

An introduction to explainable artificial intelligence with LIME and SHAP

Aleix Nieto Juscafresa

Treball Final de Grau
Universitat de Barcelona
GRAU DE MATEMÀTIQUES

Dirigit per Dr. Albert Clapés i Dr. Sergio Escalera

29 de juny de 2022



UNIVERSITAT DE
BARCELONA

Facultat de Matemàtiques
i Informàtica

- 1 Introduction
- 2 Machine learning
- 3 Random forest
- 4 Regression
- 5 Explainable artificial intelligence
- 6 Conclusions

Introduction

Today



Confusion with Today's AI Black Box

- Why did you do that?
- Why did you not do that?
- When do you succeed or fail?
- How do I correct an error?

Introduction

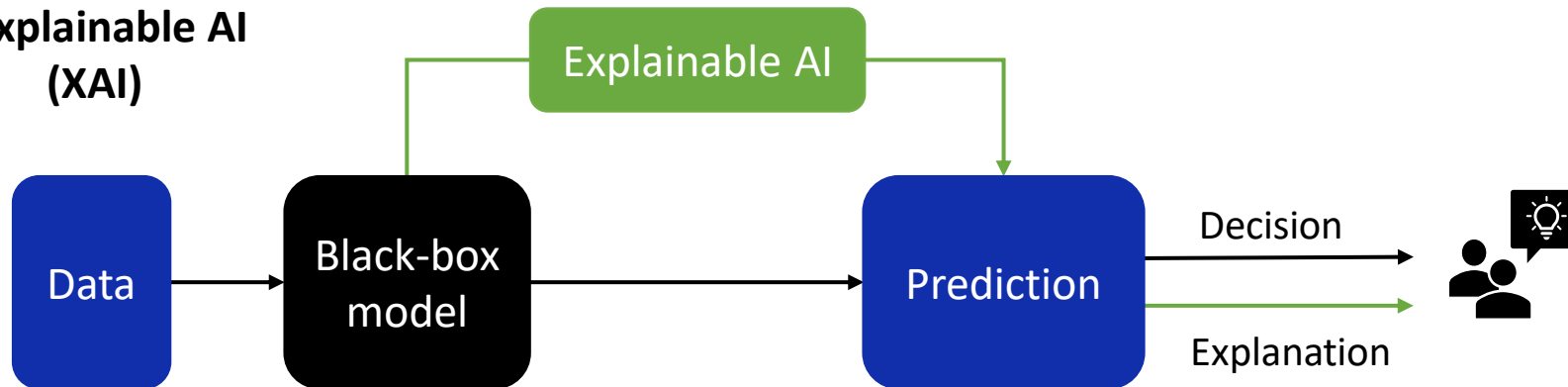
Today



Confusion with Today's AI Black Box

- Why did you do that?
- Why did you not do that?
- When do you succeed or fail?
- How do I correct an error?

Explainable AI (XAI)



Clear & Transparent Decisions

- I understand why
- I understand why not
- I know why you succeed or fail
- I understand, I trust you more

- 1 Introduction
- 2 Machine learning
- 3 Random forest
- 4 Regression
- 5 Explainable artificial intelligence
- 6 Conclusions

Machine learning

Machine learning

Application of artificial intelligence dedicated to the creation of algorithms that allow systems to learn without human intervention.

SUPERVISED
LEARNING

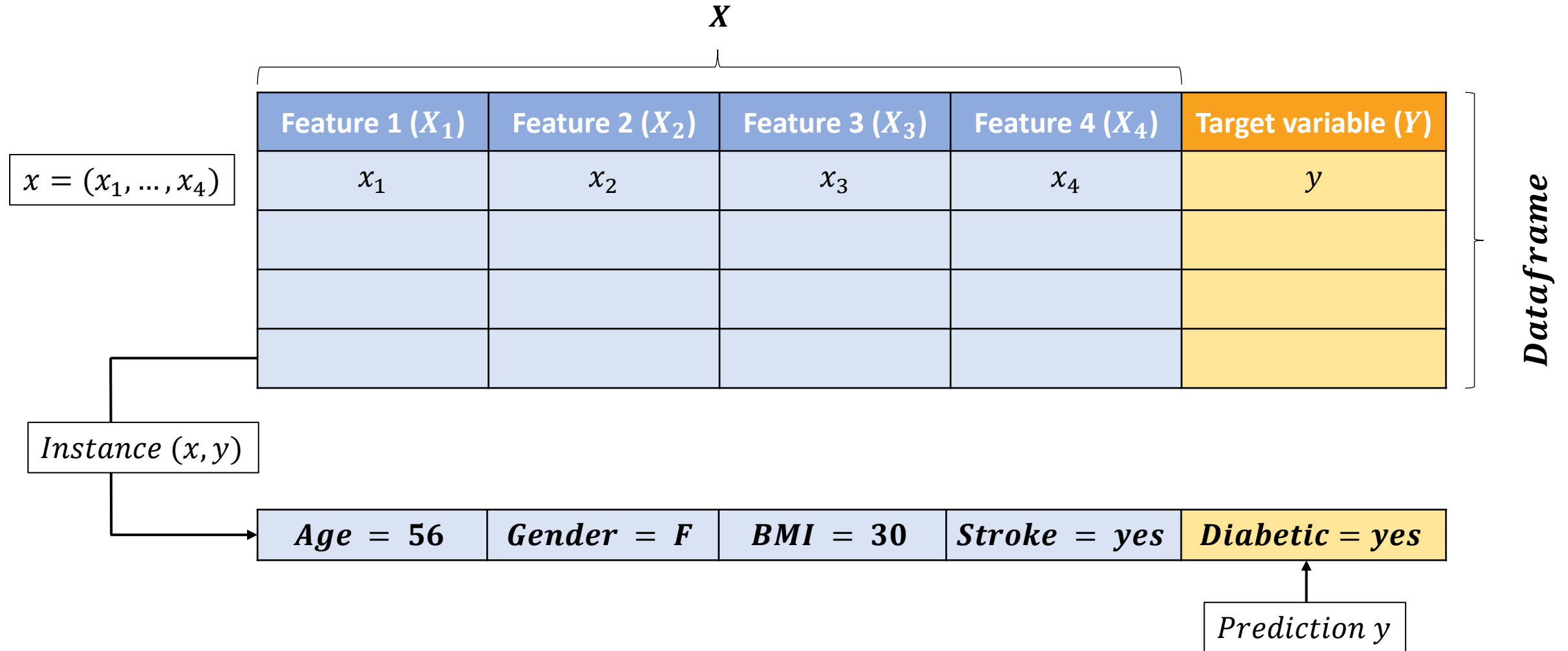
UNSUPERVISED
LEARNING

SEMI-
SUPERVISED
LEARNING

REINFORCEMENT
LEARNING

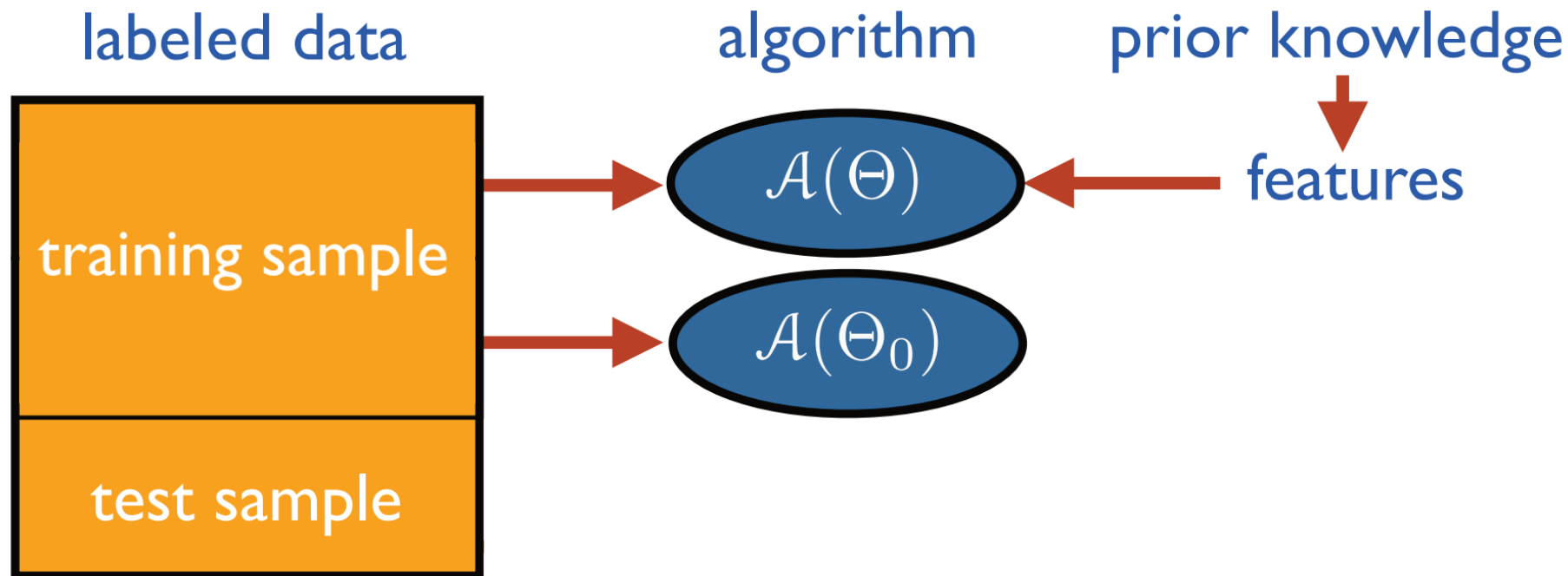
...

