# COMP6211I:
# Trustworthy Machine Learning

**Interpretability (XAI) part 2**

# Introduction

- Model Specific vs. Model Agnostic

  Can it explain a particular model or many models?

- Global Methods vs. Local Methods

  Does it explain a particular sample or entire model?

- Pre-Model vs. In-Model vs. Post-Model

  When does it occur?

- Surrogate Methods vs. Visualization Methods

  Does it work separately from the model, or does it visualize the model?

**The categories are non-exclusive. There is no universally accepted taxonomy of XAI techniques!**

# Introduction

- **Model Specific vs. Model Agnostic**

Model-specific interpretation methods are based on the parameters of the individual models.

Model Agnostic methods are mainly applicable in post-hoc analysis and not limited to specified model architecture.

# Introduction

- Global Methods vs. Local Methods

Global methods concentrate on the inside of a model by exploiting the overall knowledge about the model, the training, and the associated data.

Local interpretable methods are applicable to a single outcome of the model. This can be done by designing methods that can explain the reason for a particular prediction or outcome.

# Introduction

**LIME**

- Title: "Why Should I Trust You?" Explaining the Predictions of Any Classifier
- Conference: KDD2016
- Authors: Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin (University Of Washington)

**SHAP**

- Title: A unified approach to interpreting model predictions
- Conference: NIPS2017
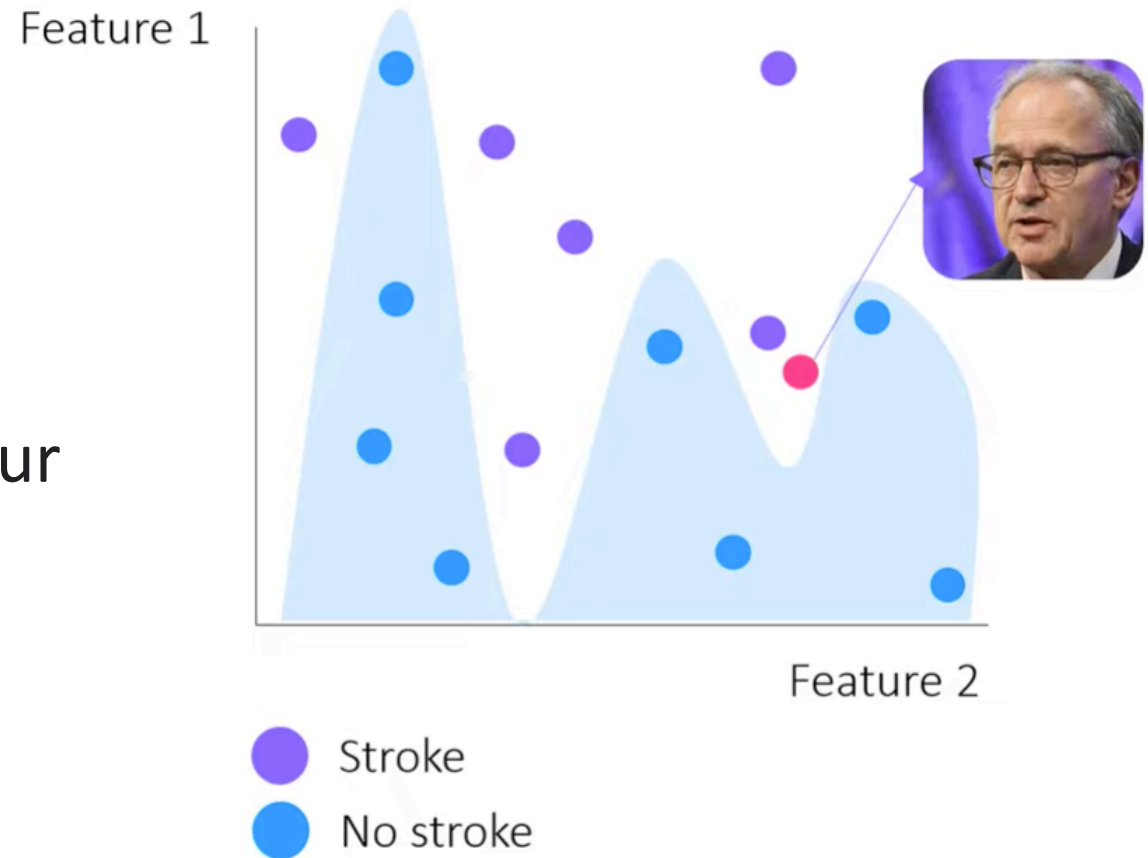- Authors: Scott M. Lundberg, Su-In Lee (University Of Washington)

