

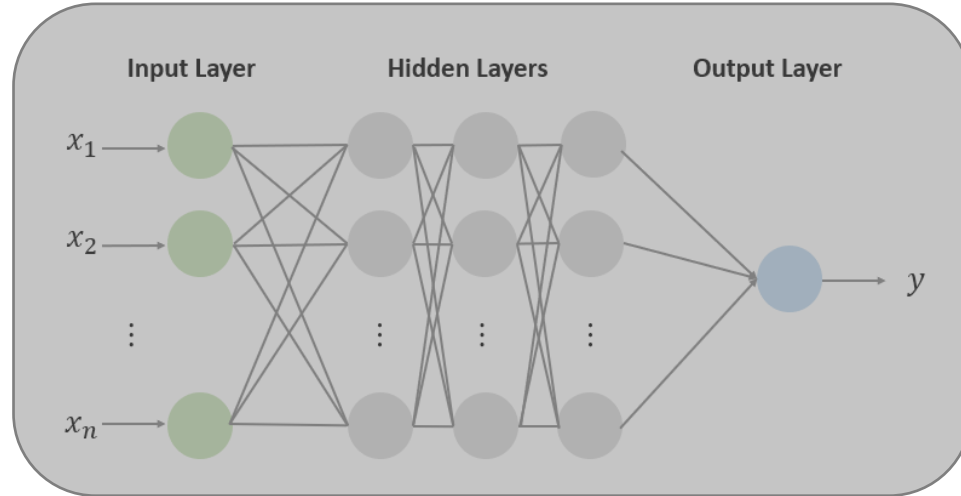
CS 4501/6501 Interpretable Machine Learning

Post-hoc explanations: perturbation-based methods

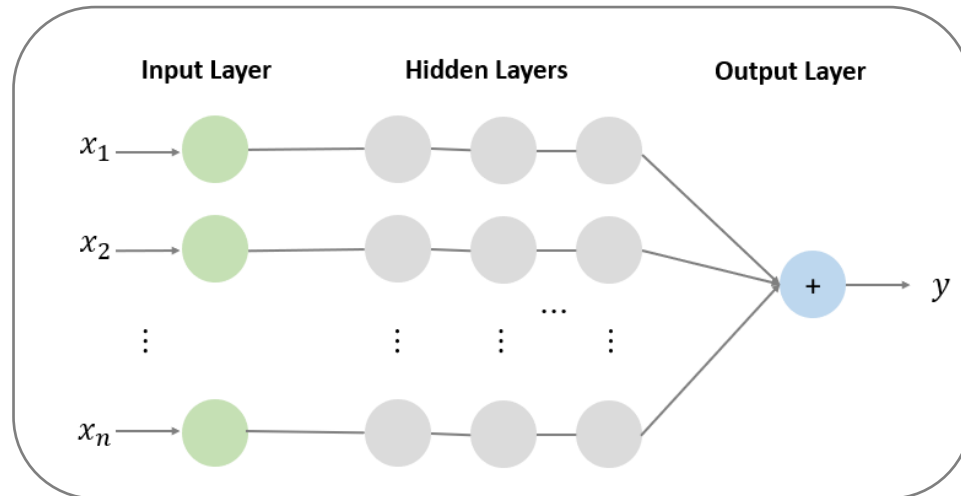
Hanjie Chen, Yangfeng Ji
Department of Computer Science
University of Virginia
{hc9mx, yangfeng}@virginia.edu

Trade-off

Black-box
Neural Network



Interpretable GAM



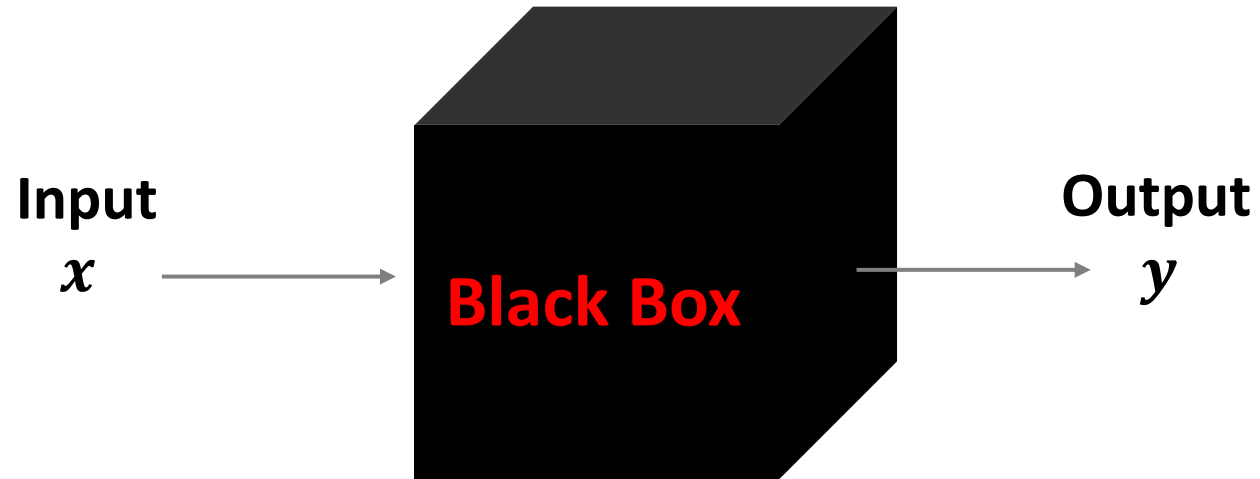
Limitations

- Ignoring complex feature interactions
- Performance drop

Explaining Black-box Model

How to improve model interpretability?

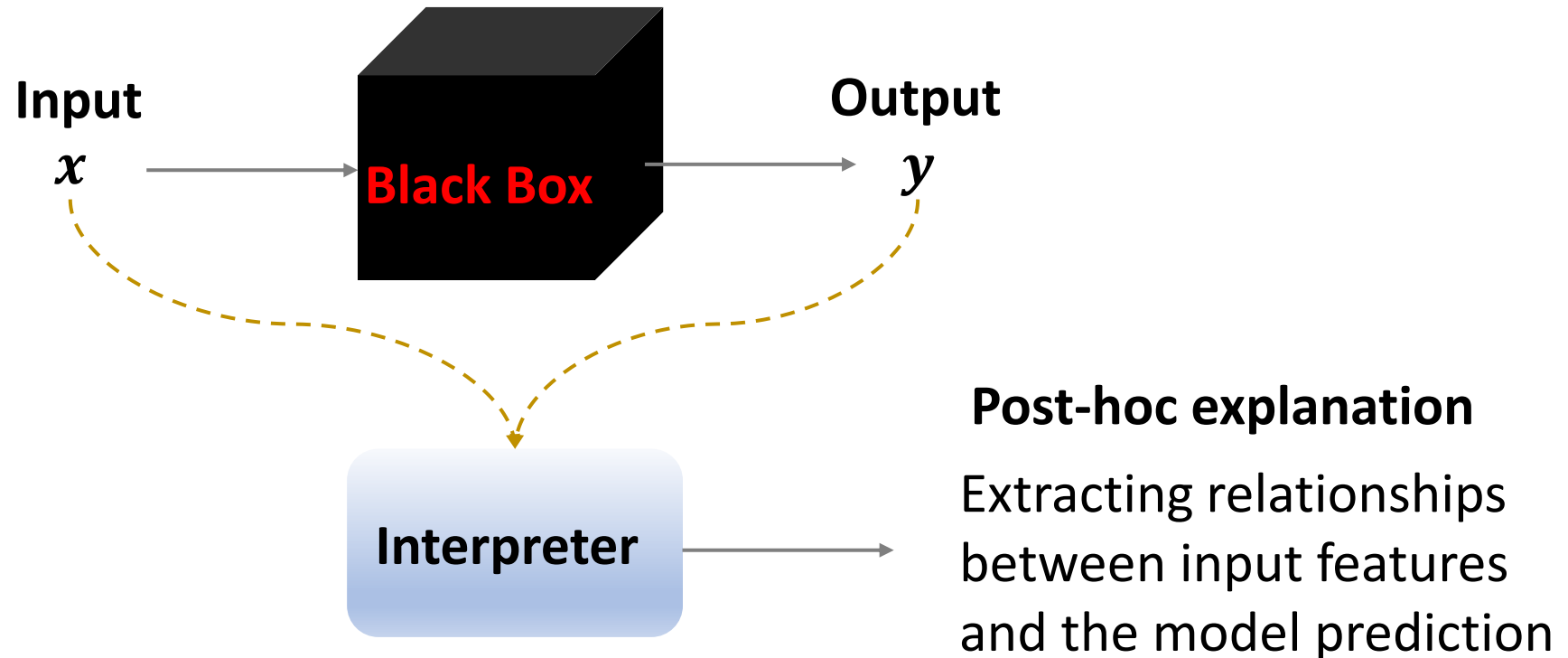
Model's inner working and decision making are hidden in the black box



Explaining Black-box Model

How to improve model interpretability?

Explaining model predictions from the post-hoc manner

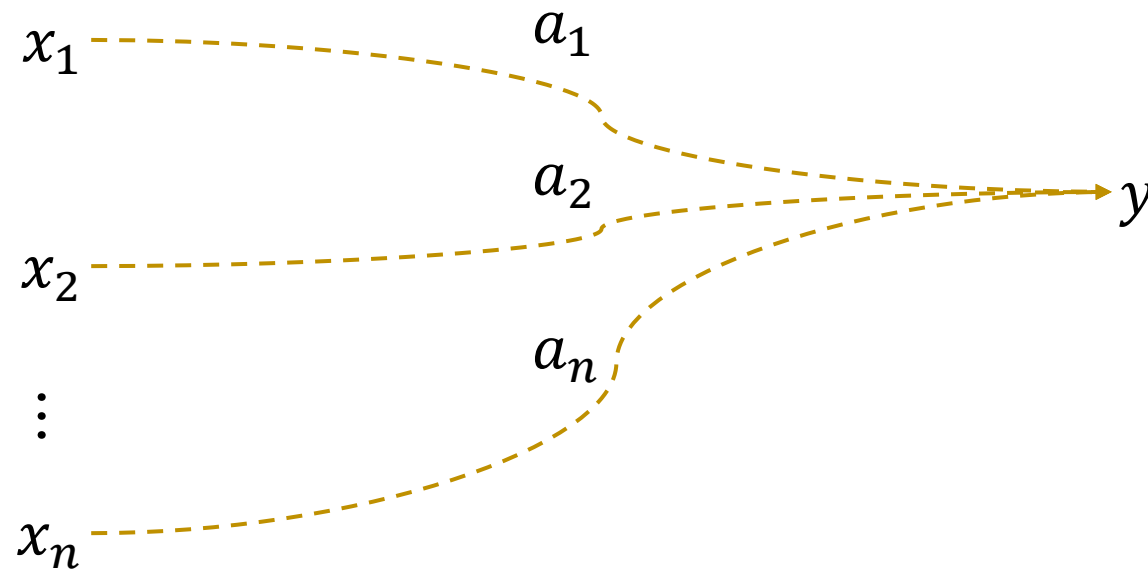


Post-hoc Explanation

Input features

Importance

Model prediction



Identifying important features

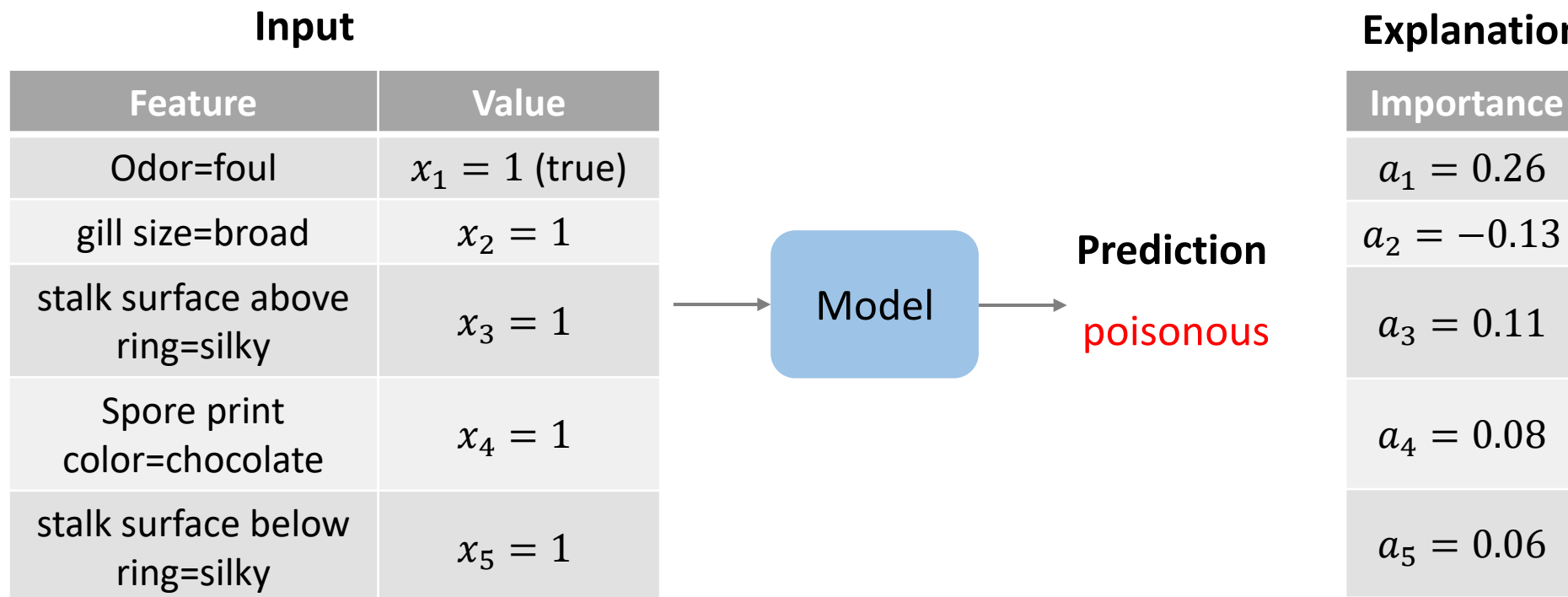
Post-hoc Explanation

- Example 1: tabular data
 - Mushroom dataset
 - **Task:** predicting if a mushroom is **edible** or **poisonous**

