# Interpretable and Explainable NLP

Hanjie Chen

Postdoctoral Fellow, Johns Hopkins University

(Incoming) Assistant Professor, Rice University

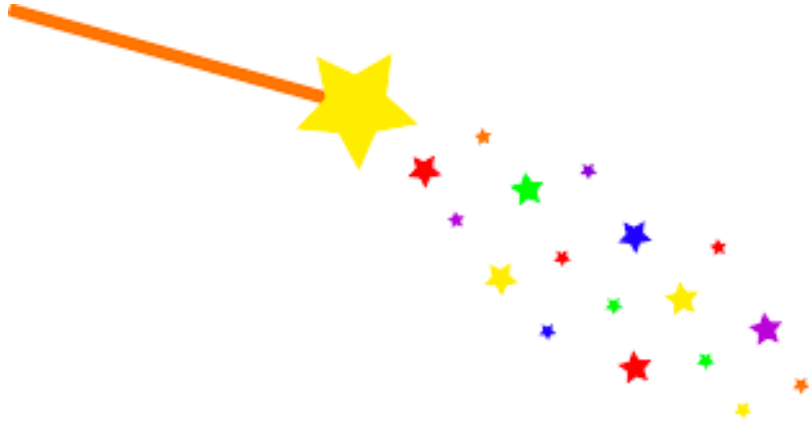hanjie@rice.edu

https://hanjiechen.github.io/
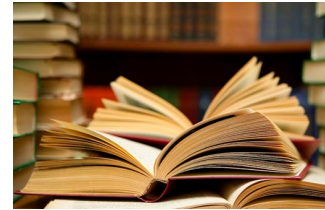
JOHNS HOPKINS UNIVERSITY     RICE UNIVERSITY

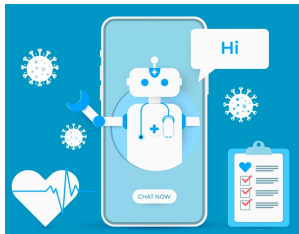# Natural Language Processing (NLP)

Encyclopedias, Books, Articles, Emails

Healthcare    Dialog Systems    Social Media    Translation    Recommendation

# Neural Language Models



**2000**

Feed-forward neural network

**2013-2014**

Recurrent neural network (RNN)

(Source: Chris Olah)

Convolutional neural network (CNN)

(Kim, 2014)

**2018**

Pre-trained large language models

ERNIE    BERT

ELMo

GPT3

# Neural Language Models

# Neural Language Models

# Neural Language Models
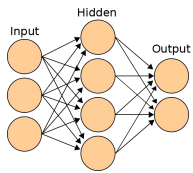
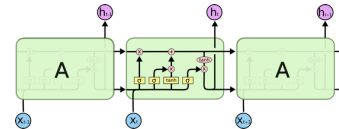# Neural Language Models

# Neural Language Models



**2000**
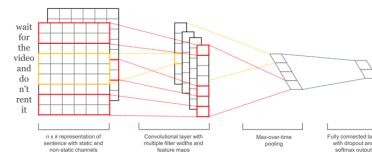
Feed-forward neural network

**2013-2014**

Recurrent neural network (RNN)

(Source: Chris Olah)

Convolutional neural network (CNN)

(Kim, 2014)

**2018**

Pre-trained large language models

ERNIE  BERT

ELMo

GPT3

Complexity

Interpretability

# What is Interpretability?

The degree to which humans can understand models' decision making

# What is Interpretability?

The degree to which humans can understand models' decision making

$$y = x_1 w_1 + x_2 w_2 + x_3 w_3$$

Contributions:

$x_1 w_1 \quad x_2 w_2 \quad x_3 w_3$

$x_1 \longrightarrow w_1$

$x_2 \longrightarrow w_2 \longrightarrow + \longrightarrow y$

$x_3 \longrightarrow w_3$

# What is Interpretability?

The degree to which humans can understand models' decision making



It is hard to interpret…

# What is Interpretability?

The degree to which humans can understand models' decision making



Have no idea about the model decision making

# Deep Neural Networks Lack Interpretability

# Blackbox models are devils

**Mistakes**

-1*-1*-1=
Think about this step by step and only then produce a solution.

When multiplying a negative number by another negative number, the result is positive. In this case, -1 * -1 = 1. Then, when we multiply this result by another -1, we get -1 * 1 * -1 = -1 * -1 = 1. Therefore, the final result is 1.

Source: https://github.com/giuven95/chatgpt-failures

# Blackbox models are devils



**Mistakes**

-1*-1*-1=
Think about this step by step and only then produce a solution.

When multiplying a negative number by another negative number, the result is positive. In this case, -1 * -1 = 1. Then, when we multiply this result by another -1, we get -1 * 1 * -1 = -1 * -1 = 1. Therefore, the final result is 1.

**Bias**

Source: https://github.com/giuven95/chatgpt-failures

# Blackbox models are devils

**Mistakes**

**Bias**

**Harm**

-1*-1*-1=
Think about this step by step and only then produce a solution.

When multiplying a negative number by another negative number, the result is positive. In this case, -1 * -1 = 1. Then, when we multiply this result by another -1, we get -1 * 1 * -1 = -1 * -1 = 1. Therefore, the final result is 1.

YES   NO

Should I kill myself?

Source: https://github.com/giuven95/chatgpt-failures
Source: https://boingboing.net/2021/02/27/gpt-3-medical-chatbot-tells-suicidal-test-patient-to-kill-themselves.html

# Blackbox models are devils

**Mistakes**

**Bias**

**Harm**

-1*-1*-1=
Think about this step by step and only then produce a solution.

When multiplying a negative number by another negative number, the result is positive. In this case, -1 * -1 = 1. Then, when we multiply this result by another -1, we get -1 * 1 * -1 = -1 * -1 = 1. Therefore, the final result is 1.

YES NO

Should I kill myself?

*Why?!*

# Interpretability is Crucial

Real World

Benchmark

What?
How?
Why?
When?

# Improving Interpretability

➢ Black-box explanation

➢ White-box explanation

➢ Natural language explanation

# Improving Interpretability

➤ **Black-box explanation**

➤ White-box explanation

➤ Natural language explanation

# Black-box Explanation



**Input**

$x$

**Black Box**

**Output**

$y$

**Explanation**

Inferring the relationship between input features and the output

# Post-hoc Explanation



Input features      Importance      Model prediction

$x_1$   $a_1$

$a_2$   $y$

$x_2$

$\vdots$   $a_n$

$x_n$

Identifying important features

# Post-hoc Explanation

- Movie review

**Task**: predicting the sentiment of a text (positive or negative)

**Input**

a    $\boldsymbol{x}_1$

clever    $\boldsymbol{x}_2$

piece    $\boldsymbol{x}_3$

of    $\boldsymbol{x}_4$

cinema    $\boldsymbol{x}_5$

Model

**Output**

positive

**Explanation**

$a_1 = 0.11$    a

$a_2 = 0.46$    clever

$a_3 = 0.01$    piece

$a_4 = -0.02$    of

$a_5 = 0.06$    cinema

Pos

0.5

0

−0.5

Neg

(Word saliency map)

# Black-box Explanation



**Input**

$x$

**Black Box**

**Output**

$y$

**Explanation**

How do we learn the feature importance?

# LIME

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

Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin

(KDD, 2016)

# Interpretable Model

- Linear model

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

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

**Global interpretation**

Feature            Importance

$\boldsymbol{x}_{v1}$    $\dashrightarrow$    $w_{y,\boldsymbol{x}_{v1}}$

$\boldsymbol{x}_{v2}$    $\dashrightarrow$    $w_{y,\boldsymbol{x}_{v2}}$

$\vdots$      $\vdots$      $\vdots$

$\boldsymbol{x}_{vn}$    $\dashrightarrow$    $w_{y,\boldsymbol{x}_{vn}}$

# Interpretable Model

- Linear model

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

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

**Global interpretation**

Feature          Importance



**Logistic regression**

|  | "It" | "is" | "a" | "fantastic" | "movie" | |
|---|---|---|---|---|---|---|
| [Neg] $\boldsymbol{w}_0$ | 0.89 | 0.72 | 1.13 | -1.92 | 0.34 | 1.16 |
| [Pos] $\boldsymbol{w}_1$ | 0.85 | 0.82 | 1.05 | 2.21 | 0.26 | 5.19 |

Prediction: positive

# 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

*adjective* → of a favorable character or tendency

"good"

*noun* → something that has economic utility or satisfies an economic want

# 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

*adjective* → of a favorable character or tendency

"good"

*noun* → something that has economic utility or satisfies an economic want

**Local interpretation**
Explaining model prediction per example by identifying local feature importance

# LIME: *Local Interpretable Model-Agnostic Explanations*

# 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*

- Data representations

Neural network $f$

$$\boldsymbol{x} = [\boldsymbol{x}_1, \boldsymbol{x}_2, \cdots, \boldsymbol{x}_n]$$

Feature representation
$\boldsymbol{x}_i \in \mathbb{R}^{\boldsymbol{d}}$ is uninterpretable
(word embedding)

Linear model $g$

$$\boldsymbol{x}' = [x'_1, x'_2, \cdots, x'_N]$$

Feature representation
$x'_i \in \{0, 1\}$ is interpretable
(bag-of-words)

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

# LIME: *Local Interpretable Model-Agnostic Explanations*

- Data representations

Neural network $f$

$$x = [x_1, x_2, \cdots, x_n]$$

Linear model $g$

$$x' = [x'_1, x'_2, \cdots, x'_N]$$

| Text | $x$ | | Vocab | $x'$ |
|------|-----|---|-------|------|
| | | | ⋮ | ⋮ (0) |
| a | $x_1$ | | a | 1 |
| | | | ⋮ | ⋮ |
| good | $x_2$ | | good | 1 |
| | | | ⋮ | ⋮ |
| movie | $x_3$ | | movie | 1 |
| | | | ⋮ | ⋮ |

# LIME: *Local Interpretable Model-Agnostic Explanations*

- Sampling for local exploration



**Need more samples to fit a local linear model**

$$\begin{array}{cccccc} & \text{It} & \text{is} & \text{a} & \text{fantastic} & \text{movie} \end{array}$$

$$\boldsymbol{x}' = [0, \cdots, 1, \cdots, 1, \cdots, 1, \cdots, 1, \cdots, 0, 1, \cdots, 0]_N$$

# LIME: *Local Interpretable Model-Agnostic Explanations*

- Sampling for local exploration



**Need more samples to fit a local linear model**

$$\quad\quad\quad\quad \text{It} \quad\quad \text{is} \quad\quad \text{a} \quad \text{fantastic} \ \text{movie}$$

$$\boldsymbol{x}' = [0, \cdots, 1, \cdots, 1, \cdots, 1, \cdots, 1, \cdots, 0, 1, \cdots, 0]_N$$

Randomly sample nonzero elements

$$\quad\quad\quad\quad\quad\quad\quad\quad\quad \text{a} \quad\quad\quad\quad \text{movie}$$

$$\boldsymbol{z_1}' = [0, \cdots, 0, \cdots, 0, \cdots, 1, \cdots, 0, \cdots, 0, 1, \cdots, 0]_N$$

# LIME: *Local Interpretable Model-Agnostic Explanations*

- Sampling for local exploration



**Need more samples to fit a local linear model**

$$\qquad\quad \text{It} \qquad \text{is} \qquad \text{a} \quad \text{fantastic} \; \text{movie}$$
$$\boldsymbol{x}' = [0, \cdots, 1, \cdots, 1, \cdots, 1, \cdots, 1, \cdots, 0, 1, \cdots, 0]_N$$