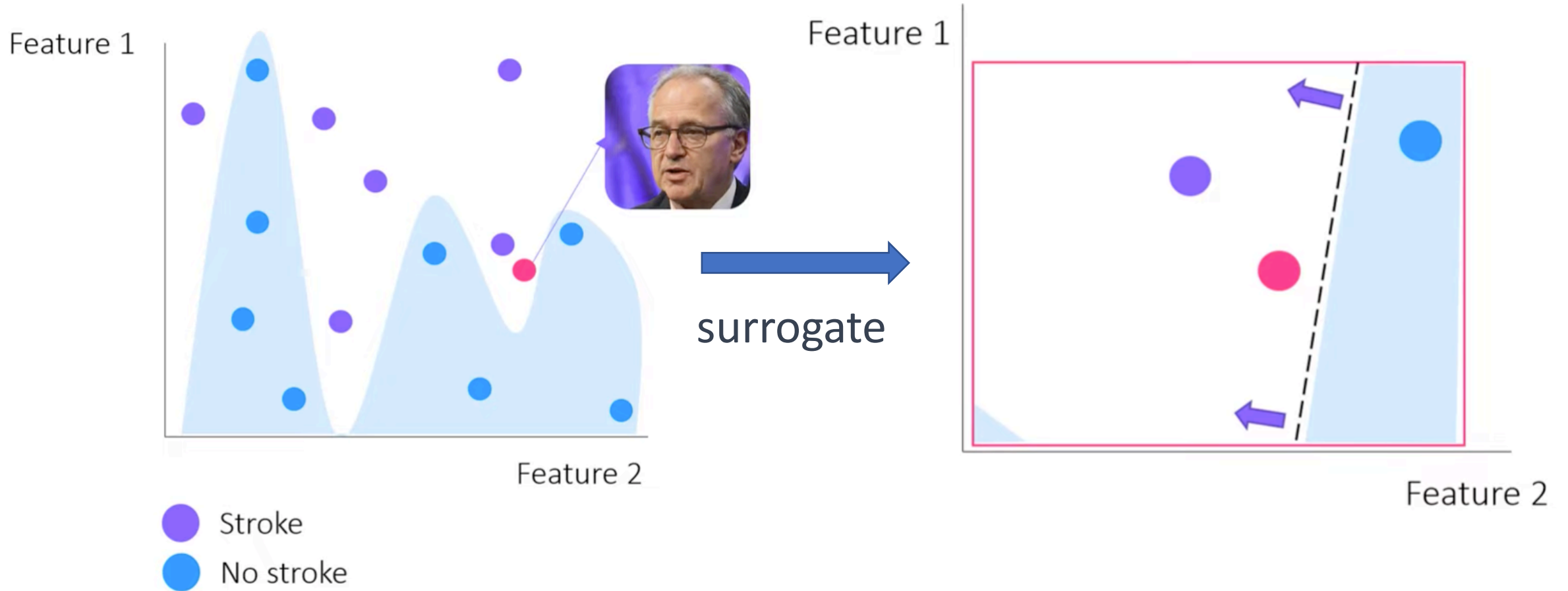# LIME: **L**ocal **i**nterpretable **m**odel-agnostic **e**xplanations

Task: Stroke Prediction

Feature 1: age
Feature 2: body mass index

How could we explain to him why our model outputs stroke?



Feature 1

Feature 2

● Stroke

● No stroke

# LIME: **L**ocal **i**nterpretable **m**odel-agnostic **e**xplanations



surrogate

Feature 1

Feature 2

Stroke

No stroke

# LIME: Local interpretable model-agnostic explanations

- Works on any black-box model

- Model internals are "hidden"

