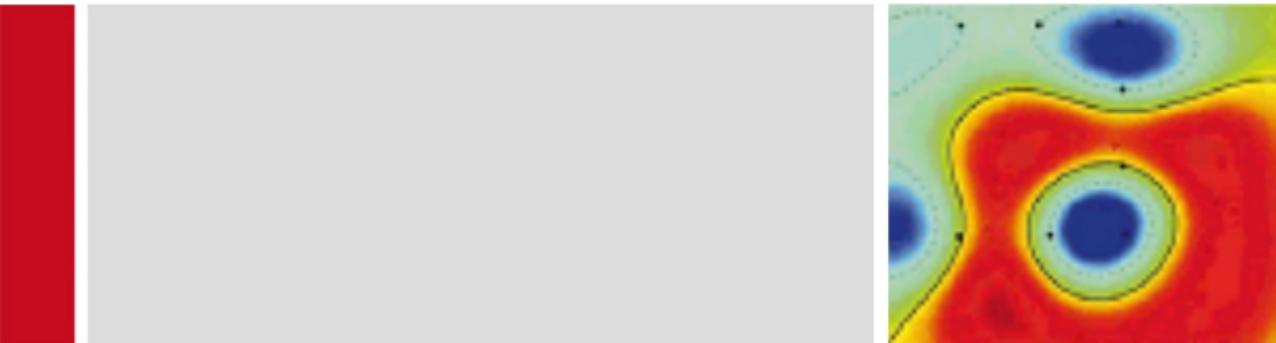


WiSe 2024/25

Machine Learning 1/1-X



## Lecture 14 | Explainable AI

# Outline

---

- ▶ What is Explainable AI?
- ▶ Desiderata of an Explainable AI technique
- ▶ Uses of Explainable AI
- ▶ Methods for Explainable AI
  - ▶ Activation Maximization
  - ▶ Shapley Values
  - ▶ Taylor Expansions
  - ▶ Layer-wise Relevance Propagation

# What is Explainable AI?

◆ 函数  $f$  通常被认为是一个“黑箱”，其参数是通过数据学习得到的，例如使用梯度下降法。优化目标是使预测值  $f(x)$  尽可能接近训练数据和测试数据的真实值。

## Standard machine learning:

- ▶ The function  $f$  is typically considered to be a “black-box” whose parameters are learned from the data using e.g. gradient descent. The objective to minimize encourages the predictions  $f(x)$  to coincide with the ground truth on the training and test data.



## Machine learning + Explainable AI:

- ▶ We do not only look at the outcome  $f(x)$  of the prediction but also at the way the prediction is produced by the ML model, e.g. which features are used, how these features are combined, or to what input pattern the model responds the most.
  - 模型对哪些输入模式最敏感

# What is Explainable AI?

合成一个输入模式，使其能够最强烈地激活机器学习模型的某个特定类别的输出。

**Example 1:** Synthesize an input pattern that most strongly activates the output of the ML model associated to a particular class.



Image source: Nguyen et al. (2016) Multifaceted Feature Visualization: Uncovering the Different Types of Features Learned By Each Neuron in Deep Neural Networks

在神经网络中，我们可以使用反向传播优化输入数据，使其尽可能激活目标类别的神经元。例如：

- 在图像分类任务中，找到一张合成的图片，它能够最大程度地激活“猫”这个类别的神经元。
- 在文本分类任务中，找到一段文本，使其能最大程度地被模型识别为“正面情感”类别。

# What is Explainable AI?

**Example 2:** Highlight features that have contributed for a given data point to the ML prediction.

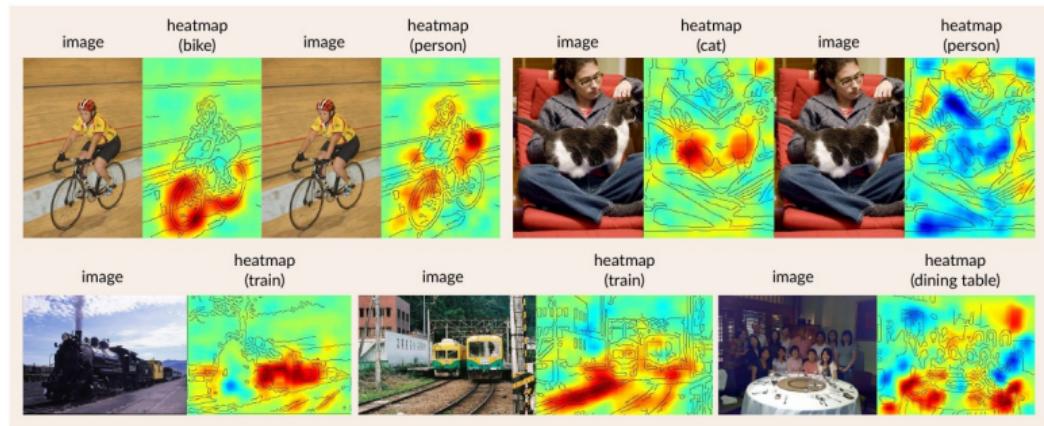
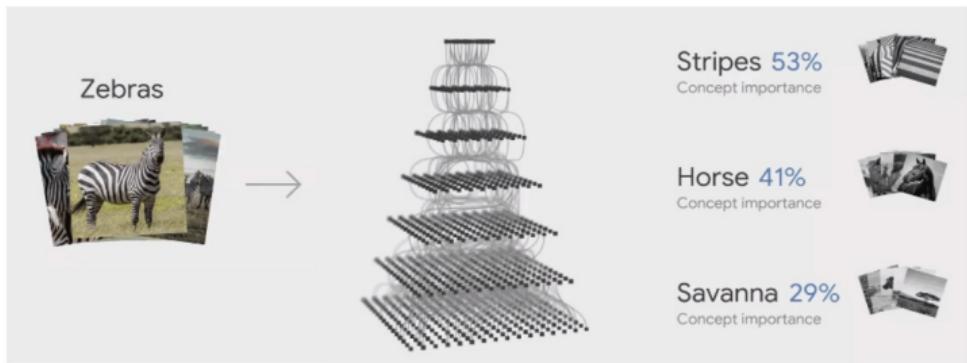