Terminology



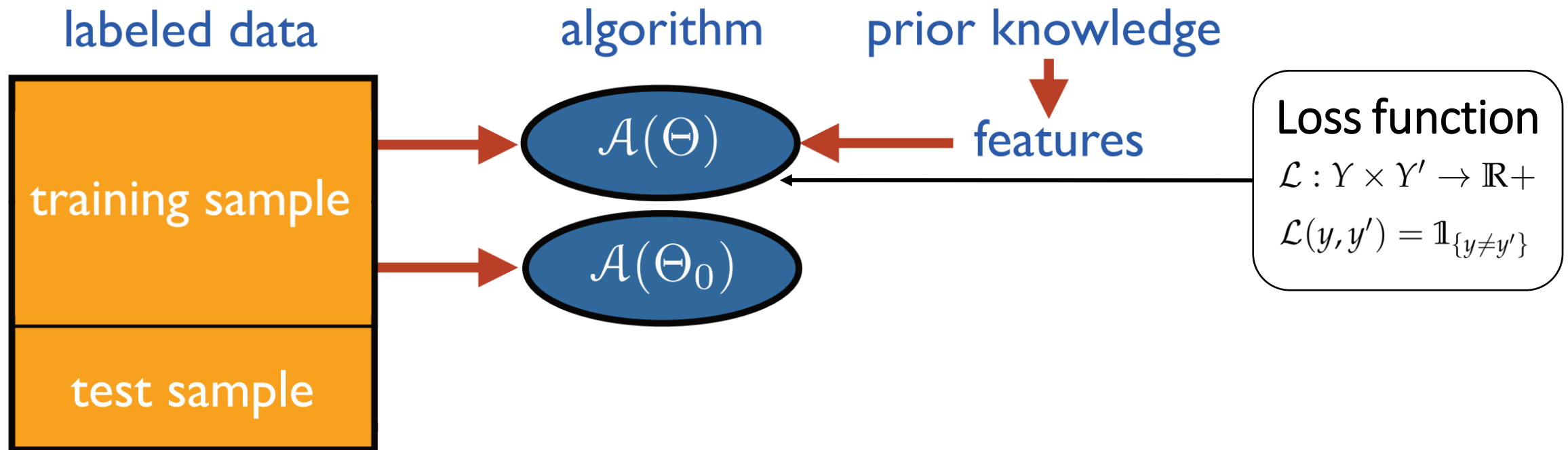


Machine learning model learning stages



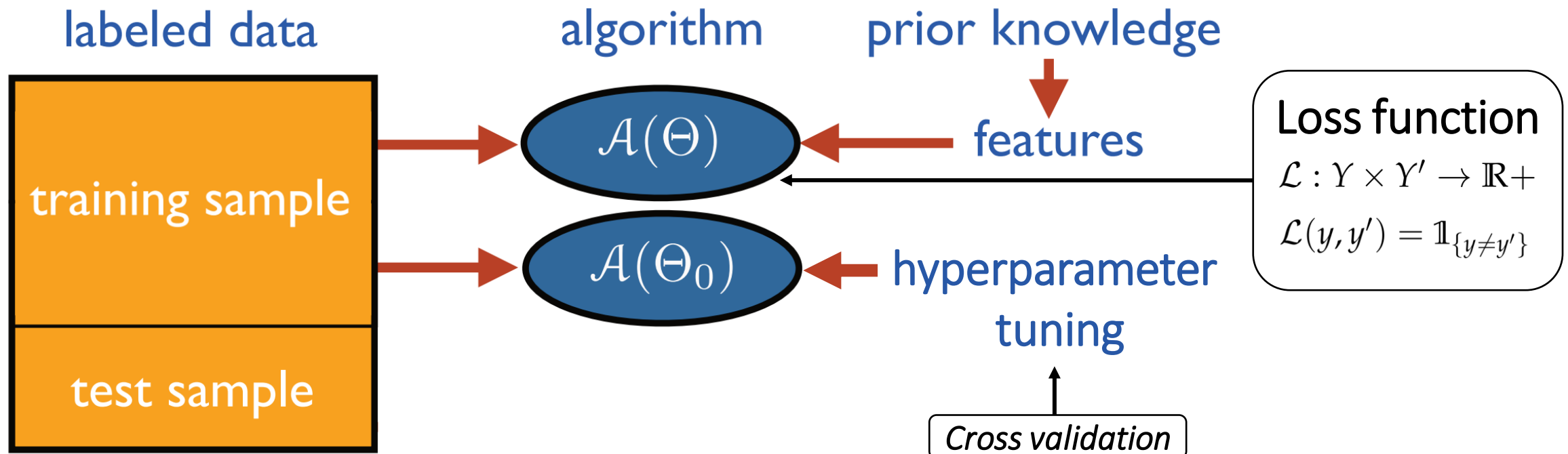


Machine learning model learning stages



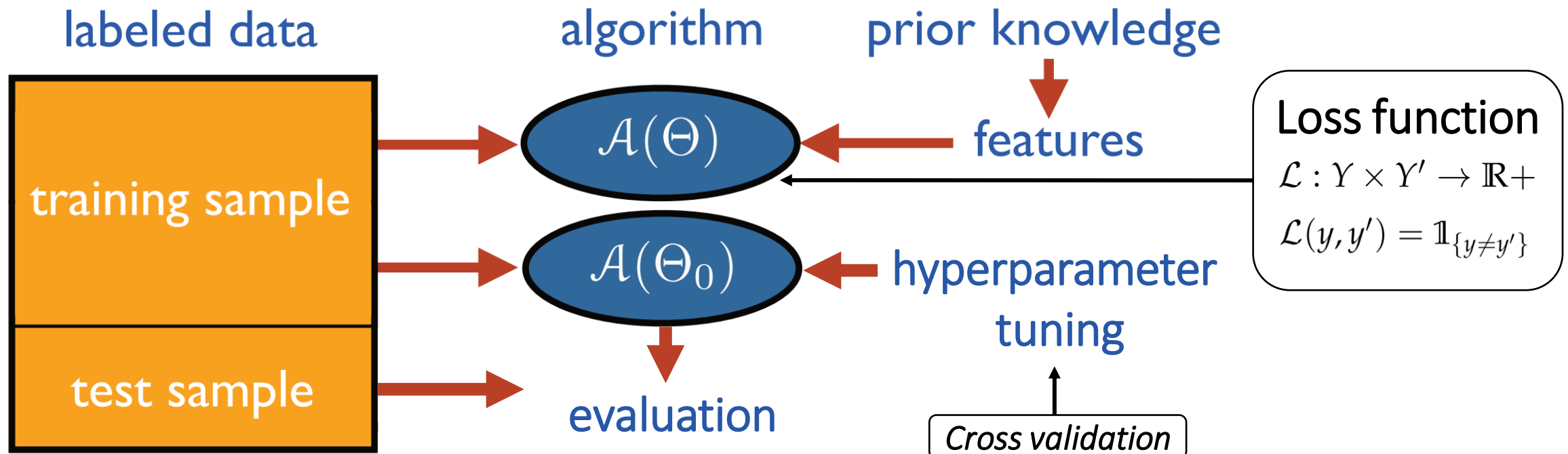


Machine learning model learning stages



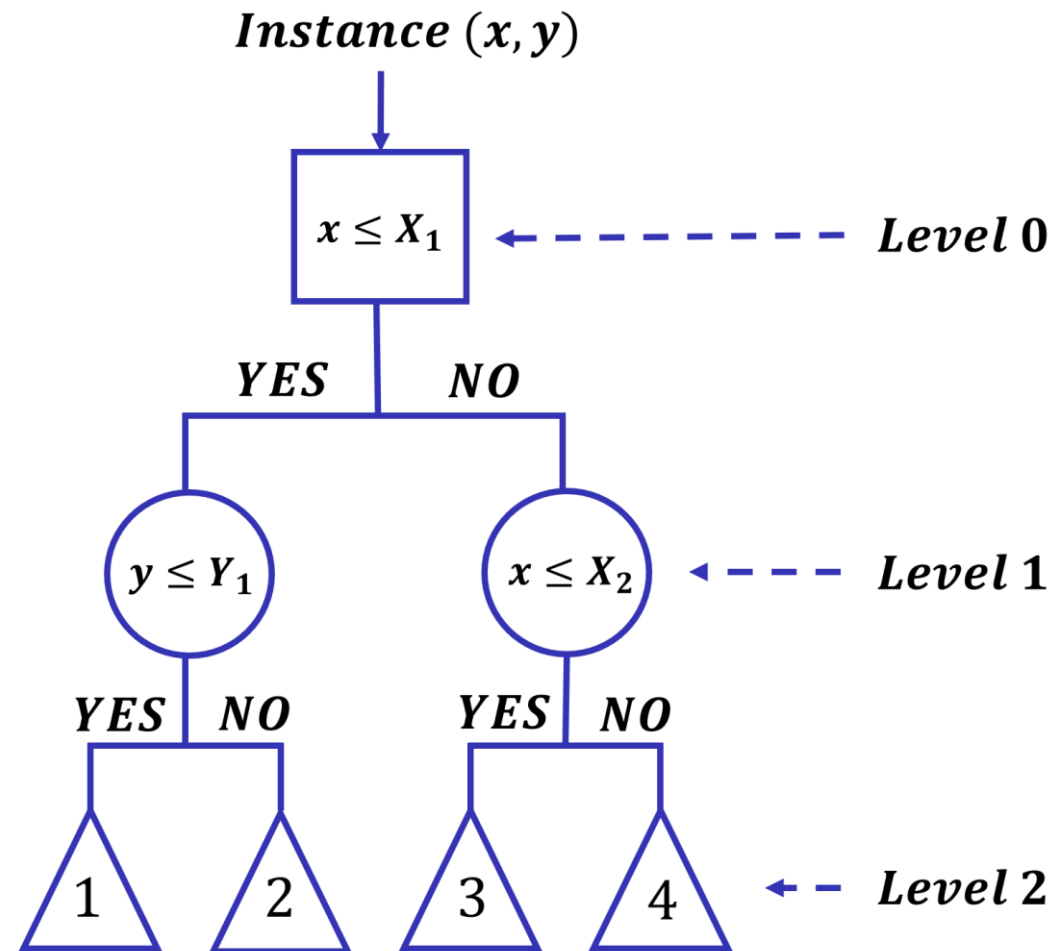


Machine learning model learning stages

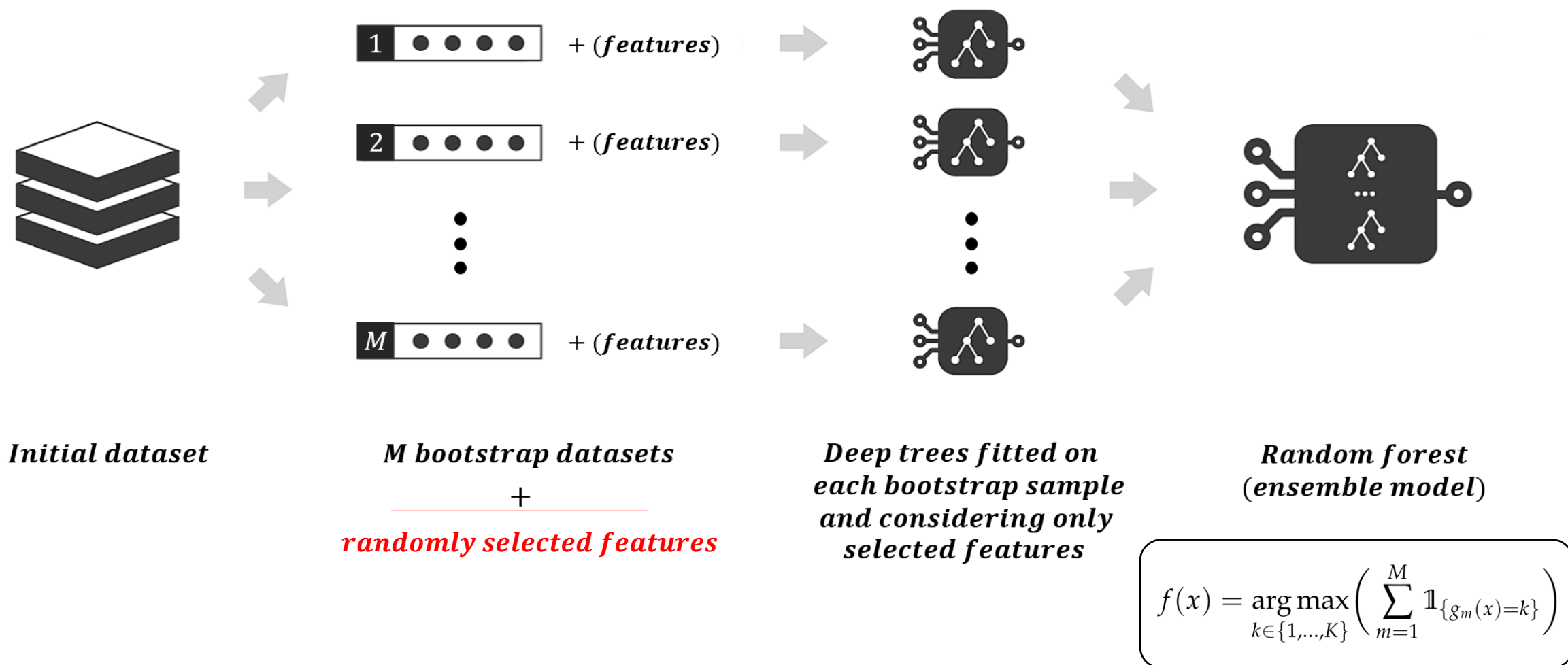


- 1 Introduction
- 2 Machine learning
- 3 Random forest
- 4 Regression
- 5 Explainable artificial intelligence
- 6 Conclusions

Decision tree



Random forest

UNIVERSITAT DE
BARCELONA

- 1 Introduction
- 2 Machine learning
- 3 Random forest
- 4 Regression
- 5 Explainable artificial intelligence
- 6 Conclusions

Linear regression

UNIVERSITAT DE
BARCELONA

$$Y = f(X) + \epsilon = \beta_0 + \sum_{j=1}^p \beta_j X_j + \epsilon$$



Linear regression

$$Y = f(X) + \epsilon = \beta_0 + \sum_{j=1}^p \beta_j X_j + \epsilon$$

Optimisation problem

$$\arg \min_{\beta_0, \dots, \beta_p} \sum_{i=1}^n (y_i - (\beta_0 + \sum_{j=1}^p \beta_j x_{ij}))^2 = \arg \min_{\beta_0, \dots, \beta_p} \sum_{i=1}^n \epsilon_i^2$$



Linear regression

$$Y = f(X) + \epsilon = \beta_0 + \sum_{j=1}^p \beta_j X_j + \epsilon$$

Optimisation problem

$$\arg \min_{\beta_0, \dots, \beta_p} \sum_{i=1}^n (y_i - (\beta_0 + \sum_{j=1}^p \beta_j x_{ij}))^2 = \arg \min_{\beta_0, \dots, \beta_p} \sum_{i=1}^n \epsilon_i^2$$

$$C = \mathbb{E}[\epsilon \epsilon^T] = \sigma^2 I$$
$$C = \begin{pmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{pmatrix}$$

Model assumptions

- *Linearity*
- *Normality*
- *Independence*
- *Absence of multicollinearity*
- *Fixed features*
- *Homoscedasticity*



Weighted linear regression

$$Y = f(X) + \epsilon^* = \beta_0 + \sum_{j=1}^p X_j \beta_j + \epsilon^*$$



Weighted linear regression

$$Y = f(X) + \epsilon^* = \beta_0 + \sum_{j=1}^p X_j \beta_j + \epsilon^*$$

Optimisation problem

$$\arg \min_{\beta_0, \dots, \beta_p} \sum_{i=1}^n w_i (y_i - (\beta_0 + \sum_{j=1}^p \beta_j x_{ij}))^2 = \arg \min_{\beta_0, \dots, \beta_p} \sum_{i=1}^n \epsilon_i^{*2}$$



Weighted linear regression

$$Y = f(X) + \epsilon^* = \beta_0 + \sum_{j=1}^p X_j \beta_j + \epsilon^*$$

Optimisation problem

$$\arg \min_{\beta_0, \dots, \beta_p} \sum_{i=1}^n w_i (y_i - (\beta_0 + \sum_{j=1}^p \beta_j x_{ij}))^2 = \arg \min_{\beta_0, \dots, \beta_p} \sum_{i=1}^n \epsilon_i^{*2}$$

$$C = \begin{pmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{pmatrix}$$

$$W = \begin{pmatrix} w_1 & 0 & \dots & 0 \\ 0 & w_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & w_n \end{pmatrix}$$

$$w_i = 1/\sigma_i^2$$

Violation of homoscedasticity

Weighted linear regression uses different weights for each observation based on their variance. A small error variance observation has a large weight since it includes more information than a large error variance observation, which has a small weight.



L1 and L2 regularisation

$$\text{RSS}(\beta) = \sum_{i=1}^n (y_i - (\beta_0 + \sum_{j=1}^p \beta_j x_{ij}))^2$$

Lasso regression (L1 regularisation)

$$\hat{\beta}_{\text{lasso}} = \arg \min_{(\beta_0, \dots, \beta_p) \in \mathbb{R}^n} \left\{ \text{RSS}(\beta) + \underbrace{\lambda \sum_{j=1}^p |\beta_j|}_{\text{Penalty term}} \right\}$$

It tends far more to drive small weights to 0.