- Works with many data types

- Using prior knowledge we can validate the explanations and create trust

- Explanations are locally faithful, but not necessarily globally

# The Math in LIME

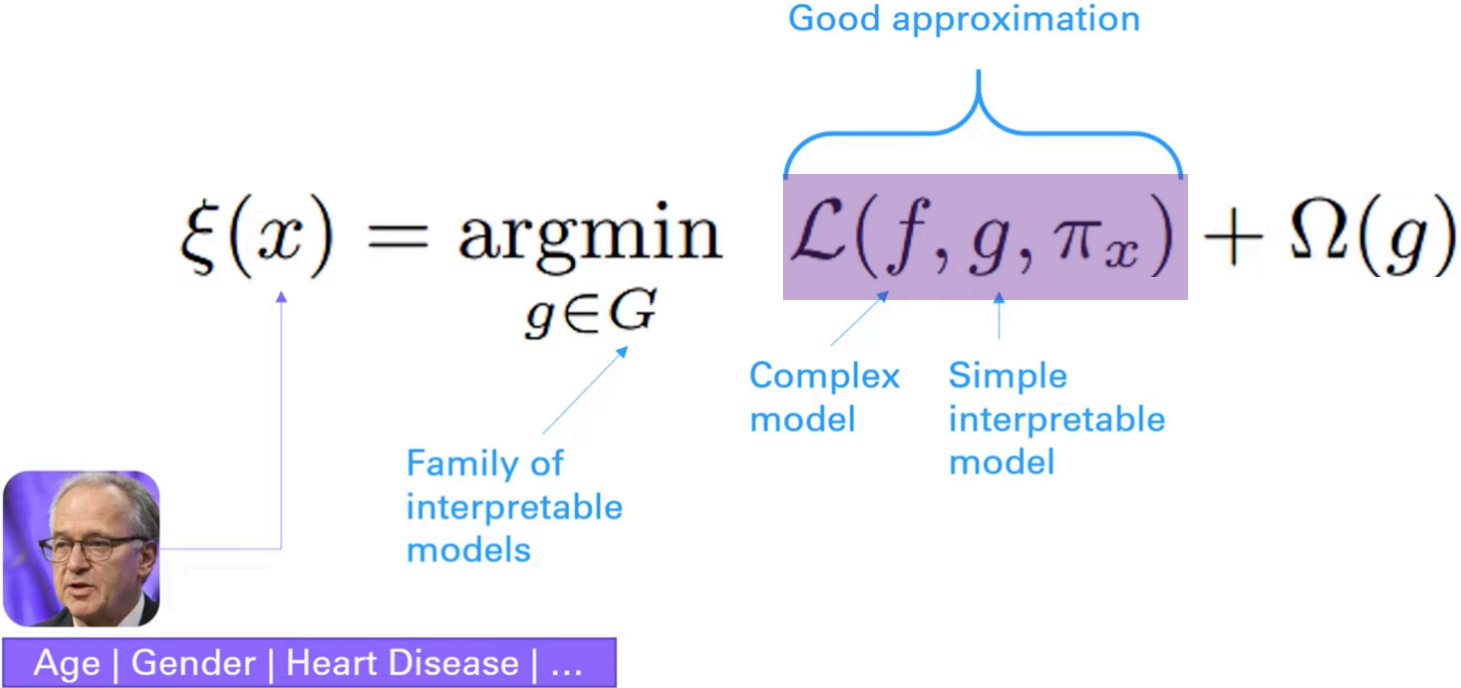$$\xi(x) = \operatorname*{argmin}_{g \in G} \; \mathcal{L}(f, g, \pi_x) + \Omega(g)$$
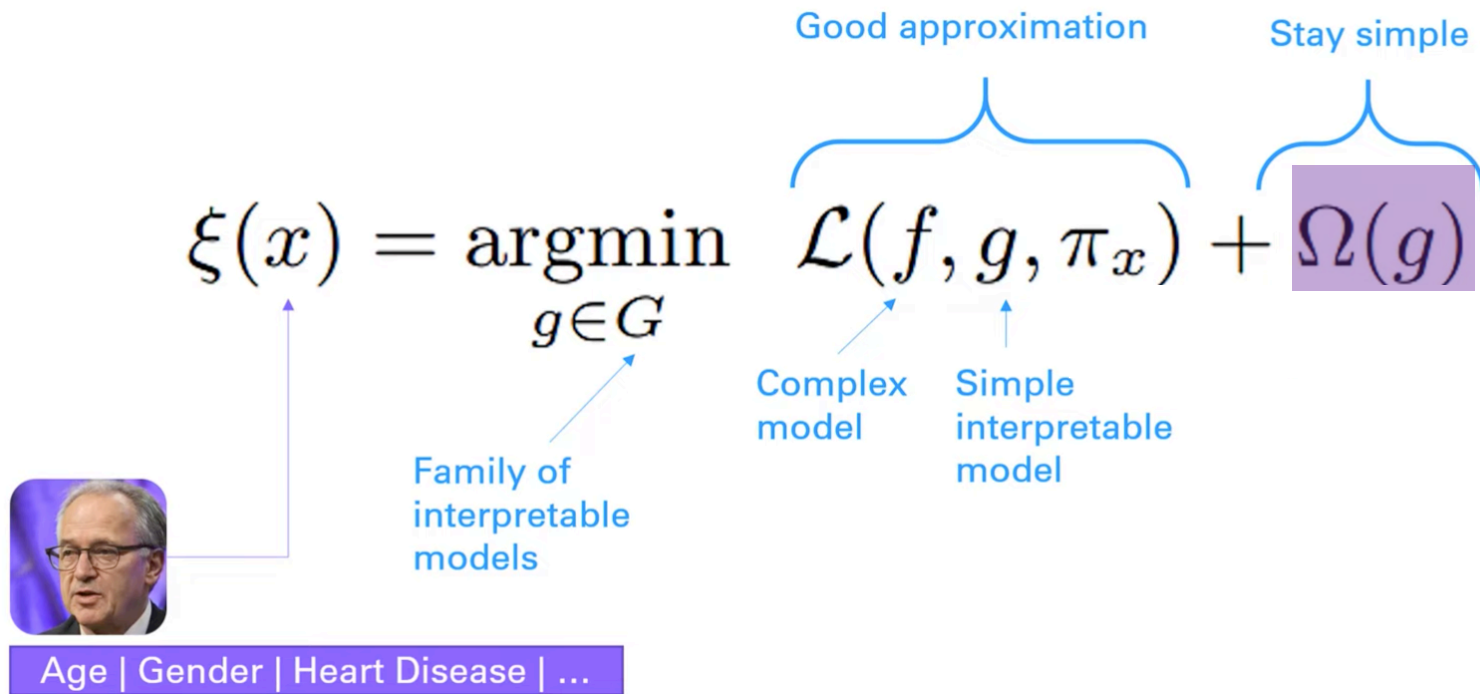
Family of interpretable models

Complex model

Simple interpretable model

Age | Gender | Heart Disease | …

# The Math in LIME



$$\xi(x) = \underset{g \in G}{\text{argmin}} \quad \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

Good approximation

Family of interpretable models

Complex model

Simple interpretable model

Age | Gender | Heart Disease | …

# The Math in LIME



$$\xi(x) = \underset{g \in G}{\text{argmin}} \; \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

Good approximation

Stay simple

Complex model

Simple interpretable model

Family of interpretable models

Age | Gender | Heart Disease | ...