Feature
Odor
gill size
stalk surface above ring
Spore print color
stalk surface below ring

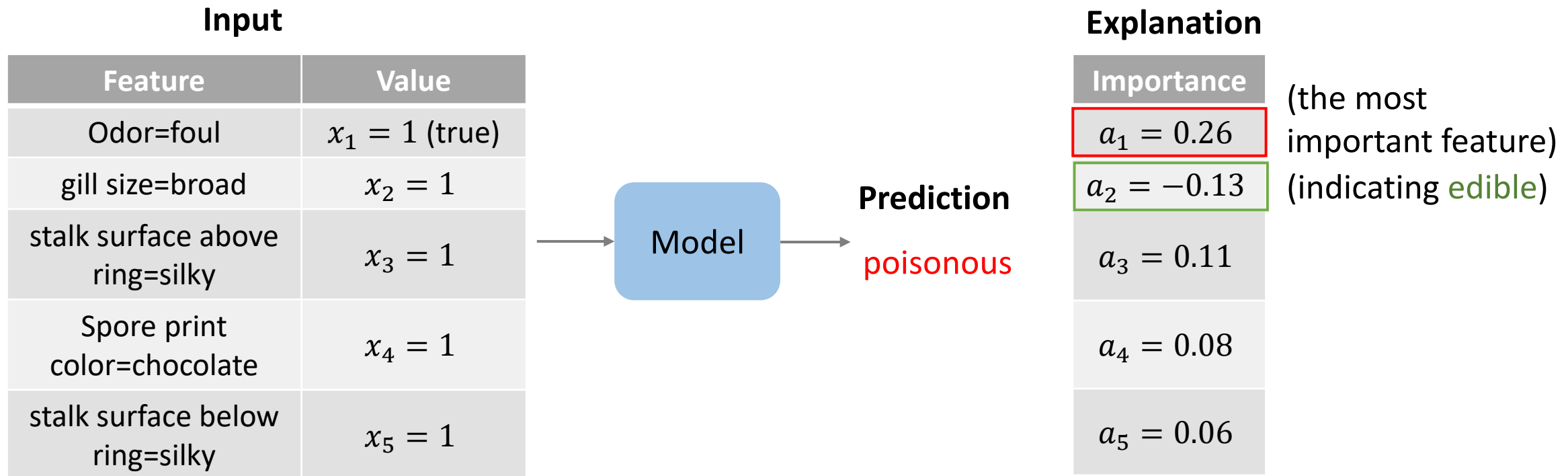
Post-hoc Explanation

- Example 1: tabular data
 - Mushroom dataset
 - **Task:** predicting if a mushroom is **edible** or **poisonous**



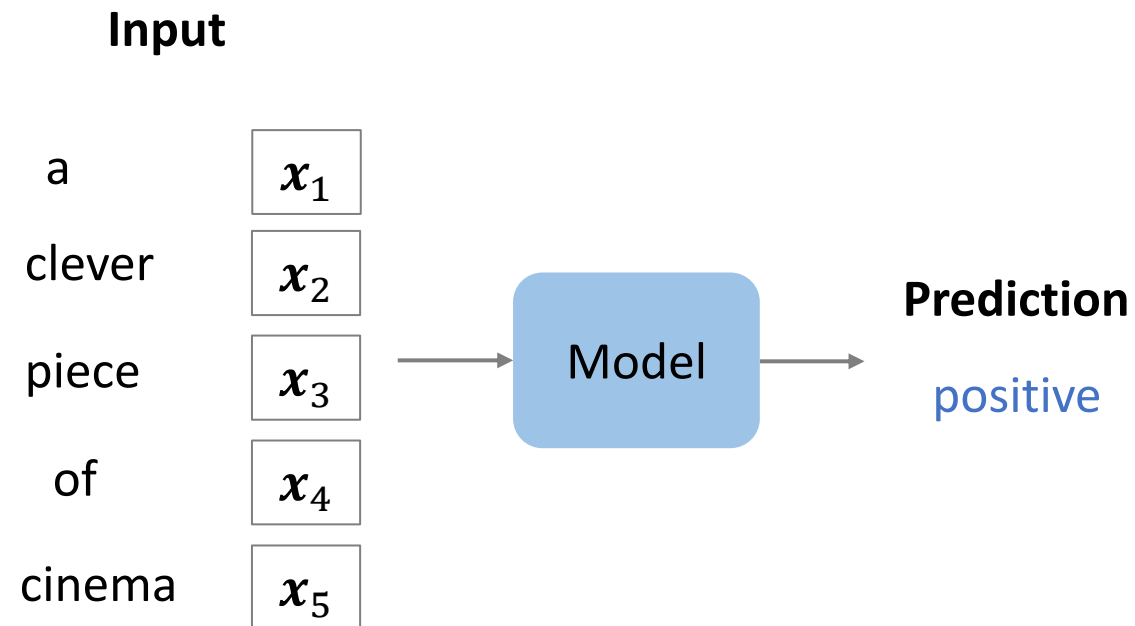
Post-hoc Explanation

- Example 1: tabular data
 - Mushroom dataset
 - **Task:** predicting if a mushroom is **edible** or **poisonous**



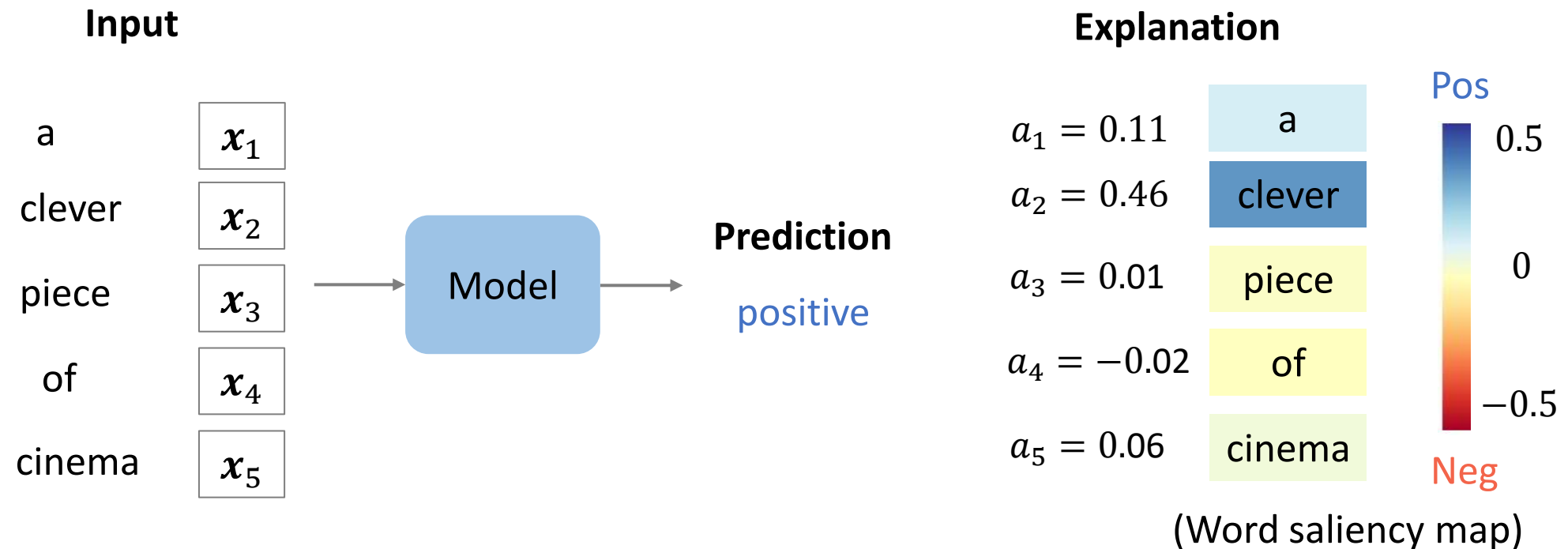
Post-hoc Explanation

- Example 2: text data
 - Movie review
 - **Task:** predicting the sentiment of a text (**positive** or **negative**)



Post-hoc Explanation

- Example 2: text data
 - Movie review
 - **Task:** predicting the sentiment of a text (positive or negative)



Post-hoc Explanation

- Example 3: image data

Task: Object recognition



Feature: a pixel x_{ij}
(color, intensity...)

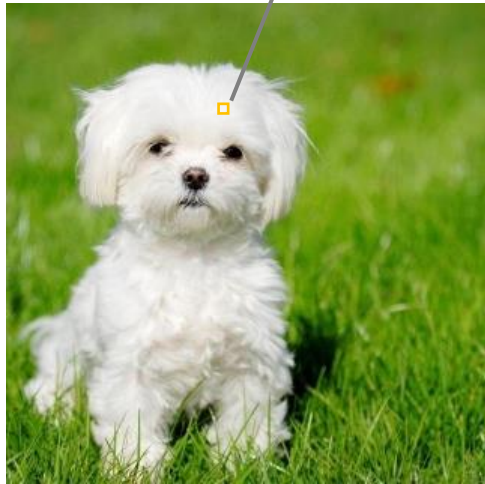
Post-hoc Explanation

- Example 3: image data

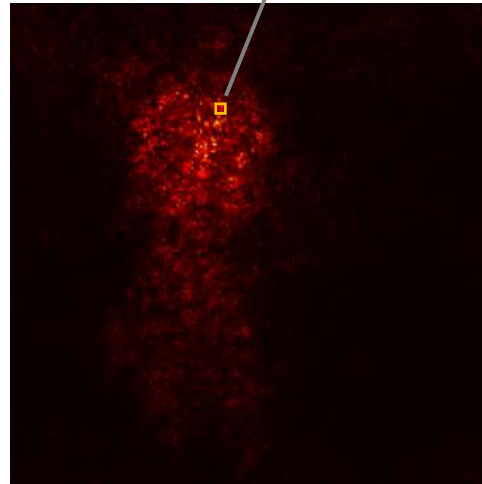
Task: Object recognition

Prediction: Dog

Input



Explanation



Saliency map: the lighter color, the larger value

Post-hoc Explanation

How to learn feature importance?

Perturbation-based methods

- Model-agnostic (black-box): not requiring access to model inner working
- Local: explaining model prediction per example

Perturbation-based methods

- LIME (Ribeiro et al., KDD, 2016)
- SHAP (Lundberg and Lee, NIPS, 2017)

"Why Should I Trust You?" Explaining the Predictions of Any Classifier

Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin

(KDD, 2016)

Interpretable Model

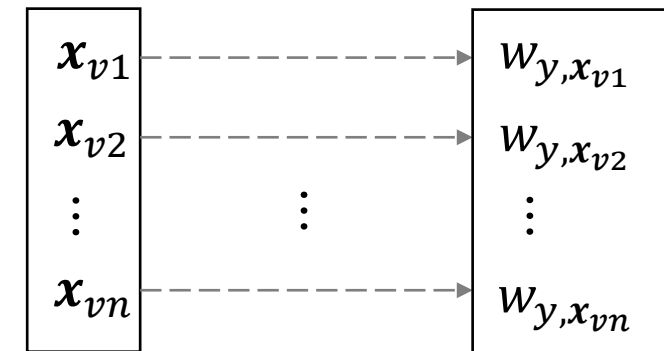
- Linear model

$$h_y(\mathbf{x}) = \mathbf{w}_y^T \mathbf{x} \quad \mathbf{x} \in \{0, 1\}^n$$

- $w_{y,j}$: the contribution of x_j
- Higher weights indicate more important features

Global interpretation

Feature Importance



Interpretable Model

- Linear model

$$h_y(\mathbf{x}) = \mathbf{w}_y^T \mathbf{x} \quad \mathbf{x} \in \{0, 1\}^n$$

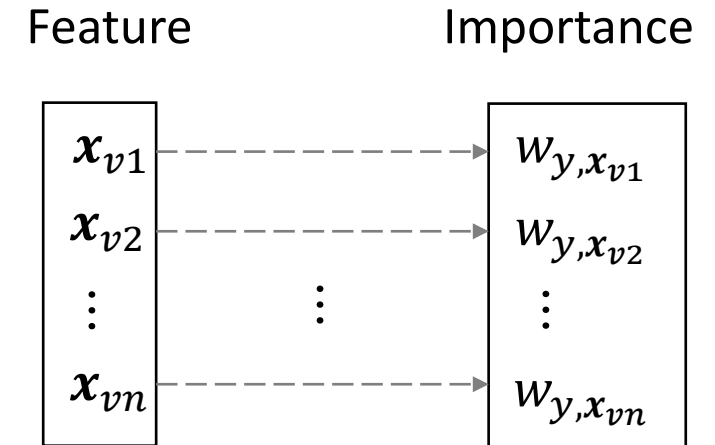
- $w_{y,j}$: the contribution of x_j
- Higher weights indicate more important features

Logistic regression

		“It”	“is”	“a”	<u>“fantastic”</u>	“movie”	
[Neg]	w_0	0.89	0.72	1.13	-1.92	0.34	1.16
[Pos]	w_1	0.85	0.82	1.05	2.21	0.26	5.19

Prediction: positive

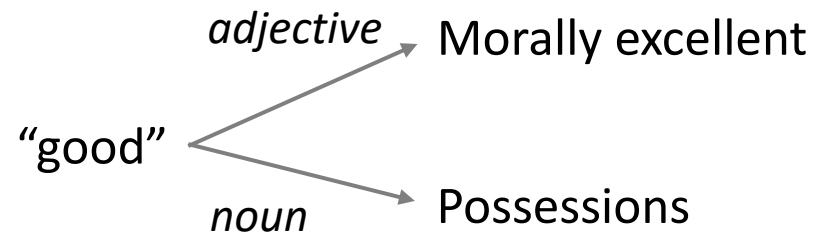
Global interpretation



Neural Networks

Global interpretation is not capable of explaining each specific model prediction