Ridge regression (L2 regularisation)

$$\hat{\beta}_{\text{ridge}} = \arg \min_{(\beta_0, \dots, \beta_p)} \left\{ \text{RSS}(\beta) + \underbrace{\lambda \sum_{j=1}^p \beta_j^2}_{\text{Penalty term}} \right\}$$

It pushes down big weights than tiny ones.

- 1 Introduction
- 2 Machine learning
- 3 Random forest
- 4 Regression
- 5 Explainable artificial intelligence
- 6 Conclusions

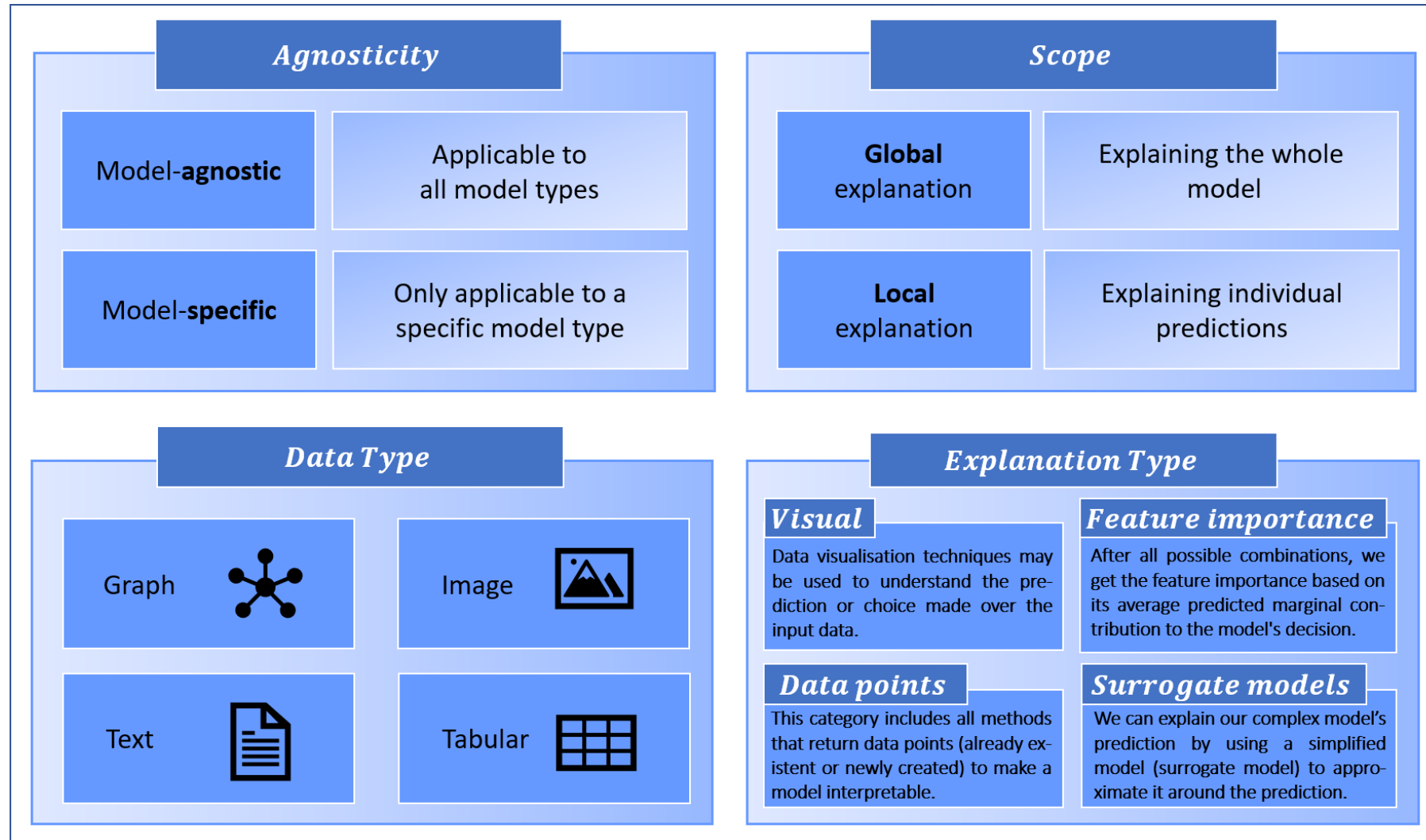
What is explainable AI?

Explainable artificial intelligence

Set of techniques that either produce more understandable models keeping high levels of performance or provide external tools to better understand the models that are inherently not interpretable.



XAI taxonomy



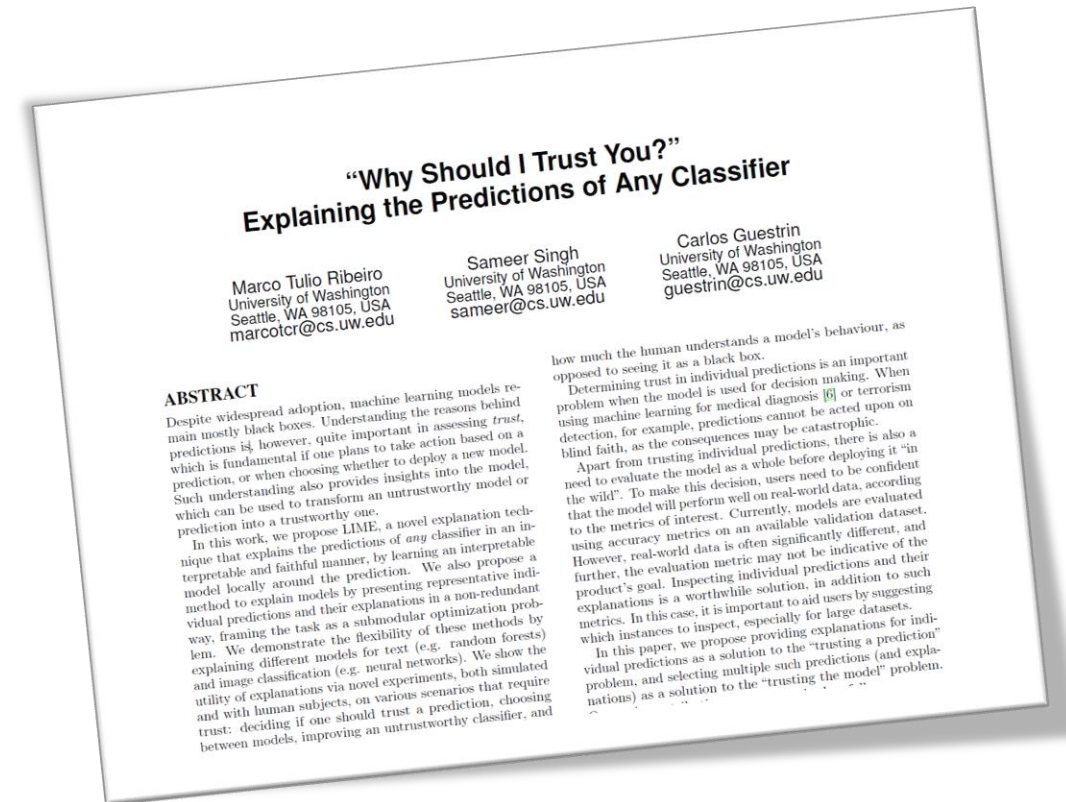
LIME

Local

Interpretable

Model-agnostic

Explanations



LIME optimisation problem

$$\zeta(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$



LIME optimisation problem

$$\zeta(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$




$x \in \mathbb{R}^d \longrightarrow$ number of features

<i>Age</i> = 56	<i>Gender</i> = F	<i>BMI</i> = 30	<i>Stroke</i> = yes
-----------------	-------------------	-----------------	---------------------



LIME optimisation problem

$$\zeta(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$


 $x \in \mathbb{R}^d \longrightarrow$ number of features

Complex model
 $f : \mathbb{R}^d \rightarrow \mathbb{R}$


Simple interpretable model

<i>Age</i> = 56	<i>Gender</i> = F	<i>BMI</i> = 30	<i>Stroke</i> = yes
-----------------	-------------------	-----------------	---------------------



LIME optimisation problem

$$\zeta(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$


 $x \in \mathbb{R}^d \rightarrow$ number of features


$g \in G$ → Family of interpretable models
 $f : \mathbb{R}^d \rightarrow \mathbb{R}$ → Complex model
 $\Omega(g)$ → Simple interpretable model

<i>Age</i> = 56	<i>Gender</i> = <i>F</i>	<i>BMI</i> = 30	<i>Stroke</i> = <i>yes</i>
-----------------	--------------------------	-----------------	----------------------------



LIME optimisation problem

$$\zeta(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$


 A stick figure icon representing a user or patient.

$x \in \mathbb{R}^d \rightarrow$ number of features

$g \in G$ Family of interpretable models

$f : \mathbb{R}^d \rightarrow \mathbb{R}$ Complex model

π_x Neighbourhood of x


$\Omega(g)$ Simple interpretable model

<i>Age</i> = 56	<i>Gender</i> = <i>F</i>	<i>BMI</i> = 30	<i>Stroke</i> = <i>yes</i>
-----------------	--------------------------	-----------------	----------------------------



LIME optimisation problem

$$\zeta(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$


 $x \in \mathbb{R}^d \rightarrow$ number of features

Family of interpretable models
 $g \in G$

Complex model
 $f : \mathbb{R}^d \rightarrow \mathbb{R}$

Simple interpretable model
 g

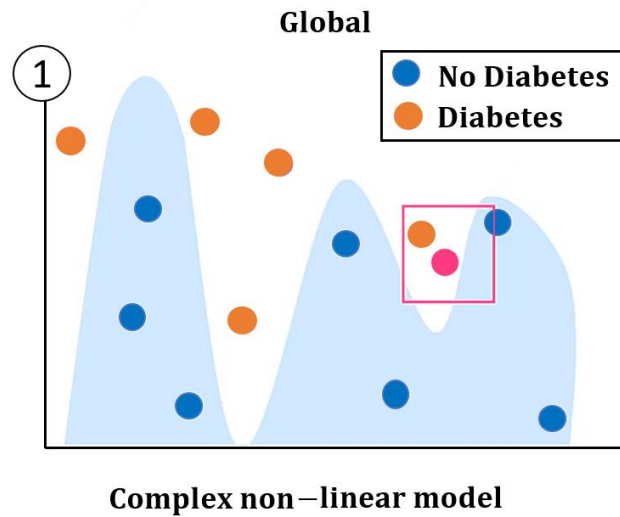
Neighbourhood of x
 π_x

(1) $\mathcal{L}(f, g, \pi_x)$
 (2) $\Omega(g)$

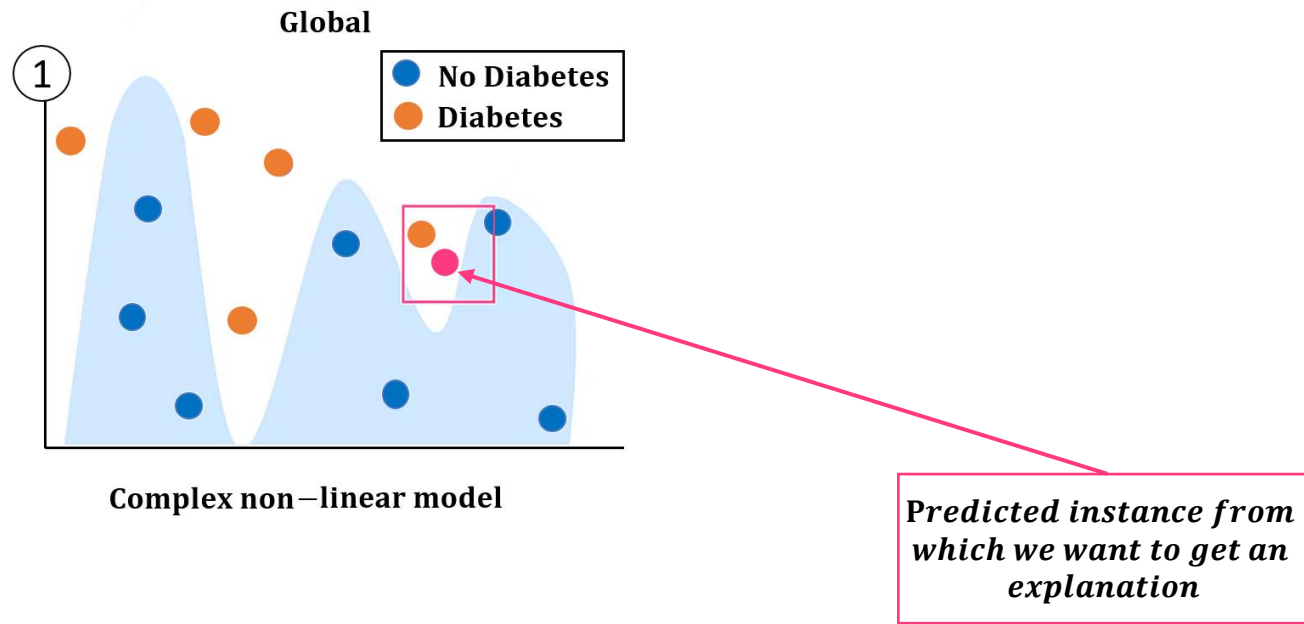
<i>Age</i> = 56	<i>Gender</i> = F	<i>BMI</i> = 30	<i>Stroke</i> = yes
-----------------	-------------------	-----------------	---------------------