# How to train the surrogate

$$\xi(x) = \underset{g \in G}{\text{argmin}} \; \boxed{\mathcal{L}(f, g, \pi_x)} + \Omega(g)$$

# How to train the surrogate

$$\xi(x) = \underset{g \in G}{\mathrm{argmin}} \ \boxed{\mathcal{L}(f, g, \pi_x)} + \Omega(g)$$

# How to train the surrogate

$$\xi(x) = \underset{g \in G}{\mathrm{argmin}} \; \boxed{\mathcal{L}(f, g, \pi_x)} + \Omega(g)$$
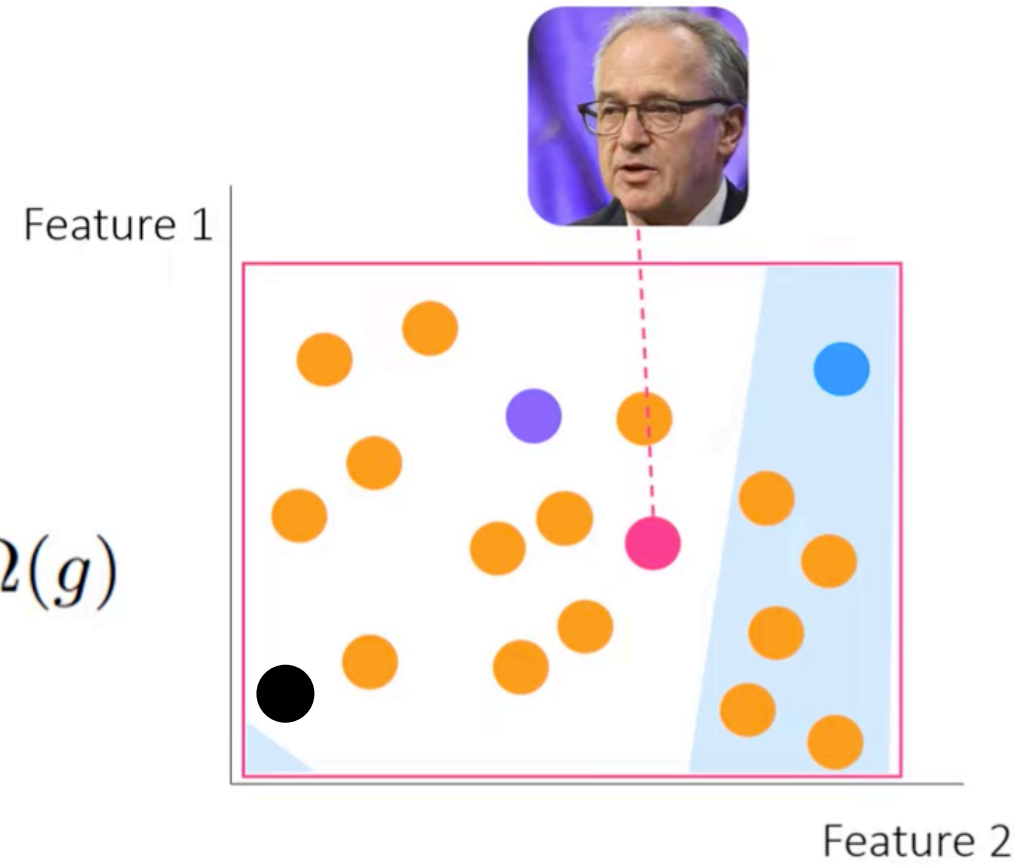
New dataset
Labels: Prediction of complex model
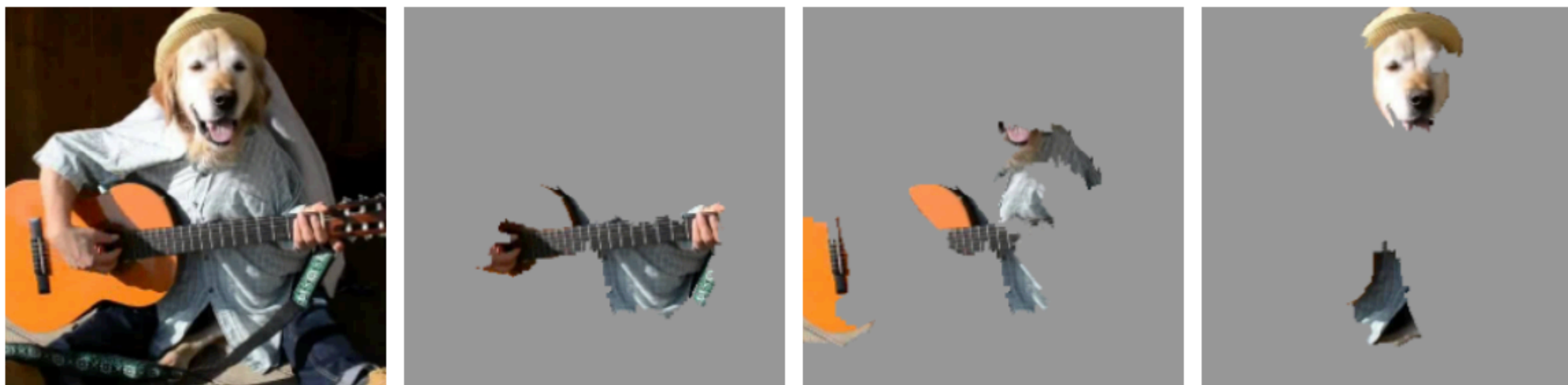Features: Newly generated datapoints