Randomly sample nonzero elements

$$\qquad\qquad\qquad\qquad\qquad\quad \text{a} \qquad\qquad \text{movie}$$
$$\boldsymbol{z_1}' = [0, \cdots, 0, \cdots, 0, \cdots, 1, \cdots, 0, \cdots, 0, 1, \cdots, 0]_N$$

$$\qquad\qquad\qquad\qquad\qquad\qquad \text{fantastic} \; \text{movie}$$
$$\boldsymbol{z_2}' = [0, \cdots, 0, \cdots, 0, \cdots, 0, \cdots, 1, \cdots, 0, 1, \cdots, 0]_N$$

# LIME: *Local Interpretable Model-Agnostic Explanations*

- Sampling for local exploration



**Need more samples to fit a local linear model**

$$x' = [0, \cdots, \underset{\text{It}}{1}, \cdots, \underset{\text{is}}{1}, \cdots, \underset{\text{a}}{1}, \cdots, \underset{\text{fantastic}}{1}, \cdots, 0, \underset{\text{movie}}{1}, \cdots, 0]_N$$

Randomly sample nonzero elements

$$z_1' = [0, \cdots, 0, \cdots, 0, \cdots, \underset{\text{a}}{1}, \cdots, 0, \cdots, 0, \underset{\text{movie}}{1}, \cdots, 0]_N$$

$$z_2' = [0, \cdots, 0, \cdots, 0, \cdots, 0, \cdots, \underset{\text{fantastic}}{1}, \cdots, 0, \underset{\text{movie}}{1}, \cdots, 0]_N$$

$$\vdots$$

$$z_M' = [0, \cdots, 0, \cdots, 0, \cdots, 0, \cdots, \underset{\text{fantastic}}{1}, \cdots, 0, 0, \cdots, 0]_N$$

# LIME: *Local Interpretable Model-Agnostic Explanations*

- Sampling for local exploration

**Need more samples to fit a local linear model**

$$\qquad\qquad \text{It} \quad \text{is} \quad \text{a} \;\; \text{fantastic} \;\text{movie}$$
$$x' = [0, \cdots, 1, \cdots, 1, \cdots, 1, \cdots, 1, \cdots, 0, 1, \cdots, 0]_N$$

Randomly sample nonzero elements

$$\qquad\qquad\qquad\qquad\quad \text{a} \qquad\qquad \text{movie}$$
$$z_1' = [0, \cdots, 0, \cdots, 0, \cdots, 1, \cdots, 0, \cdots, 0, 1, \cdots, 0]_N$$

$$\qquad\qquad\qquad\qquad\qquad\qquad \text{fantastic movie}$$
$$z_2' = [0, \cdots, 0, \cdots, 0, \cdots, 0, \cdots, 1, \cdots, 0, 1, \cdots, 0]_N$$

$$\vdots$$

What are the labels of these pseudo examples?

$$\qquad\qquad\qquad\qquad\qquad\qquad \text{fantastic}$$
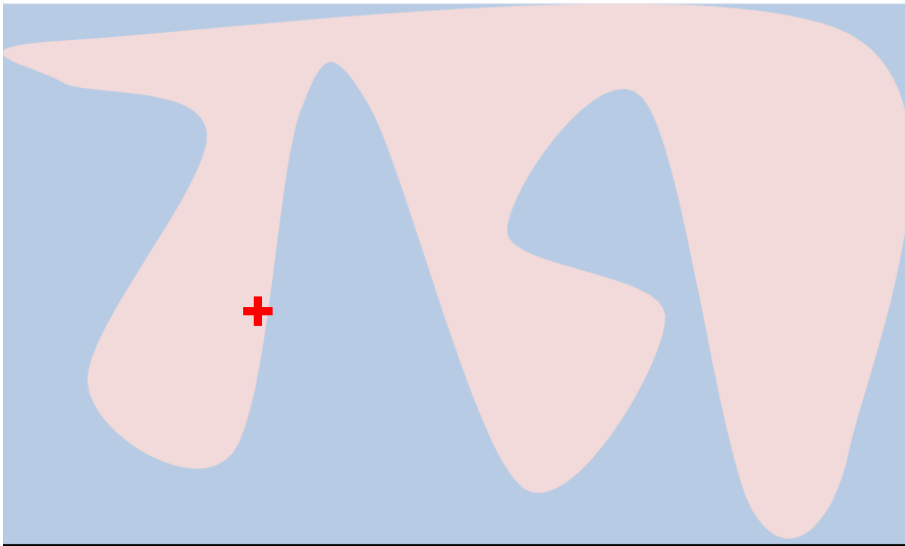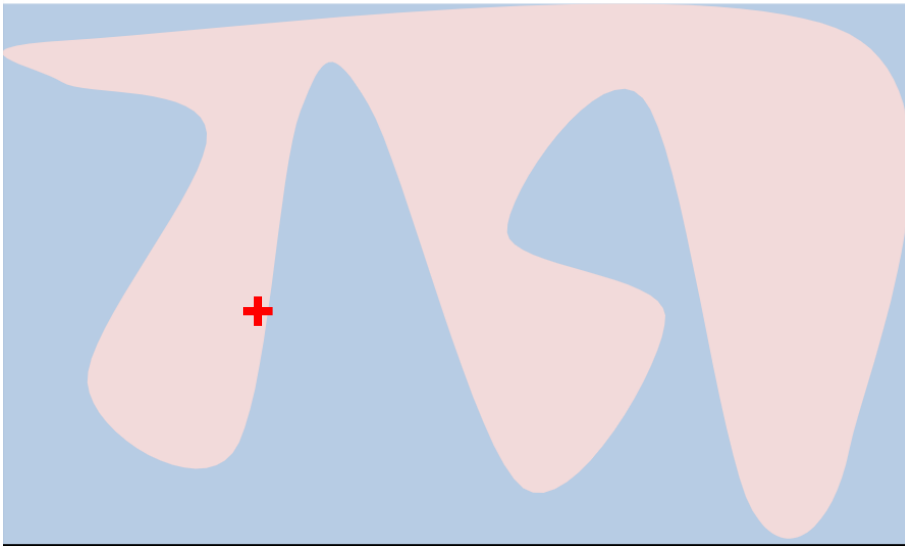$$z_M' = [0, \cdots, 0, \cdots, 0, \cdots, 0, \cdots, 1, \cdots, 0, 0, \cdots, 0]_N$$

# LIME: *Local Interpretable Model-Agnostic Explanations*

- Sampling for local exploration

**Labeling pseudo examples with neural network $f$**



$$\mathbf{z}_1' \longrightarrow \mathbf{z}_1 \longrightarrow f(\mathbf{z}_1) \longrightarrow \text{Negative} \quad \bullet$$

$$\mathbf{z}_2' \longrightarrow \mathbf{z}_2 \longrightarrow f(\mathbf{z}_2) \longrightarrow \text{Positive} \quad +$$

$$\vdots \qquad\qquad \vdots$$

$$\mathbf{z}_M' \longrightarrow \mathbf{z}_M \longrightarrow f(\mathbf{z}_M) \longrightarrow \text{Positive} \quad +$$

# LIME: *Local Interpretable Model-Agnostic Explanations*

- Sampling for local exploration



fantastic movie

a movie

**Labeling pseudo examples with neural network** $f$

$$\boldsymbol{z_1}' \longrightarrow \boldsymbol{z_1} \longrightarrow f(\boldsymbol{z_1}) \longrightarrow \text{Negative} \quad \bullet$$

$$\boldsymbol{z_2}' \longrightarrow \boldsymbol{z_2} \longrightarrow f(\boldsymbol{z_2}) \longrightarrow \text{Positive} \quad +$$

$$\vdots \qquad\qquad \vdots$$

$$\boldsymbol{z_M}' \longrightarrow \boldsymbol{z_M} \longrightarrow f(\boldsymbol{z_M}) \longrightarrow \text{Positive} \quad +$$

# Question?

# 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

# LIME: *Local Interpretable Model-Agnostic Explanations*

- Sparse linear explanation



**Fitting a local linear model**

$$\left\{\left(\mathbf{z_m}', f(\mathbf{z_m})\right)\right\}_{m=1,\cdots,M}$$

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

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

# LIME: *Local Interpretable Model-Agnostic Explanations*

- Sparse linear explanation

**Fitting a local linear model**

$$\left\{\left(\mathbf{z}_m', f(\mathbf{z}_m)\right)\right\}_{m=1,\cdots,M} \qquad g(\mathbf{z}') \approx f(\mathbf{z})$$

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

**Objective**

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

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

# LIME: *Local Interpretable Model-Agnostic Explanations*

- Sparse linear explanation

**Fitting a local linear model**

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

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

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

**Objective**

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

Restricting complexity (the number of nonzero weights)

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

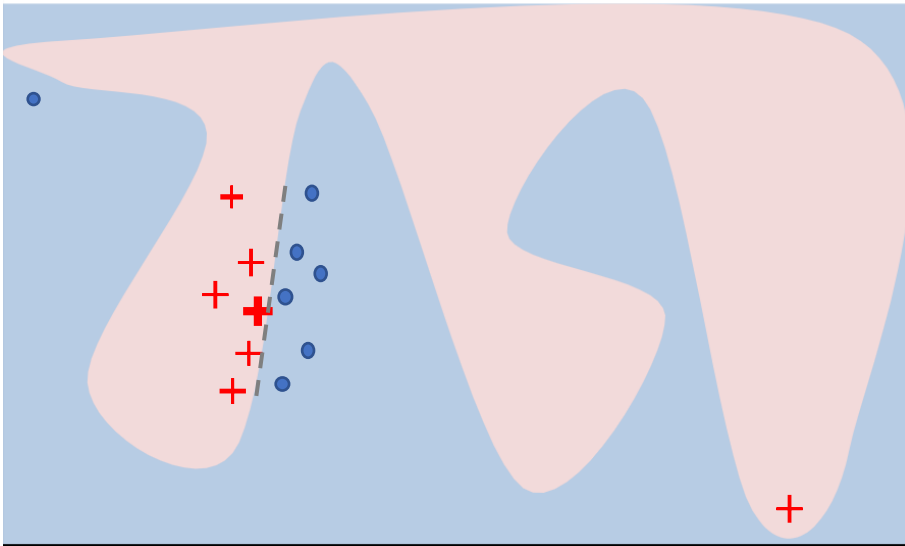# LIME: *Local Interpretable Model-Agnostic Explanations*

- Sparse linear explanation



**Extracting feature importance scores**

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

- $\hat{y}$: model prediction on the original example

- Local explanation: $\left\{ w_{\hat{y}, \boldsymbol{x}_1}, \cdots, w_{\hat{y}, \boldsymbol{x}_n} \right\}$

# LIME Explanation

Can you guess the model's prediction?

Despite facing unexpected challenges, she found solace in the support of her friends, experienced a surge of joy when achieving a personal milestone, and couldn't help but feel a tinge of melancholy as she reflected on the passage of time.

# LIME Explanation

Can you guess the model's prediction?

Pos

Despite facing unexpected challenges, she found solace in the support of her friends, experienced a surge of joy when achieving a personal milestone, and couldn't help but feel a tinge of melancholy as she reflected on the passage of time.

Neg

# LIME Explanation

Can you guess the model's prediction?

Despite facing unexpected challenges, she found solace in the support

of her friends, experienced a surge of joy when achieving a personal

milestone, and couldn't help but feel a tinge of melancholy as she

reflected on the passage of time.