2 LOSS TERMS

LIME step by step

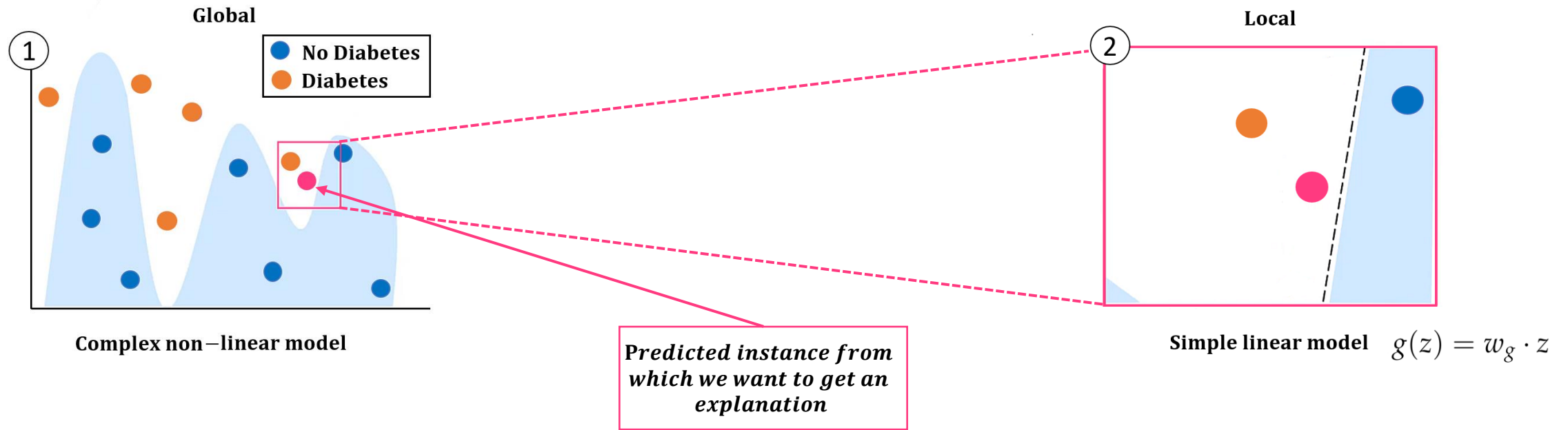


LIME step by step



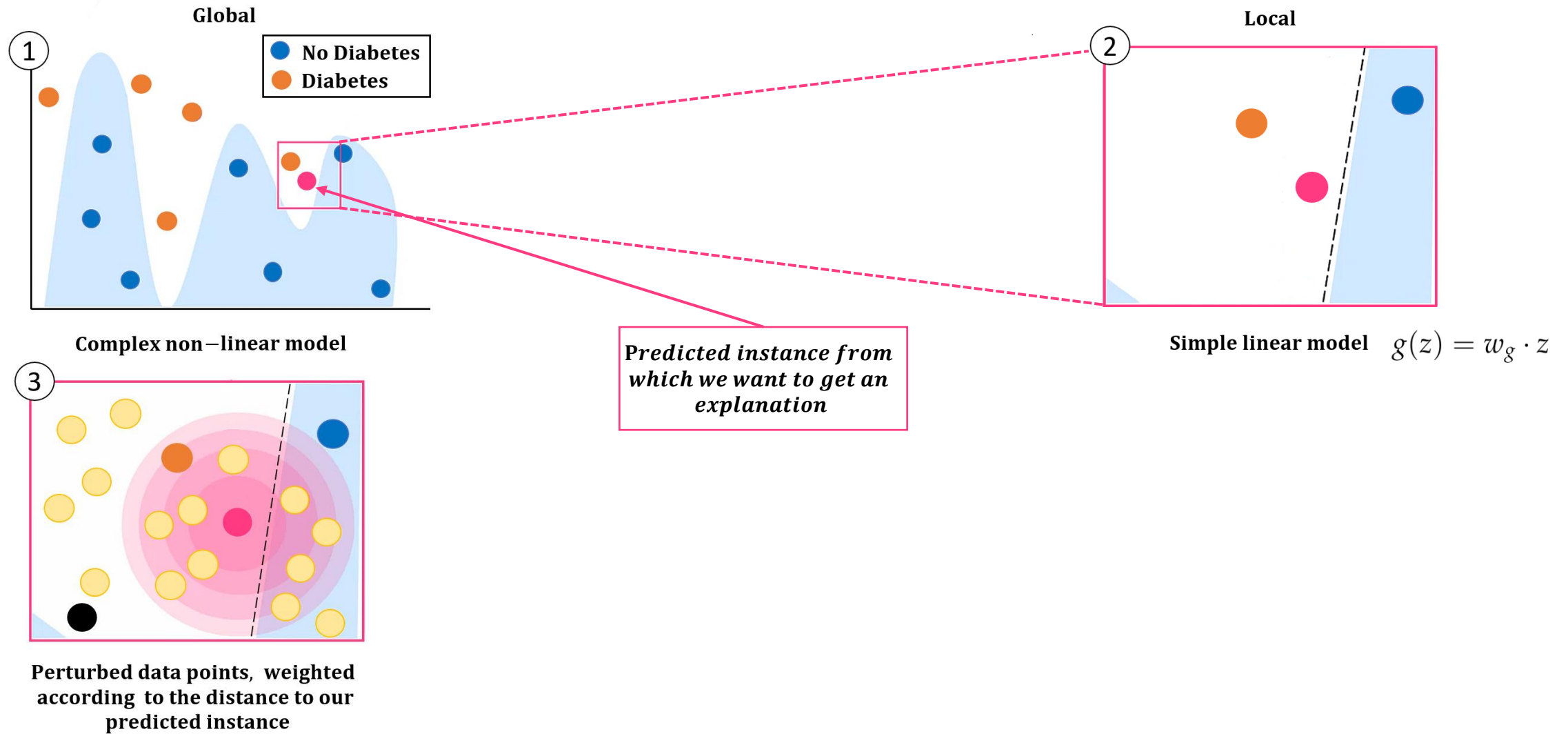


LIME step by step



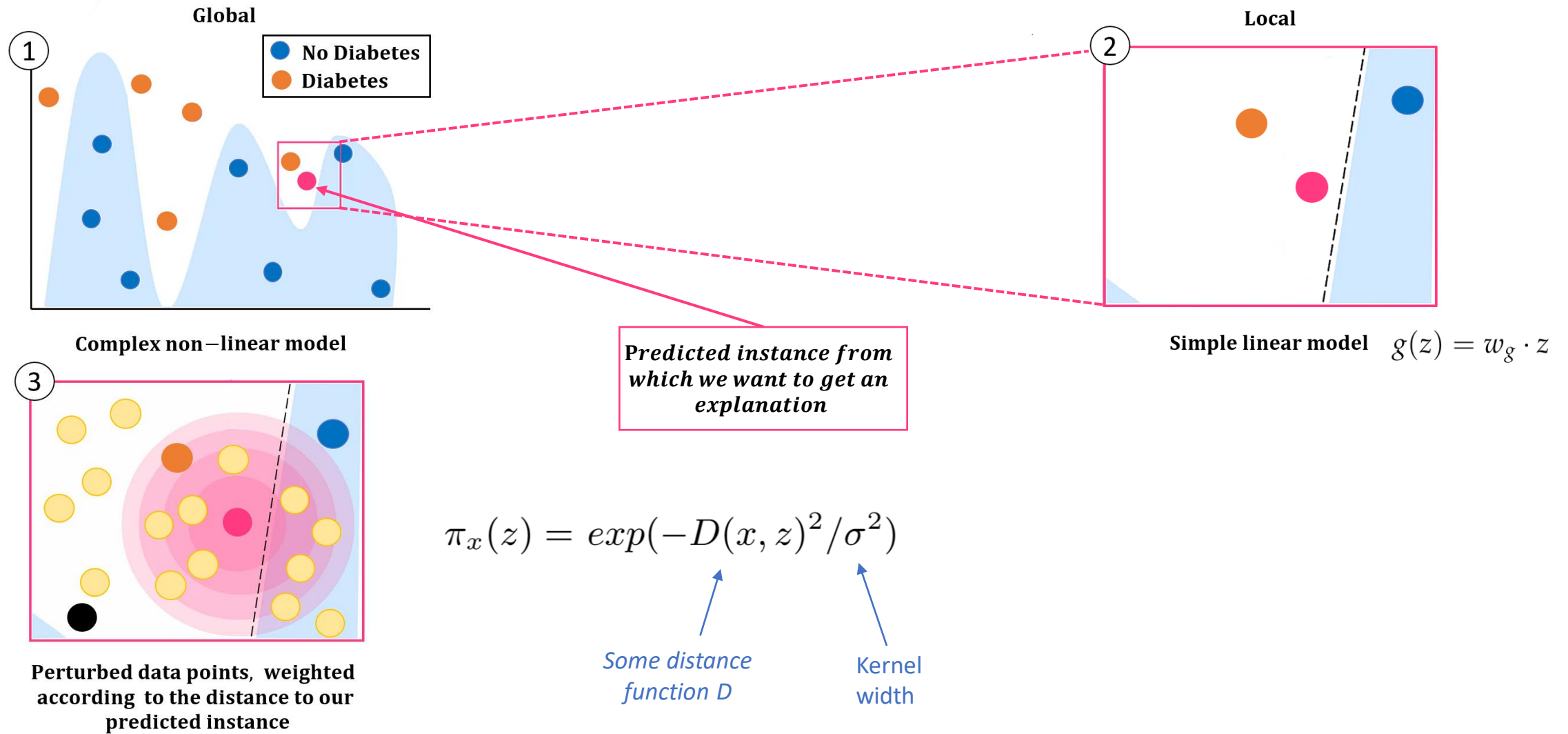


LIME step by step



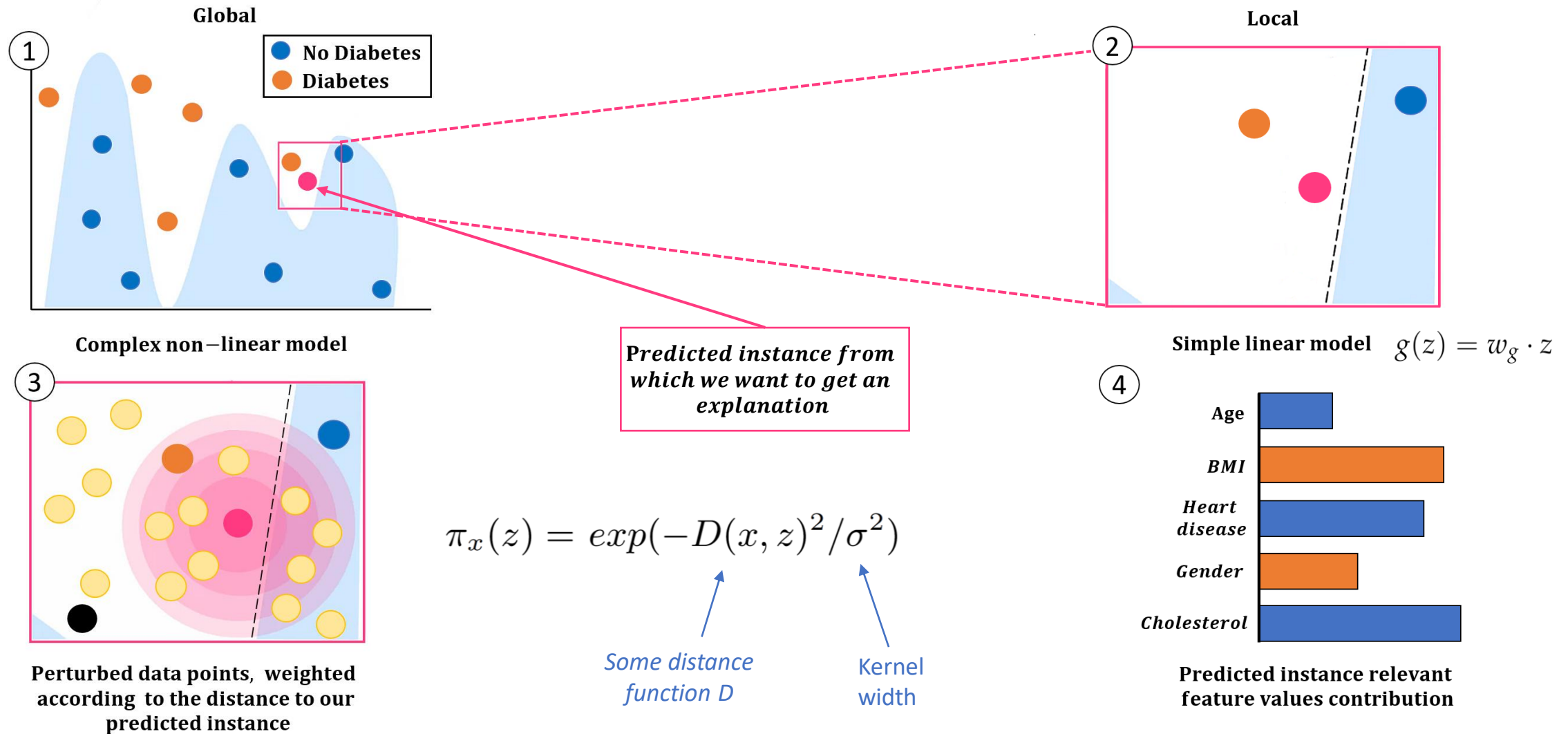


LIME step by step





LIME step by step






Loss terms

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

Train a weighted, interpretable model on the dataset with the perturbed instances

$$(1) \quad \mathcal{L}(f, g, \pi_x) = \sum_{z \in \mathcal{Z}} \pi_x(z) \left(f(z) - g(z) \right)^2$$



