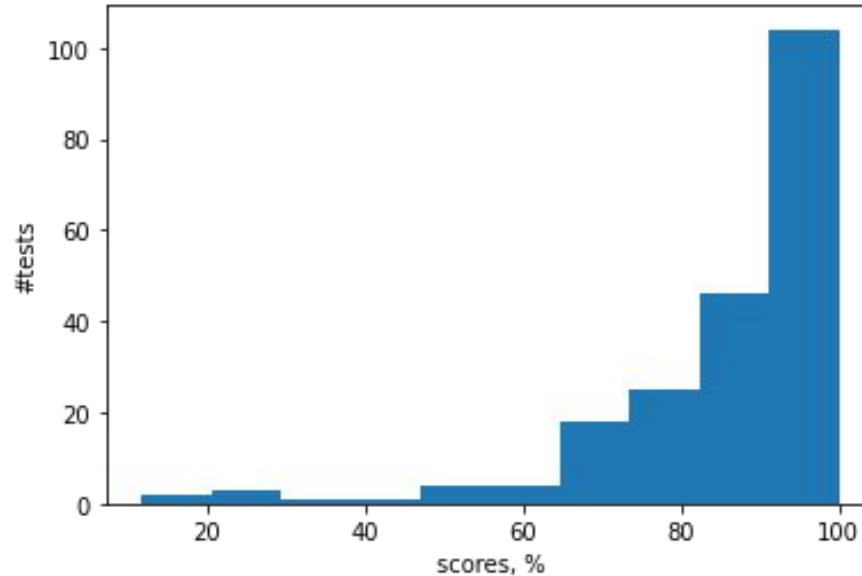


# Test results

Total number of test results : 208

Average score : 87.19 %

Score distribution :



# Preliminary results. Dynamics

Taking the currently maximum scores, we have the following results:

“5” (from 80%) :	159	76%
------------------	-----	-----

“4” (between 60 and 80%) :	37	18%
----------------------------	----	-----

“3” (between 40 and 60%) :	6	3%
----------------------------	---	----

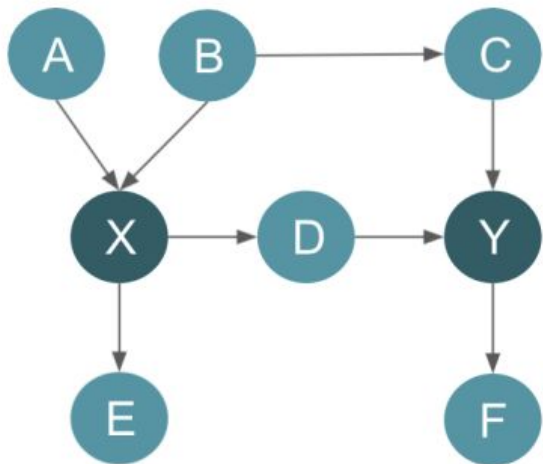
“failed” (less 40%) :	6	3%
-----------------------	---	----

# Some common mistakes

В двухступенчатом MDL (а именно,  $L(D, H) = L(D|H) + L(H)$ ) отсутствует компонент, учитывающий сложность модели  $H$

# Some common mistakes

Примените backdoor критерий для оценки влияния  $X$  на  $Y$ . Выберите наименьшее по размеру множество вершин, необходимых для блокирования всех backdoor-путей