Feature 1

Feature 2

# How to train the surrogate

$$\xi(x) = \operatorname*{argmin}_{g \in G} \boxed{\mathcal{L}(f, g, \pi_x)} + \Omega(g)$$

$$\mathcal{L}(f, g, \pi_x) = \sum_{z, z' \in \mathcal{Z}} \pi_x(z) \left(f(z) - g(z')\right)^2$$



Feature 1

Feature 2

14

# Example for LIME



(a) Original Image  (b) Explaining *Electric guitar*  (c) Explaining *Acoustic guitar*  (d) Explaining *Labrador*

Figure 4: Explaining an image classification prediction made by Google's Inception neural network. The top 3 classes predicted are "Electric Guitar" ($p = 0.32$), "Acoustic guitar" ($p = 0.24$) and "Labrador" ($p = 0.21$)

# Example for LIME



Prediction probabilities

| | |
|---|---|
| atheism | 0.58 |
| christian | 0.42 |

atheism          christian

Posting
0.15
Host
0.14
NNTP
0.11
edu
0.04
have
0.01
There
0.01

**Text with highlighted words**
From: johnchad@triton.unm.edu (jchadwic)
Subject: Another request for Darwin Fish
Organization: University of New Mexico, Albuquerque
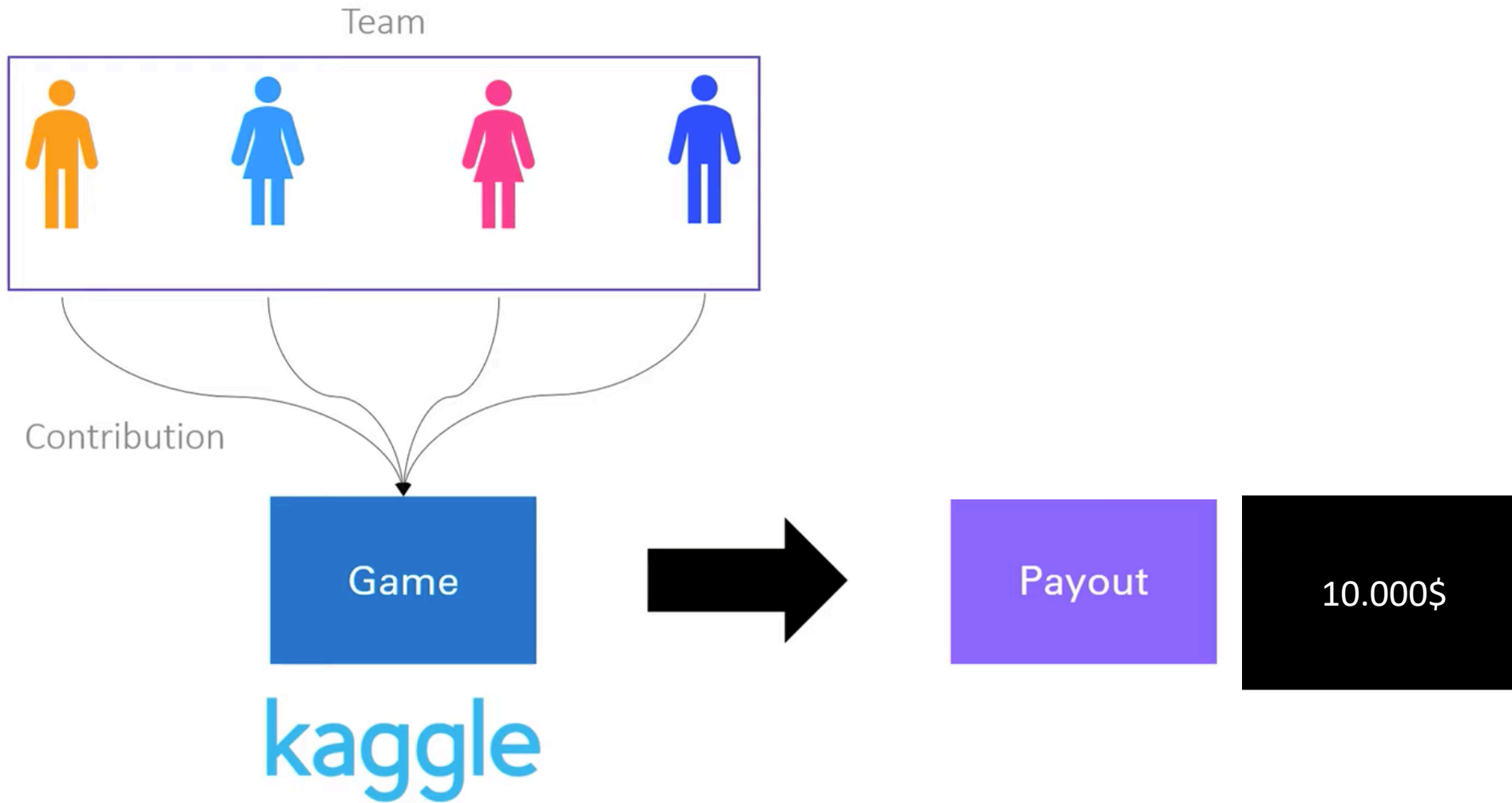Lines: 11
NNTP-Posting-Host: triton.unm.edu