Complex model prediction Simple model prediction

Loss terms

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

Train a weighted, interpretable model on the dataset with the perturbed instances

$$(1) \quad \mathcal{L}(f, g, \pi_x) = \sum_{z \in \mathcal{Z}} \pi_x(z) \left(\underset{\substack{\uparrow \\ \text{Complex model} \\ \text{prediction}}}{f(z)} - \underset{\substack{\uparrow \\ \text{Simple model} \\ \text{prediction}}}{g(z)} \right)^2$$

$$(2) \quad \Omega(g) \quad ?$$

Loss terms

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

Train a weighted, interpretable model on the dataset with the perturbed instances

$$(1) \quad \mathcal{L}(f, g, \pi_x) = \sum_{z \in \mathcal{Z}} \pi_x(z) \left(\underset{\substack{\uparrow \\ \text{Complex model} \\ \text{prediction}}}{f(z)} - \underset{\substack{\uparrow \\ \text{Simple model} \\ \text{prediction}}}{g(z)} \right)^2$$

$$(2) \quad \Omega(g) \quad \text{LIME uses sparse linear models (K - LASSO)}$$

$$\hat{\beta}_{lasso} = \operatorname{argmin}_{(\beta_0, \dots, \beta_p) \in \mathbb{R}^n} \left\{ \operatorname{RSS}(\beta) + \underbrace{\lambda \sum_{j=1}^p |\beta_j|}_{\text{Penalty term}} \right\}$$