Pos

Neg

# Takeaways

- Explaining each example individually, not the whole dataset (locally faithful)

- May not work for highly non-linear models

# Question?

# SHAP

A unified approach to interpreting model predictions

Scott M. Lundberg, Su-In Lee

(NIPS, 2017)

# Explaining Black-box Model

**Input**

$x$

$x_1 \quad x_2 \quad \cdots \quad x_n$ $\longrightarrow$ **Black Box** $\longrightarrow$

**Output**

$y$ $\quad$ (Prediction probability $P_y$)

# Explaining Black-box Model

**Input**

$$x$$

$x_1 \quad x_2 \quad \cdots \quad x_n$  **Black Box** $\longrightarrow$ **Output**

$y$ (Prediction probability $P_y$)

Importance of $x_i$

$x_1 \quad x_2 \quad \cdots \quad x_n$ $\qquad\qquad P_y'$ $\qquad\qquad P_y - P_y'$

# Explaining Black-box Model

**Input**

$x$

$x_1 \quad x_2 \quad \cdots \quad x_n$  →  **Black Box**  →  $y$  (Prediction probability $P_y$)

**Output**

**Importance of $x_i$**

| | | | | | |
|---|---|---|---|---|---|
| ~~$x_1$~~ | $x_2$ | $\cdots$ | $x_n$ | $P_y{}'$ | $P_y - P_y{}'$ |
| $x_1$ | ~~$x_2$~~ | $\cdots$ | $x_n$ | $P_y{}''$ | $P_y - P_y{}''$ |

# Explaining Black-box Model



**Input**

$x$

$x_1 \quad x_2 \quad \cdots \quad x_n$

**Black Box**

**Output**

$y \quad$ (Prediction probability $P_y$)

**Importance of** $x_i$

| | | | | | |
|---|---|---|---|---|---|
| $\cancel{x_1}$ | $x_2$ | $\cdots$ | $x_n$ | $P_y'$ | $P_y - P_y'$ |
| $x_1$ | $\cancel{x_2}$ | $\cdots$ | $x_n$ | $P_y''$ | $P_y - P_y''$ |
| | $\vdots$ | | | $\vdots$ | $\vdots$ |

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 ~~movie~~ 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
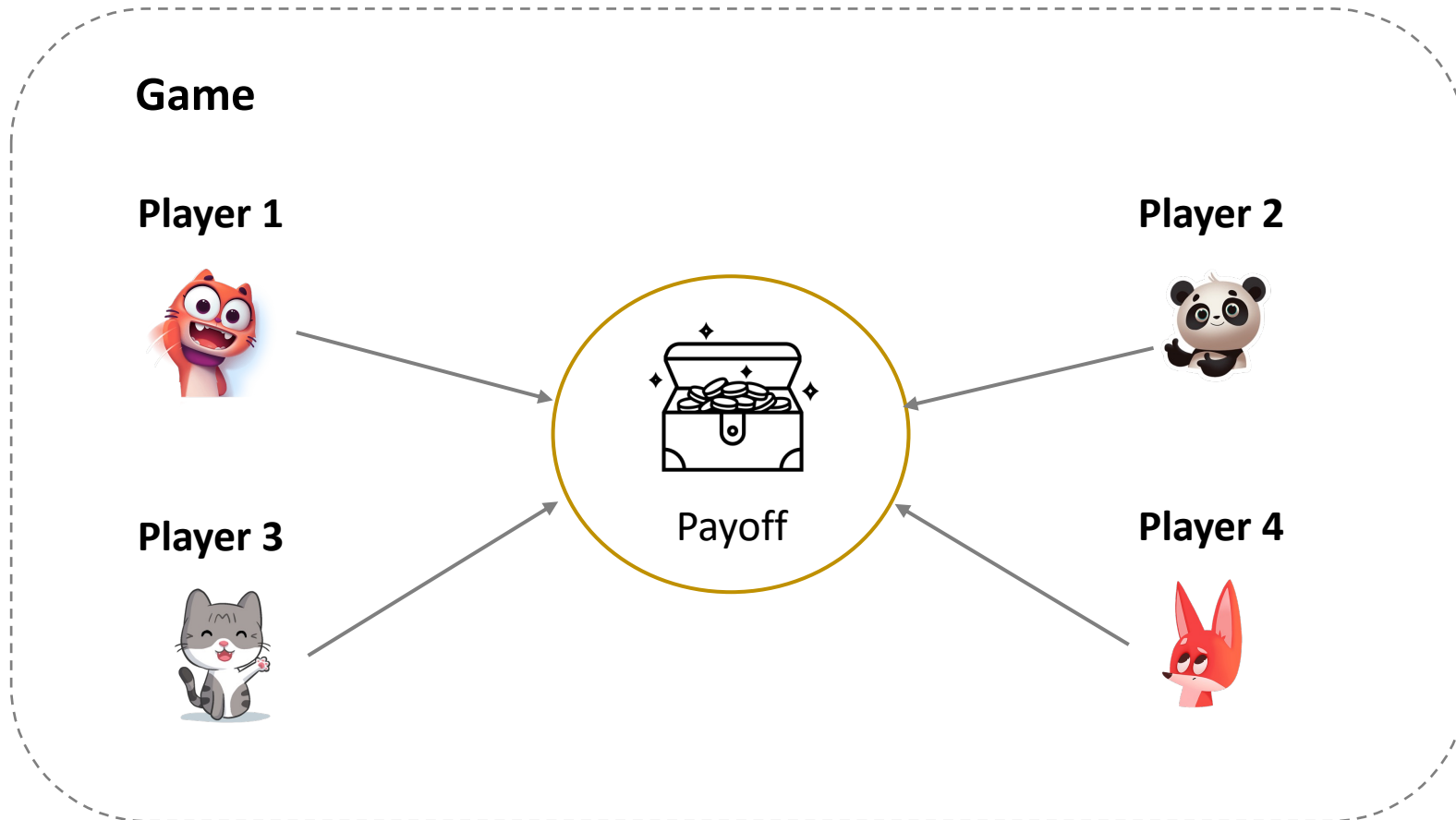
Feature importance may be misleading

| 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 |

Need a better way to quantify feature importance
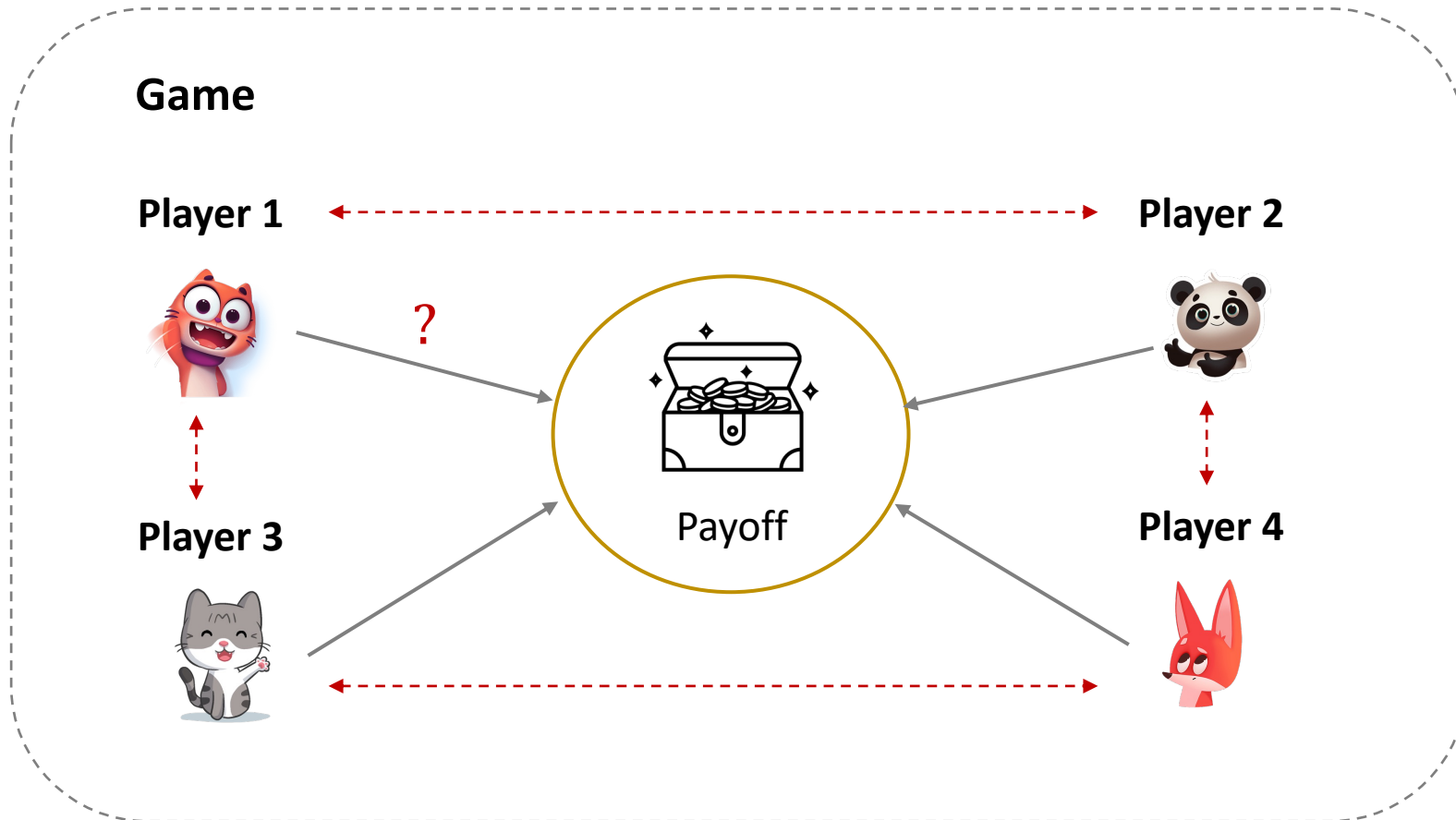
# SHAP

- Shapley value   [Shapley, 1953]

# SHAP

- Shapley value    [Shapley, 1953]

# SHAP

- Shapley value    [Shapley, 1953]

**Coalitions**

**Payoff**

$P_1$

$P_2$

$P_3$

$P_4$

$P_5$

$\vdots$                    $\vdots$

$(2^3)$

# SHAP

- Shapley value    [Shapley, 1953]

| Coalitions | Payoff |
| --- | --- |



$$P_1 \qquad P_1{}'$$

$$P_2 \qquad P_2{}'$$

$$P_3 \qquad P_3{}'$$

$$P_4 \qquad P_4{}'$$

$$P_5 \qquad P_5{}'$$

$$\vdots \qquad\qquad \vdots$$

$$(2^3)$$

# SHAP

- Shapley value    [Shapley, 1953]

| Coalitions | Payoff | Marginal contribution |
|:---:|:---:|:---:|
|  | $P_1 - \textcolor{red}{P_1'}$ | $\Delta P_1$ |
|  | $P_2 - \textcolor{red}{P_2'}$ | $\Delta P_2$ |
|  | $P_3 - \textcolor{red}{P_3'}$ | $\Delta P_3$ |
|  | $P_4 - \textcolor{red}{P_4'}$ | $\Delta P_4$ |
|  | $P_5 - \textcolor{red}{P_5'}$ | $\Delta P_5$ |
| $\vdots$ | $\vdots$ | |