Hello Gang,

There have been some notes recently asking where to obtain the
DARWIN fish.
This is the same question I have and I have not seen an answer on
the
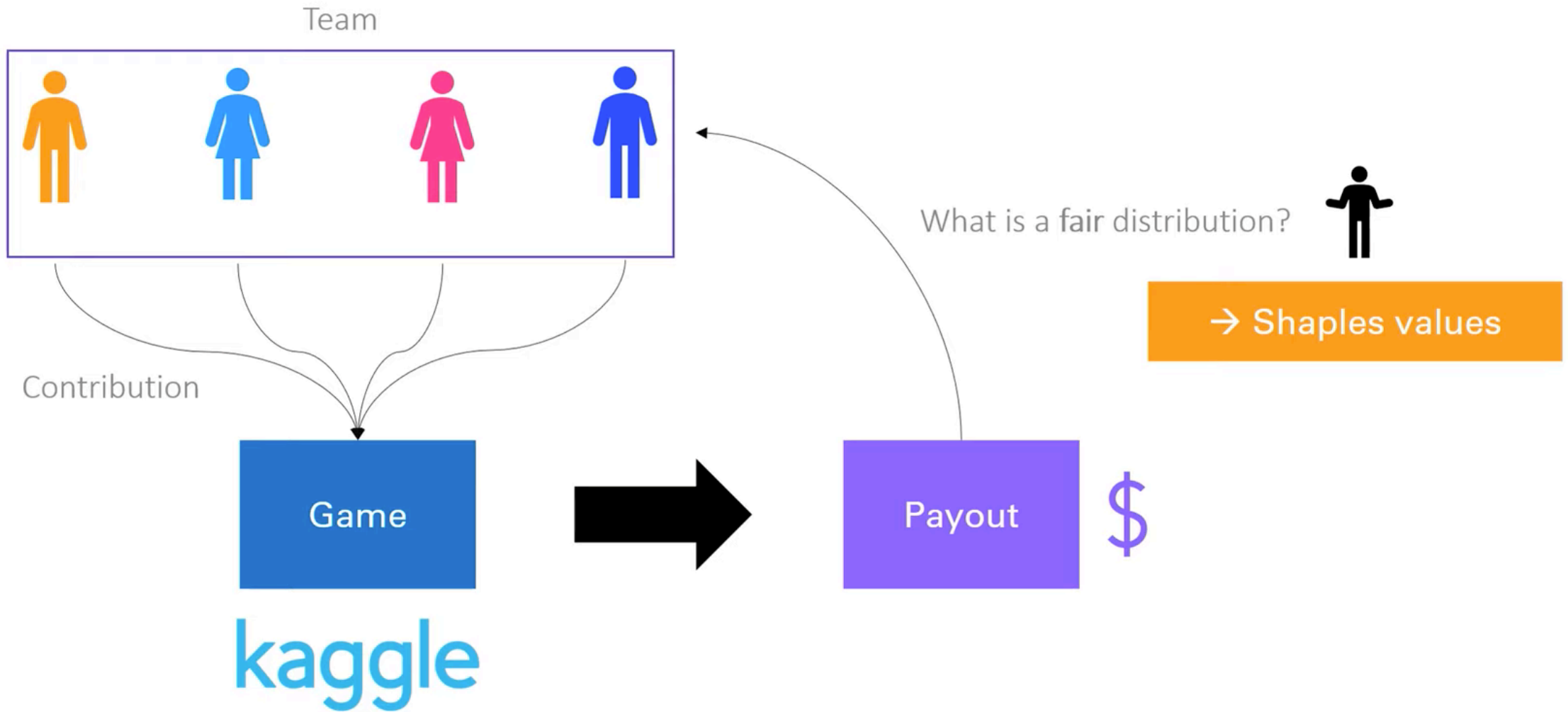net. If anyone has a contact please post on the net or email me.