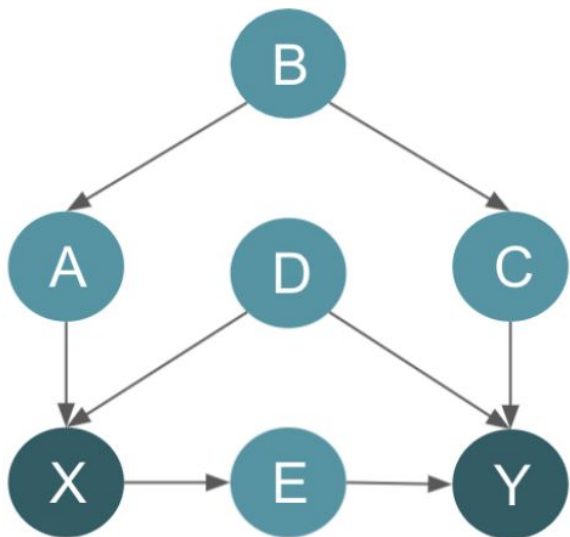


- A)  $\{B\}$
- Б)  $\{A, C\}$
- В)  $\{B, C\}$
- Г)  $\{A, B, C\}$
- Д)  $\{D\}$
- Е)  $\{B, D\}$

**Правильные ответы: А**

# Some common mistakes

Примените backdoor критерий для оценки влияния  $X$  на  $Y$ . Выберите наименьшие по размеру множество вершин, необходимых для блокирования всех backdoor-путей



- A)  $\{B\}$
- Б)  $\{D\}$
- В)  $\{B, D\}$
- Г)  $\{C, D\}$
- Д)  $\{A, C\}$
- Е)  $\{E\}$
- Ж)  $\{A, E\}$

**Правильные ответы: В, Г**

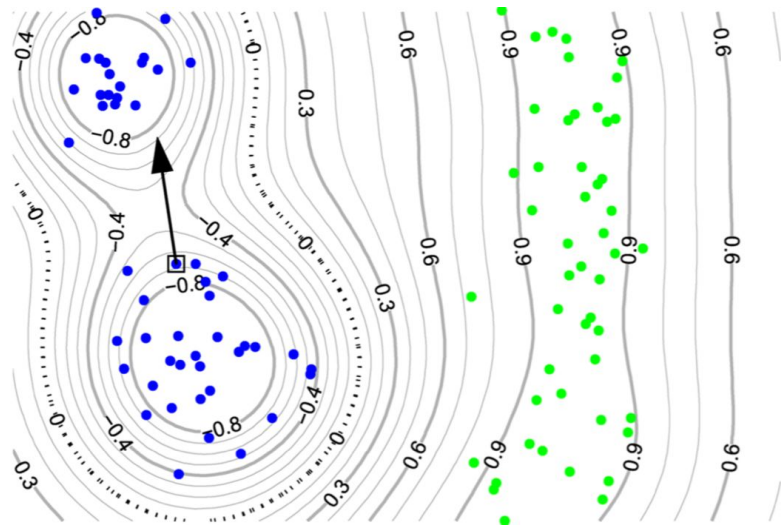
# Lecture 7. Layer-wise Relevance Propagation and Deep Taylor Decomposition

# Part 1. Layer-wise Relevance Propagation

# First order Taylor approximation

The first order Taylor approximation  $f(x) \approx f(x_0) + \sum_{d=1}^V f'(x_0)(x_{(d)} - x_{0(d)})$

Saliency (sensitivity) maps: an image  $x$  is explained by  $f'(x)$



The **blue** dots are labeled **negatively**, the **green** dots are labeled **positively**

Local gradient of the classification function at the prediction point

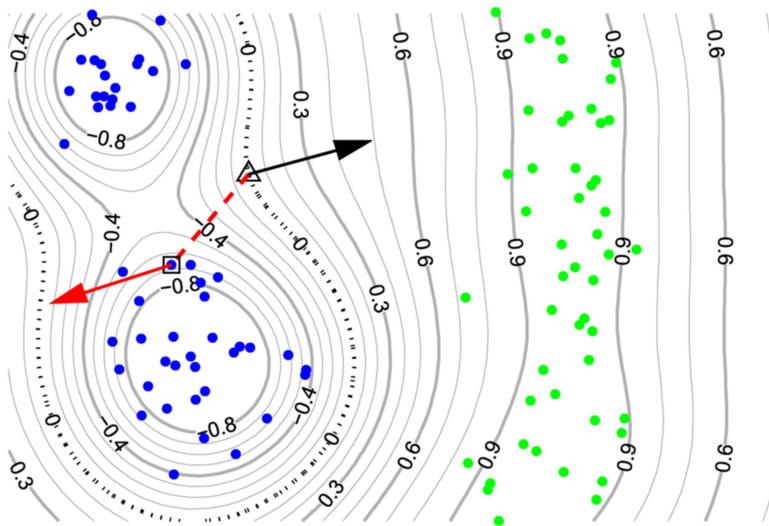
The closest neighbors of the other class can be found at a very different angle. Thus, the local gradient at the prediction point  $x$  may not be a good explanation for the contributions of single dimensions to the function value  $f(x)$