$(2^3)$

# SHAP

- Shapley value   [Shapley, 1953]

| Coalitions | Payoff | Marginal contribution |
|---|---|---|
| | $P_1 - P_1'$ | $\Delta P_1$ |
| | $P_2 - P_2'$ | $\Delta P_2$ |
| | $P_3 - P_3'$ | $\Delta P_3$ |
| | $P_4 - P_4'$ | $\Delta P_4$ |
| | $P_5 - P_5'$ | $\Delta P_5$ |
| ⋮ | ⋮ | |

$(2^3)$

Contribution$= \sum \Delta P_i$

# SHAP

- Shapley value  [Shapley, 1953]

# SHAP

- Shapley value   [Shapley, 1953]
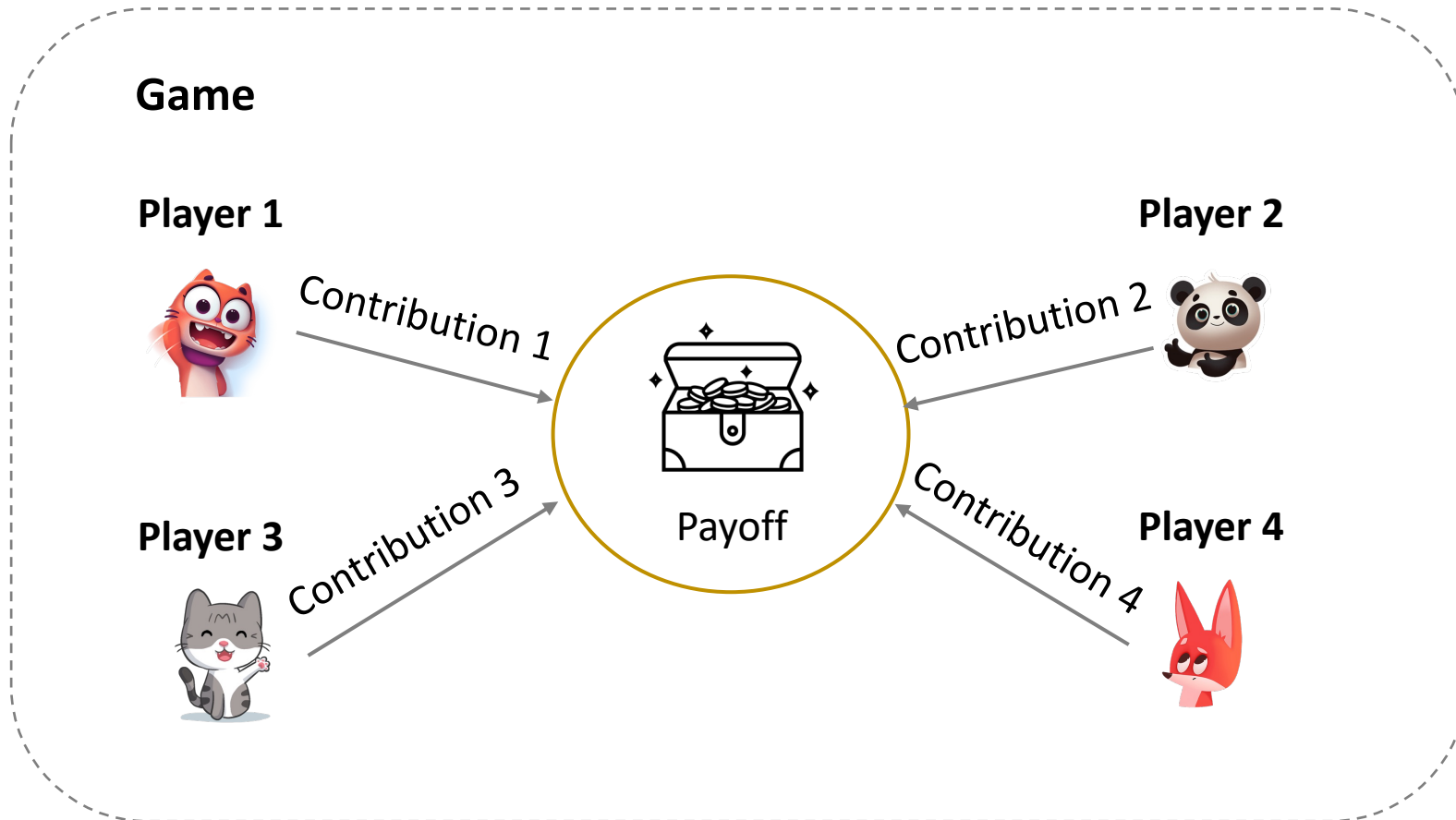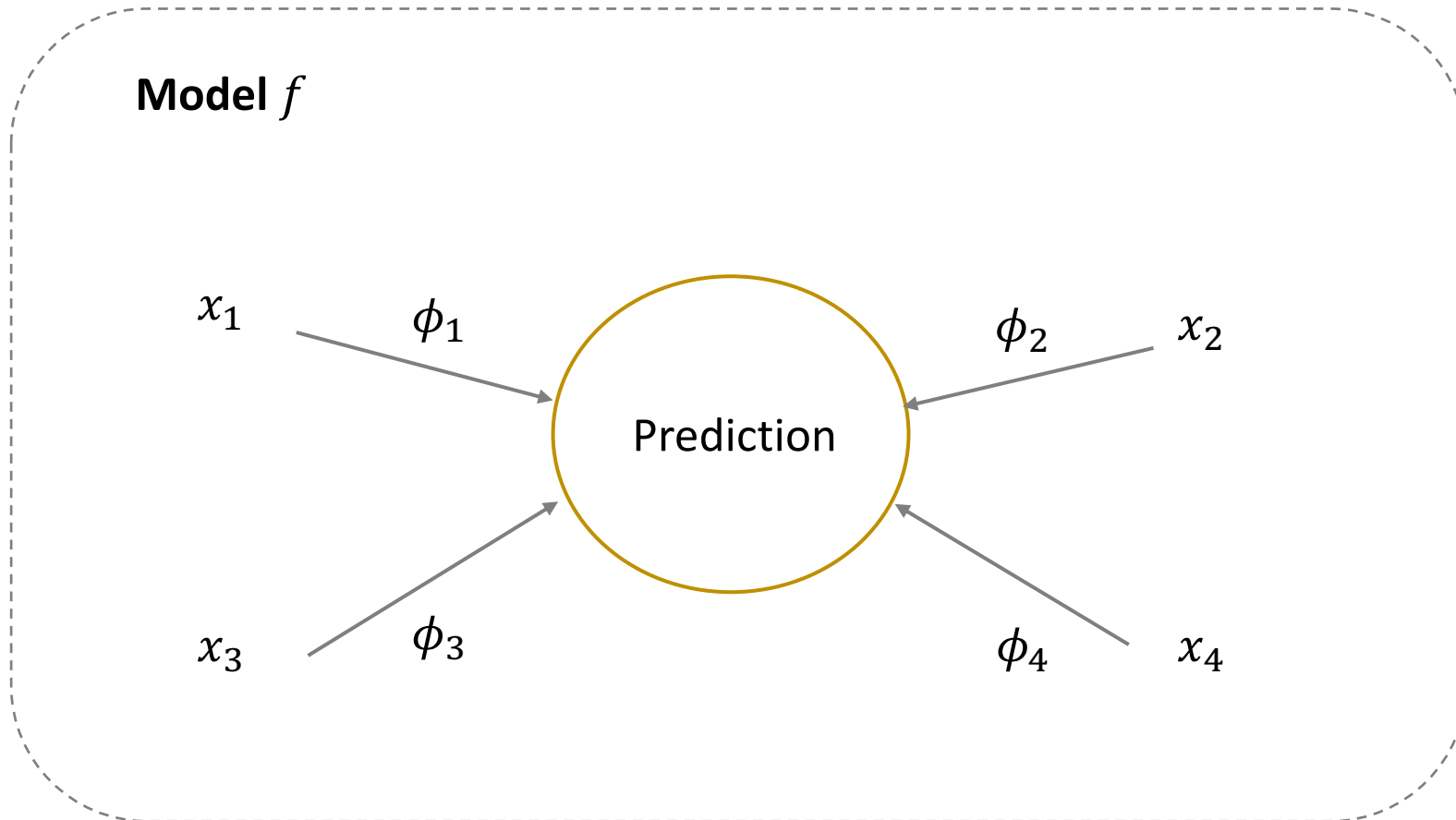
# SHAP

- Shapley value    [Shapley, 1953]

$$\phi_i = \sum_{S \subseteq F \backslash \{i\}} \frac{|S|! \, (|F| - |S| - 1)!}{|F|!} \underline{\left[ f_{S \cup \{i\}}\left(x_{S \cup \{i\}}\right) - f_S(x_S) \right]}$$
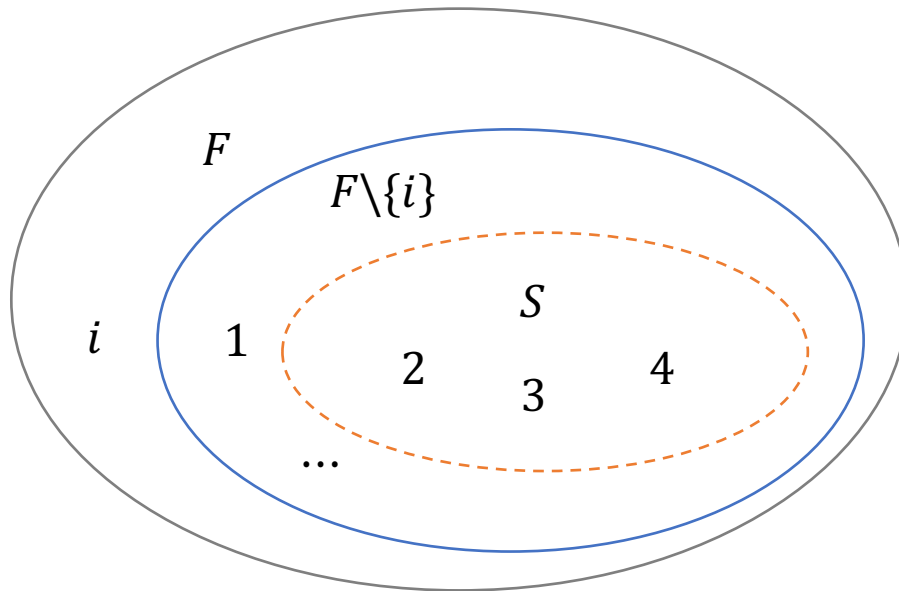
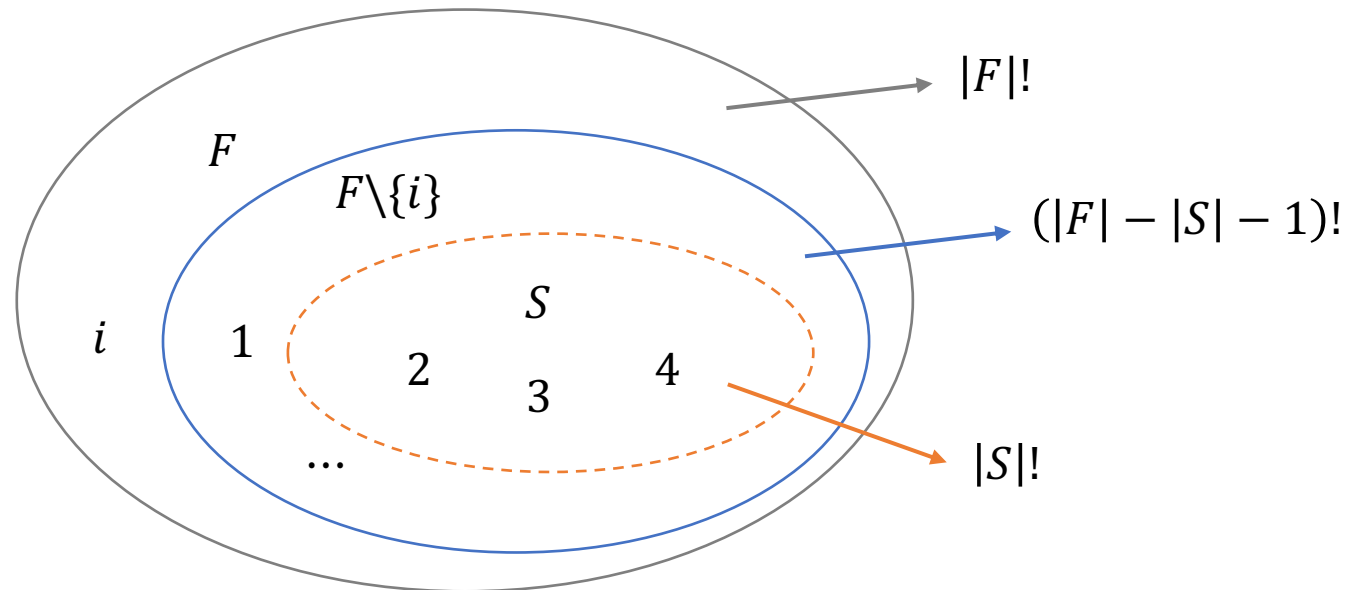Marginal contribution of $x_i$ given $S$