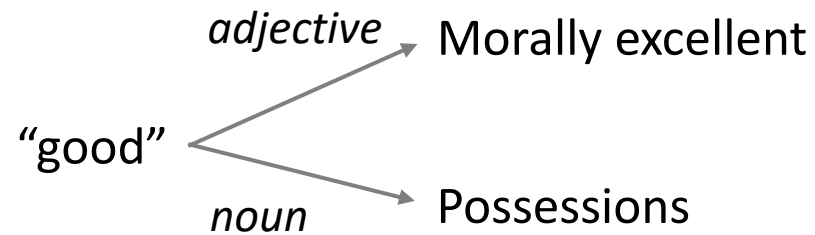
- Neural networks can capture complex relationships between features and the response
- The meaning of a feature may vary across different examples



Neural Networks

Global interpretation is not capable of explaining each specific model prediction

- Neural networks can capture complex relationships between features and the response
- The meaning of a feature may vary across different examples



Local interpretation

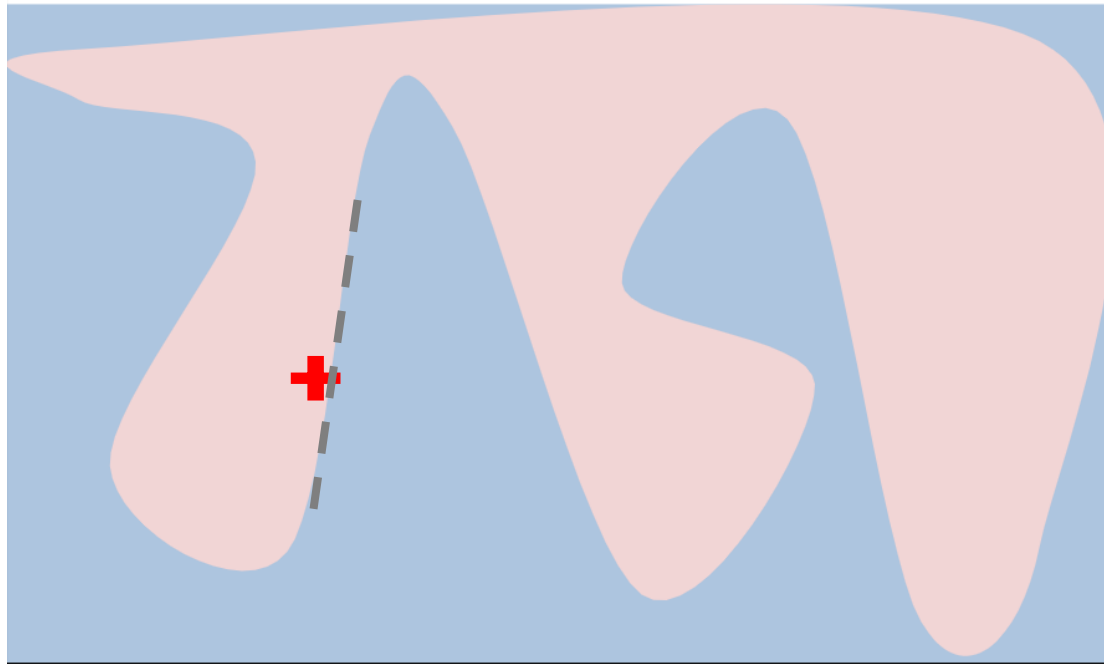
Explaining model prediction per example by identifying local feature importance

LIME: *Local Interpretable Model-Agnostic Explanations*

The way that explains model predictions or the generated explanations are understandable to humans

LIME: *Local Interpretable Model-Agnostic Explanations*

Idea: using local linear model to approximate neural network for each example



- Decision boundary of a neural network f
- Blue/pink background represents negative (-) /positive (+) class
- Bold red cross: the instance x being explained
- Dashed line: local linear model g

$$g \approx f$$

LIME: *Local Interpretable Model-Agnostic Explanations*

- Interpretable data representations

Neural network f

$$\mathbf{x} = [x_1, x_2, \dots, x_n]$$

Feature representation

$x_i \in \mathbb{R}^d$ is uninterpretable

Linear model g

$$\mathbf{x}' = [x'_1, x'_2, \dots, x'_N]$$

Feature representation

$x'_i \in \{0, 1\}$ is interpretable

- n : the number of features in the example
- N : the number of all features

LIME: Local Interpretable Model-Agnostic Explanations

- Interpretable data representations

Neural network f

$$\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$$

Feature representation

$\mathbf{x}_i \in \mathbb{R}^d$ is uninterpretable

Image

\mathbf{x}_i : a tensor with three color channels per pixel

Text

\mathbf{x}_i : a high-dimensional vector (word embedding)

Linear model g

$$\mathbf{x}' = [x'_1, x'_2, \dots, x'_N]$$

Feature representation

$x'_i \in \{0, 1\}$ is interpretable

Image

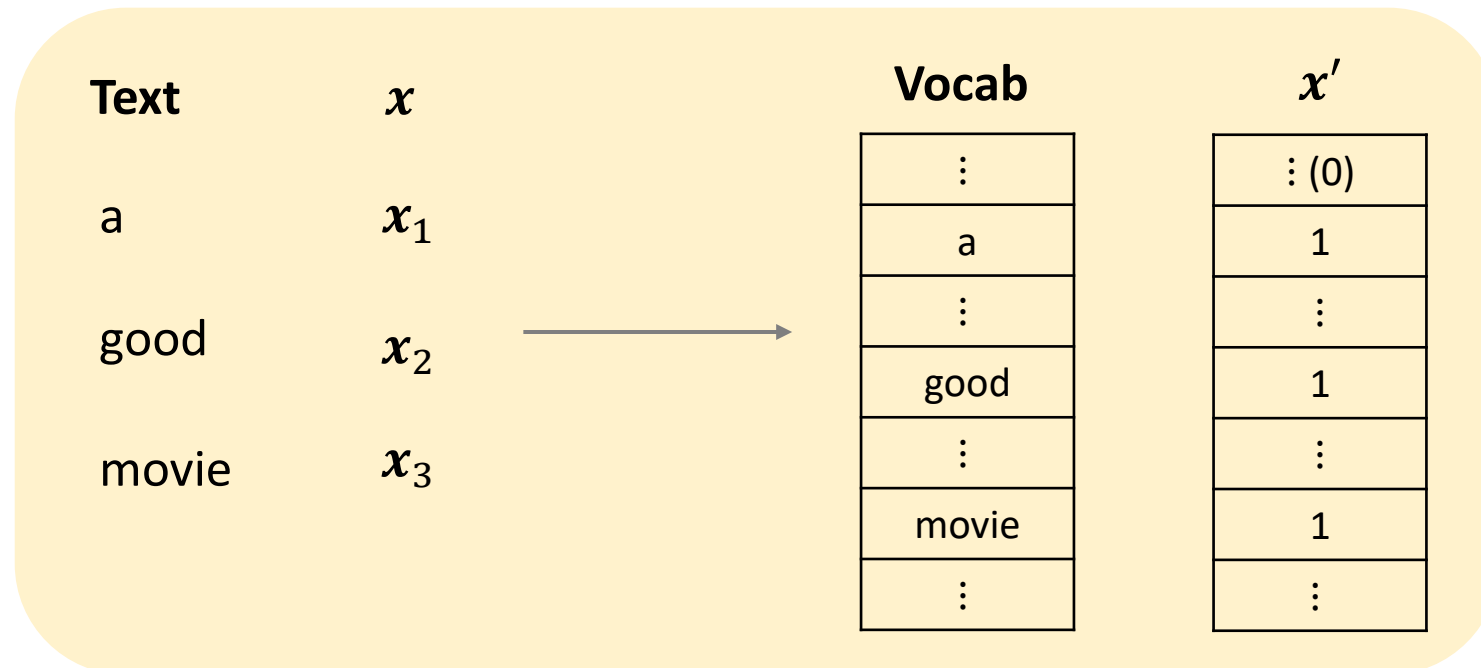
0/1 indicates the absence/presence of a patch of pixels

Text

0/1 indicates the absence/presence of a word
(bag-of-words representation)

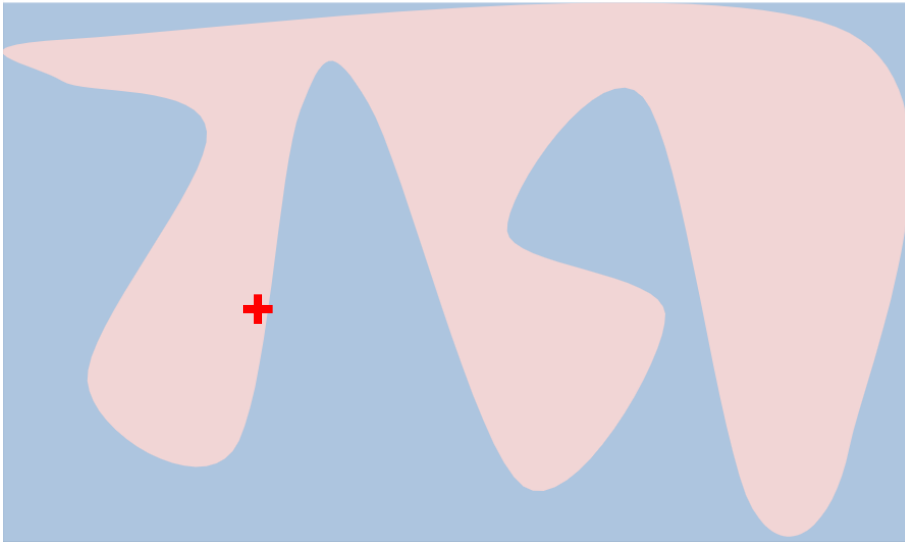
LIME: Local Interpretable Model-Agnostic Explanations

- Interpretable data representations



LIME: *Local Interpretable Model-Agnostic Explanations*

- Sampling for local exploration



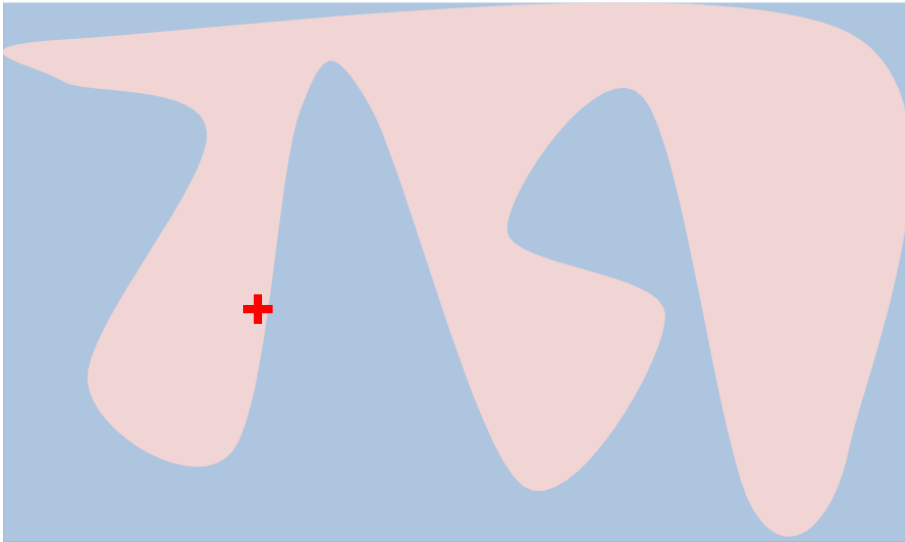
Need more samples to fit a local linear model

It is a fantastic movie

$$\mathbf{x}' = [0, \dots, 1, \dots, 1, \dots, 1, \dots, 1, \dots, 0, 1, \dots, 0]_N$$

LIME: Local Interpretable Model-Agnostic Explanations

- Sampling for local exploration



Need more samples to fit a local linear model

It is a fantastic movie

$$\mathbf{x}' = [0, \dots, 1, \dots, 1, \dots, 1, \dots, 1, \dots, 0, 1, \dots, 0]_N$$



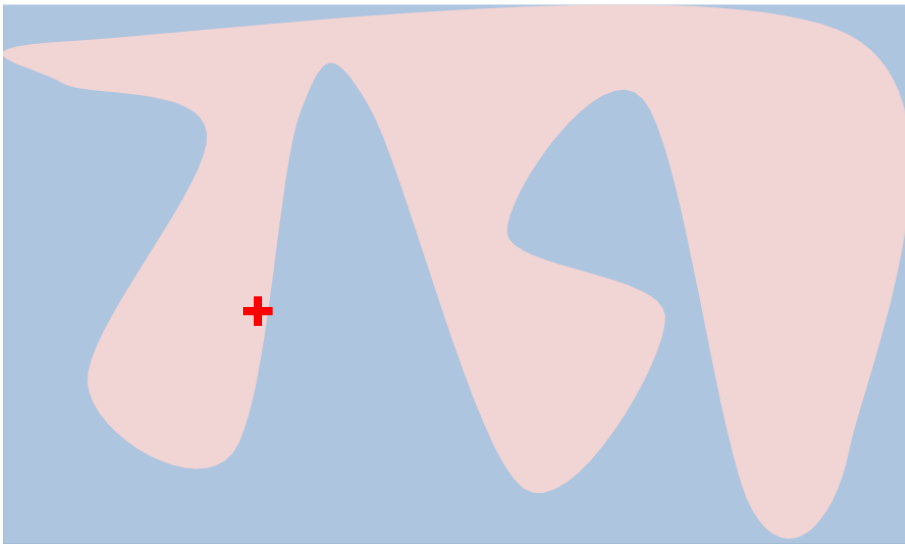
Randomly sample nonzero elements

a movie

$$\mathbf{z}_1' = [0, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 0, 1, \dots, 0]_N$$

LIME: Local Interpretable Model-Agnostic Explanations

- Sampling for local exploration



Need more samples to fit a local linear model

It is a fantastic movie

$$\mathbf{x}' = [0, \dots, 1, \dots, 1, \dots, 1, \dots, 1, \dots, 0, 1, \dots, 0]_N$$