# First order Taylor approximation

The first order Taylor approximation  $f(x) \approx f(x_0) + \sum_{d=1}^V f'(x_0)(x_{(d)} - x_{0(d)})$

We are interested to find out the contribution of each pixel relative to the state of maximal uncertainty of the prediction, i.e.,  $f(x_0) = 0$ , i.e.,  $f(x) \approx \sum_{d=1}^V f'(x_0)(x_{(d)} - x_{0(d)})$



$\Delta$  the nearest root point  $x_0$  on the decision boundary

$\rightarrow f'(x_0)$

---  $x - x_0$

$\rightarrow$  the approximation of  $f(x)$  by Taylor expansion around  $x_0$  (equivalent to the diagonal of the outer product between  $f'(x_0)$  and  $x - x_0$ )

# Layer-wise relevance propagation

**Goal:** to find out the contribution of each input pixel to a particular prediction

**Idea:** In case of classification, to find out the contribution of each pixel relative to the state of maximal uncertainty of the prediction which is given by the set of points  $f(x_0) = 0$ , since  $f(x) > 0$  denotes presence and  $f(x) < 0$  absence of the learned structure,  $x_0$  is called a reference point

**Remark:**  $f(x) < 0$  is less desirable values, since it is difficult to interpret negative evidence for a class

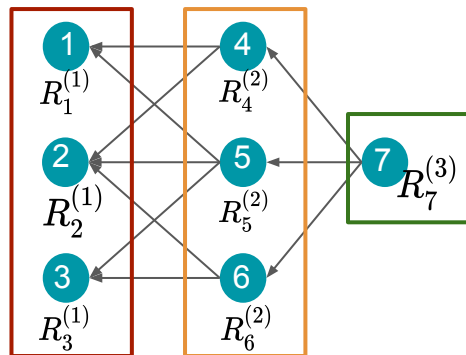
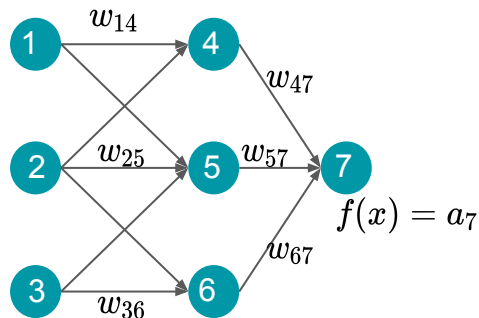
**Basic hypothesis** (the conservation property): relevance is constant throughout the layers, i.e.,

$$f(x) = \dots = \sum_{d \in l+1} R_d^{(l+1)} = \sum_{d \in l} R_d^{(l)} = \dots = \sum_{d \in 1} R_d^{(1)}$$

where  $R_d^{(l)}$  a relevance score for each dimension  $d$  at the level  $l$

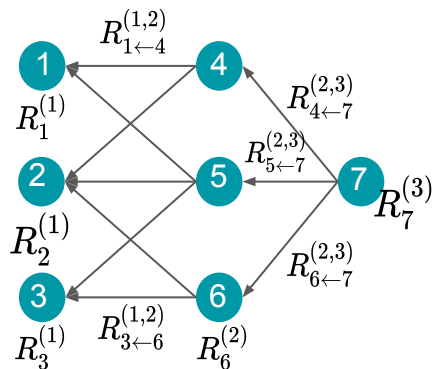
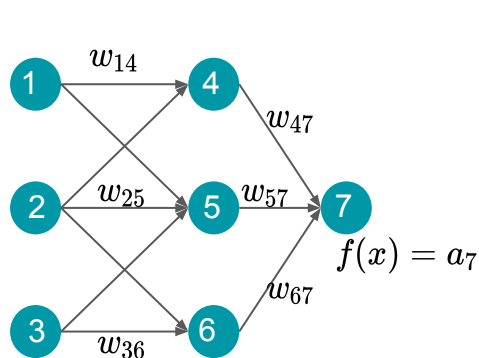
**Remark:** the decomposition satisfying the conservation property is not unique and there is not guarantee that it yields a meaningful interpretation