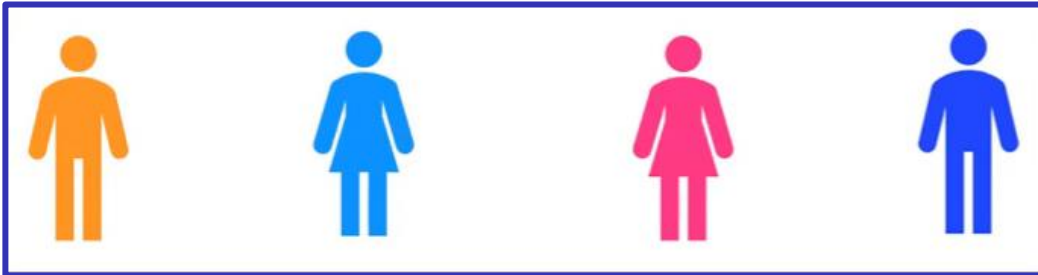
LIME limitations

- Neighbourhood
- Non-linearity
- Improbable instances
- Instability



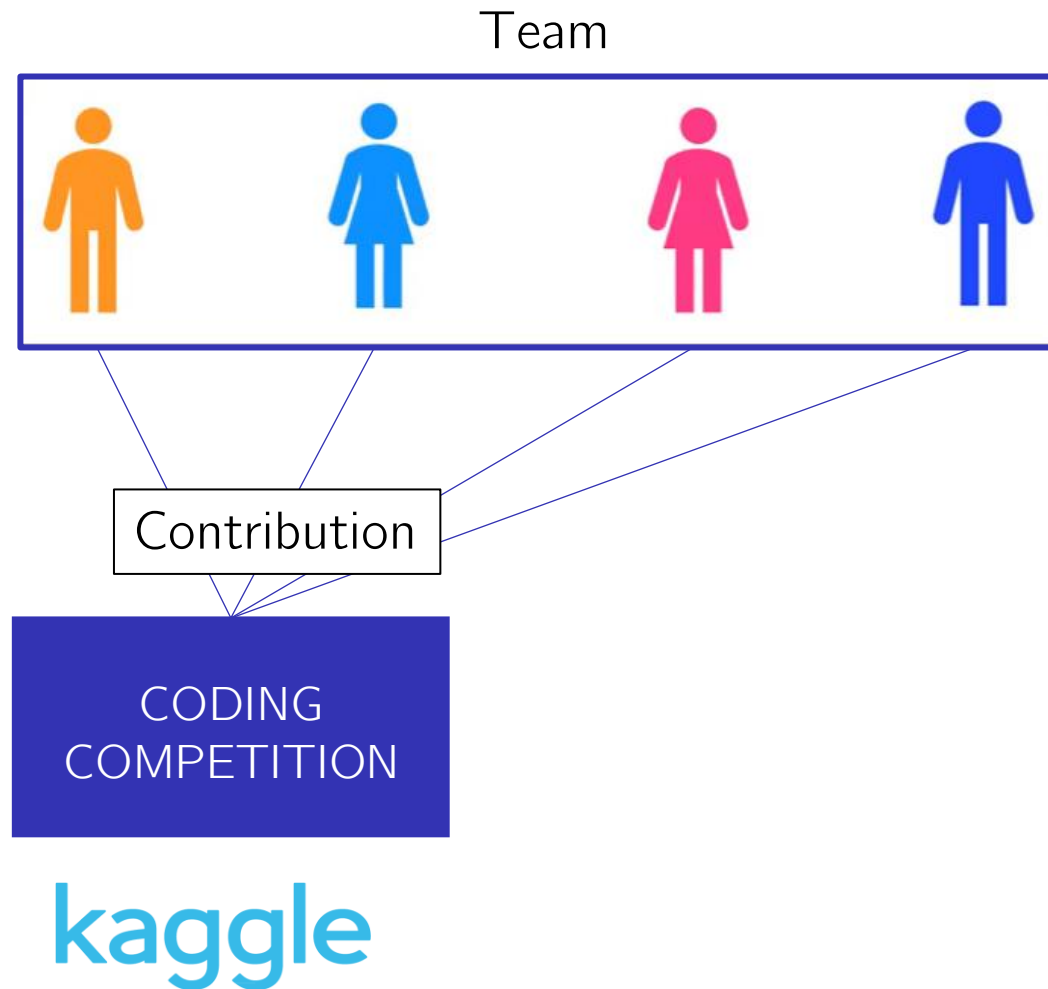
Shapley values

Team



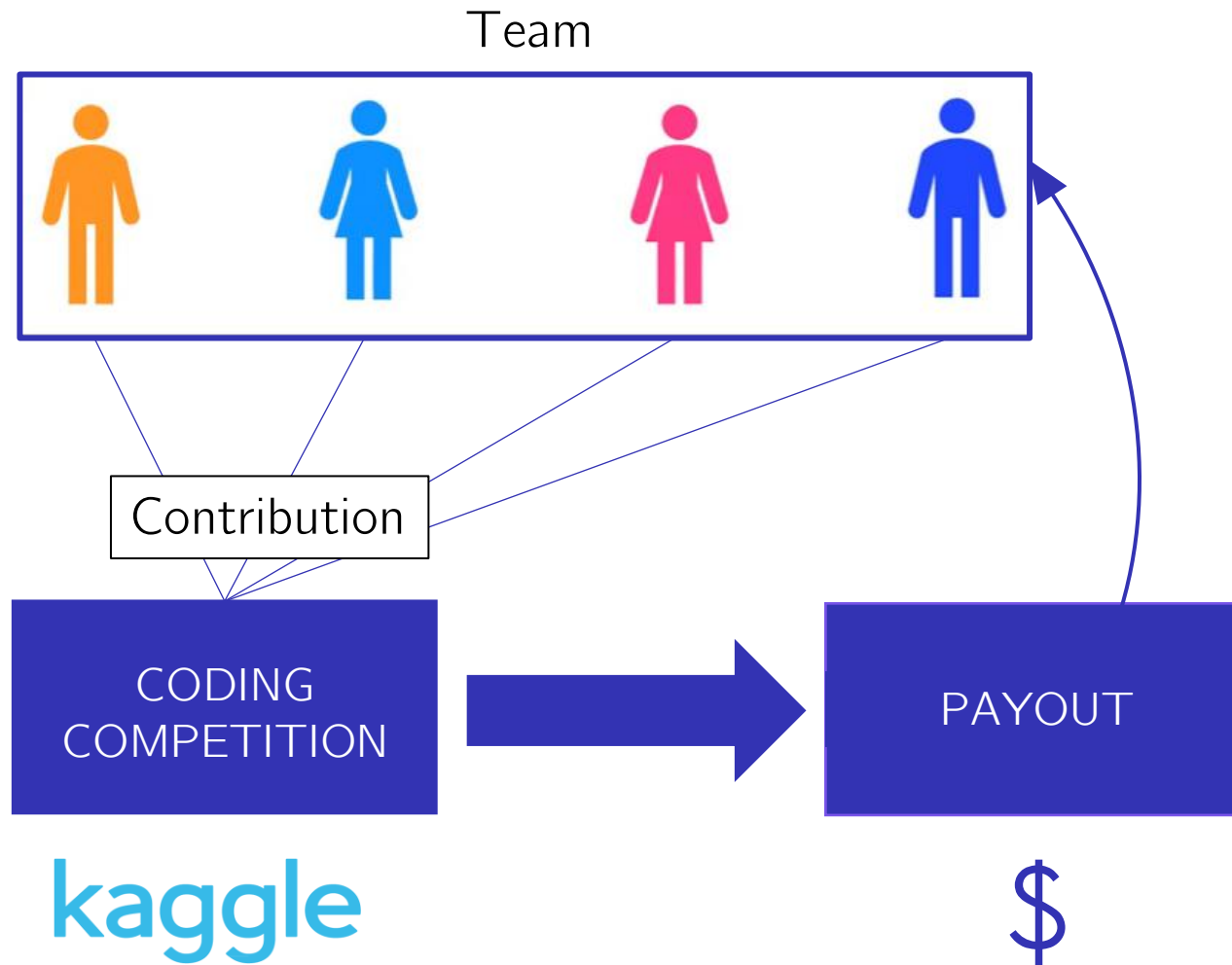


Shapley values

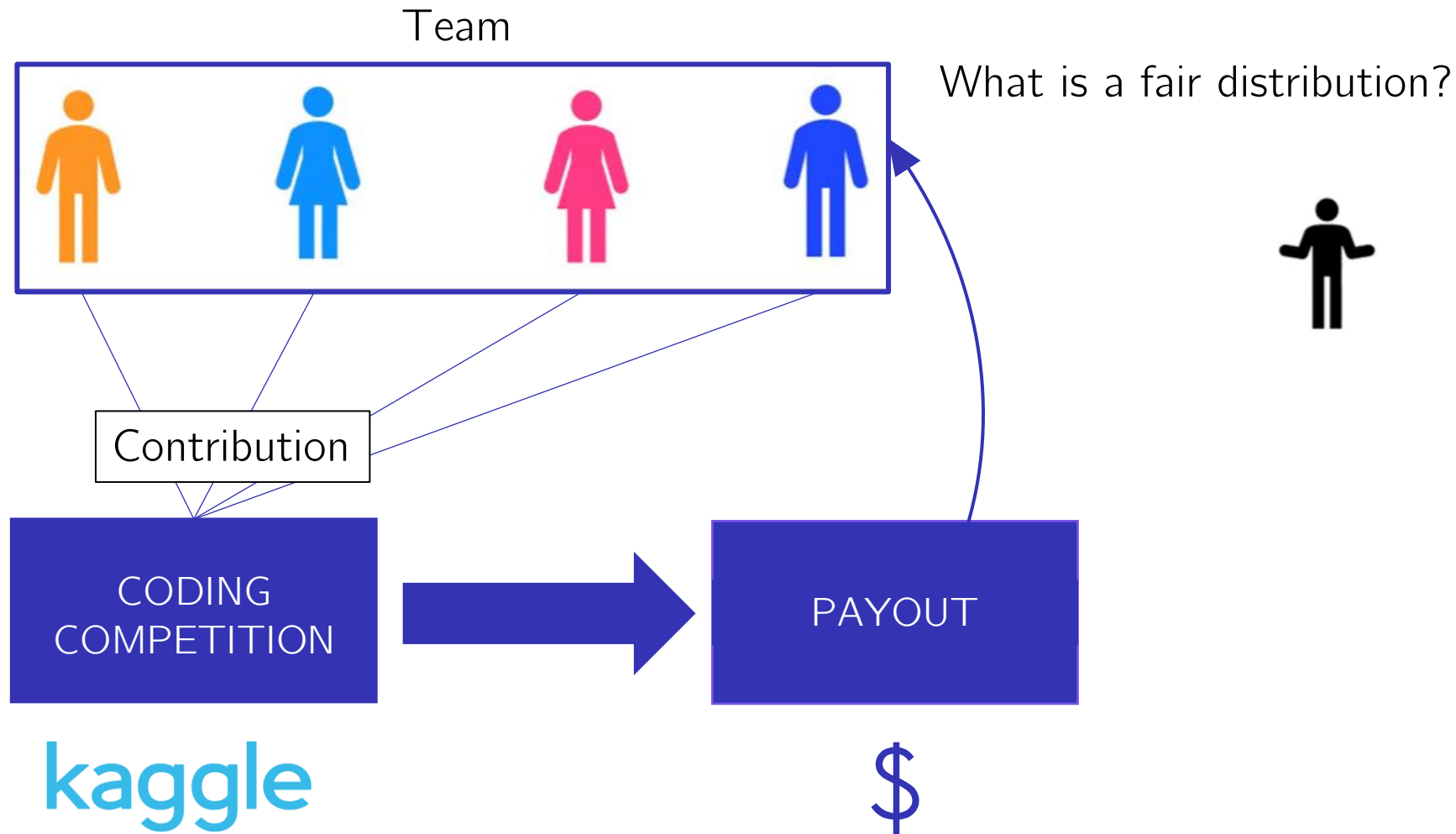




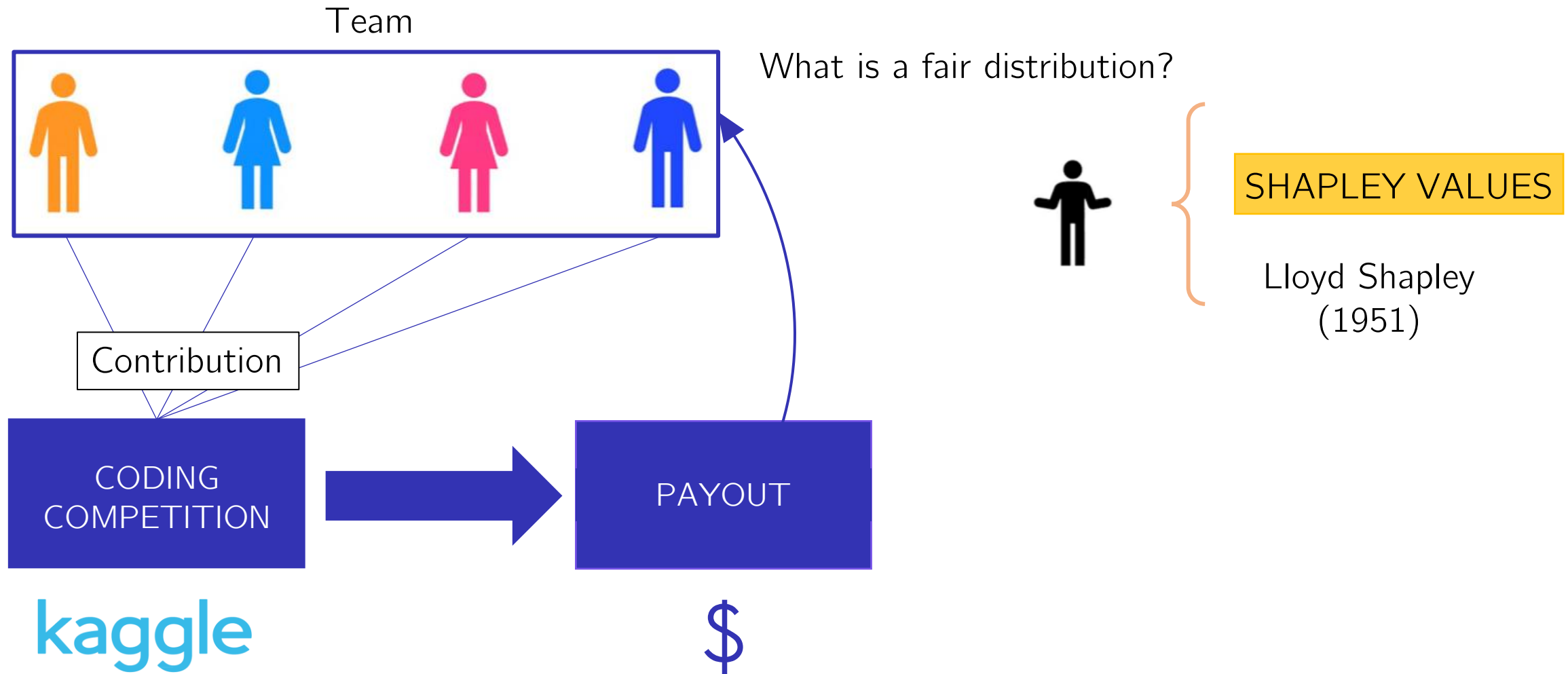
Shapley values



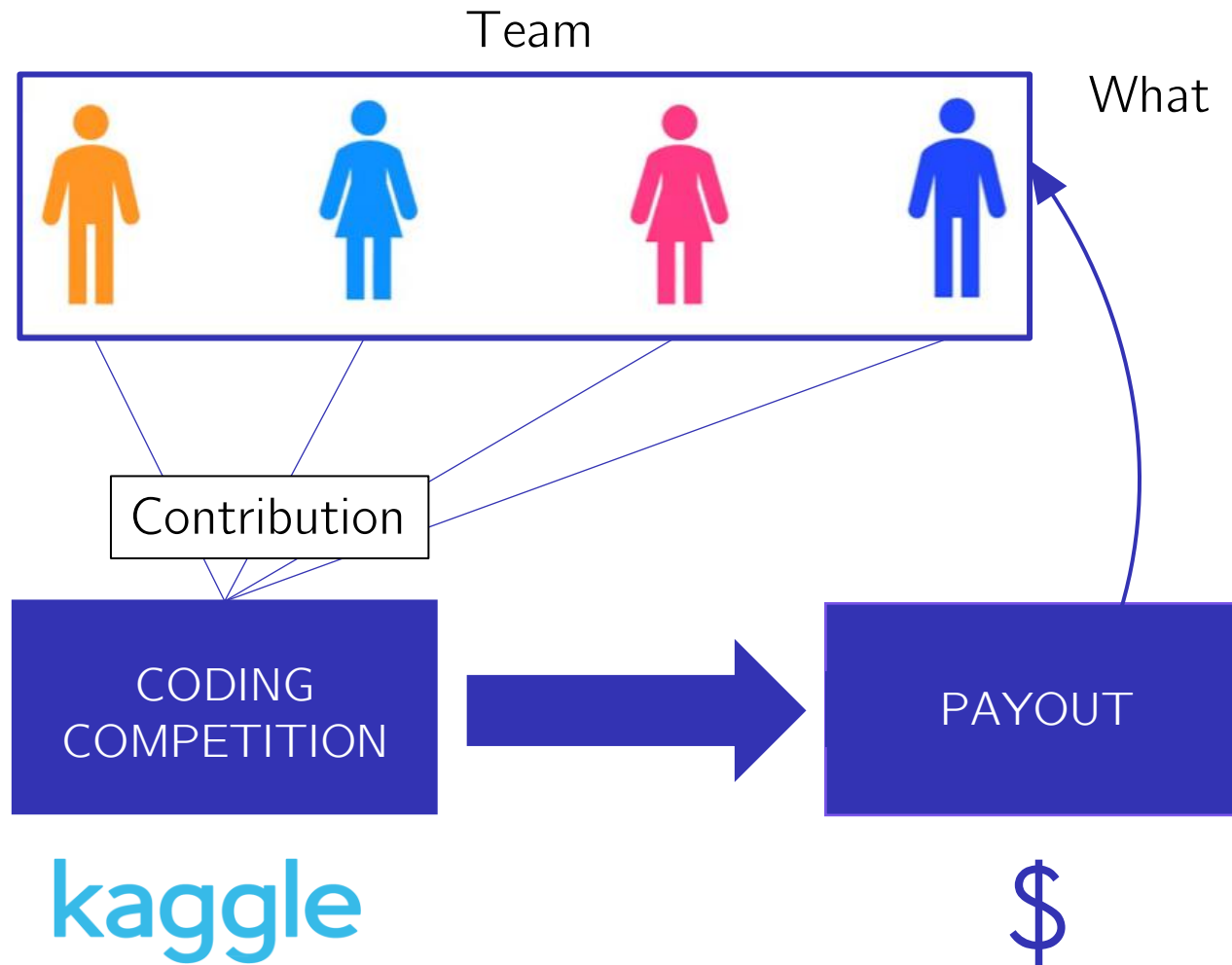
Shapley values



Shapley values



Shapley values



SHAPLEY VALUES

Lloyd Shapley
(1951)

$$\phi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} (v(S \cup \{i\}) - v(S))$$