Randomly sample nonzero elements

a movie

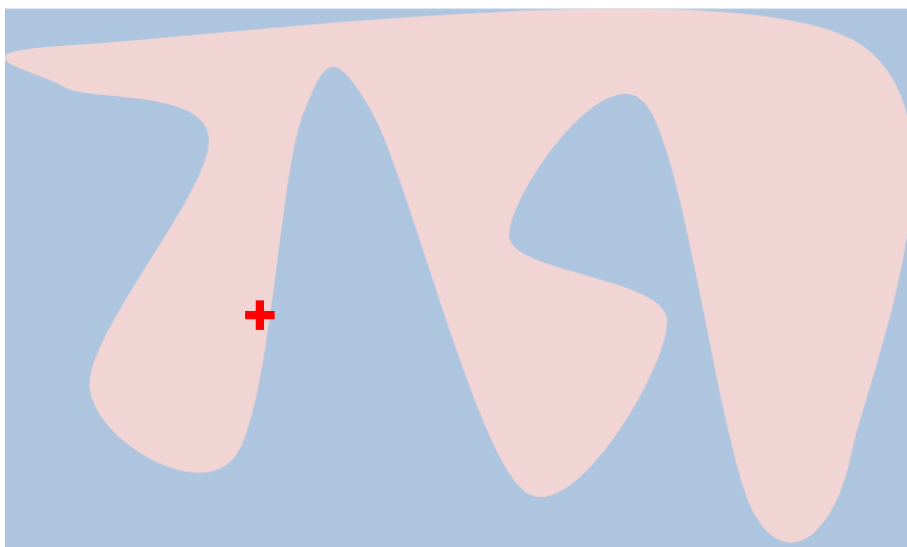
$$\mathbf{z}_1' = [0, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 0, 1, \dots, 0]_N$$

fantastic movie

$$\mathbf{z}_2' = [0, \dots, 0, \dots, 0, \dots, 0, \dots, 1, \dots, 0, 1, \dots, 0]_N$$

LIME: Local Interpretable Model-Agnostic Explanations

- Sampling for local exploration



Need more samples to fit a local linear model

It is a fantastic movie

$$\mathbf{x}' = [0, \dots, 1, \dots, 1, \dots, 1, \dots, 1, \dots, 0, 1, \dots, 0]_N$$



Randomly sample nonzero elements

a movie

$$\mathbf{z}_1' = [0, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 0, 1, \dots, 0]_N$$

fantastic movie

$$\mathbf{z}_2' = [0, \dots, 0, \dots, 0, \dots, 0, \dots, 1, \dots, 0, 1, \dots, 0]_N$$

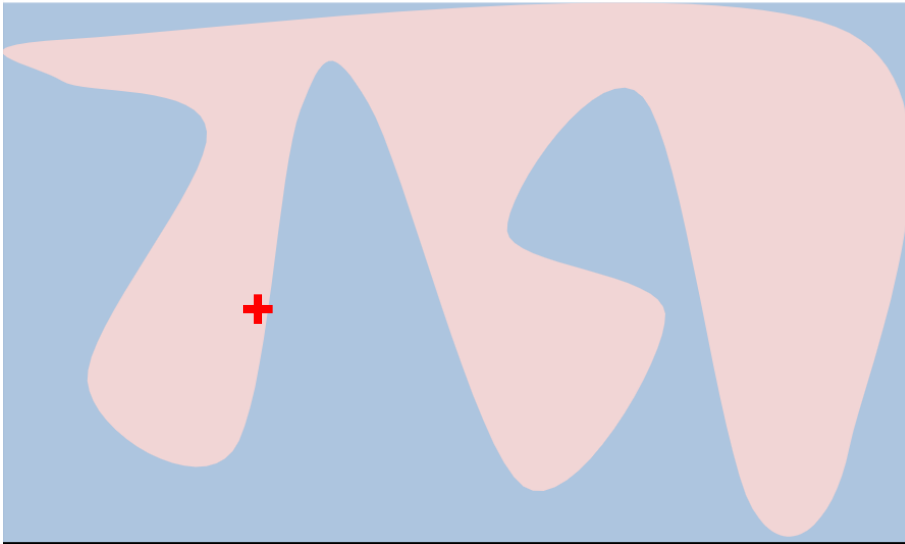
⋮

fantastic

$$\mathbf{z}_M' = [0, \dots, 0, \dots, 0, \dots, 0, \dots, 1, \dots, 0, 0, \dots, 0]_N$$

LIME: Local Interpretable Model-Agnostic Explanations

- Sampling for local exploration



What are the labels of these pseudo examples?

Need more samples to fit a local linear model

It is a fantastic movie

$$\mathbf{x}' = [0, \dots, 1, \dots, 1, \dots, 1, \dots, 1, \dots, 0, 1, \dots, 0]_N$$



Randomly sample nonzero elements

a movie

$$\mathbf{z}_1' = [0, \dots, 0, \dots, 0, \dots, 1, \dots, 0, \dots, 0, 1, \dots, 0]_N$$

fantastic movie

$$\mathbf{z}_2' = [0, \dots, 0, \dots, 0, \dots, 0, \dots, 1, \dots, 0, 1, \dots, 0]_N$$

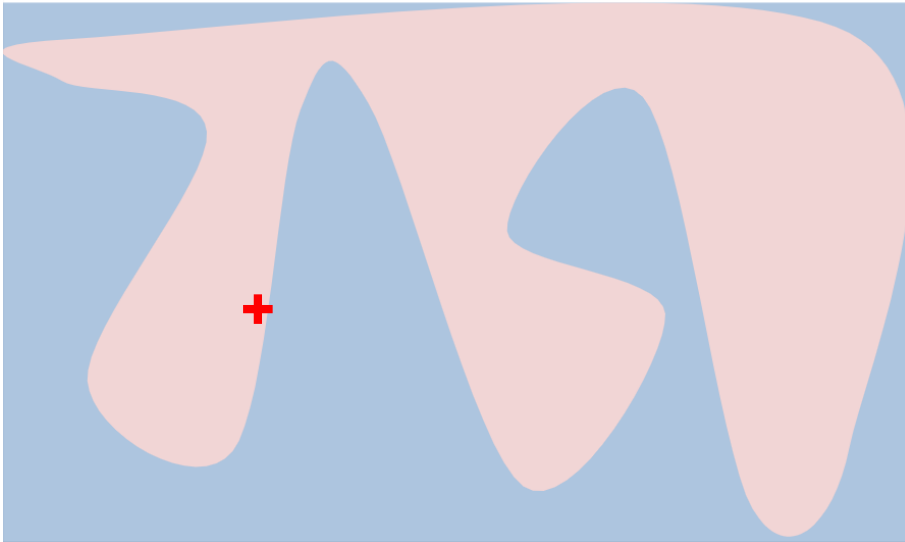
⋮

fantastic

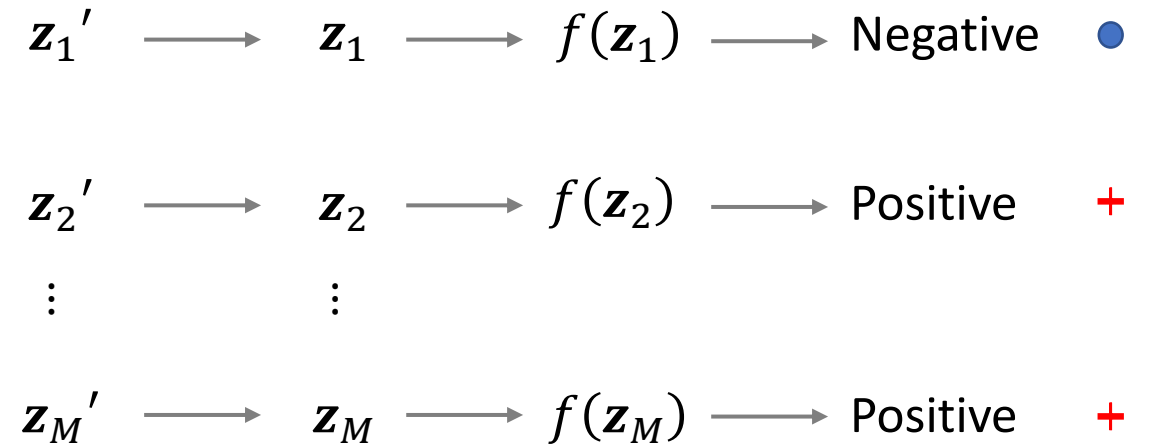
$$\mathbf{z}_M' = [0, \dots, 0, \dots, 0, \dots, 0, \dots, 1, \dots, 0, 0, \dots, 0]_N$$

LIME: Local Interpretable Model-Agnostic Explanations

- Sampling for local exploration

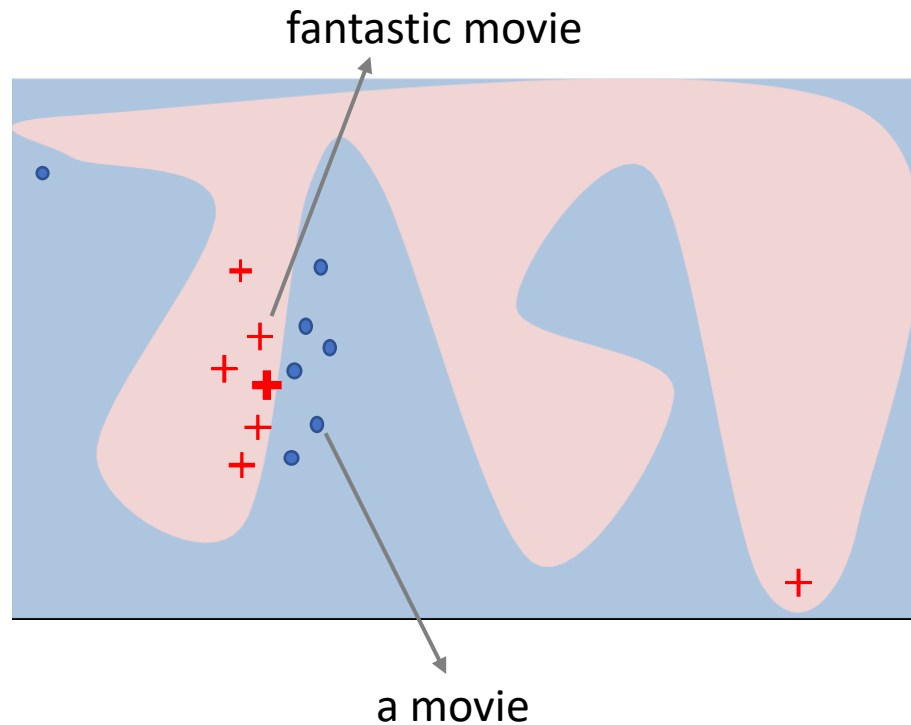


Labeling pseudo examples with neural network f

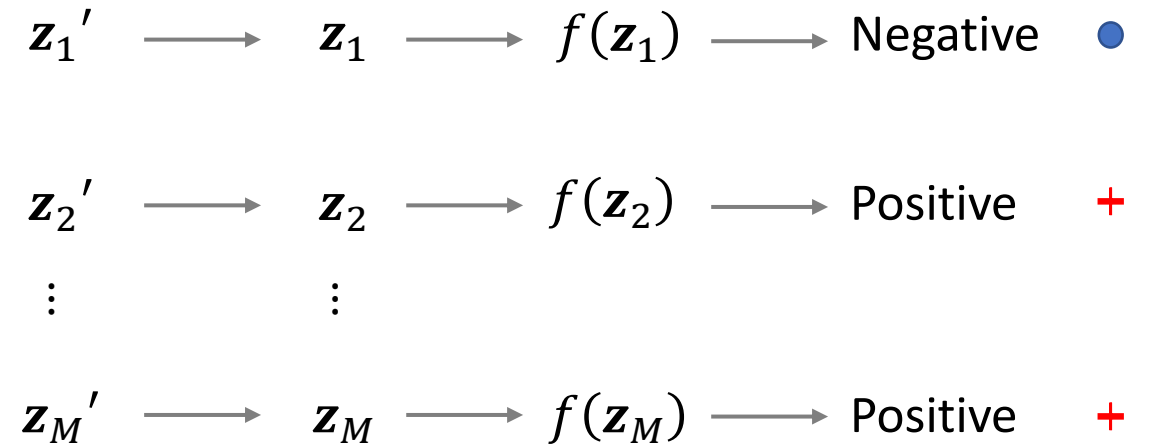


LIME: *Local Interpretable Model-Agnostic Explanations*

- Sampling for local exploration

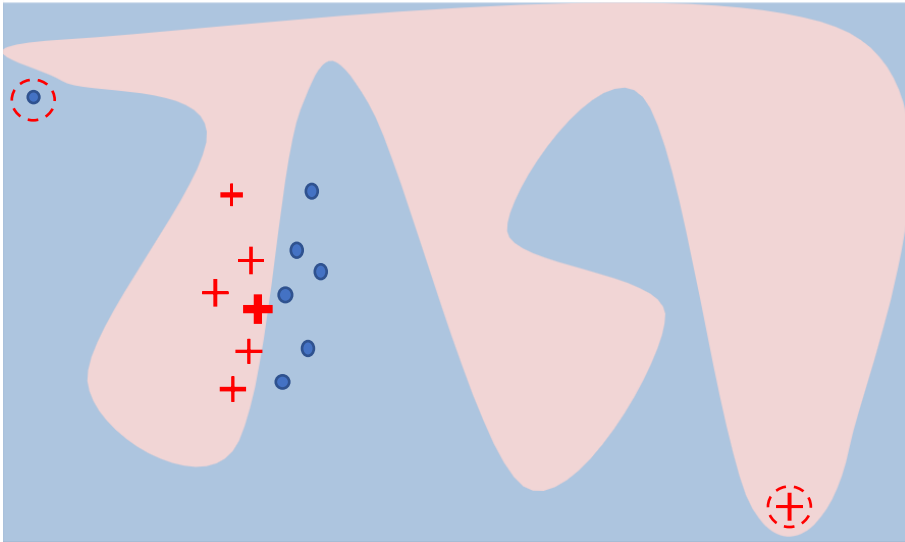


Labeling pseudo examples with neural network f



LIME: Local Interpretable Model-Agnostic Explanations

- Sampling for local exploration



Penalize noisy examples

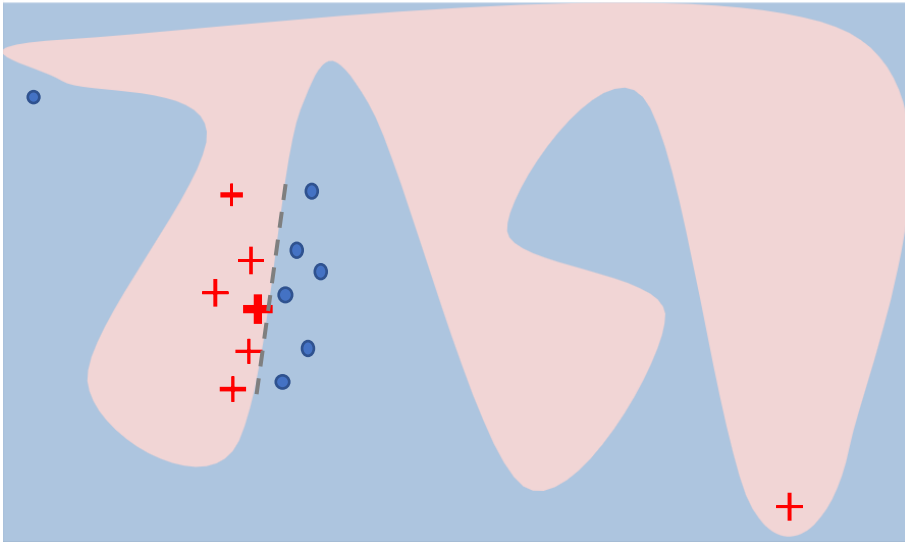
Distance between x and z_m

$$\pi_x(z_m) = e^{(-D(x,z_m)^2/\sigma^2)}$$

D : cosine distance (for text), L_2 distance (for image)