# Relevance propagation: global conservation



$$f(x) = \boxed{R_7^{(3)}} = \boxed{R_4^{(2)} + R_5^{(2)} + R_6^{(2)}} = \boxed{R_1^{(1)} + R_2^{(1)} + R_3^{(1)}}$$

# Relevance propagation: local conservation



$$R_7^{(3)} = R_{4 \leftarrow 7}^{(2,3)} + R_{5 \leftarrow 7}^{(2,3)} + R_{6 \leftarrow 7}^{(2,3)}$$

$$R_k^{(l+1)} = \sum_{i: i \text{ is input for neuron } k} R_{i \leftarrow k}^{(l, l+1)}$$

$$R_1^{(1)} = R_{1 \leftarrow 4}^{(1,2)} + R_{1 \leftarrow 5}^{(1,2)}$$

$$R_i^{(l)} = \sum_{k: i \text{ is input for neuron } k} R_{i \leftarrow k}^{(l, l+1)}$$

Connection between global and local relevance:

$$\sum_k R_k^{(l+1)} = \sum_k \sum_{i: i \text{ is input for neuron } k} R_{i \leftarrow k}^{(l, l+1)} = \sum_i \sum_{k: i \text{ is input for neuron } k} R_{i \leftarrow k}^{(l, l+1)} = \sum_i R_i^{(l)}$$

# Relevance propagation. main properties

Global conservation:

$$f(x) = \dots = \sum_{d \in l+1} R_d^{(l+1)} = \sum_{d \in l} R_d^{(l)} = \dots = \sum_{d \in 1} R_d^{(1)}$$

Local conservation:

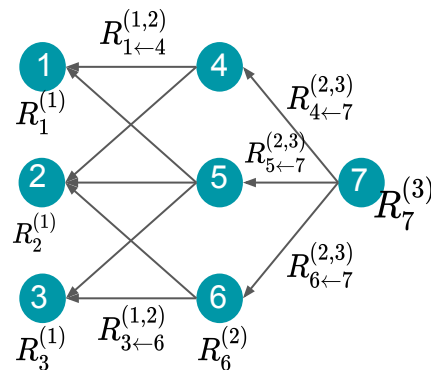
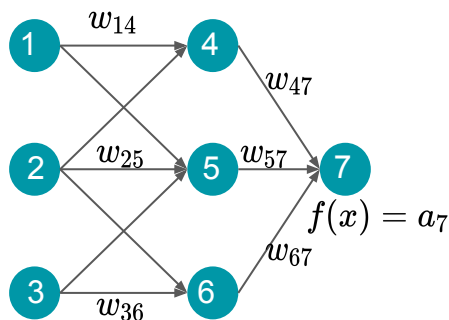
$$R_k^{(l+1)} = \sum_{i : i \text{ is input for neuron } k} R_{i \leftarrow k}^{(l, l+1)} \qquad R_i^{(l)} = \sum_{k : i \text{ is input for neuron } k} R_{i \leftarrow k}^{(l, l+1)}$$

# Example of relevance propagation

Let  $a_j$  be an activation of the  $j$ -th neuron

Then the relevance can be distributed as follows:  $R_j = \sum_k \frac{a_j w_{jk}}{\sum_j a_j w_{jk}} R_k$