SHAP

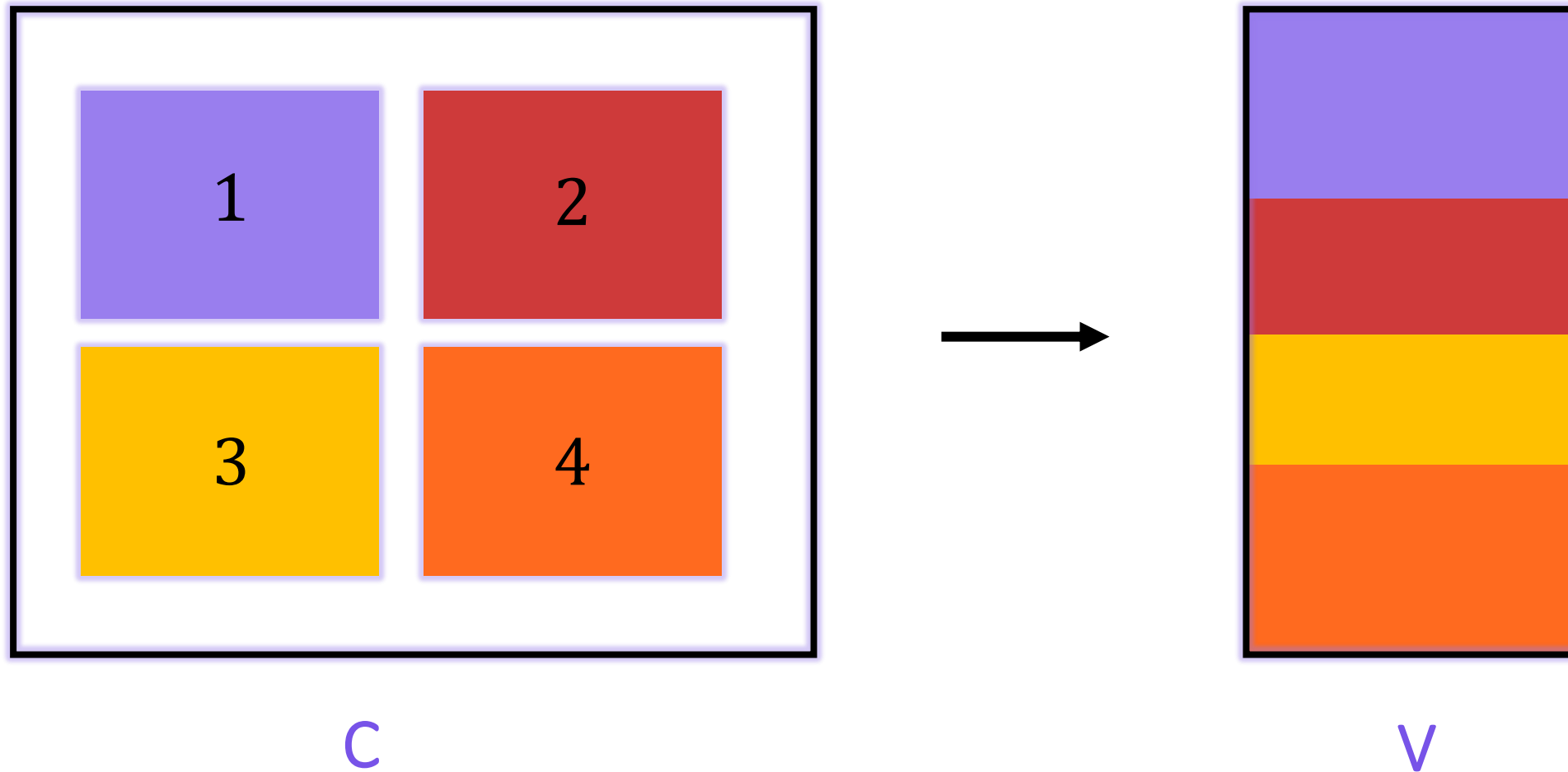
SHapley
Additive
ex**P**lanations



UNIVERSITAT DE
BARCELONA

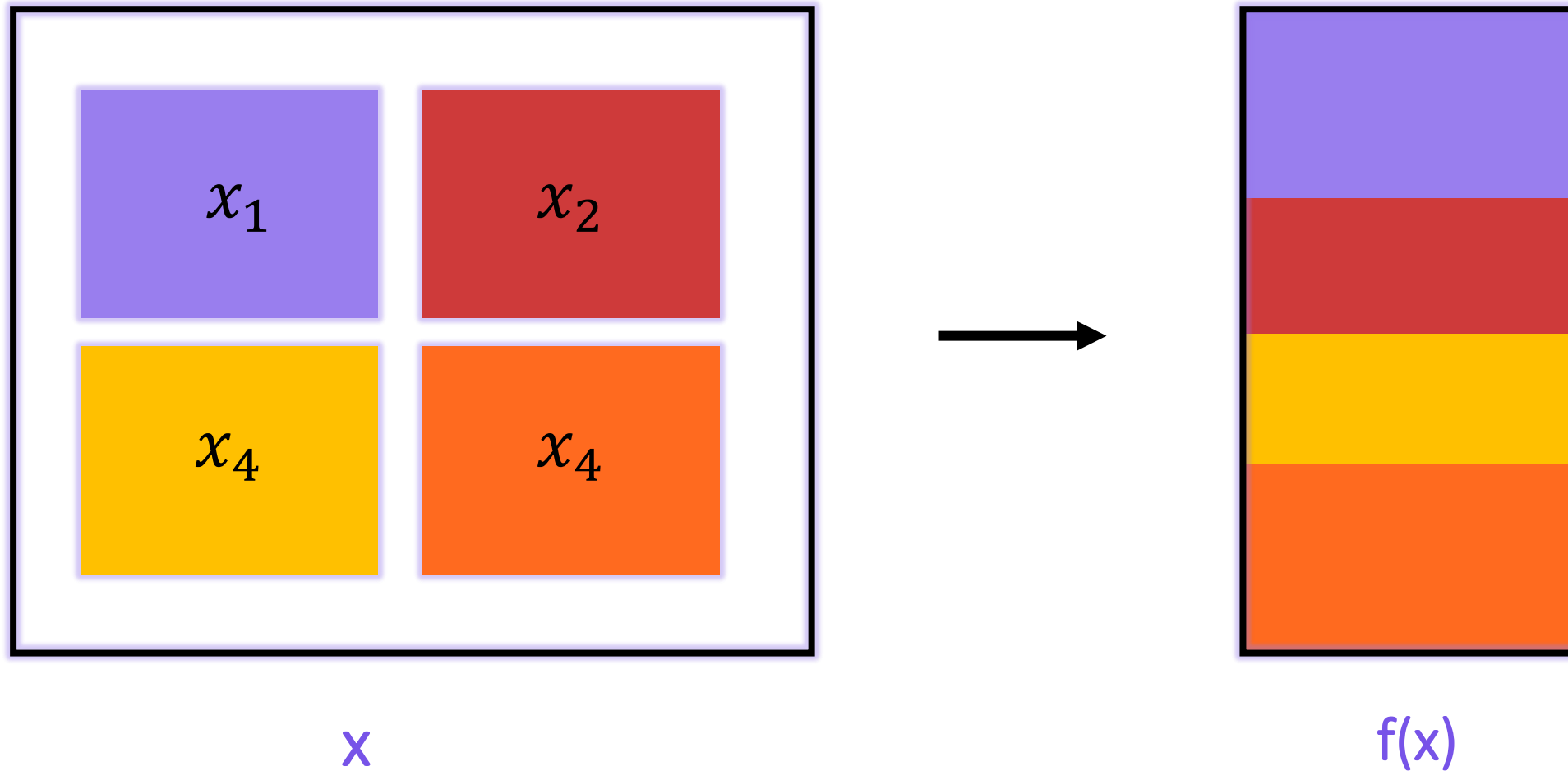


SHAP





SHAP



SHAP

SHappley **A**dditive ex**P**lanations

Additive feature attribution methods

1 *if $x \approx x'$ then $f(x) \approx g(x')$*



Additive feature attribution methods

1 *if $x \approx x'$ then $f(x) \approx g(x')$*

2
$$g(x') = \phi_0 + \sum_{i=1}^p \phi_i x'_i$$

Additive feature attribution methods properties

Properties



Additive feature attribution methods properties

Properties

1

Local accuracy

$$f(x) = g(x') = \phi_0 + \sum_{i=1}^p \phi_i x'_i$$

Additive feature attribution methods properties

Properties

2

Missingness

$$x'_i = 0 \implies \phi_i = 0$$



Additive feature attribution methods properties

Properties

3

Consistency

Let $f_x(z') = f(h_x(z'))$ and $z' \setminus i$ denote setting $z'_i = 0$. For any two models f and f' :

$$\forall z' \in \{0, 1\}^p, f'_x(z') - f'_x(z' \setminus i) \geq f_x(z') - f_x(z' \setminus i) \implies \phi_i(f', x) \geq \phi_i(f, x).$$

Attribution methods satisfying properties 1, 2, 3

$$g(x') = \phi_0 + \sum_{i=1}^p \phi_i x'_i$$



Attribution methods satisfying properties 1, 2, 3

$$g(x') = \phi_0 + \sum_{i=1}^p \phi_i x'_i$$



Attribution methods satisfying properties 1, 2, 3

$$g(x') = \phi_0 + \sum_{i=1}^p \phi_i x'_i$$

Theorem *Only one possible explanation model g follows additive feature attribution methods definition and satisfies Properties 1, 2, and 3:*

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



Calculating Shapley values

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



Calculating Shapley values

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

Shapley value
for feature i



Calculating Shapley values

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



Calculating Shapley values

Black Box model

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



Calculating Shapley values

Black Box model Input data point

AGE

Shapley value for feature i

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



Calculating Shapley values

Black Box model Input data point

AGE

Shapley value for feature i

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

$x =$	<i>Age = 56</i>	<i>Gender = F</i>	<i>BMI = 30</i>	<i>Stroke = yes</i>	<i>...</i>
-------	------------------------	--------------------------	------------------------	----------------------------	-------------------



Calculating Shapley values

Black Box model Input data point

AGE

Shapley value for feature i

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

$x =$	<i>Age = 56</i>	<i>Gender = F</i>	<i>BMI = 30</i>	<i>Stroke = yes</i>	<i>...</i>
-------	------------------------	--------------------------	------------------------	----------------------------	-------------------

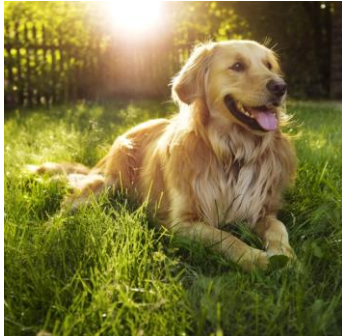


Calculating Shapley values

Black Box model Input data point

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 Simplified data input



$x =$	Age = 56	Gender = F	BMI = 30	Stroke = yes	...
-------	-----------------	-------------------	-----------------	---------------------	-----



Calculating Shapley values

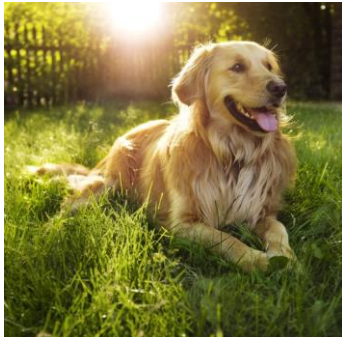
Black Box model Input data point

AGE

Shapley value for feature i

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

Subset Simplified data input



$x =$	Age = 56	Gender = F	BMI = 30	Stroke = yes	...
-------	-----------------	-------------------	-----------------	---------------------	-----



Calculating Shapley values

Black Box model Input data point

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	BMI = 30	Stroke = yes	...
-----------------	------------	-----------------	--------------	-----



Calculating Shapley values

Black Box model Input data point

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	BMI = 30	Stroke = yes	...
-----------------	-------------------	-----------------	---------------------	-----



Calculating Shapley values

Black Box model Input data point

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	<i>Gender = F</i>	BMI = 30	<i>Stroke = yes</i>	...
-----------------	-------------------	-----------------	---------------------	-----



Calculating Shapley values

Black Box model Input data point

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

Age = 56 Body Mass Index = 30

Body Mass Index = 30

$x =$

Age = 56	<i>Gender = F</i>	BMI = 30	<i>Stroke = yes</i>	...
-----------------	-------------------	-----------------	---------------------	-----



Calculating Shapley values

Black Box model Input data point

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

Age = 56 Body Mass Index = 30

Body Mass Index = 30

$x =$

Age = 56	<i>Gender = F</i>	BMI = 30	<i>Stroke = yes</i>	...
-----------------	-------------------	-----------------	---------------------	-----



Calculating Shapley values

Black Box model Input data point

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

Age = 56 Body Mass Index = 30

70% stroke

Body Mass Index = 30

$x =$

Age = 56	<i>Gender = F</i>	BMI = 30	<i>Stroke = yes</i>	...
-----------------	-------------------	-----------------	---------------------	-----



Calculating Shapley values

Black Box model Input data point

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

Age = 56 Body Mass Index = 30

70% stroke

10% stroke

Body Mass Index = 30

$x =$

Age = 56	Gender = F	BMI = 30	Stroke = yes	...
-----------------	-------------------	-----------------	---------------------	-----



Calculating Shapley values

Black Box model Input data point

AGE

Shapley value for feature i

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

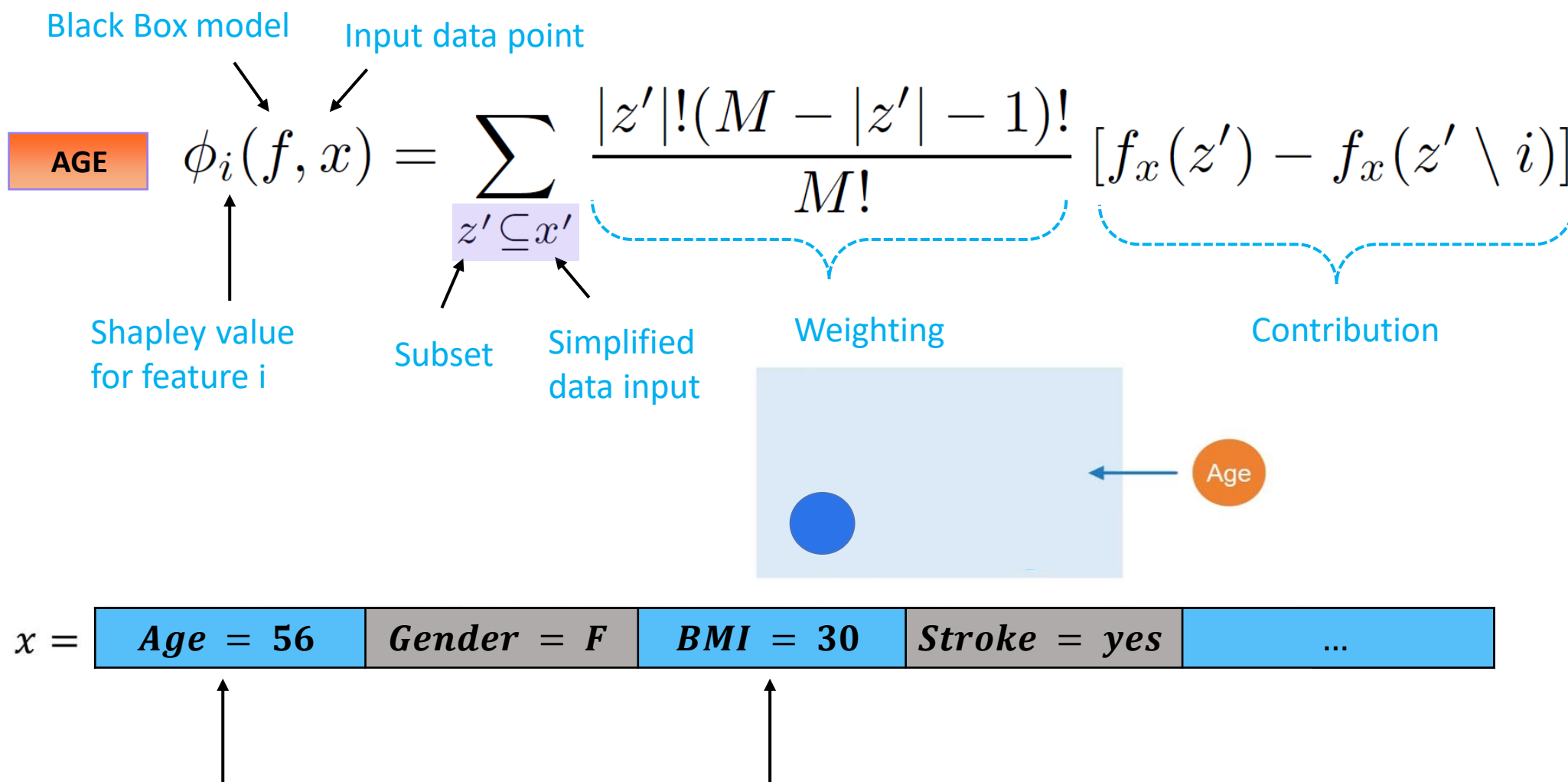
Subset Simplified data input Weighting

$x =$

Age = 56	Gender = F	BMI = 30	Stroke = yes	...
-----------------	------------	-----------------	--------------	-----