Image source: Lapuschkin et al. (2016) Analyzing Classifiers: Fisher Vectors and Deep Neural Networks

# What is Explainable AI?

概念激活向量 (TCAV)。突出显示用于解释某个数据点的机器学习预测的中层概念。

**Example 3:** Concept activation vectors (TCAV). Highlight the mid-level concepts that explain, for a given data point, the ML prediction.



帮助我们理解模型的决策是基于哪些概念。

[www.youtube.com/watch?v=lyRPyRKH08M&t=2279s](https://www.youtube.com/watch?v=lyRPyRKH08M&t=2279s))

示例分析（斑马分类）：

- 假设一个深度神经网络被训练用于识别动物，并正确分类了一张斑马的图片。
- 但是，模型的决定是基于什么？ TCAV 可以分析中层特征，发现：
  - 条纹 (Stripes) 对于识别“斑马”最重要，占比 53%。
  - 马 (Horse) 的特征也有一定影响，占比 41%（因为斑马和马形态相似）。
  - 稀树草原 (Savanna) 也影响预测，占比 29%（斑马的自然栖息地）。



# Desiderata of an Explanation

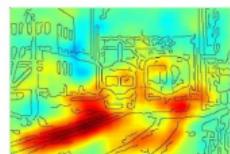
In practice, we would like the explanation technique to satisfy a number of properties:

1. **Fidelity:** The explanation should reflect the quantity being explained and not something else.
2. **Understandability:** The explanation must be easily understandable by its receiver.
3. **Sufficiency:** The explanation should provide sufficient information on how the model came up with its prediction.
4. **Low Overhead:** The explanation should not cause the prediction model to become less accurate or less efficient.
5. **Runtime Efficiency:** Explanations should be computable in reasonable time.

image



heatmap (train)



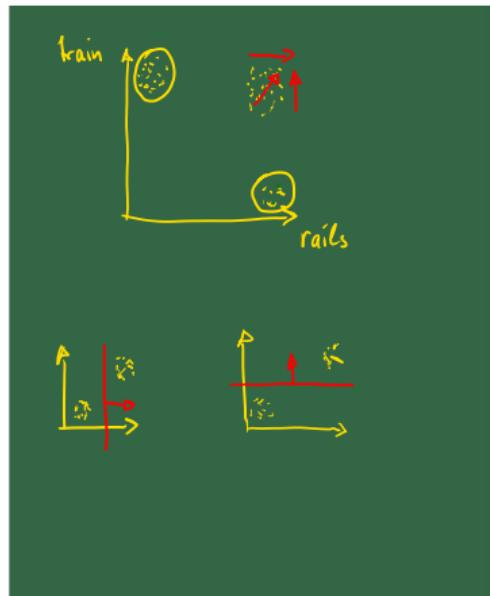
0.0	0.0
0.8	0.2

see also Swartout & Moore (1993), *Explanation in Second Generation Expert Systems*.

# Uses of an Explanation

## Verify (and improve?) a ML model

- ▶ Verify that the model is based on features which *generalize well* to examples outside the current data distribution (this cannot be done with standard validation techniques!).
- ▶ Reliance of the ML models on wrong features is often encountered when there are *spurious correlation* in the data.
- ▶ From the explanation, the model's trustworthiness can be reevaluated, and the flawed ML model can be potentially *retrained* based on the user feedback.