$$R_{j \leftarrow k}^{(l, l+1)} = \frac{a_j w_{jk}}{\sum_j a_j w_{jk}}$$

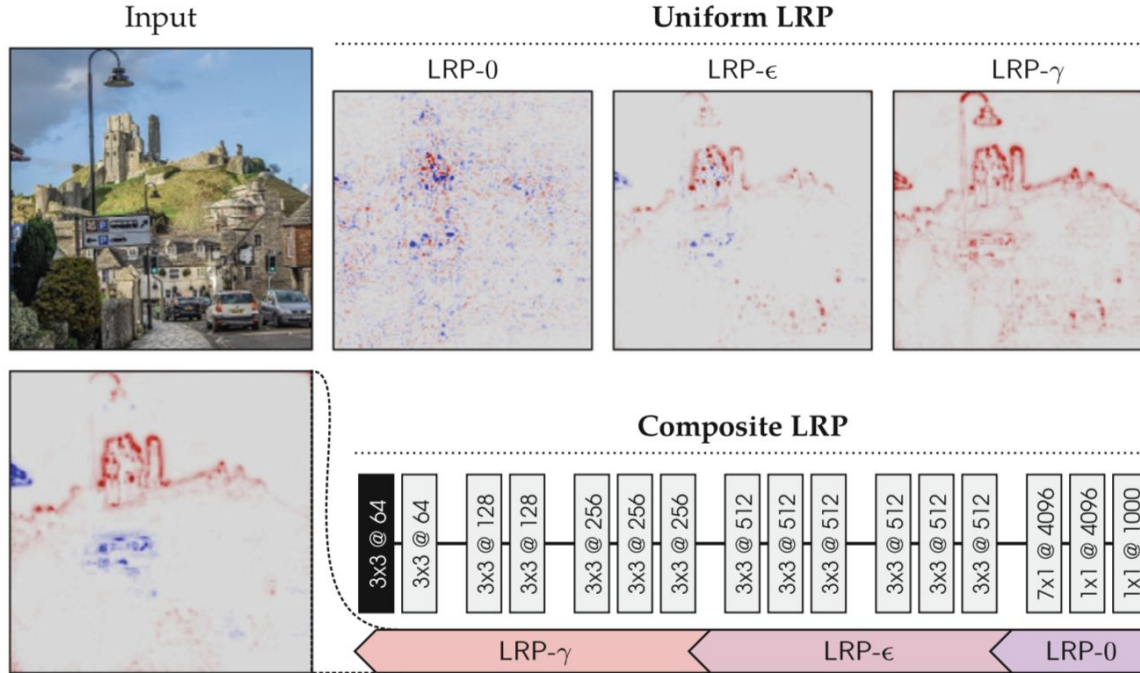


# LRP for Deep Rectifier Networks

DRN is composed of neurons  $a_j = \max(0, \sum_i a_i w_{ij} + b_i)$

	propagation property	approximation of the conservation property
LRP-0	$R_j = \sum_k \frac{a_j w_{jk}}{\sum_j a_j w_{jk} + b_j} R_k$	$\sum_j R_{j \leftarrow k}^{(l, l+1)} = R_k^{(l+1)} \left( 1 - \frac{b_k}{\sum_j a_j w_{jk} + b_k} \right)$
LRP- $\varepsilon$	$R_j = \sum_k \frac{a_j w_{jk}}{\sum_j \varepsilon + a_j w_{jk} + b_j} R_k$	$\sum_j R_{j \leftarrow k}^{(l, l+1)} = R_k^{(l+1)} \left( 1 - \frac{b_k + \varepsilon}{\sum_j a_j w_{jk} + b_k + \varepsilon} \right)$
LRP- $\gamma$	$R_j = \sum_k \frac{a_j (w_{jk} + \gamma w_{jk}^+)}{\sum_j \varepsilon + a_j (w_{jk} + \gamma w_{jk}^+) + b_j} R_k$	$\sum_j R_{j \leftarrow k}^{(l, l+1)} = R_k^{(l+1)} \left( 1 - \gamma \frac{b_k + b_k^+}{\sum_j a_j w_{jk} + b_k + \sum_j (a_j w_{jk})^+ + b_k^+} \right)$
LRP- $\alpha\beta$	$R_j = \sum_k \left( \alpha \frac{(a_j w_{jk})^+}{\sum_j (a_j w_{jk}^+) + b_j} + \beta \frac{(a_j w_{jk})^-}{\sum_j (a_j w_{jk}^-) + b_j} \right) R_k$ $\alpha + \beta = 1$	$\sum_j R_{j \leftarrow k}^{(l, l+1)} = R_k^{(l+1)} \left( 1 - \alpha \frac{(a_j w_{jk})^+}{\sum_j (a_j w_{jk})^+ + b_k^+} - \beta \frac{(a_j w_{jk})^+}{\sum_j (a_j w_{jk})^- + b_k^-} \right)$

# An example of a combination of rules



Input image and pixel-wise explanations of the output neuron 'castle' obtained with various LRP procedures.

Parameters are  $\epsilon = 0.25$  std and  $\gamma = 0.25$ .