LIME: *Local Interpretable Model-Agnostic Explanations*

- Sparse linear explanation



Fitting a local linear model

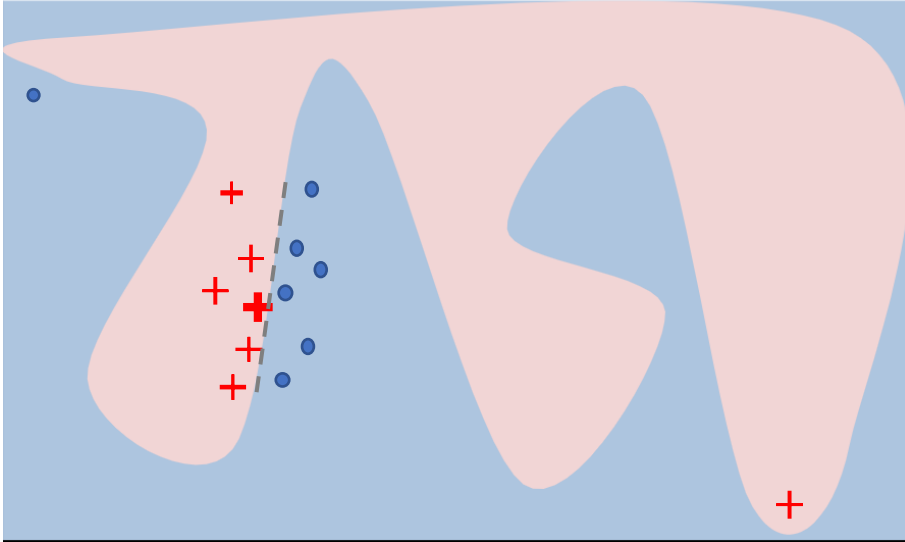
$$\{(\mathbf{z}_m', f(\mathbf{z}_m))\}_{m=1, \dots, M}$$

$$g(\mathbf{z}') \approx f(\mathbf{z})$$

$$g(\mathbf{z}') = \mathbf{w}^T \mathbf{z}'$$

LIME: Local Interpretable Model-Agnostic Explanations

- Sparse linear explanation



Fitting a local linear model

$$\{(\mathbf{z}_m', f(\mathbf{z}_m))\}_{m=1, \dots, M}$$

$$g(\mathbf{z}') \approx f(\mathbf{z})$$

$$g(\mathbf{z}') = \mathbf{w}^T \mathbf{z}'$$

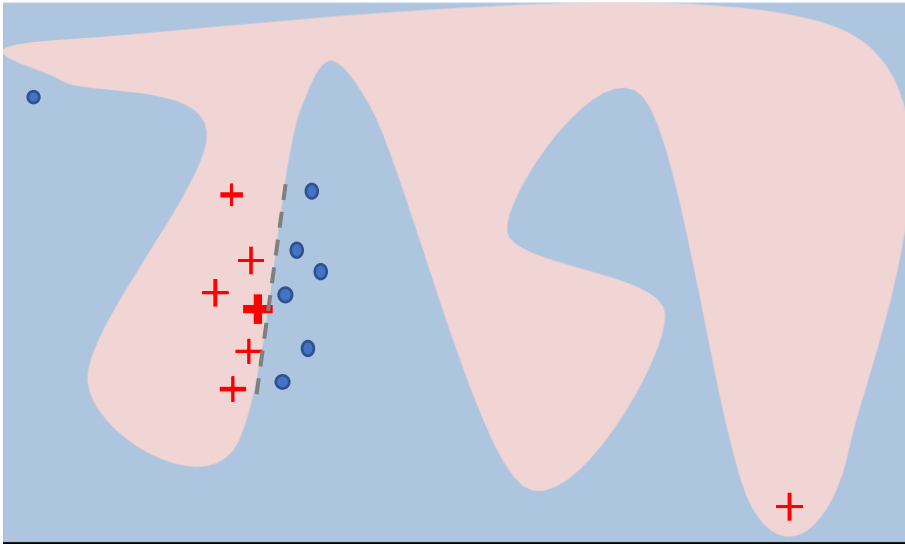
Objective

$$\min \mathcal{L}(f, g)$$

$$\mathcal{L}(f, g) = \sum \pi_x(\mathbf{z})(f(\mathbf{z}) - g(\mathbf{z}'))^2$$

LIME: Local Interpretable Model-Agnostic Explanations

- Sparse linear explanation



Fitting a local linear model

$$\{(\mathbf{z}_m', f(\mathbf{z}_m))\}_{m=1, \dots, M}$$

$$g(\mathbf{z}') \approx f(\mathbf{z})$$

$$g(\mathbf{z}') = \mathbf{w}^T \mathbf{z}'$$

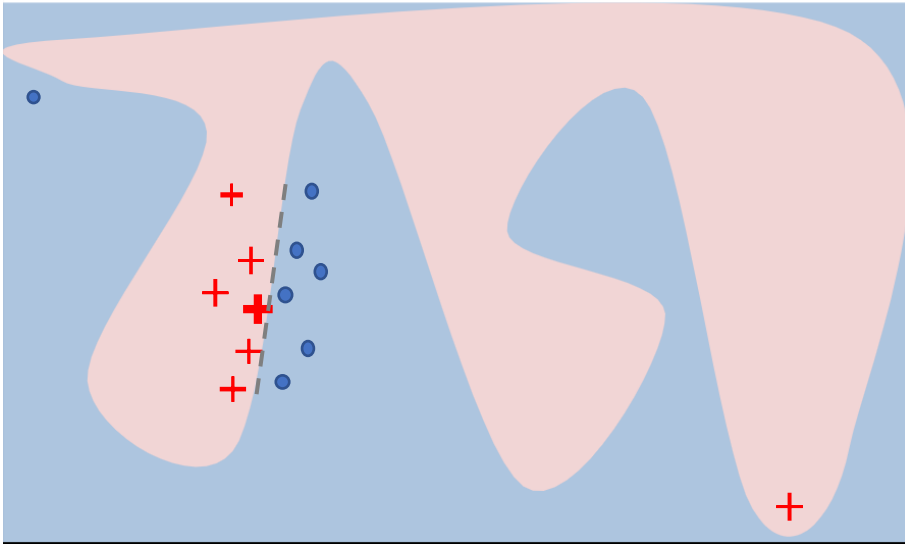
Objective

$$\min \mathcal{L}(f, g) + \Omega(g) \quad \text{Restricting complexity (the number of nonzero weights)}$$

$$\mathcal{L}(f, g) = \sum \pi_x(\mathbf{z})(f(\mathbf{z}) - g(\mathbf{z}'))^2$$

LIME: *Local Interpretable Model-Agnostic Explanations*

- Sparse linear explanation



Extracting feature importance scores

$$\mathbf{w}_{\hat{y}}^T$$

- \hat{y} : model prediction on the original example
- Local explanation: $\{w_{\hat{y},x_1}, \dots, w_{\hat{y},x_n}\}$

Question?

LIME: *Local Interpretable Model-Agnostic Explanations*

- Explaining each example individually, not the whole dataset (**locally faithful**)
- May not work for highly non-linear models

LIME: *Local Interpretable Model-Agnostic Explanations*

- Example: Deep networks for image classification

Model: Google's pre-trained Inception neural network

Top 3 predicted classes

Electric guitar

Acoustic guitar

Labrador



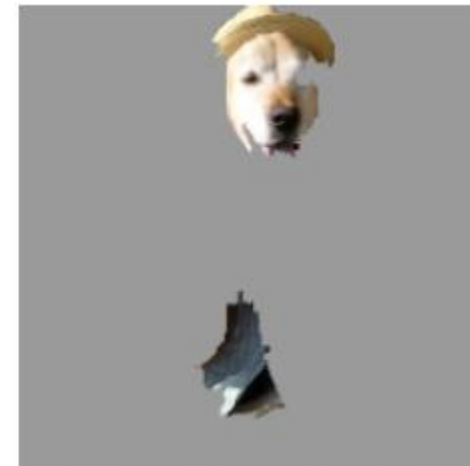
(a) Original Image



(b) Explaining *Electric guitar*



(c) Explaining *Acoustic guitar*



(d) Explaining *Labrador*

The explanations enhance trust in the model, as it acts in a reasonable manner

Question?

Submodular Pick for Explaining Models

Single explanation is not sufficient to evaluate and assess trust in the model as a whole



Providing a global understanding of the model by explaining a set of individual instances

Submodular Pick for Explaining Models

Single explanation is not sufficient to evaluate and assess trust in the model as a whole



Providing a global understanding of the model by explaining a set of individual instances

How to select these instances judiciously?

Submodular Pick for Explaining Models

- Budget B : the number of explanations users are willing to look at in order to understand a model

A set of instances X $\xrightarrow{\text{select}}$ B instances (diverse, representative)

Submodular Pick for Explaining Models

- Budget B : the number of explanations users are willing to look at in order to understand a model

A set of instances X $\xrightarrow{\text{select}}$ B instances (diverse, representative)

➤ Submodular pick (SP) algorithm

$x^{(1)}$ contains two features

Explanation matrix W

	x_1	x_2	x_3	x_4	x_5
$x^{(1)}$	$w_{1,1}$	$w_{1,2}$			
$x^{(2)}$		$w_{2,2}$	$w_{2,3}$		
$x^{(3)}$		$w_{3,2}$	$w_{3,3}$		
$x^{(4)}$		$w_{4,2}$		$w_{4,4}$	
$x^{(5)}$				$w_{5,4}$	$w_{5,5}$

- Each row represents an instance
- Each column represents a feature
- Each value represents a local importance

Submodular Pick for Explaining Models

- Budget B : the number of explanations users are willing to look at in order to understand a model

A set of instances X $\xrightarrow{\text{select}}$ B instances (diverse, representative)

➤ Submodular pick (SP) algorithm

Explanation matrix W

	x_1	x_2	x_3	x_4	x_5
$x^{(1)}$	$w_{1,1}$	$w_{1,2}$			
$x^{(2)}$		$w_{2,2}$	$w_{2,3}$		
$x^{(3)}$		$w_{3,2}$	$w_{3,3}$		
$x^{(4)}$		$w_{4,2}$		$w_{4,4}$	
$x^{(5)}$				$w_{5,4}$	$w_{5,5}$

- ❑ Select instances that cover important features (x_2)

I_2 : global importance of x_2

Submodular Pick for Explaining Models

- Budget B : the number of explanations users are willing to look at in order to understand a model

A set of instances X $\xrightarrow{\text{select}}$ B instances (diverse, representative)

➤ Submodular pick (SP) algorithm

Explanation matrix W

	x_1	x_2	x_3	x_4	x_5
$x^{(1)}$	$w_{1,1}$	$w_{1,2}$			
$x^{(2)}$		$w_{2,2}$	$w_{2,3}$		
$x^{(3)}$		$w_{3,2}$	$w_{3,3}$		
$x^{(4)}$		$w_{4,2}$		$w_{4,4}$	
$x^{(5)}$				$w_{5,4}$	$w_{5,5}$

If $x^{(2)}$ is selected,
there is no need to
select $x^{(3)}$

- ❑ Select instances that cover important features (x_2)
- ❑ Avoid selecting instances with similar explanations (redundant features)

Submodular Pick for Explaining Models

- Budget B : the number of explanations users are willing to look at in order to understand a model

A set of instances X $\xrightarrow{\text{select}}$ B instances (diverse, representative)

➤ Submodular pick (SP) algorithm

Explanation matrix W

	x_1	x_2	x_3	x_4	x_5
$x^{(1)}$	$w_{1,1}$	$w_{1,2}$			
$x^{(2)}$		$w_{2,2}$	$w_{2,3}$		
$x^{(3)}$		$w_{3,2}$	$w_{3,3}$		
$x^{(4)}$		$w_{4,2}$		$w_{4,4}$	
$x^{(5)}$				$w_{5,4}$	$w_{5,5}$

- Select instances that cover important features (x_2)
- Avoid selecting instances with similar explanations (redundant features)
- Select less instances, while covering more features ($x^{(2)}$ and $x^{(5)}$)