验证（并改进？）机器学习模型

- ◆ 验证 模型是否基于能够很好泛化到当前数据分布之外的示例的特征（标准验证技术无法做到这一点！）。
- ◆ 机器学习模型依赖错误特征 的情况通常出现在数据中存在“虚假相关”（spurious correlation）\*\* 时。
- ◆ 通过解释，我们可以重新评估模型的可信度，并基于用户反馈对有缺陷的机器学习模型进行再训练（retrain）。

# Uses of an Explanation

- 图片（左侧）：一个骑马的人，图片左下角带有版权标记。

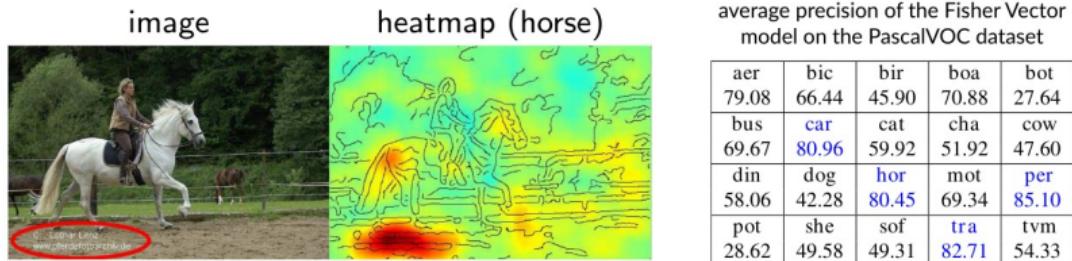
- 热力图（中间）：机器学习模型用于识别“马”(horse)的特征重要性图。

- 表格（右侧）：在Pascal VOC数据集上的Fisher Vector模型的类别平均精度。

分析：◆ 在这个示例中，分类器准确预测了“马”类别，但它基于错误的特征（如左下角的版权标记）进行判断。

◆ 这种错误的决策方式无法通过测试误差检测到，只有使用解释技术才能发现。

## Example: The classifier is right for the wrong reasons



- In this example, the classifier accurately predicts the horse class, but based on the wrong features (some copyright tag in the corner).
- This incorrect decision strategy cannot be detected by just looking at the test error.

cf. Lapuschkin et al. (2019) Unmasking Clever Hans Predictors and Assessing What Machines Really Learn. Nature Communications

# Uses of an Explanation

Learn something about the data  
(or about the system that produced the data)

- ▶ Step 1: Train a ML model that predicts well the data.
- ▶ Step 2: Apply XAI to the trained ML model to produce explanations of the ML decision strategy.
- ▶ Step 3: Based on the XAI explanations, the user can compare his reasoning with that of the ML model, and can potentially refine his own domain knowledge.

学习数据（或生产数据的系统）的信息

步骤 1：训练一个能很好预测数据的机器学习模型。

步骤 2：将可解释性人工智能（XAI）应用于训练好的机器学习模型，以生成模型决策策略的解释。

步骤 3：基于 XAI 生成的解释，用户可以比较自己的推理与模型的推理，并潜在地改进自己的领域知识。

图片右侧：

- 这是一个神经影像分析的示例，展示了如何使用 XAI 解释深度学习模型的决策过程。

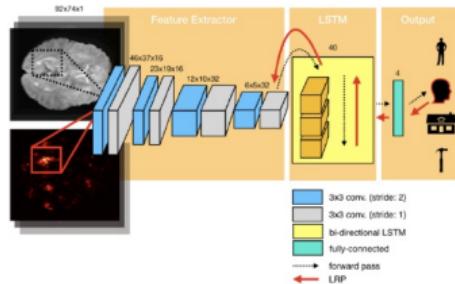


Image source: Thomas et al. (2019) Analyzing Neuroimaging Data Through Recurrent Deep Learning Models

# Part II: Methods of XAI

---

## Presented methods

- ▶ Activation maximization
- ▶ Shapley values
- ▶ Taylor expansion
- ▶ Layer-wise relevance propagation

## Other methods

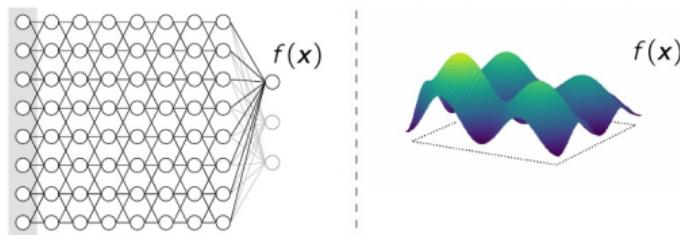
- ▶ Surrogate models (LIME)
- ▶ Integrated gradients / expected gradients / SmoothGrad
- ▶ Influence functions
- ▶ ...