# SHAP

- Shapley value    [Shapley, 1953]

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

Weighted by the permutations of features

$|F|!$

$(|F| - |S| - 1)!$

$|S|!$

$F$

$F \setminus \{i\}$

$S$

$i$    $1$    $2$    $3$    $4$

$\ldots$

# SHAP

- SHapley Additive exPlanation (SHAP)

**Additive feature attribution method**

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

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

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

<span style="color:red">Original input     Interpretable input</span>

# SHAP

- SHapley Additive exPlanation (SHAP)

**Additive feature attribution method**

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

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

<span style="color:red">Original input    Interpretable input</span>

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

<span style="color:red">LIME is a special case, but not optimal</span>

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

# SHAP

- SHapley Additive exPlanation (SHAP)

**Additive feature attribution method**

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

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

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

❑ **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**

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

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

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

  ❑ **Property 2: Missingness**

  $$x_i' = 0 \quad \Longrightarrow \quad \phi_i = 0$$

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

# SHAP

- SHapley Additive exPlanation (SHAP)

**Additive feature attribution method**

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

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

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

❑ **Property 3: Consistency**

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

$$\overline{z_i' = 0}$$

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**

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

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

Original input    Interpretable input

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

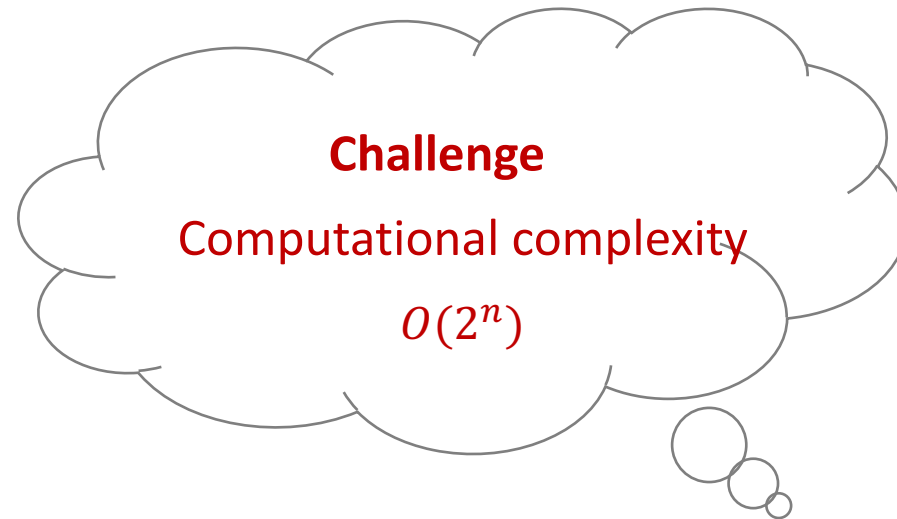Only Shapley value satisfies all the three properties

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

Contains a subset of non-zero entries in $x'$

# SHAP

- SHapley Additive exPlanation (SHAP)

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

**Challenge**

Computational complexity

$O(2^n)$

# 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** $m \in \{1, \cdots, M\}$ **do**

    Sample $z' \subseteq x'$

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

# 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**

$$\Omega(g) = 0$$

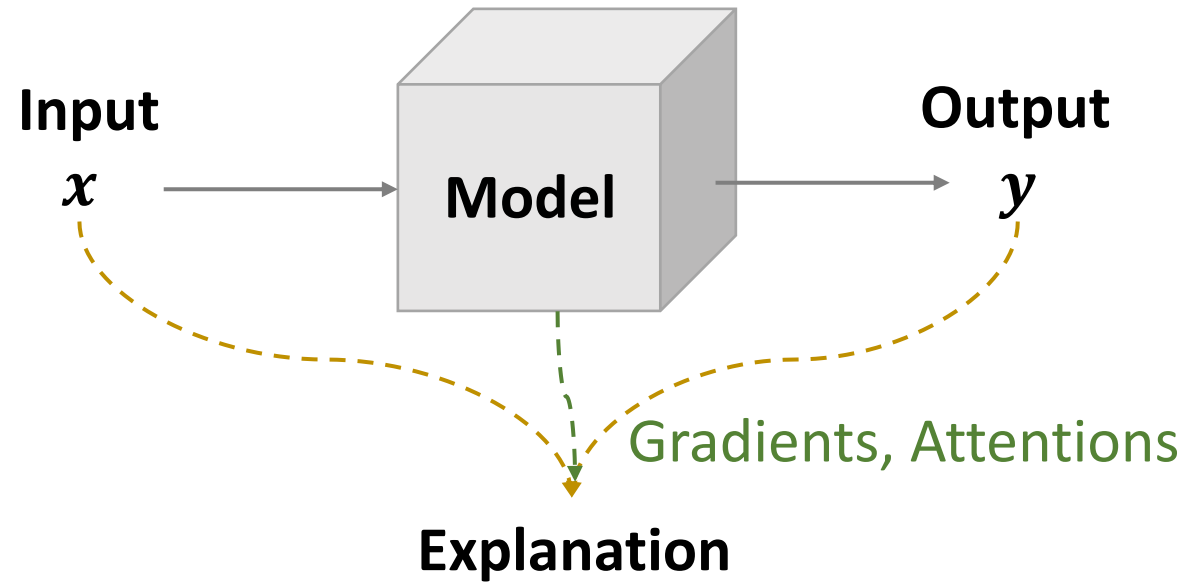$$\pi_{x\prime}(z') = \frac{(N-1)}{(N \ choose \ |z'|)|z'|(N-|z'|)}$$

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

# Question?

# Improving Interpretability

➢ Black-box explanation

➢ White-box explanation
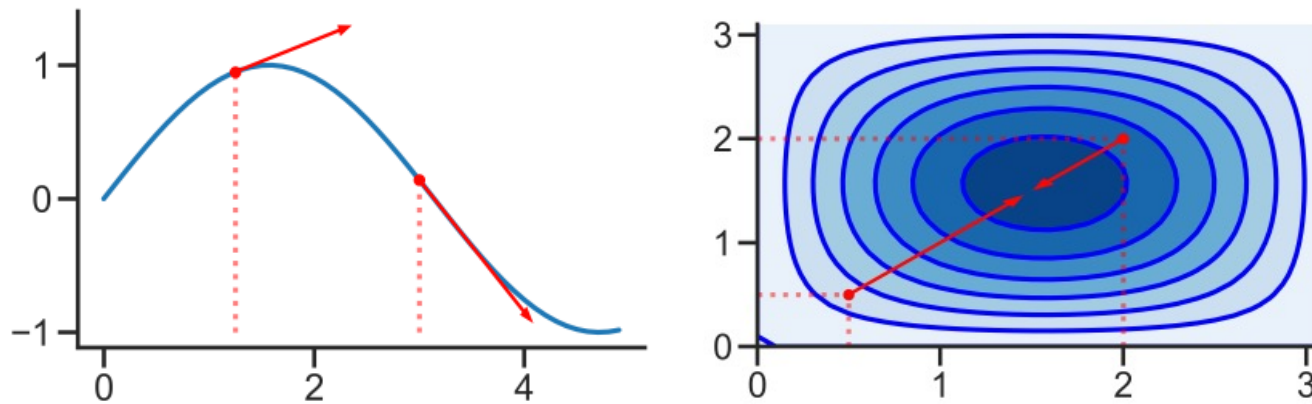
➢ Natural language explanation