Submodular Pick for Explaining Models

➤ Submodular pick (SP) algorithm

Algorithm Submodular pick (SP) algorithm

Require: Instances X , Budget B

for all $x^{(i)}$ in X **do**

$w^{(i)} \leftarrow LIME(x^{(i)})$

Construct the explanation matrix

Submodular Pick for Explaining Models

➤ Submodular pick (SP) algorithm

Algorithm Submodular pick (SP) algorithm

Require: Instances X , Budget B

for all $x^{(i)}$ in X **do**

$$w^{(i)} \leftarrow LIME(x^{(i)})$$

Construct the explanation matrix

for $j \in \{1, \dots, N\}$ **do**

$$I_j \leftarrow \sqrt{\sum_{i=1}^{|X|} w_{i,j}}$$

Compute global feature importance

Submodular Pick for Explaining Models

➤ Submodular pick (SP) algorithm

Algorithm Submodular pick (SP) algorithm

Require: Instances X , Budget B

for all $\mathbf{x}^{(i)}$ in X **do**

$$\mathbf{w}^{(i)} \leftarrow \text{LIME}(\mathbf{x}^{(i)})$$

Construct the explanation matrix

for $j \in \{1, \dots, N\}$ **do**

$$I_j \leftarrow \sqrt{\sum_{i=1}^{|\mathcal{X}|} w_{i,j}}$$

Compute global feature importance

$V \leftarrow \{\}$

while $|V| < B$ **do**

Greedly add examples that maximize the coverage gain

$$i^* = \underset{i}{\operatorname{argmax}} c(V \cup \mathbf{x}^{(i)}, W, I) \quad c(V, W, I) = \sum_{j=1}^N \mathbf{1}_{[\exists i \in V: w_{i,j} > 0]} I_j$$

$$V \leftarrow V \cup \mathbf{x}^{(i^*)}$$

The coverage computes the total global importance of the features that appear in at least one instance in a set V

Return V

SP-LIME

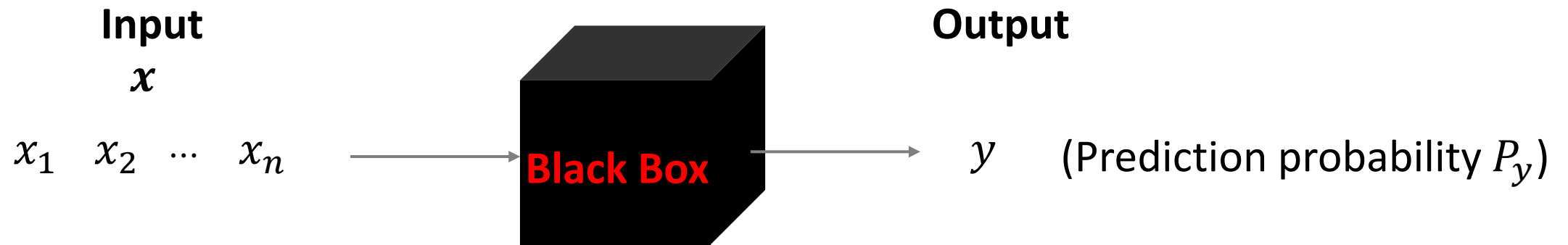
- Compute global feature importance based on local feature importance from LIME
- Provide a global understanding of the model by selecting a set of representative instances

Question?

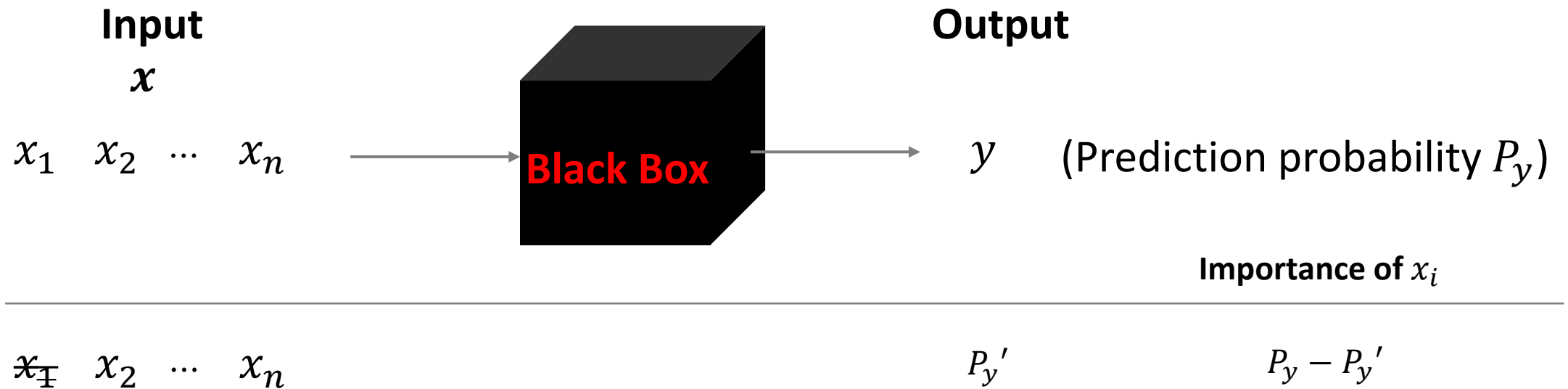
Perturbation-based methods

- LIME (Ribeiro et al., KDD, 2016)
- SHAP (Lundberg and Lee, NIPS, 2017)

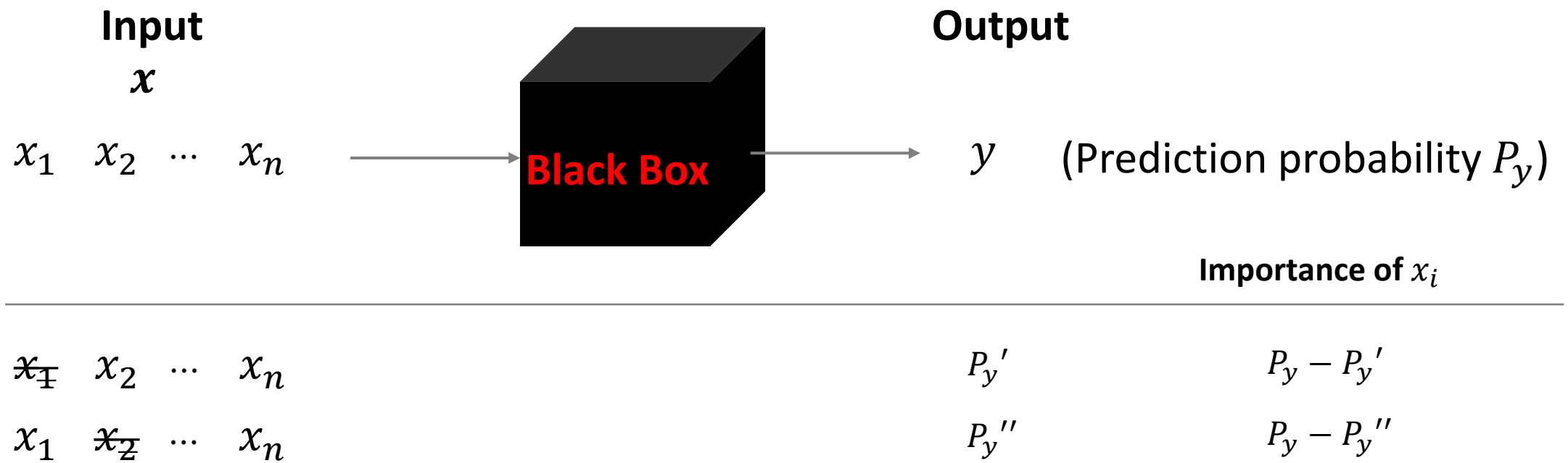
Explaining Black-box Model



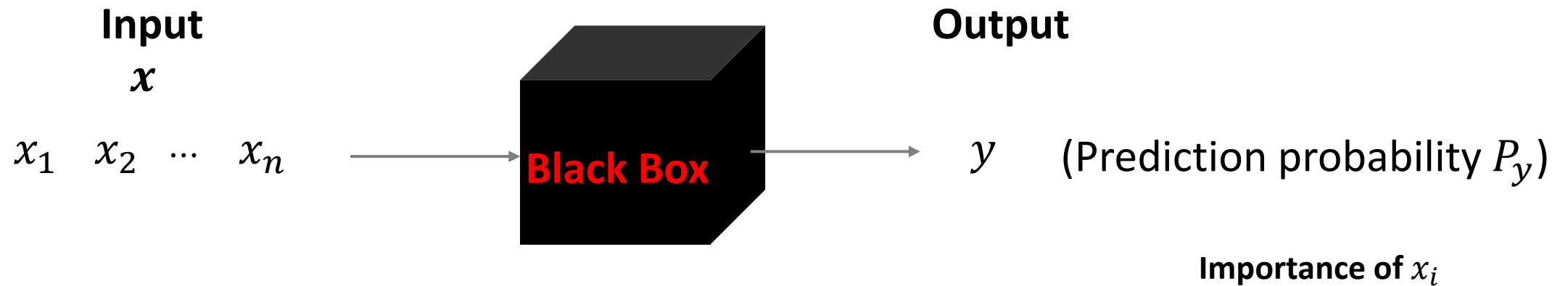
Explaining Black-box Model



Explaining Black-box Model



Explaining Black-box Model



				Importance of x_i	
x_1	x_2	\dots	x_n	$P_{y'}$	$P_y - P_{y'}$
x_1	x_2	\dots	x_n	$P_{y''}$	$P_y - P_{y''}$
	\vdots			\vdots	\vdots

[Leave-one-out, (Li et al., 2016)]

Leave-one-out

- Sentiment classification

Model prediction: positive

Text	Confidence	Word importance
The movie is interesting	0.98	
The movie is interesting	0.95	The 0.03
The mov ie is interesting	0.87	movie 0.11
The movie is interesting	0.96	is 0.02
The movie is interesting	0.61	interesting 0.37

Leave-one-out

- Leave **ONE** feature out at each step

Feature importance may be misleading

Text	Confidence	Word importance
The movie is interesting and impressive	0.97	
The movie is interesting and impressive	0.95	interesting 0.02
The movie is interesting and impressive	0.96	impressive 0.01

Leave-one-out

- Leave **ONE** feature out at each step

Feature importance may be misleading

Text	Confidence	Word importance
The movie is interesting and impressive	0.97	
The movie is interesting and impressive	0.95	interesting 0.02
The movie is interesting and impressive	0.96	impressive 0.01



SHAP

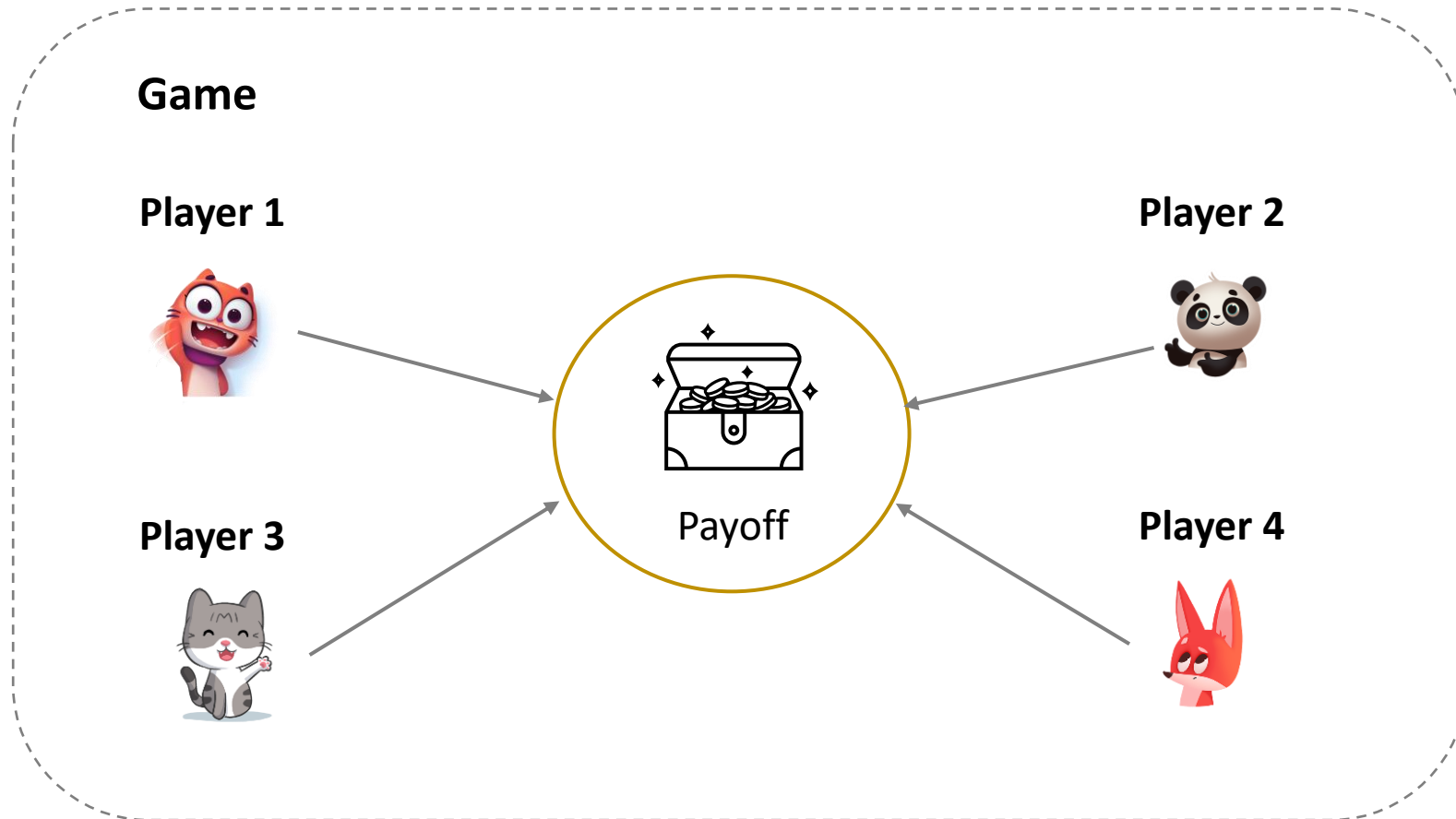
A unified approach to interpreting model predictions

Scott M. Lundberg, Su-In Lee

(NIPS, 2017)