# SHAP: **SH**apley **A**dditive ex**P**lanations

**Cooperative Game Theory**

# SHAP

# SHAP

# Shapley Values



Payout
10.000 $

kaggle

# Shapley Values

3.000$

kaggle

Domain expert

# Shapley Values



Domain expert

Payout 10.000 $

7.000$

# Shapley Values

# Shapley Values

# Shapley Values

# Shapley Values



Marginal
Contribution

Payout
10.000 $

kaggle

# Shapley Values



3.000$    1.500$    2.500$    3.000$

Payout
10.000 $

kaggle

# Calculating shapley value

Black model    Input datapoint

Age  $\phi_i(f, x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$

Shapley value
for feature i

x =  | Age = 56 | Gender = F | Body Mass Index = 30 | Hear disease = yes | ... |

# Calculating shapley value

Black model · Input datapoint

$$\phi_i(f, x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$$

Age

Shapley value
for feature i

subset

Simplified
data input

x = | Age = 56 | Gender = F | Body Mass Index = 30 | Hear disease = yes | ... |

# Calculating shapley value

Black model    Input datapoint

Age

$$\phi_i(f, x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$$

Shapley value
for feature i

subset

Simplified
data input

x = | Age = 56 | Gender = F | Body Mass Index = 30 | Hear disease = yes | … |

# Calculating shapley value

Black model    Input datapoint

Age = 56 | Body Mass Index = 30

Age

$$\phi_i(f, x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$$

Shapley value
for feature i

Body Mass Index = 30

subset

Simplified
data input

x = | Age = 56 | Gender = F | Body Mass Index = 30 | Hear disease = yes | ... |

31

# Calculating shapley value

Black model    Input datapoint

Age = 56 | Body Mass Index = 30

Age $\phi_i(f, x) = \sum_{z' \subseteq x'} \dfrac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$

70% Stroke

Shapley value
for feature i

subset

Body Mass Index = 30

Simplified
data input

x = | Age = 56 | Gender = F | Body Mass Index = 30 | Hear disease = yes | ... |

# Calculating shapley value

Black model   Input datapoint

Age = 56 | Body Mass Index = 30

10% Stroke

Age
$$\phi_i(f, x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$$

Shapley value
for feature i

subset

Simplified
data input

Body Mass Index = 30

x =  | Age = 56 | Gender = F | Body Mass Index = 30 | Hear disease = yes | ... |

# Calculating shapley value

Black model   Input datapoint

Age = 56 | Body Mass Index = 30

Age  $\phi_i(f, x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$

Shapley value
for feature i

subset

Simplified
data input

Body Mass Index = 30

x =  Age = 56 | Gender = F | Body Mass Index = 30 | Hear disease = yes | ...

# Calculating shapley value

Black model   Input datapoint

Shapley value
for feature i

$$\phi_i(f, x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$$

Age

subset

Simplified
data input

Weighting

Contribution

x =  | Age = 56 | Gender = F | Body Mass Index = 30 | Hear disease = yes | ... |

# Calculating shapley value



Black model   Input datapoint

$$\phi_i(f, x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$$
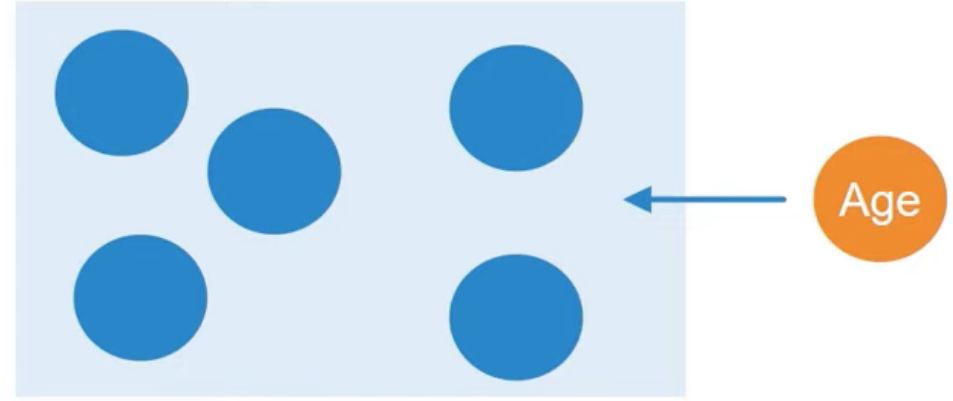
Age