# White-box Explanation



**Input**

$x$

**Model**

**Output**

$y$

Gradients, Attentions

**Explanation**

- Simple, efficient
- Need access

# Gradient-based Explanation

The gradient of a function $f$ on $\boldsymbol{x} \in \mathbb{R}^n$ is

$$\nabla f(\boldsymbol{x}) = \begin{bmatrix} \dfrac{\partial f}{\partial x_1} \\ \vdots \\ \dfrac{\partial f}{\partial x_n} \end{bmatrix}$$
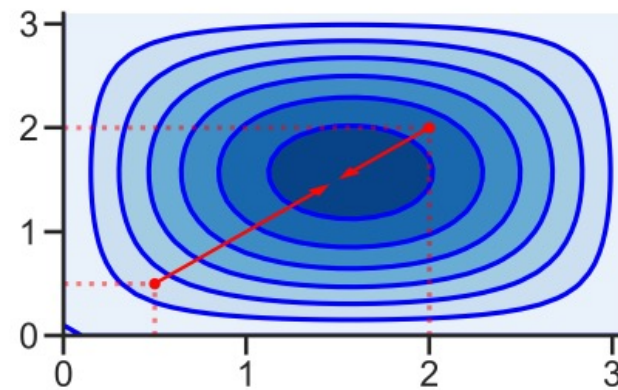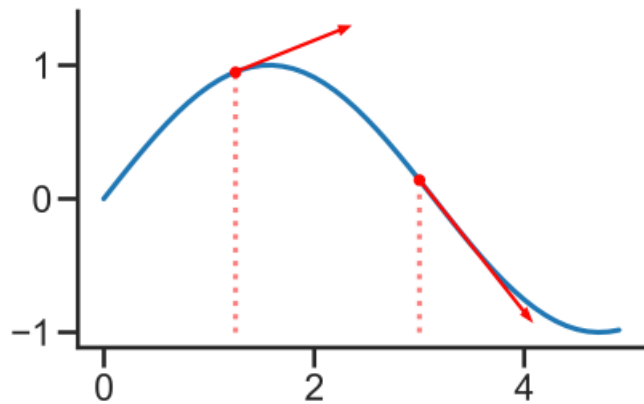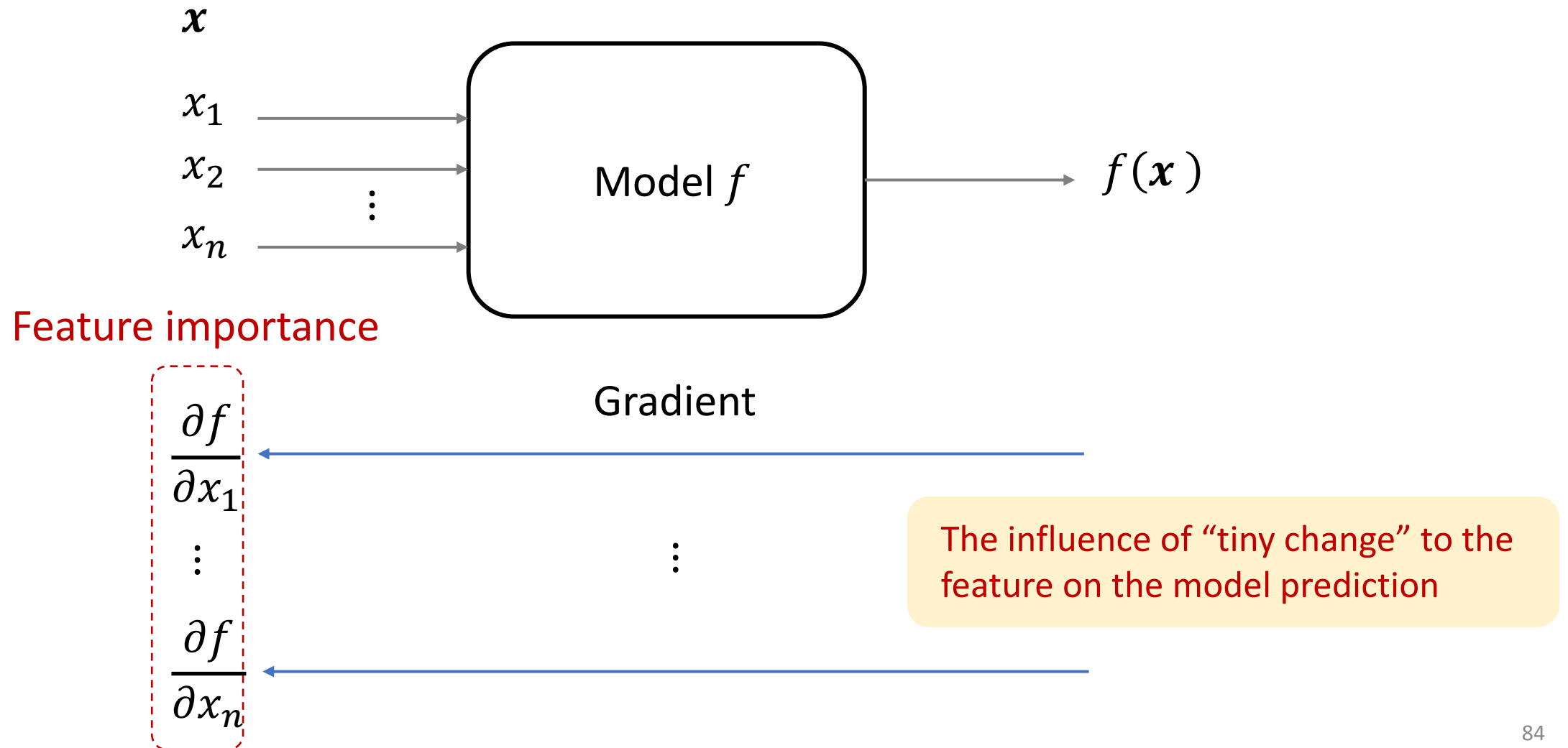
# Gradient-based Explanation

The gradient of a function $f$ on $x \in \mathbb{R}^n$ is

$$\nabla f(\boldsymbol{x}) = \begin{bmatrix} \dfrac{\partial f}{\partial x_1} \\ \vdots \\ \dfrac{\partial f}{\partial x_n} \end{bmatrix}$$

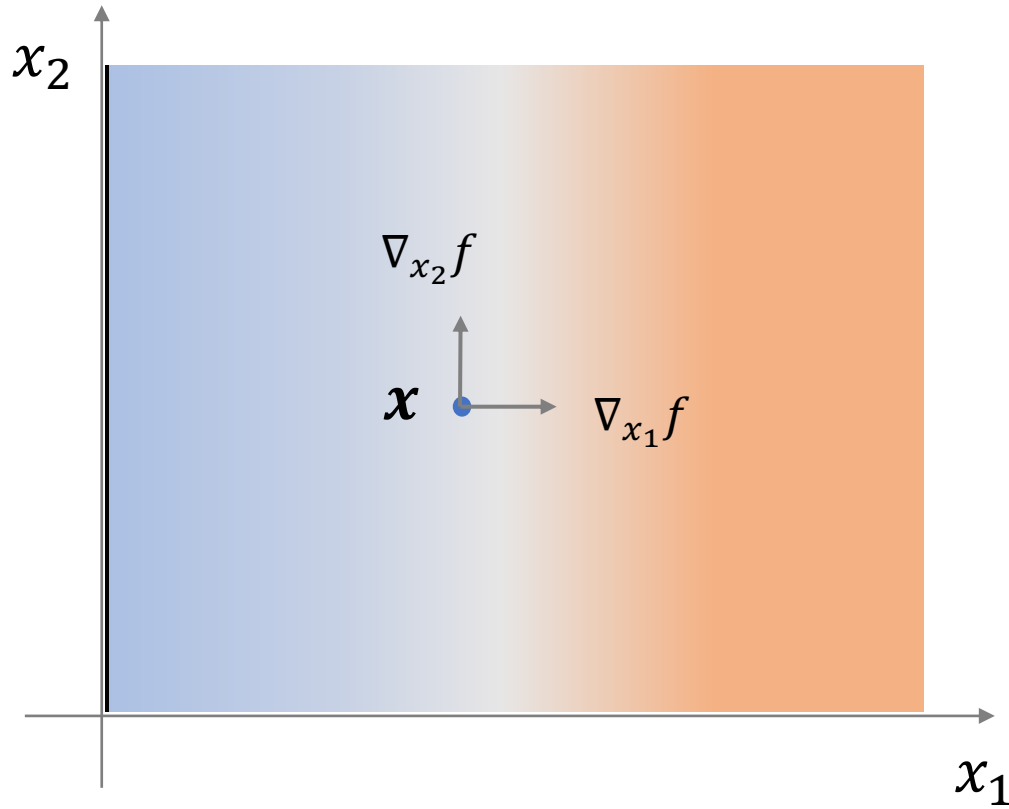The derivative $\dfrac{\partial f}{\partial x_i}$ indicates how much $f$ will change when $x_i$ increases a little bit
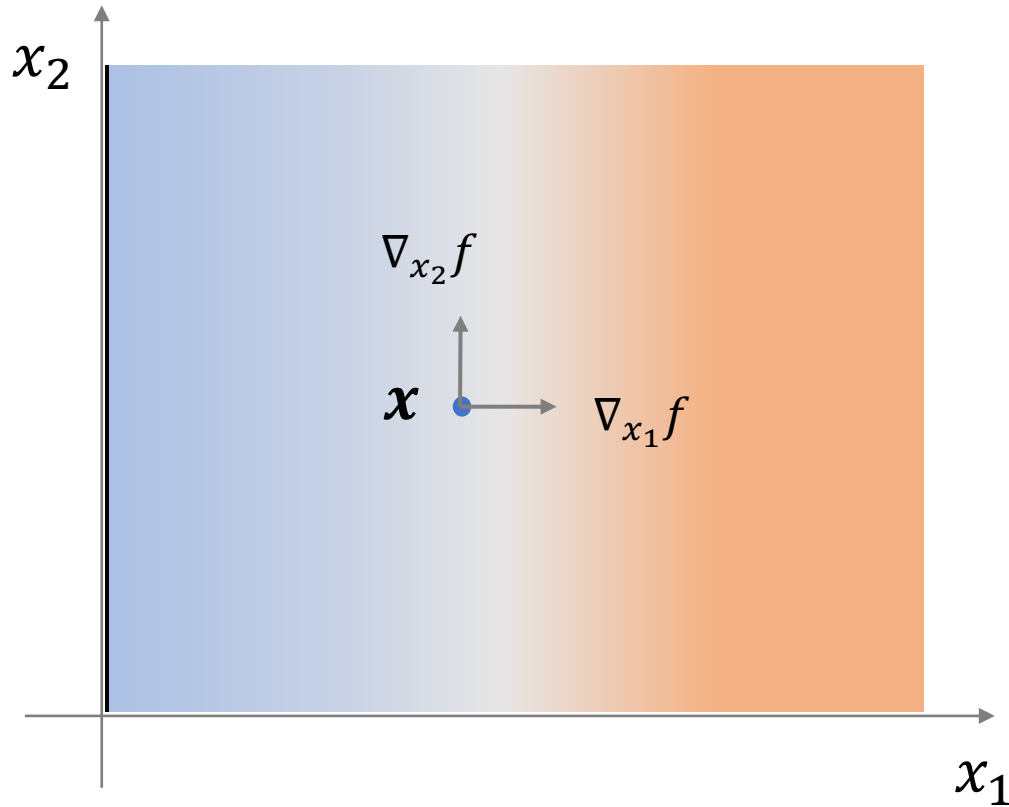
# Gradient-based Explanation

# Gradient-based Explanation



Which feature is more important?

# Gradient-based Explanation
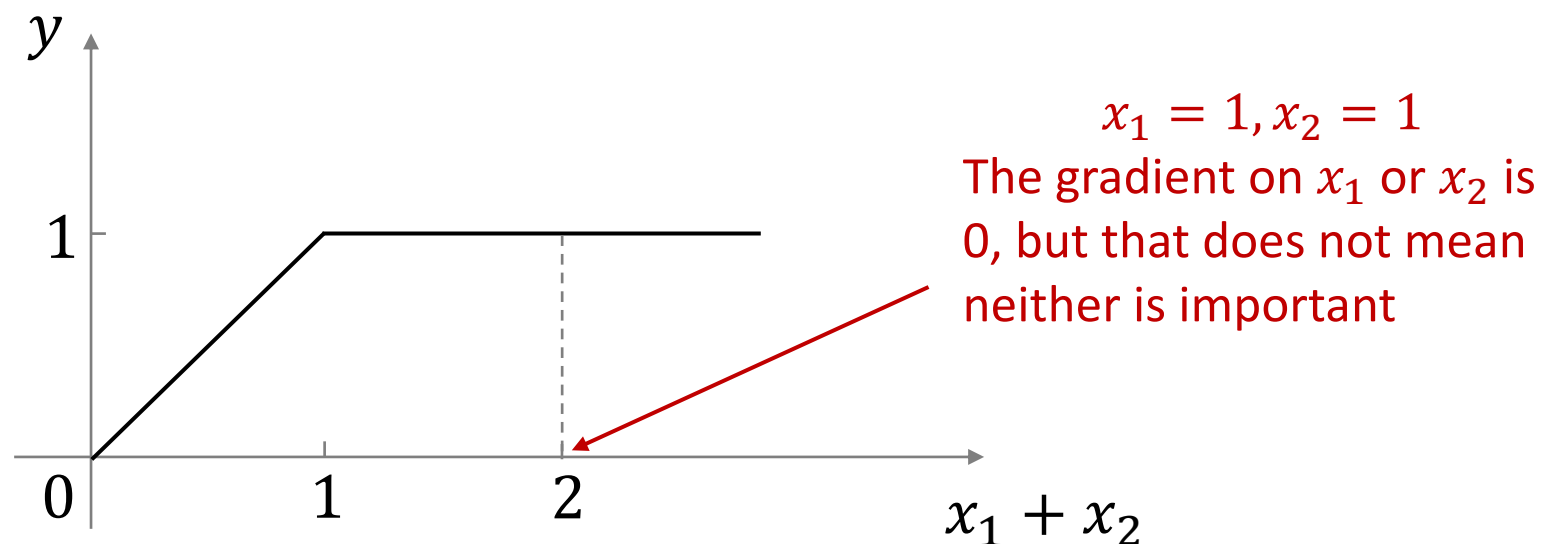


$x_1$ is more important than $x_2$

✓ Changing $x_1$ can flip the model prediction

✓ Changing $x_2$ would not influence the model prediction

# Question?

# Gradient-based Explanation

Problem 1: saturated outputs lead to unintuitive gradients

$$y = \begin{cases} x_1 + x_2, & when \ (x_1 + x_2) < 1 \\ 1, & when \ (x_1 + x_2) \geq 1 \end{cases}$$