# When and which rules should be used?

**LRP-0** picks many local artifacts of the function. Thus, the explanation is **overly complex** and **does not focus sufficiently** on the “actual explanation”, the explanation is neither faithful nor understandable

**LRP- $\epsilon$**  removes **noise elements** in the explanation to keep only a limited number features for the explanation. It provides a **faithful** explanation, but **too sparse to be easily understandable**

**LRP- $\gamma$**  is easier for a human to understand because **features are more densely highlighted**, but it also picks unrelated concepts for the explanation, thus it is rather unfaithful

Composite LRP overcomes the disadvantages of the approaches above

# Deep Taylor Decomposition. Motivation

## How to justify LRP rules theoretically?

In LRP we express relevance of a neuron  $k$  using the relevance of the neurons from the upper layer

## How to interpret negative values of relevance?

Classifying a ball, a dark ball on a bright background would have negative gradient, while white ball on darker background would have a positive gradient\*.

\* Smilkov, Daniel, et al. "SmoothGrad: removing noise by adding noise." arXiv preprint arXiv:1706.03825 (2017).

## Part 2. Deep Taylor Decomposition

# Positive relevance propagation

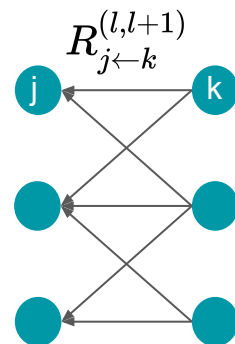
Let  $a_j$  be a **non-negative** activation of the  $j$ -th neuron

We do not use negative values since it is difficult to interpret negative evidence for a class

Then the relevance can be distributed as follows:

$$R_j = \sum_k \frac{a_j w_{jk}}{\sum_j a_j w_{jk}} R_k \quad R_j = \sum_k \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} R_k$$

$$R_{i \leftarrow j}^{(l,l+1)} = \frac{a_i w_{ij}}{\sum_i a_i w_{ij}} \quad R_{j \leftarrow k}^{(l,l+1)} = \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} R_k$$



# Positive relevance propagation

$$R_j = \sum_k \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} \boxed{R_k}$$

assume  $R_k = a_k c_k$ , where  $c_k$  is a positive constant

$$R_j = a_j \boxed{\sum_k w_{jk}^+ \frac{\max(0, \sum_j a_j w_{jk})}{\sum_j a_j w_{jk}^+} c_k}$$

$R_j = a_j c_j$  where  $c_j$  is positive and approximately constant

# Deep Taylor Decomposition

In DTD we suppose that the relevance of the current neuron  $k$  is a function of the lower-level neuron activations  $\{a_i\}$ , i.e.,  $R_j(\{a_i\})$ , where  $\{\}$  denotes a vector.

$$R_j = a_j c_j$$

$$\begin{aligned} R_j(\{a_i\}) &= a_j(\{a_i\}) \cdot c_j = \max(0, \sum_i a_i w_{ij} + b_j) \cdot c_j \\ &= \max(0, \sum_i a_i w_{ij} c_j + b_j c_j) \\ &= \max(0, \sum_i a_i w'_{ij} + b'_j) \end{aligned}$$

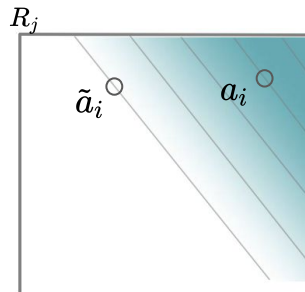
# Deep Taylor Decomposition

$$R_j(\{a_i\}) = \max(0, \sum_i a_i w'_{ij} + b'_j)$$

First-order Taylor expansion  $R_j(\{a_i\}) = R_j(\{\tilde{a}_i^{(j)}\}) + \sum_i \left. \frac{\partial R_j}{\partial a_i} \right|_{\tilde{a}_i^{(j)}} \cdot (a_i - \tilde{a}_i^{(j)}) + \varepsilon$

