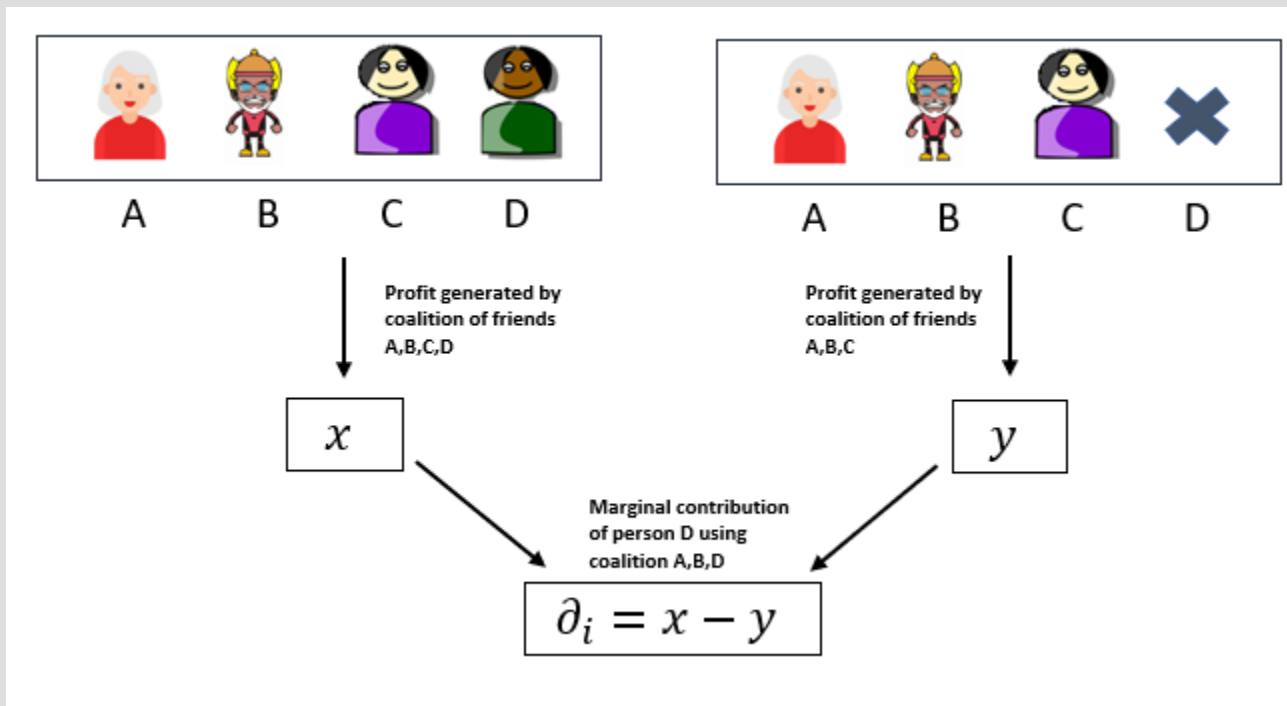......
$y_{5000}$

Weight based on distance

$w_1$
$w_2$
......
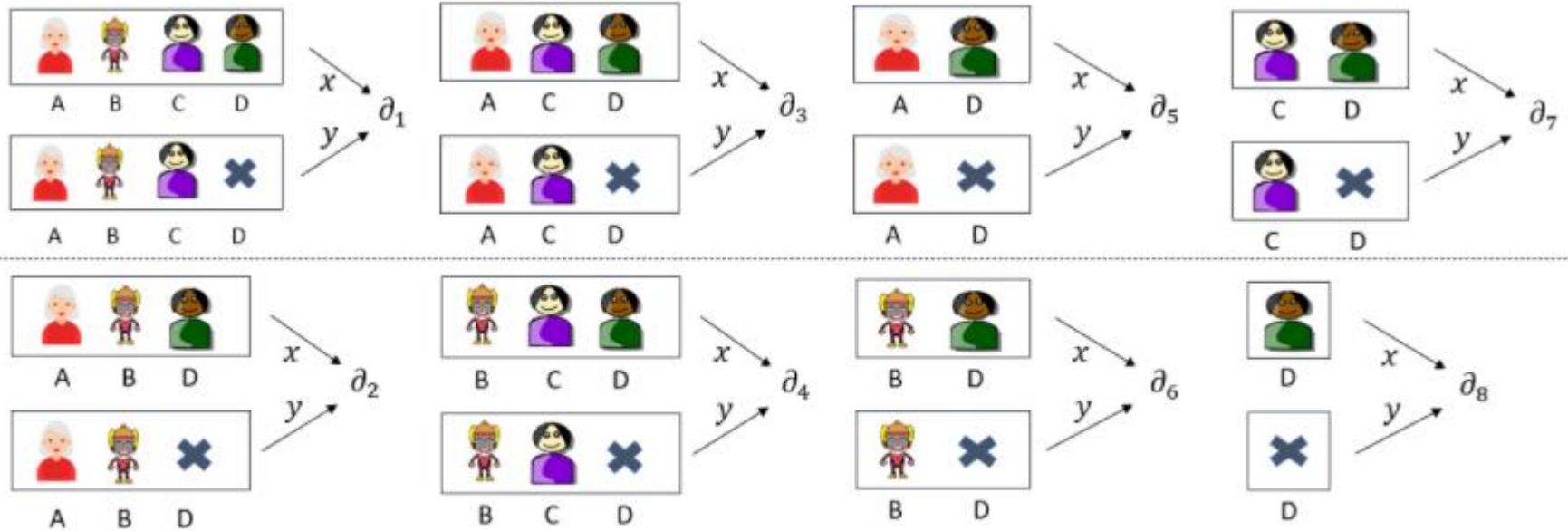$w_{5000}$

Feature importance on output

**Explain results for candidate i:**
- Sample 5,000 times feature vector i using a normal distribution
- Predict output
- Assign weight base on distance
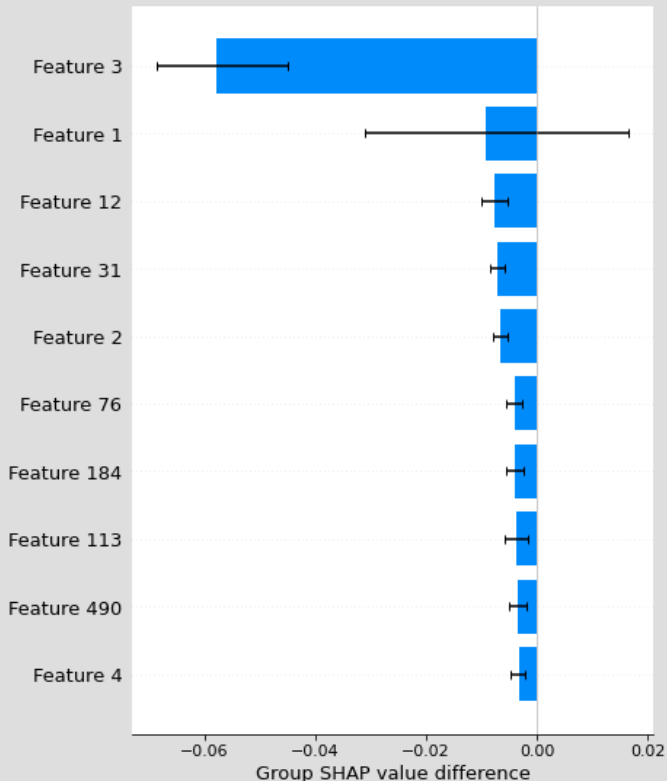- Feature selection (Lasso)

# SHAPLEY Values

# SHAPLEY Values



The shapley value for person D is therefore: $\Phi_D = \frac{\delta_1 + \delta_2 + \delta_3 + \delta_4 + \delta_5 + \delta_6 + \delta_7 + \delta_8}{8}$

# 2. Interactions with Fairness

# Questions one can answer

- Are the most influential factors reasonable? Are they proxy for a protected characteristics?

- Is the model relying too much on one feature?

- Are they the influential factors the same across different groups ?

# Adaptation to fairness

- Instead of explaining output → explain fairness metric

- Example: effect of permutation importance on Disparate Impact metric.

- Answers the question: what are the features most responsible for the observed bias (if any) ?

# Further readings

- **Interpretable Machine Learning.** *A Guide for Making Black Box Models Explainable* (https://christophm.github.io/interpretable-ml-book/)