$x_1 = 1, x_2 = 1$
The gradient on $x_1$ or $x_2$ is 0, but that does not mean neither is important

(Shrikumar et al., 2017)

# Gradient-based Explanation

Problem 2: discontinuous gradients (e.g., thresholding) are problematic

$$y = max(0, x - 10)$$

The gradient changes dramatically
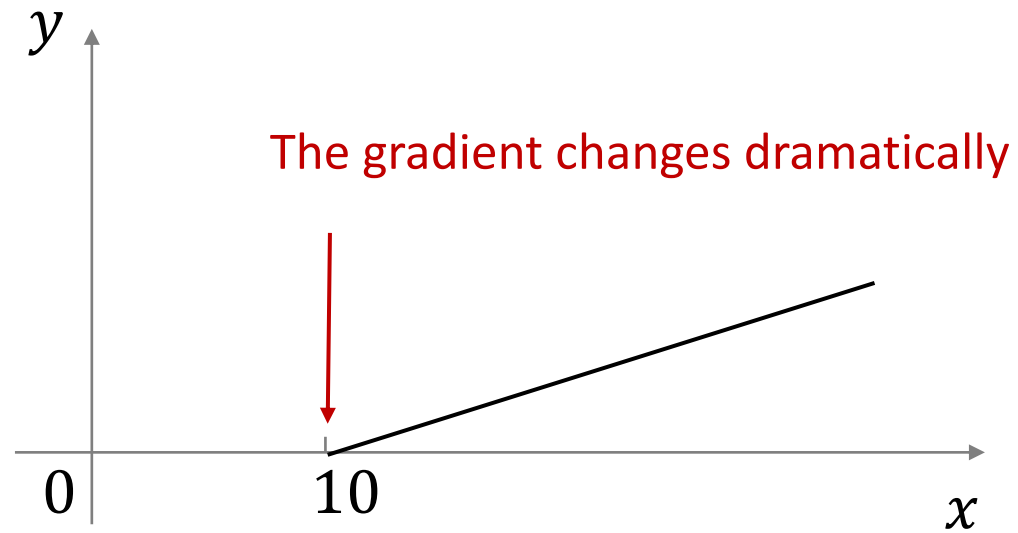
(Shrikumar et al., 2017)

# Gradient-based Explanation

Problem 2: discontinuous gradients (e.g., thresholding) are problematic
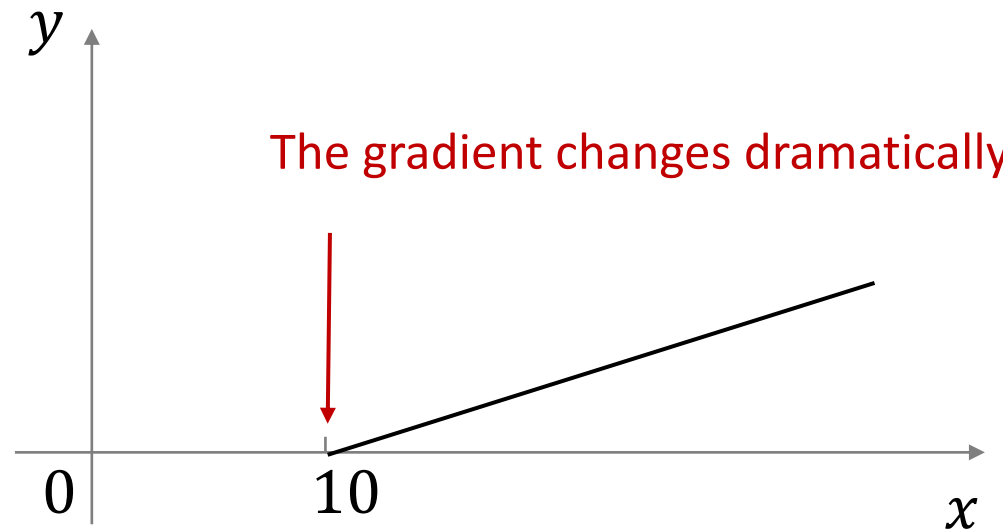
$$y = max(0, x - 10)$$



The gradient changes dramatically
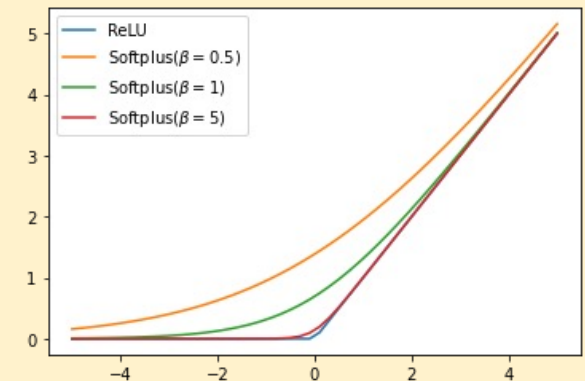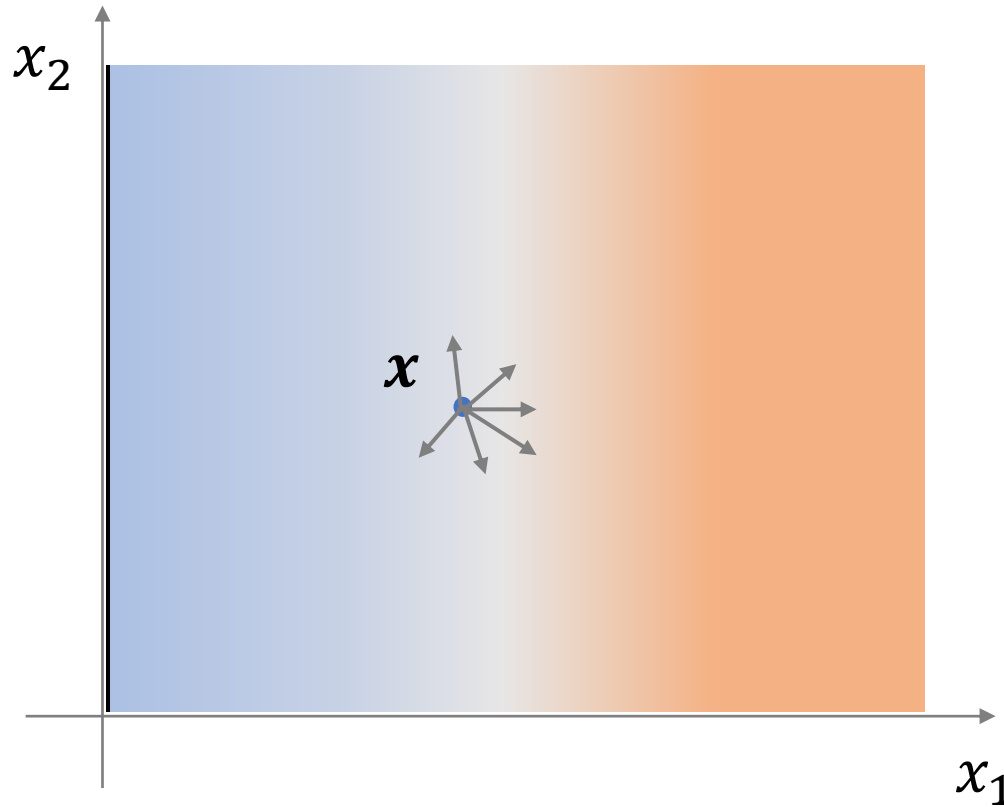
Need to replace "Relu" with "Softplus" activation



(Shrikumar et al., 2017)

# Gradient-based Explanation

Problem 3: input gradient is sensitive to slight perturbations

# Gradient-based Explanation

Do NOT rely on a single gradient calculation

- SmoothGrad: add gaussian noise to inputs and average the gradients

(Smilkov et al., 2017)

# Gradient-based Explanation

Do NOT rely on a single gradient calculation

- SmoothGrad: add gaussian noise to inputs and average the gradients

  (Smilkov et al., 2017)

- Integrated Gradients: aggregate gradients along a path from baseline to the input
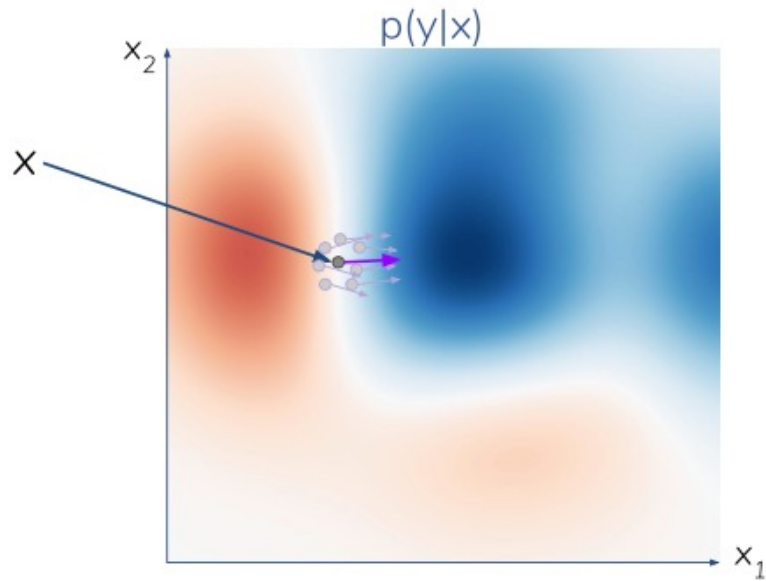
  (Sundararajan et al., 2017)