# Activation Maximization

---

Assume we have trained a ML model (e.g. a neural network), and we would like to understand what concept is associated to some particular output neuron of the ML model, e.g. the output neuron that codes for the class 'cat'. Activation maximization proceeds in two steps:

- ▶ **Step 1:** Think of the ML model as a function of the input



- ▶ **Step 2:** Explain the function  $f$  by generating a maximally activating input pattern:

$$x^* = \arg \max_x f(x, \theta)$$

# Activation Maximization

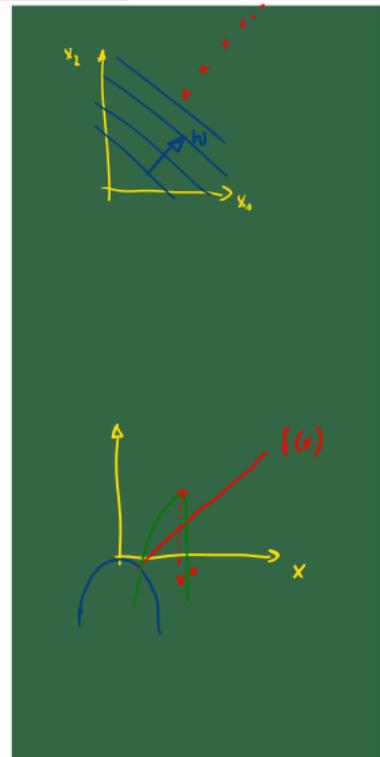
**Problem:** In most cases  $f(x)$  does not have single point corresponding to the maximum.

*E.g. in linear models,  $f(x) = \mathbf{w}^\top \mathbf{x} + b$ , we can keep moving the point  $x$  further along the direction  $\mathbf{w}$ , and the output continues to grow).*

Therefore, we would like to apply a preference for 'regular' regions of the input domain, i.e.

$$\mathbf{x}^* = \arg \max_{\mathbf{x}} f(\mathbf{x}) + \Omega(\mathbf{x})$$

In practice, the preference can be for data points with small norm (i.e. we set  $\Omega(\mathbf{x}) = -\lambda \|\mathbf{x}\|^2$  so that points with large norm are penalized.)



# Activation Maximization: Examples

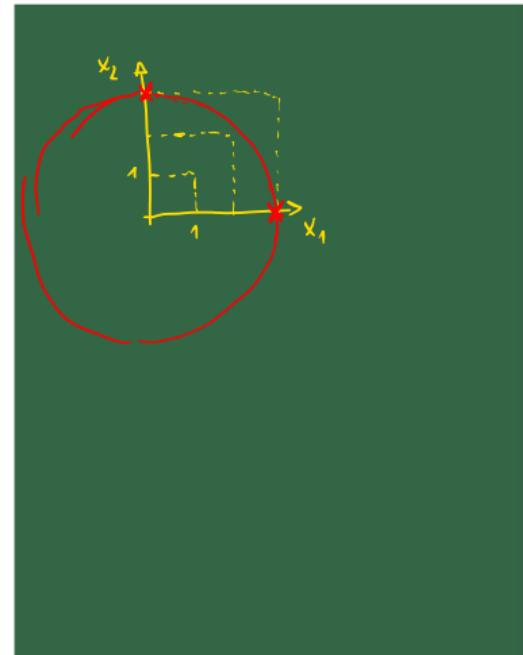
$$f(x) = \mathbf{w}^\top \mathbf{x} + b \text{ and}$$
$$\Omega(x) = -\lambda \|\mathbf{x}\|^2$$

$$f(x) = \max(x_1, x_2) \text{ and}$$
$$\Omega(x) = -\lambda \|\mathbf{x}\|^2$$

$$\underset{\mathbf{x}}{\operatorname{argmax}} \underbrace{\mathbf{w}^\top \mathbf{x} + b}_{J(\mathbf{x})} - \lambda \|\mathbf{x}\|^2$$

$$\frac{\partial J}{\partial \mathbf{x}} = \mathbf{w} - 2\lambda \mathbf{x} = \mathbf{0}$$

$$\mathbf{x}^* = \frac{1}{2\lambda} \cdot \mathbf{w}$$



# Activation Maximization: Probability View

Assume the model produces a log-probability for class  $\omega_c$ :

$$f(\mathbf{x}) = \log p(\omega_c | \mathbf{x})$$

The input  $\mathbf{x}^*$  that maximizes this function can be interpreted as the point where the classifier is the most *sure* about class  $\omega_c$ .

Choose the regularizer  $\Omega(\mathbf{x}) = \log p(\mathbf{x})$ , i.e. favor points that are *likely*.