Conservation property:

$$\begin{aligned} \sum_j R_j &= \left( \left. \frac{\partial (\sum_j R_j)}{\partial \{x_i\}} \right|_{\{\tilde{x}_i\}} \right)^\top \cdot (\{x_i\} - \{\tilde{x}_i\}) + \varepsilon \\ &= \sum_i \underbrace{\sum_j \left. \frac{\partial R_j}{\partial x_i} \right|_{\{\tilde{x}_i\}}}_{R_i} \cdot (x_i - \tilde{x}_i) + \varepsilon, \end{aligned}$$



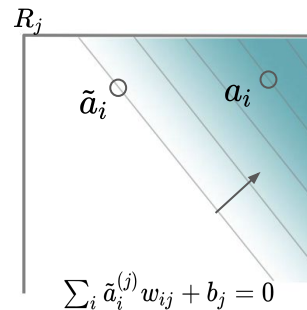
**Remark:** due to the potentially complex relation between  $a_j$  and  $R_k$ , finding an appropriate reference point and computing the gradient locally is difficult

# Deep Taylor Decomposition

First-order Taylor expansion  $R_j(\{a_i\}) = R_j(\{\tilde{a}_i^{(j)}\}) + \sum_i \frac{\partial R_j}{\partial a_i} \Big|_{\tilde{a}_i^{(j)}} \cdot (a_i - \tilde{a}_i^{(j)}) + \varepsilon$

We search for a root point such that  $R_j(\{a_i\}) = 0 + \sum_i \frac{\partial R_j}{\partial a_i} \Big|_{\tilde{a}_i^{(j)}} \cdot (a_i - \tilde{a}_i^{(j)}) + 0$

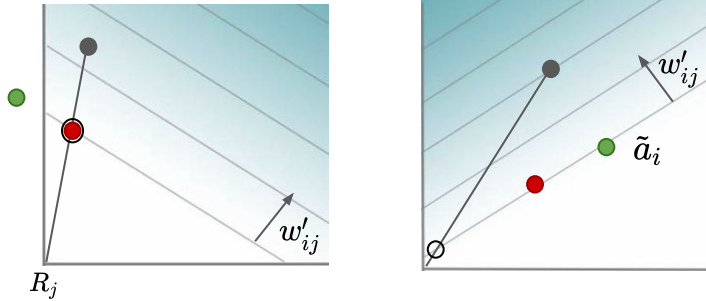
$R_{i \leftarrow j}$



Closed form solution:  $R_{i \leftarrow j} = \frac{v_{ij} w_{ij}}{\sum_i v_{ij} w_{ij}} R_j$  where  $v_{ij} \propto a_i - \tilde{a}_i^{(j)}$  is the root search direction



# Choosing search directions for a root point



root point	search direction
nearest root	$v_{ij} = w_{ij}$
origin	$v_{ij} = a_i$
LRP equivalent	$v_{ij} = a_i \mathbf{1}_{w'_{ij} > 0}$

**Remark:** Only the last one guarantees that the root point 1) belongs to the input domain, 2) has non-negative relevance scores

# Commonly used LRP rules

Name	Formula	Usage	DTD
LRP-0 [7]	$R_j = \sum_k \frac{a_j w_{jk}}{\sum_{0,j} a_j w_{jk}} R_k$	Upper layers	✓
LRP- $\epsilon$ [7]	$R_j = \sum_k \frac{a_j w_{jk}}{\epsilon + \sum_{0,j} a_j w_{jk}} R_k$	Middle layers	✓
LRP- $\gamma$	$R_j = \sum_k \frac{a_j (w_{jk} + \gamma w_{jk}^+)}{\sum_{0,j} a_j (w_{jk} + \gamma w_{jk}^+)} R_k$	Lower layers	✓
LRP- $\alpha\beta$ [7]	$R_j = \sum_k \left( \alpha \frac{(a_j w_{jk})^+}{\sum_{0,j} (a_j w_{jk})^+} - \beta \frac{(a_j w_{jk})^-}{\sum_{0,j} (a_j w_{jk})^-} \right) R_k$	Lower layers	$\times^a$
flat [30]	$R_j = \sum_k \frac{1}{\sum_j 1} R_k$	Lower layers	$\times$
$w^2$ -rule [36]	$R_i = \sum_j \frac{w_{ij}^2}{\sum_i w_{ij}^2} R_j$	First layer ( $\mathbb{R}^d$ )	✓
$z^{\mathcal{B}}$ -rule [36]	$R_i = \sum_j \frac{x_i w_{ij} - l_i w_{ij}^+ - h_i w_{ij}^-}{\sum_i x_i w_{ij} - l_i w_{ij}^+ - h_i w_{ij}^-} R_j$	First layer (pixels)	✓

(<sup>a</sup>DTD interpretation only for the case  $\alpha = 1, \beta = 0$ .)