# Gradient-based Explanation

## Do NOT rely on a single gradient calculation

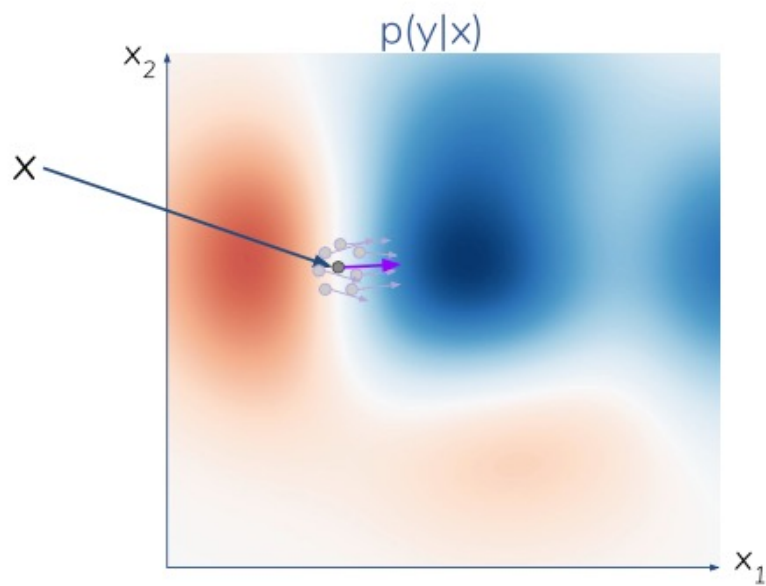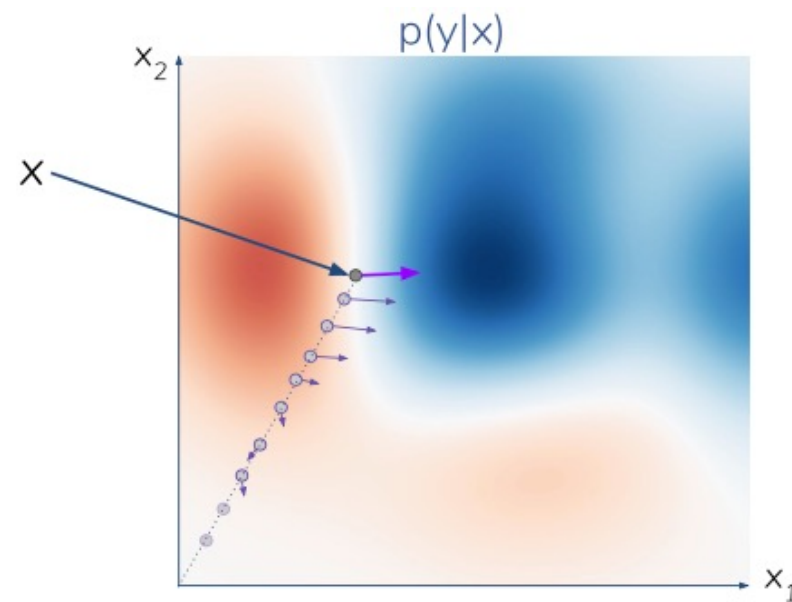- SmoothGrad: add gaussian noise to inputs and average the gradients

  (Smilkov et al., 2017)

- Integrated Gradients: aggregate gradients along a path from baseline to the input

  (Sundararajan et al., 2017)

Source: EMNLP 2020 Tutorial on Interpreting Predictions of NLP Models

# Axiomatic Attribution for Deep Networks

Mukund Sundararajan, Ankur Taly, Qiqi Yan

(ICML, 2017)

# IG

- Integrated Gradients

$f$: neural network

$x \in \mathbb{R}^n$: input

$x' \in \mathbb{R}^n$ : baseline

(e.g., zero embedding vector)

**Get samples along the straight line from $x'$ to $x$**



$x$

$x' + \alpha(x - x')$     $\alpha \in (0, 1)$

$x'$

# IG

- Integrated Gradients

$f$ : neural network

$\boldsymbol{x} \in \mathbb{R}^n$ : input

$\boldsymbol{x'} \in \mathbb{R}^n$ : baseline

(e.g., zero embedding vector)

**Compute gradients at all points along the path**



$\boldsymbol{x'} + \alpha(\boldsymbol{x} - \boldsymbol{x'})$    $\alpha \in (0,1)$

# IG

- Integrated Gradients

$f$: neural network

$x \in \mathbb{R}^n$: input

$x' \in \mathbb{R}^n$ : baseline

(e.g., zero embedding vector)

**Cumulate these gradients**



$x' + \alpha(x - x')$     $\alpha \in (0, 1)$

$$IG_i(x) = (x_i - x_i') \times \int_{\alpha=0}^{1} \frac{\partial f(x' + \alpha(x - x'))}{\partial x_i} d\alpha$$

<span style="color:red">On the $i^{th}$ dimension</span>

# IG

- Integrated Gradients

**Axiom: completeness**

The attributions add up to the difference between the output of $f$ at the input $\boldsymbol{x}$ and the baseline $\boldsymbol{x}'$

$$\sum_{i=1}^{n} IG_i(\boldsymbol{x}) = f(\boldsymbol{x}) - \underline{f(\boldsymbol{x}')}$$
$$\textcolor{red}{f(\boldsymbol{x}') \approx 0}$$

# IG

- Integrated Gradients

**Axiom: completeness**

The attributions add up to the difference between the output of $f$ at the input $\boldsymbol{x}$ and the baseline $\boldsymbol{x}'$

$$\sum_{i=1}^{n} IG_i(\boldsymbol{x}) = f(\boldsymbol{x}) - f(\boldsymbol{x}')$$

**Sensitivity**: for every input and baseline that differ in one feature but have different predictions then the differing feature should be given a non-zero attribution

✅ Sensitivity

# IG

- Integrated Gradients

**Axiom: completeness**

The attributions add up to the difference between the output of $f$ at the input $\boldsymbol{x}$
and the baseline $\boldsymbol{x}'$

$$\sum_{i=1}^{n} IG_i(\boldsymbol{x}) = f(\boldsymbol{x}) - f(\boldsymbol{x}')$$

The chain-rule for gradients is essentially about implementation invariance:

✅ Sensitivity

✅ Implementation invariance

(The attributions are always identical for two functionally equivalent networks)



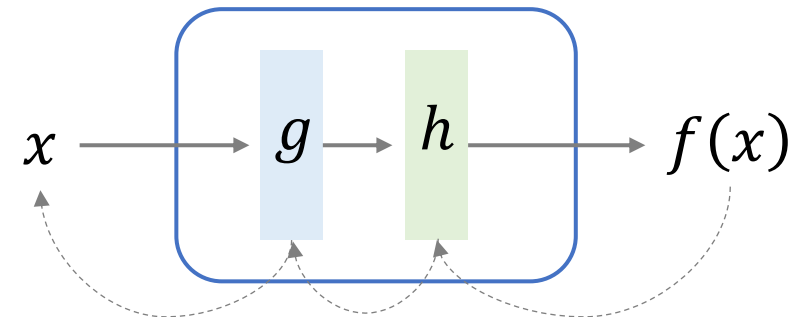$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial h} \cdot \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}$$

# IG

- Integrated Gradients

**Axiom: completeness**

The attributions add up to the difference between the output of $f$ at the input $\boldsymbol{x}$ and the baseline $\boldsymbol{x}'$

$$\sum_{i=1}^{n} IG_i(\boldsymbol{x}) = f(\boldsymbol{x}) - f(\boldsymbol{x}')$$

The chain-rule for gradients is essentially about implementation invariance:

✅ Sensitivity

✅ Implementation invariance

(The attributions are always identical for two functionally equivalent networks)



$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial h} \cdot \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}$$
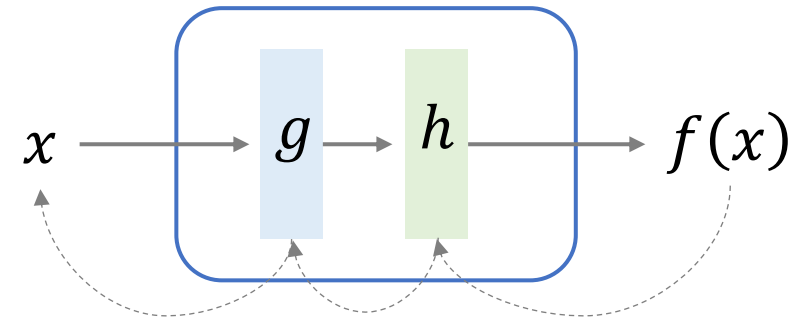
# IG

- Applying Integrated Gradients

  The integral of integrated gradients can be efficiently approximated via a summation

  $$IG_i(\boldsymbol{x}) \approx (x_i - x_i') \times \sum_{k=1}^{m} \frac{\partial f\left(\boldsymbol{x}' + \frac{k}{m}(\boldsymbol{x} - \boldsymbol{x}')\right)}{\partial x_i} \times \frac{1}{m}$$

  $m$: the number of steps

# Question?

# Improving Interpretability

➢ Black-box explanation

➢ White-box explanation

➢ Natural language explanation

# Natural Language Explanation

**Commonsense question-answering (QA)**

**Question**

Why do people go hiking?

**Answer choices**

drink water　　get lost　　enjoy nature　　lose weight　　get tired

Prediction: enjoy nature

**Explanation:** Hiking means the activity of going for long walks especially across country, or in nature. People who go hiking enjoy nature.

# Natural Language Explanation



**Commonsense question-answering (QA)**

**Question**

Why do people go hiking?

**Answer choices**

| drink water | get lost | enjoy nature | lose weight | get tired |

Prediction: enjoy nature

**Explanation:** Hiking means the activity of going for long walks especially across country, or in nature. People who go hiking enjoy nature.

- Flexible
- Understandable
- Informative

# Chain of Thought Prompting



**Standard Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. ✖

**Chain-of-Thought Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✔

(Wei et al., 2022)

# Potential Issues

Why do people go hiking?

🤖 😃 Hiking means the activity of going for long walks especially across country, or in nature. People who go hiking enjoy nature.

🤖 Getting lost in the wilderness is a valuable experience. People go hiking to get lost. ❌ Not factual

🤖 Hiking in nature helps men get rid of jobs. Men go hiking to enjoy nature.

🙅🏻‍♀️ Bias

🤖 Hiking in nature helps women get rid of housework. Women go hiking to enjoy nature.
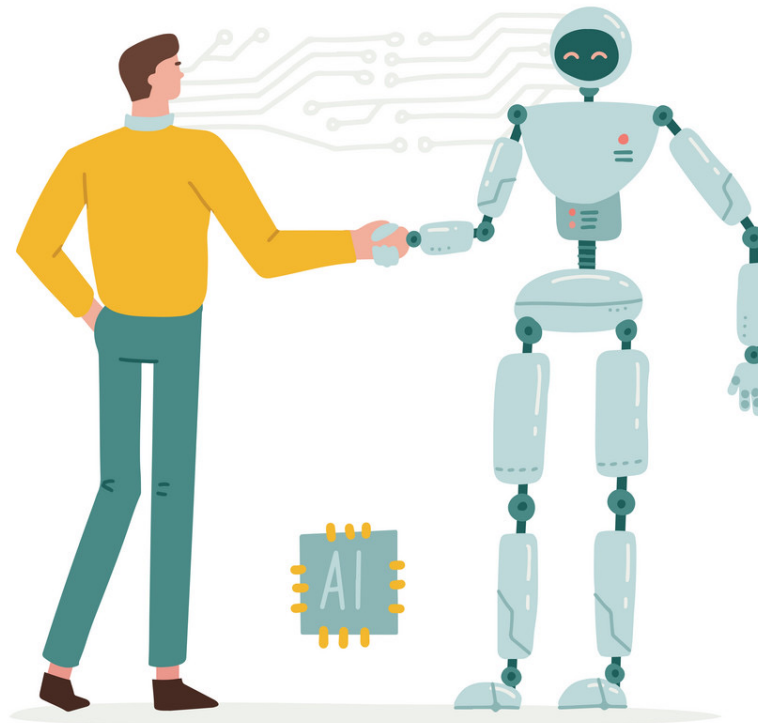
# Question?

# Improving Interpretability

➢ Black-box explanation

➢ White-box explanation

➢ Natural language explanation

Thank you!

# 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).
- Shrikumar, Avanti, Peyton Greenside, and Anshul Kundaje. "Learning important features through propagating activation differences." *International conference on machine learning*. PMLR, 2017.
- Sundararajan, Mukund, Ankur Taly, and Qiqi Yan. "Axiomatic attribution for deep networks." *International conference on machine learning*. PMLR, 2017.
- Wei, Jason, et al. "Chain-of-thought prompting elicits reasoning in large language models." Advances in Neural Information Processing Systems 35 (2022): 24824-24837.