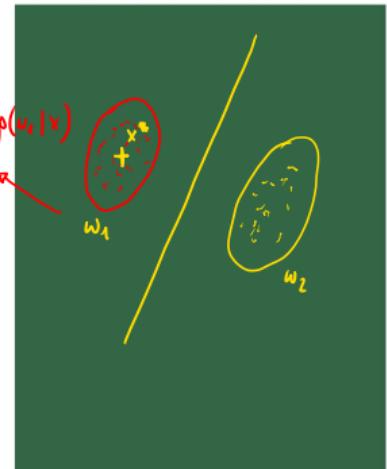
The optimization problem becomes:

$$\begin{aligned}\mathbf{x}^* &= \arg \max_{\mathbf{x}} \log p(\omega_c | \mathbf{x}) + \log p(\mathbf{x}) \\ &= \arg \max_{\mathbf{x}} \log p(\mathbf{x} | \omega_c)\end{aligned}$$

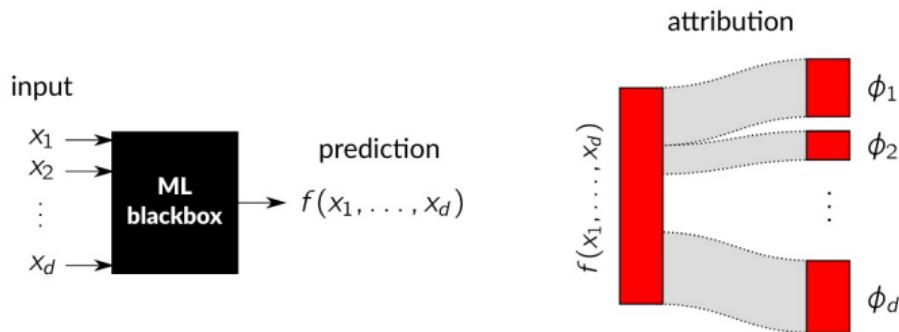
$$p(\omega_c | \mathbf{x}) = \frac{p(\mathbf{x} | \omega_c) p(\omega_c)}{p(\mathbf{x})}$$

$$\log p(\omega_c | \mathbf{x}) = \log p(\mathbf{x} | \omega_c) + \log p(\omega_c) - \log p(\mathbf{x})$$

where  $\mathbf{x}^*$  can now be interpreted as the most *typical* input for class  $\omega_c$ .



# Attribution of a Prediction to Input Features



1. The data  $x \in \mathbb{R}^d$  is fed to the ML black-box and we get a prediction  $f(x) \in \mathbb{R}$ .
2. We explain the prediction by determining the contribution of each input feature.

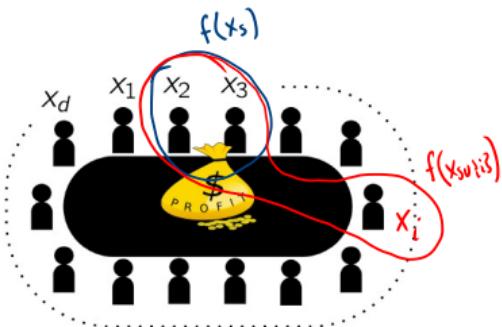
Key property of an explanation: **conservation** ( $\sum_{i=1}^d \phi_i = f(x)$ ).

# Attribution: Shapley Values

◆ 该框架最初是在博弈论 (Shapley 1951) 的背景下提出的，用于在合作博弈中分配收益，并最近被应用于机器学习模型。

◆ 每个输入变量被视为一个玩家，而函数输出则被视为合作玩家所实现的收益。

- ▶ Framework originally proposed in the context of game theory (Shapley 1951) for assigning payoffs in a cooperative game, and recently applied to ML models.
- ▶ Each input variable is viewed as a player, and the function output as the profit realized by the cooperating players.



The Shapley values  $\phi_1, \dots, \phi_d$  measuring the contribution of each feature are:

$$\phi_i = \sum_{S: i \notin S} \frac{|S|!(d-|S|-1)!}{d!} [f(x_{S \cup \{i\}}) - f(x_S)]$$

表示当加入特征  $x_i$  后，模型的预测值变化。

where  $(x_S)_S$  are all possible subsets of features contained in the input  $x$ .

- $S$  表示输入特征  $x$  的所有可能子集。
- $f(x_S)$  表示模型在特征子集  $S$  下的预测值。

# Attribution: Shapley Values

$\alpha_S$  是 Shapley 值公式中的加权因子，用于衡量子集  $S$  在所有可能排列中的重要性。

它的作用可以理解为：

- 衡量不同子集  $S$  对最终 Shapley 值的贡献权重，确保所有可能的排列都被公平计算。
- 保证 Shapley 值满足博弈论中的对称性、公平性等性质。

$$\text{Recall: } \phi_i = \sum_{S: i \notin S} \underbrace{\frac{|S|!(d-|S|-1)!}{d!}}_{\alpha_S} \underbrace{[f(\mathbf{x}_{S \cup \{i\}}) - f(\mathbf{x}_S)]}_{\Delta_S}$$

注意：  
1 · (1 + 1) = 2  
0 · (1 + 1) = 0

**Worked-through example:** Consider the function  $f(\mathbf{x}) = x_1 \cdot (x_2 + x_3)$ . Calculate the contribution of each feature to the prediction  $f(\mathbf{1}) = 1 \cdot (1+1) = 2$ .