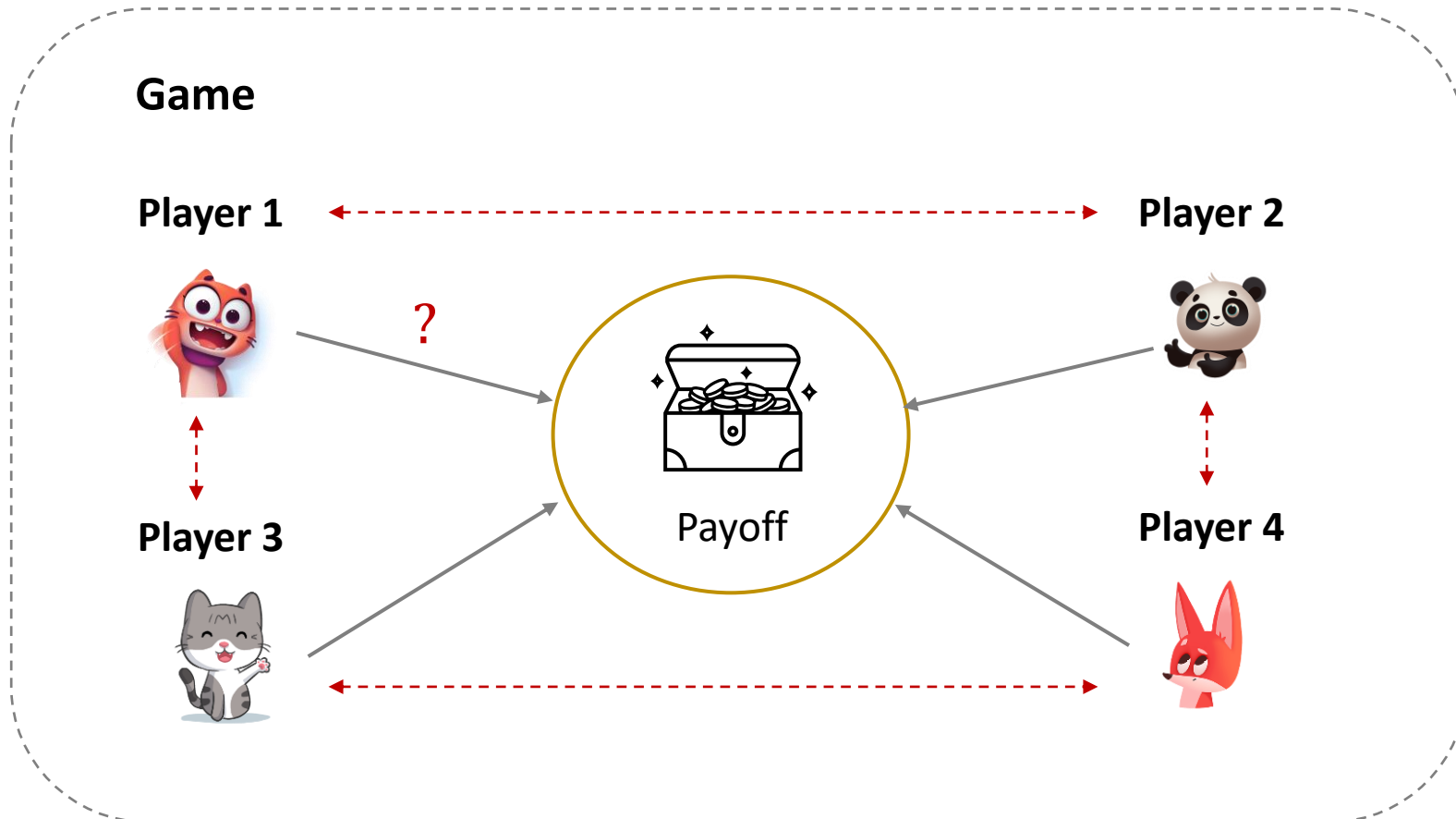
SHAP

- Shapley value [Shapley, 1953]



SHAP

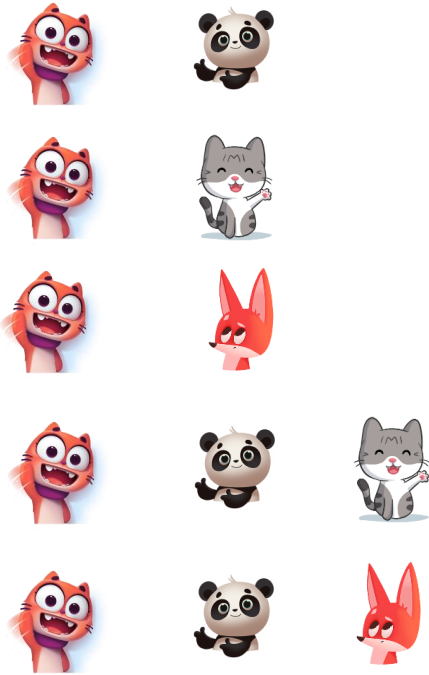
- Shapley value [Shapley, 1953]



SHAP

- Shapley value [Shapley, 1953]

Coalitions



(2^3)

Payoff

P_1

P_2

P_3

P_4

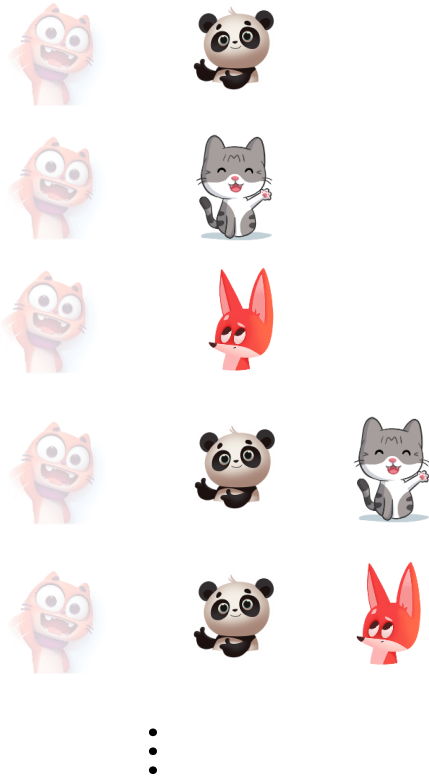
P_5

⋮

SHAP

- Shapley value [Shapley, 1953]

Coalitions



(2^3)

Payoff

P_1 P_1'

P_2 P_2'

P_3 P_3'

P_4 P_4'

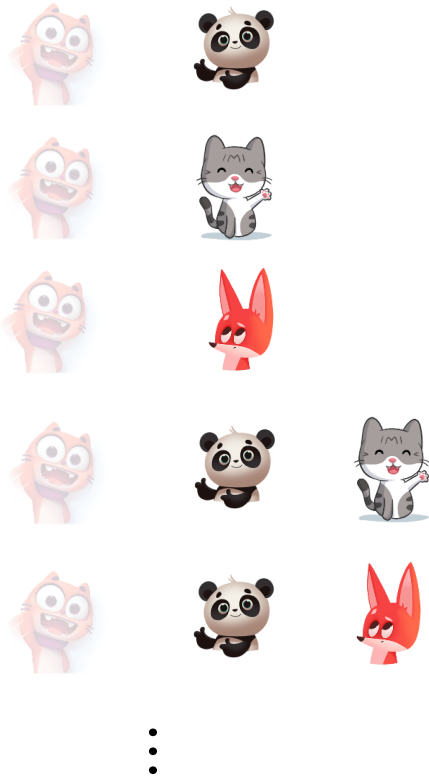
P_5 P_5'

⋮

SHAP

- Shapley value [Shapley, 1953]

Coalitions



(2^3)

Payoff

$$P_1 - P_1'$$

$$P_2 - P_2'$$

$$P_3 - P_3'$$

$$P_4 - P_4'$$

$$P_5 - P_5'$$

⋮

Marginal contribution

$$\Delta P_1$$

$$\Delta P_2$$

$$\Delta P_3$$

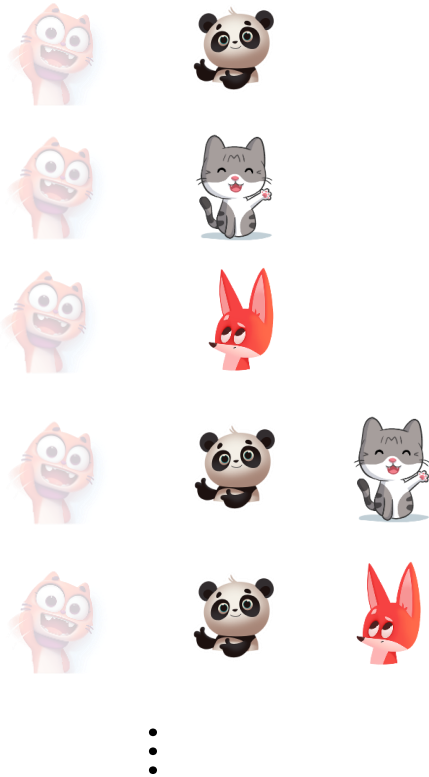
$$\Delta P_4$$

$$\Delta P_5$$

SHAP

- Shapley value [Shapley, 1953]

Coalitions



(2^3)