Shapley value
for feature i

subset

Simplified
data input

Weighting

Contribution

x = | Age = 56 | Gender = F | Body Mass Index = 30 | Hear disease = yes | ... |

# Calculating shapley value

Black model    Input datapoint

Age

$$\phi_i(f, x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$$

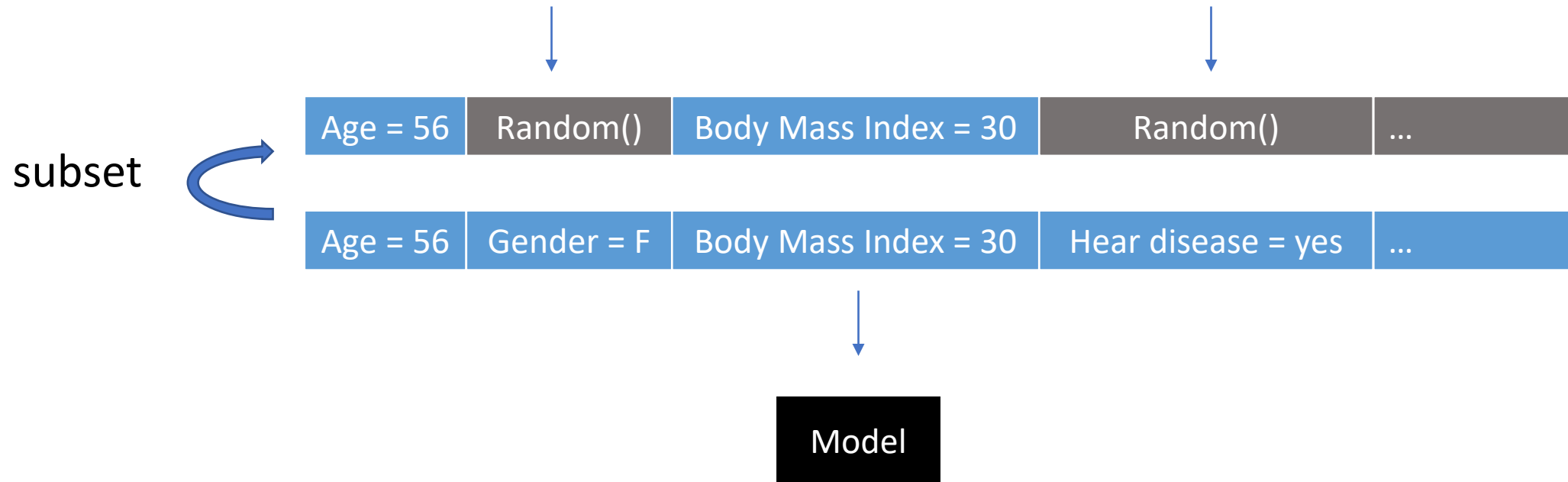Shapley value
for feature i

subset

Simplified
data input

Weighting

Contribution

x = | Age = 56 | Gender = F | Body Mass Index = 30 | Hear disease = yes | … |

# Calculating shapley value

Black model    Input datapoint

Age = 56    Body Mass Index = 30

Age

$$\phi_i(f,x) = \sum_{z' \subseteq x'} \frac{|z'|!(M - |z'| - 1)!}{M!} [f_x(z') - f_x(z' \setminus i)]$$

Shapley value
for feature i

subset

Simplified
data input

Body Mass Index = 30

x = | Age = 56 | Gender = F | Body Mass Index = 30 | Hear disease = yes | ... |

# Calculating shapley value
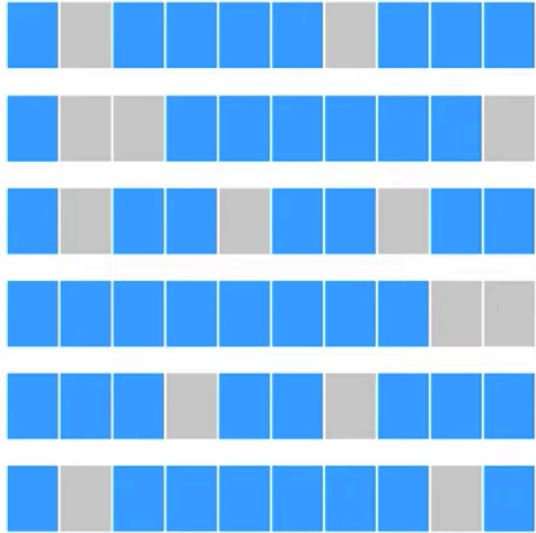
| Age = 56 | Random() | Body Mass Index = 30 | Random() | ... |
|----------|----------|----------------------|----------|-----|

subset

| Age = 56 | Gender = F | Body Mass Index = 30 | Hear disease = yes | ... |
|----------|------------|----------------------|--------------------|-----|

Model

# Calculating shapley value

$2^n$ = total number of subsets of a set

# Calculating shapley value



$2^n$ = total number of subsets of a set

...

$$2^{10} = 1024$$

# Calculating shapley value



Kernel SHAP

$2^{n}$ = total number of subsets of a set

$$Y = x_1\beta_1 + x_2\beta_2 + x_3\beta_3 \cdots$$

$$2^{10} = 1024$$

# Calculating shapley value



$2^n$ = total number of subsets of a set

$2^{10} = 1024$

Kernel SHAP

Tree SHAP

Deep SHAP