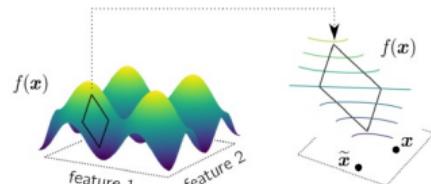
$\Phi_1: S$	$\mathbf{x}_S$	$\Delta_S$	$\Phi_2: S$	$\alpha_S$	$\Delta_S$	$\Phi_3: \text{Symmetry/conservation}$
$\{\}$	$\frac{4(1 \cdot (1-1))}{3!} = 1/3$	0	$\{\}$	$1/3$	0	
$\{2\}$	$\frac{4(1 \cdot (2-1))}{3!} = 1/6$	1	$\{1\}$	$1/6$	1	
$\{3\}$	$1/6$	1	$\{3\}$	$1/6$	0	
$\{1, 2\}$	$1/3$	2	$\{1, 3\}$	$1/3$	1	
	$\frac{4(1 \cdot (3-2-1))}{3!} = 1/6$					$\Phi_2 = 1/2$
						$\Phi_3 = 1/2$

$S$	$f(\mathbf{x}_S)$	$f(\mathbf{x}_{S \cup \{1\}})$	$\Delta_S$
$\emptyset$	$f(0, 0, 0) = 0$	$f(1, 0, 0) = 0$	0
$\{2\}$	$f(0, 1, 0) = 0$	$f(1, 1, 0) = 1$	1
$\{3\}$	$f(0, 0, 1) = 0$	$f(1, 0, 1) = 1$	1
$\{1, 2\}$	$f(0, 1, 1) = 0$	$f(1, 1, 1) = 2$	2

$S$	$f(\mathbf{x}_S)$	$f(\mathbf{x}_{S \cup \{2\}})$	$\Delta_S$
$\emptyset$	$f(0, 0, 0) = 0$	$f(0, 1, 0) = 0$	0
$\{1\}$	$f(1, 0, 0) = 0$	$f(1, 1, 0) = 1$	1
$\{3\}$	$f(0, 0, 1) = 0$	$f(0, 1, 1) = 0$	0
$\{1, 3\}$	$f(1, 0, 1) = 1$	$f(1, 1, 1) = 2$	1

# Attribution: Taylor Expansions

- ▶ Many ML models  $f(x)$  are complex and nonlinear when taken globally but are simple and linear when taken locally.
- ▶ The function can be approximated locally by some Taylor expansion:



$$f(x) = \underbrace{f(\tilde{x})}_{=0} + \sum_{i=1}^d \underbrace{[\nabla f(\tilde{x})]_i \cdot (x_i - \tilde{x}_i)}_{\phi_i} + \dots$$

- ▶ First-order terms  $\phi_i$  of the expansion can serve as an explanation.
- ▶ The explanation  $(\phi_i)_i$  depends on the choice of root point  $\tilde{x}$ .

# Attribution: Taylor Expansions

**Example:** Attribute the prediction  $f(x) = \mathbf{w}^\top \mathbf{x}$  with  $\mathbf{x} \in \mathbb{R}^d$  on the  $d$  input features.

$$\tilde{\mathbf{x}} = \mathbf{0}$$

$$f(\mathbf{x}) = \underbrace{f(\tilde{\mathbf{x}})}_{=0} + \sum_{i=1}^d \underbrace{\left[ \nabla f(\tilde{\mathbf{y}}) \right]_i}_{\omega_i} \cdot \underbrace{(\mathbf{x}_i - \tilde{\mathbf{x}}_i)}_{x_i}$$
$$\phi_i = \omega_i x_i$$

# Attribution: Taylor Expansions

局限性 (Limitations) : 梯度信息过于局部化

**Limitations:** Gradient information is too localized.

- ▶ Cannot handle *saturation effects* and *discontinuities* e.g. cannot explain the function

$$f(\mathbf{x}) = \sum_{i=1}^d x_i - \max(0, x_i - \theta)$$

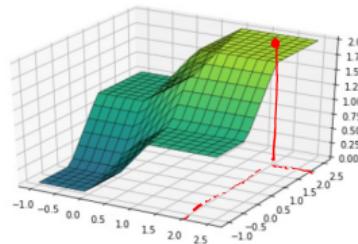
at the point  $\mathbf{x} = (2, 2)$ .

◆ 无法处理 饱和效应 (saturation effects) 和 不连续性 (discontinuities)，例如，无法解释如下函数：

$$f(\mathbf{x}) = \sum_{i=1}^d x_i - \max(0, x_i - \theta)$$

在点  $\mathbf{x} = (2, 2)$  处的行为。

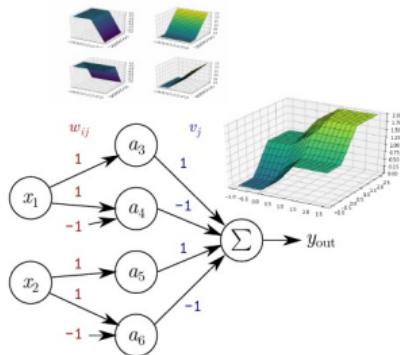
◆ 这种局限性可以通过分析模型的结构，并将解释问题分解为多个部分来克服 (→ 下一页)。



This limitation can be overcome by looking at the structure of the model and decompose the problem of explanation in multiple parts (→ next slide).