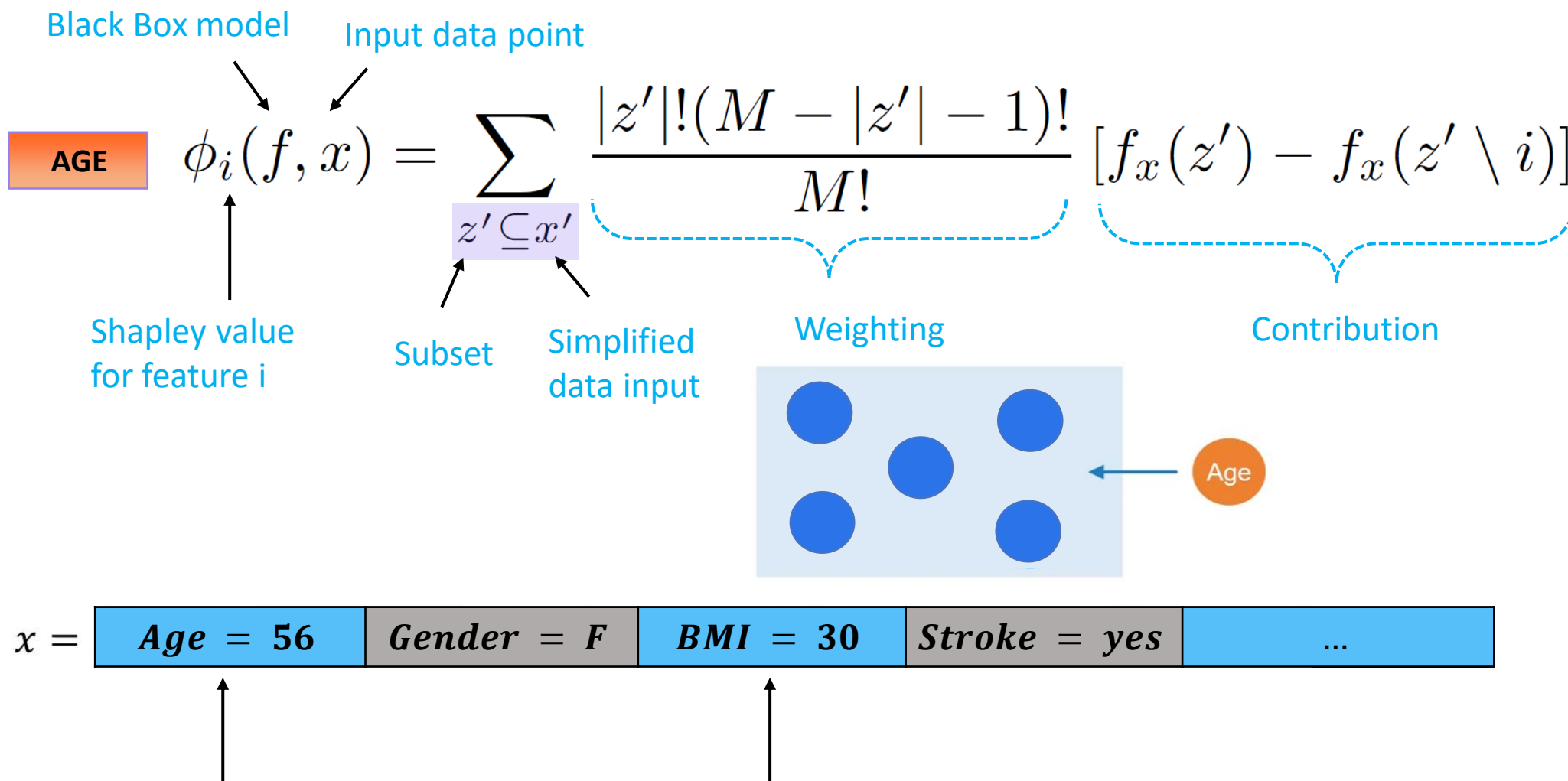


Calculating Shapley values



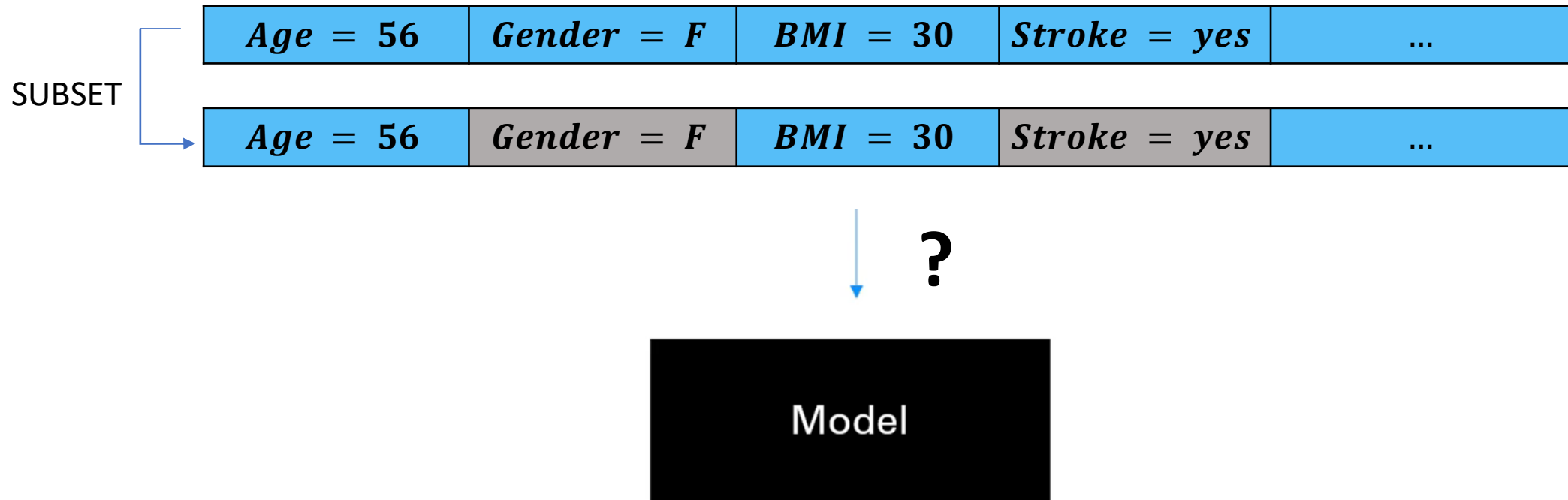


Calculating Shapley values

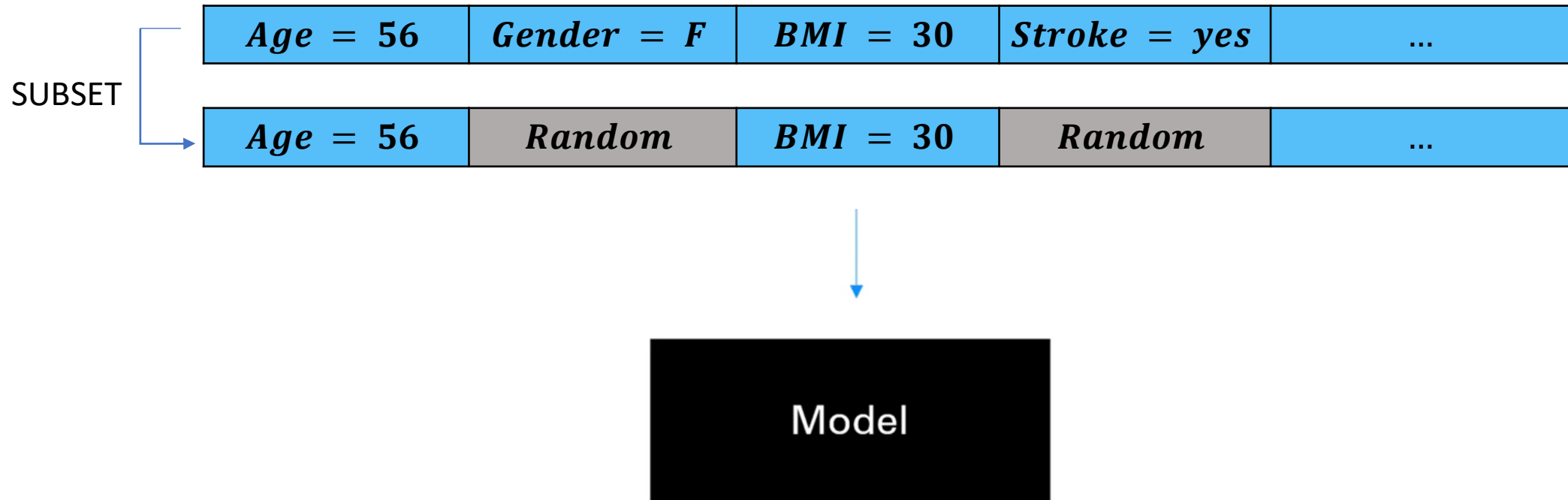




Calculating Shapley values



Calculating Shapley values



Calculating Shapley values

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

Calculating Shapley values

2^n = *total number of subsets of a set of size n*

4 features: 64 total coalitions to sample

Calculating Shapley values

2^n = *total number of subsets of a set of size n*

4 features: 64 total coalitions to sample

32 features: 17.1 billion



Shapley kernel

Shapley kernel

Shapley kernel theorem

$$\left. \begin{aligned} \Omega(g) &= 0, \\ \pi_{x'}(z') &= \frac{M-1}{\binom{M}{|z'|} |z'| (M-|z'|)}, \\ \mathcal{L}(f, g, \pi_{x'}) &= \sum_{z' \in \mathcal{Z}} [f(h_x^{-1}(z')) - g(z')]^2 \pi_{x'}(z') \end{aligned} \right\} \quad \zeta(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

Kernel SHAP = LIME + SHAPLEY VALUES

SHAP limitations

- Computational cost
- Access to data
- Feature dependencies



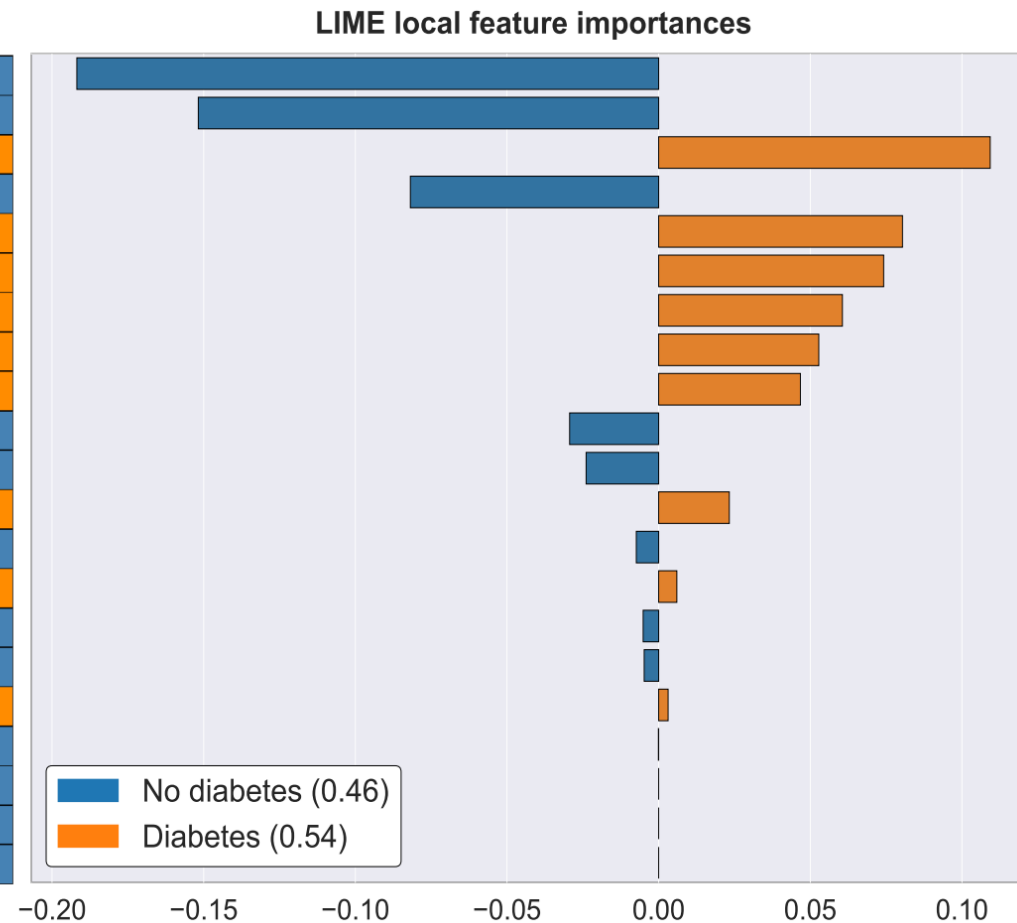
LIME and SHAP application

Step by step

- *Diabetes public tabular database*
- *Random forest fit with this database*
- *LIME and SHAP explanations*

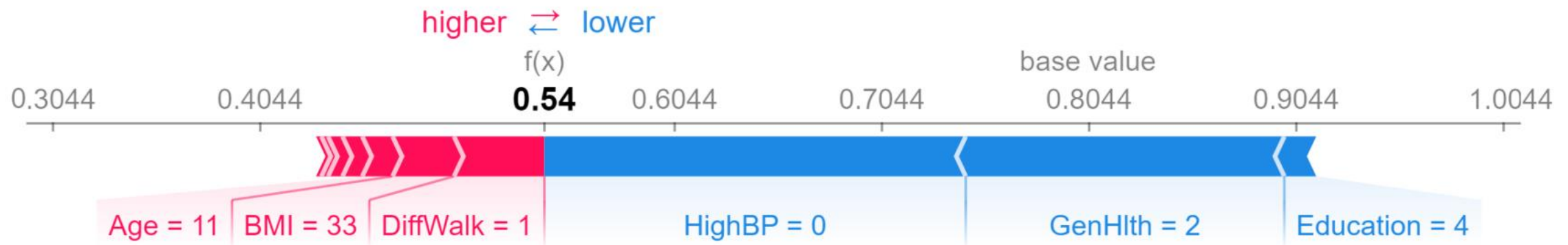
LIME explanation

Feature	Value	Weight
GenHlth	2	-0.191825
HighBP	0	-0.151779
HighChol	1	0.109313
HeartDiseaseorAttack	0	-0.081861
HvyAlcoholConsump	0	0.080394
Age	11	0.074182
BMI	33	0.060582
DiffWalk	1	0.052797
Income	4	0.046733
Stroke	0	-0.029359
Gender	0	-0.023914
Education	4	0.023293
PhysHlth	0	-0.007371
PhysActivity	0	0.005999
Fruits	1	-0.005114
MentHlth	1	-0.004775
NoDocbcCost	0	0.003120
Smoker	1	-0.000084
CholCheck	1	0.000000
Veggies	1	0.000000
AnyHealthcare	1	0.000000





SHAP explanation



- 1 Introduction
- 2 Machine learning
- 3 Random forest
- 4 Regression
- 5 Explainable artificial intelligence
- 6 Conclusions

Conclusion: LIME vs SHAP



	LIME	SHAP
Theory driven	Fails at being consistent. ✗	Supported by the Shapley values theory properties and consistency property. ✓
Time expensive	Time affordable. ✓	Computation of marginal contributions for all possible coalitions makes it time expensive. ✗
Require training data	Does not require the training set for fitting the surrogate model. ✓	Requires the training set for generating the background set that will be used to train the surrogate model. ✗
What-if explanations	Can provide what-if explanations. ✓	Cannot provide what-if explanations. ✗
Improbable instances	Improbable instances may be generated when obtaining perturbed instances. ✗	When imputing omitted features, improbable instances may be generated. ✗
Instability	Kernel width can make it unstable. ✗	Its strong theoretical properties makes it stable. ✓

Conclusion: Summary

- Detailed insight into the theory behind random forest
- Formalise and unify the theory behind LIME and SHAP
- Healthcare application for LIME and SHAP

Future work



- See how LIME explanations vary depending on the kernel width
- Expand LIME and SHAP theory and application to images
- Compare LIME and SHAP explanations