Payoff

$$P_1 - P_1'$$

$$P_2 - P_2'$$

$$P_3 - P_3'$$

$$P_4 - P_4'$$

$$P_5 - P_5'$$

⋮

Marginal contribution

$$\Delta P_1$$

$$\Delta P_2$$

$$\Delta P_3$$

$$\Delta P_4$$

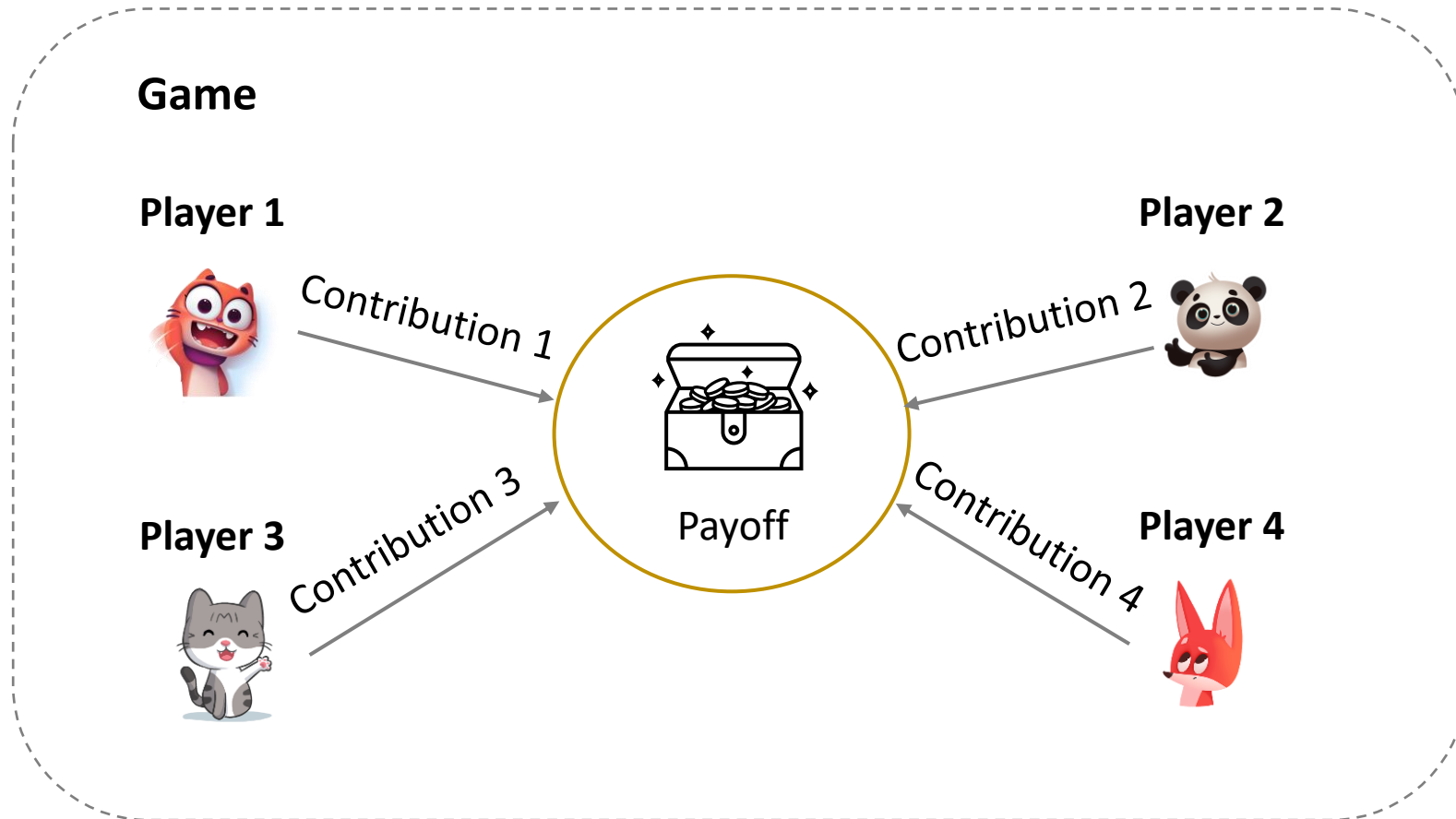
$$\Delta P_5$$



$$\text{Contribution} = \sum \Delta P_i$$

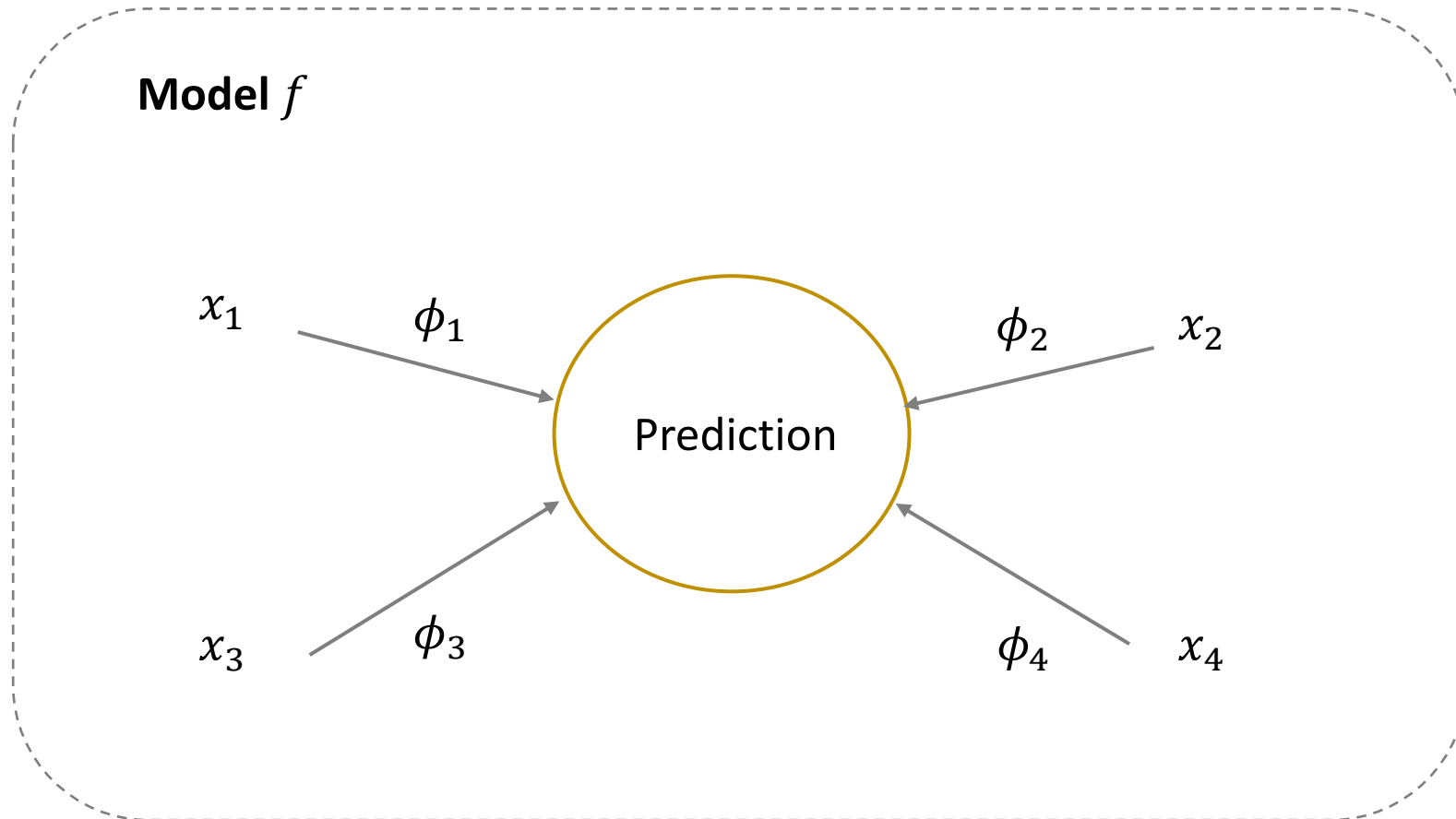
SHAP

- Shapley value [Shapley, 1953]



SHAP

- Shapley value [Shapley, 1953]

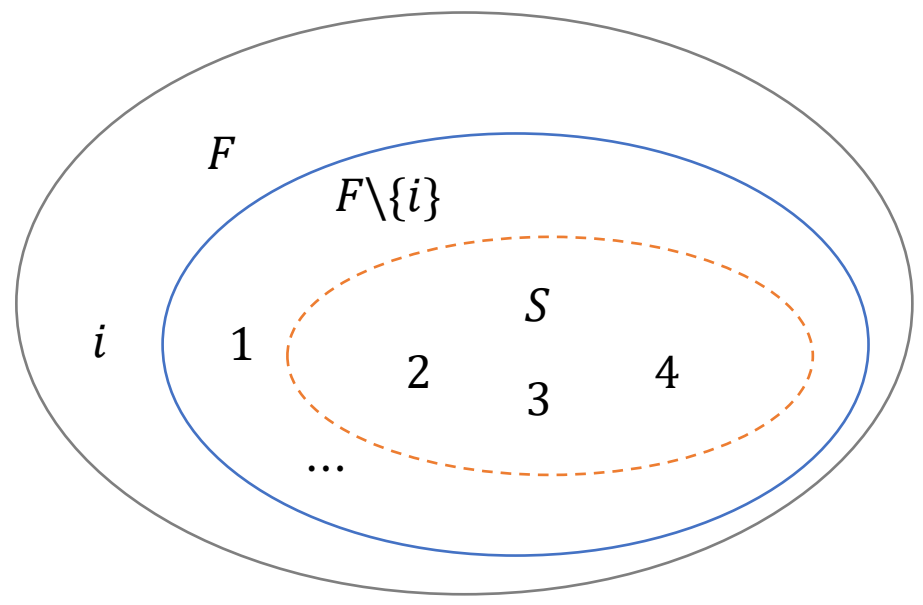


SHAP

- Shapley value [Shapley, 1953]

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|! (|F| - |S| - 1)!}{|F|!} [f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S)]$$

Marginal contribution of x_i given S

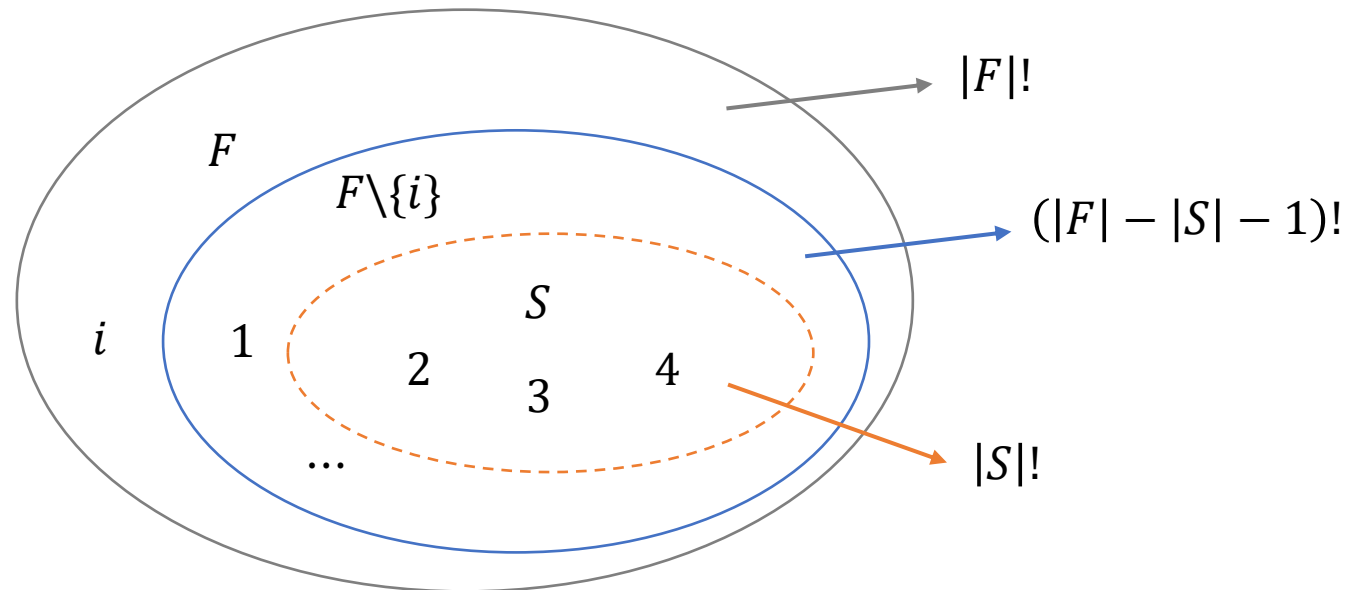


SHAP

- Shapley value [Shapley, 1953]

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|! (|F| - |S| - 1)!}{|F|!} [f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S)]$$

Weighted by the permutations of features



SHAP

- SHapley Additive exPlanation (SHAP)

Additive feature attribution method

$$g(z') \approx f(h_x(z'))$$

$$g(z') = \phi_0 + \sum_{i=1}^N \phi_i z_i'$$

$$z' \approx x' \quad \underline{x} = h_x(\underline{x'})$$

Original input Interpretable input

SHAP

- SHapley Additive exPlanation (SHAP)

Additive feature attribution method

$$g(z') \approx f(h_x(z'))$$

$$z' \approx x' \quad \underline{x} = h_x(\underline{x'})$$

Original input Interpretable input

$$g(z') = \phi_0 + \sum_{i=1}^N \phi_i z_i'$$

LIME is a special case, but not optimal

$$g(z') = \sum_{i=1}^N w_i z_i'$$

SHAP

- SHapley Additive exPlanation (SHAP)

Additive feature attribution method

$$g(z') \approx f(h_x(z'))$$

$$g(z') = \phi_0 + \sum_{i=1}^N \phi_i z_i'$$

$$z' \approx x' \quad \underline{x} = h_x(\underline{x'})$$

Original input Interpretable input

□ Property 1: Local accuracy

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

$$\phi_0 = h_x(0)$$

SHAP

- SHapley Additive exPlanation (SHAP)

Additive feature attribution method

$$g(z') \approx f(h_x(z'))$$

$$z' \approx x' \quad \underline{x} = h_x(\underline{x'})$$

Original input Interpretable input

$$g(z') = \phi_0 + \sum_{i=1}^N \phi_i z_i'$$

□ Property 2: Missingness

$$x'_i = 0 \quad \Rightarrow \quad \phi_i = 0$$

Missingness constrains features missing in the original input to have no attributed impact

SHAP

- SHapley Additive exPlanation (SHAP)

Additive feature attribution method

$$g(z') \approx f(h_x(z'))$$

$$z' \approx x' \quad \underline{x} = h_x(\underline{x'})$$

Original input Interpretable input

$$g(z') = \phi_0 + \sum_{i=1}^N \phi_i z_i'$$

□ Property 3: Consistency

For any two models f_1 and f_2 , if $f_1(h_x(z')) - f_1(h_x(z' \setminus i)) \geq f_2(h_x(z')) - f_2(h_x(z' \setminus i))$

$$\overline{z'_i = 0}$$

for all inputs $z' \in \{0, 1\}^N$, then $\phi_i(f_1, x) \geq \phi_i(f_2, x)$

SHAP

- SHapley Additive exPlanation (SHAP)

Additive feature attribution method

$$g(z') \approx f(h_x(z'))$$

$$z' \approx x' \quad \underline{x} = h_x(\underline{x'})$$

Original input Interpretable input

$$g(z') = \phi_0 + \sum_{i=1}^N \phi_i z_i'$$

Only Shapley value satisfies all the three properties

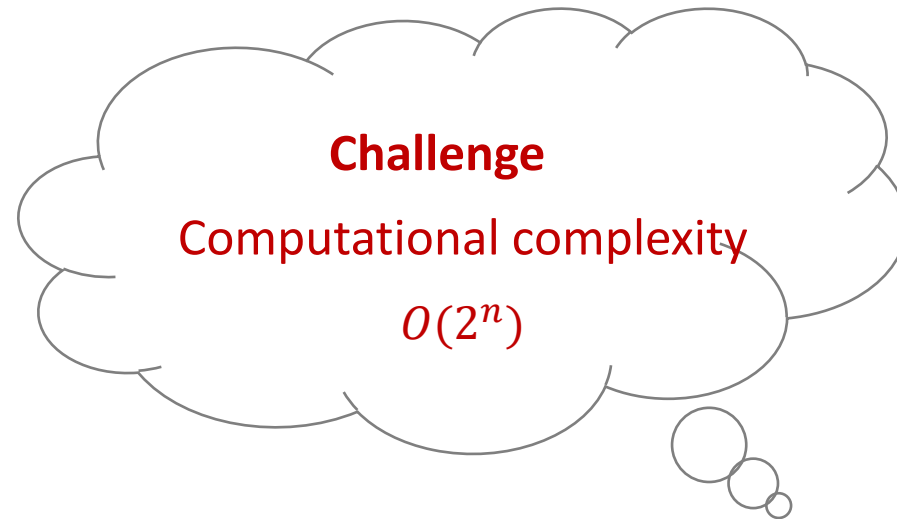
$$\phi_i(f, x) = \sum_{\underline{z}' \subseteq \underline{x}'} \frac{|\underline{z}'|! (N - |\underline{z}'| - 1)!}{N!} [f(h_x(\underline{z}')) - f(h_x(\underline{z}' \setminus i))]$$

Contains a subset of non-zero entries in x'

SHAP

- SHapley Additive exPlanation (SHAP)

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



SHAP

- SHapley Additive exPlanation (SHAP)

Model-agnostic approximations

- Shapley sampling values
- Kernel SHAP

Model-type-specific approximations

- Linear SHAP
- Low-Order SHAP
- Max SHAP
- Deep SHAP

SHAP

- SHapley Additive exPlanation (SHAP)

Model-agnostic approximations

- Shapley sampling values
- Kernel SHAP

Model-type-specific approximations

- Linear SHAP
- Low-Order SHAP
- Max SHAP
- Deep SHAP

Initialize the number of samples M

$$\phi_i \leftarrow 0$$

for $m \in \{1, \dots, M\}$ **do**

Sample $z' \subseteq x'$

$$\phi_i \leftarrow \phi_i + \frac{|z'|!(N-|z'|-1)!}{N!} [f(h_x(z')) - f(h_x(z' \setminus i))]$$

SHAP

- SHapley Additive exPlanation (SHAP)

Model-agnostic approximations

- Shapley sampling values
- Kernel SHAP **Linear LIME + Shapley values**

Model-type-specific approximations

- Linear SHAP
- Low-Order SHAP
- Max SHAP
- Deep SHAP

The solutions would be consistent with properties 1-3

$$\Omega(g) = 0$$

$$\pi_{x'}(z') = \frac{(N - 1)}{(N \text{ choose } |z'|) |z'| (N - |z'|)}$$

$$\mathcal{L}(f, g) = \sum \pi_{x'}(z') (f(h_x(z')) - g(z'))^2$$

SHAP

- SHapley Additive exPlanation (SHAP)

Model-agnostic approximations

- Shapley sampling values
- Kernel SHAP

Model-type-specific approximations

- Linear SHAP
- Low-Order SHAP
- Max SHAP
- Deep SHAP

Faster model-specific methods

SHAP values can be approximated directly from the model's weight coefficients

Question?

Reference

- Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "" Why should i trust you?" Explaining the predictions of any classifier." *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 2016.
- Lundberg, Scott M., and Su-In Lee. "A unified approach to interpreting model predictions." *Proceedings of the 31st international conference on neural information processing systems*. 2017.
- Li, Jiwei, Will Monroe, and Dan Jurafsky. "Understanding neural networks through representation erasure." *arXiv preprint arXiv:1612.08220* (2016).
- Lloyd S Shapley. 1953. A value for n-person games. *Contributions to the Theory of Games*, 2(28).