# Attribution: Look at the Structure of the Model

## Observation:



归因：观察模型的结构

观察（Observation）：

- ◆ 机器学习模型实现的函数通常是由多个简单基本函数的组合。
- ◆ 这些基本函数比整个输入-输出函数更容易分析。

思路（Idea）：

- ◆ 将解释问题视为在输入-输出图中向后传播预测值的问题。
- ◆ 分层相关传播（Layer-wise Relevance Propagation, LRP）方法可以实现该思路，并用于解释机器学习模型（→下一页）。

- ▶ The function implemented by a ML model is typically a composition of simple elementary functions.
- ▶ These functions are simpler to analyze than the whole input-output function.

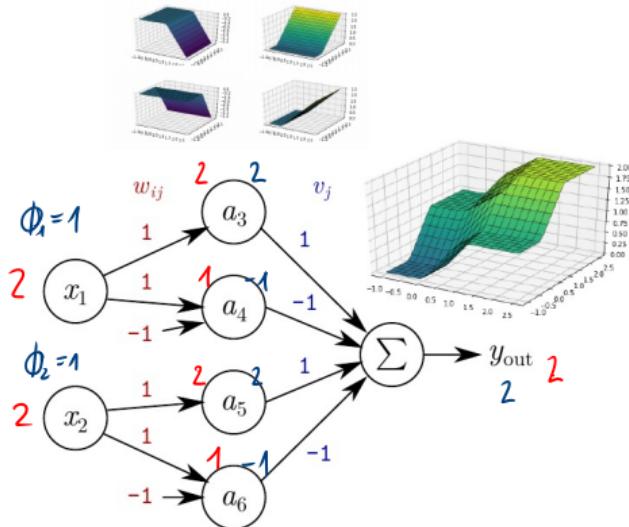
## Idea:

- ▶ Treat the problem of explanation as propagating the prediction backward in the input-output graph.
- ▶ The layer-wise relevance propagation (LRP) method implements this approach and can be used to explain ML models (→ next slide).

# Attribution: The LRP Method

**Example:** Consider  $y_{\text{out}}$  to be the quantity to explain:

$$R_{\text{out}} \leftarrow y_{\text{out}}$$



如何计算 LRP 归因?

1. 从输出层开始，分配所有预测值的贡献。
2. 逐层向后传播贡献值，计算中间神经元的贡献。
3. 最终归因到输入特征，得到每个特征对预测值的贡献程度。

示例解析

- 在隐藏层 (Step 1)，每个神经元的贡献是按比例分配的，即神经元的激活值  $a_j$  乘以它的权重  $v_j$ 。
- 在输入层 (Step 2)，每个特征的贡献是按它的连接权重  $w_{ij}$  进行分配的。

► **Step 1:** Propagate on the hidden layer

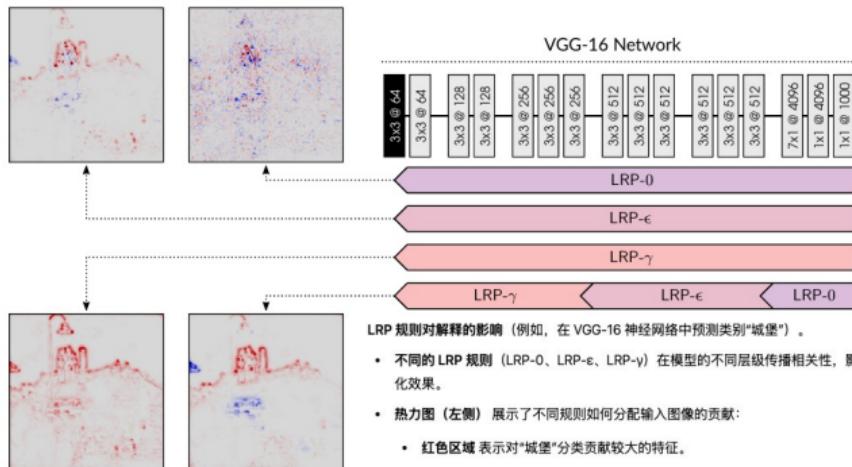
$$\forall_{j=3}^6 : R_j \leftarrow \frac{a_j v_j}{\sum_{j=3}^6 a_j v_j} R_{\text{out}}$$

► **Step 2:** Propagate on the first layer

$$\forall_{i=1}^2 : R_i \leftarrow \sum_{j=3}^6 \frac{x_i w_{ij}}{\sum_{i=1}^2 x_i w_{ij}} R_j$$

# Attribution: The LRP Method

Effect of LRP rules on the explanation (e.g. class 'castle' predicted by a VGG-16 neural network.)



## 解释

- VGG-16 网络 是一个预训练的 CNN（卷积神经网络），用于图像分类。
- LRP 规则 允许我们查看模型如何做出预测，通过不同的 LRP 规则可以调整解释的粒度：
  - LRP-0 直接反向传播。
  - LRP- $\epsilon$  (带小正则项) 可以减少噪声，提高稳定性。
  - LRP- $\gamma$  (增强重要特征) 让更重要的特征获得更高的权重。

# XAI for Unsupervised Learning

示例：聚类 (Clustering)

假设我们使用 K-means 聚类算法，我们希望解释某个数据点  $x$  的聚类分配 (cluster assignment)。

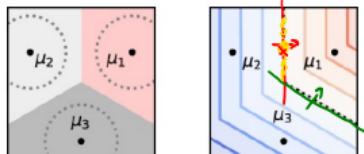
步骤 1

将 K-means 结果转换为等效的神经网络，该神经网络可对  $x$  进行聚类分配。

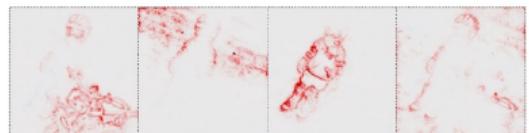
步骤 2

对神经网络应用 LRP (层级相关传播)，以解释数据点  $x$  的聚类归属。

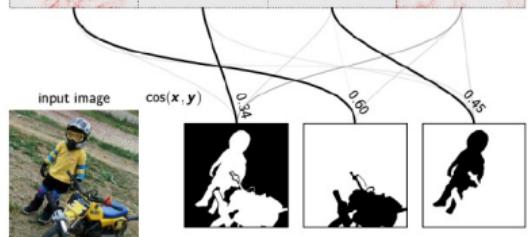
**Example Clustering:** Consider the K-means algorithm, and we would like to explain the cluster assignment of a data point  $x$ .



**Step 1:** Transform the K-means solution into an equivalent neural network that assigns a cluster to  $x$ .



**Step 2:** Apply LRP to the neural network to explain the cluster assignment.



# XAI for Clustering

---

**Step 1:** K-means as a Neural Network:

$$\begin{aligned}\text{decide cluster } c &\Leftrightarrow \min_{k \neq c} \{\|x - \mu_k\|^2\} > \|x - \mu_c\|^2 \\ &\Leftrightarrow \min_{k \neq c} \{\|x - \mu_k\|^2 - \|x - \mu_c\|^2\} > 0 \\ &\Leftrightarrow \min_{k \neq c} \underbrace{\mathbf{w}_k^\top x + b_k}_{h_k} > 0\end{aligned}$$

**Step 2:** LRP applied to the neural network:

$$\begin{aligned}R_k &= \frac{\exp(-\beta h_k)}{\sum_k \exp(-\beta h_k)} R_{\text{out}} \\ R_i &= \sum_k \frac{\mathbf{w}_{ik}(x_i - m_{ik})}{\sum_i \mathbf{w}_{ik}(x_i - m_{ik})} R_k\end{aligned}$$

# Summary

---

- ▶ Explainable AI is an important addition to classical ML models (e.g. for validating a ML model or extracting knowledge from it).
- ▶ Many XAI methods have been developed, each of them, with their strengths and limitations:
  - ▶ *Activation maximization* can be used to understand what a ML model has learned, but is unsuitable for explaining an individual prediction  $f(x)$ .
  - ▶ *Shapley value* has strong theoretical foundations, but is computationally unfeasible for high-dimensional input data.
  - ▶ *Taylor expansions* are simple and theoretically founded for simple models, but the expansion does not extrapolate well in complex nonlinear models.
  - ▶ *LRP* leverages the structure of the ML model to handle nonlinear decision functions, but requires to carefully choose the propagation rules.

# Summary

---

- ◆ 可解释性人工智能（XAI）是传统机器学习模型的重要补充（例如用于验证机器学习模型或从中提取知识）。
  - ◆ 目前已开发出多种 XAI 方法，每种方法都有其优点和局限性：
1. 激活最大化（Activation Maximization)
    - 可以用于理解机器学习模型学习到了什么，但不适用于解释单个预测结果  $f(x)$ 。
  2. Shapley 值（Shapley Value)
    - 具有强大的理论基础，但在高维输入数据情况下计算成本过高，不具备可行性。
  3. 泰勒展开（Taylor Expansions)
    - 对于简单模型来说，泰勒展开计算简单且有理论依据。
    - 但是，该方法在复杂的非线性模型中无法很好地推广（外推能力较差）。
  4. 层级相关传播（LRP, Layer-wise Relevance Propagation)
    - 利用机器学习模型的结构来处理非线性决策函数，比梯度方法更稳定。
    - 但需要谨慎选择传播规则，否则可能导致解释偏差。