# Further reading: Layer-Wise Relevance Propagation

Site with plenty of sources: <http://heatmapping.org/>

Tutorial on implementation: <https://git.tu-berlin.de/gmontavon/lrp-tutorial>

Example of implementation:

<https://github.com/atulshanbhag/Layerwise-Relevance-Propagation/blob/master/vgg/lrp.py>

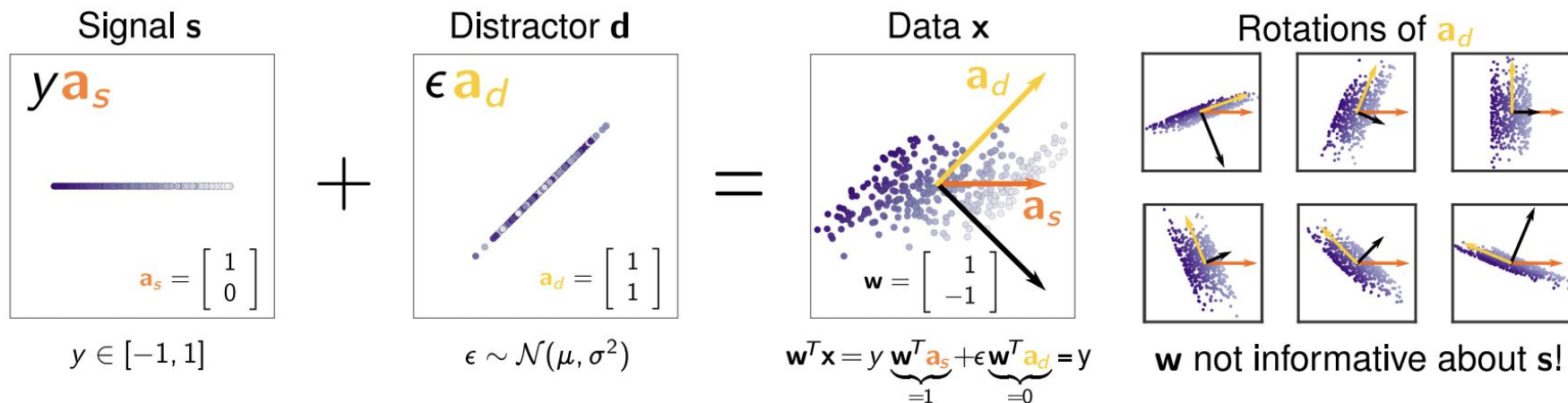
## Part 3. PatternNet and PatternAttribution

# PatternNet and PatternAttribution

## Reasoning line

- the linear model is the simplest neural network
- the explanation methods should work correctly in the limit of simplicity
- data contain the signal and some distortion
- the explanation method should work well at least for the simplest cases

# Understanding a linear model. Example



The weight vector is not aligned with the signal, its primary is to **cancel the distractor**

The weight vector is able to filter out the distractor, that is why it is called a filter

# Understanding a linear model

For a given data

$$\mathbf{x} = \underbrace{y\mathbf{a}_s}_{\text{signal}} + \underbrace{\varepsilon\mathbf{a}_d}_{\text{distractor}} = \mathbf{s} + \mathbf{d}$$

$$\mathbf{w}^T \mathbf{x} = y \underbrace{\mathbf{w}^T \mathbf{a}_s}_{=1} + \varepsilon \underbrace{\mathbf{w}^T \mathbf{a}_d}_{=0} = y$$

The filter  $\mathbf{w}$  tells us how to extract  $y$  optimally from data  $\mathbf{x}$ . The pattern  $\mathbf{a}_s$  is the direction in the data along which the desired output  $y$  varies

# Going back to the studied models

- **Functions** (gradients, saliency maps):  $\partial y / \partial \mathbf{x} = \mathbf{w}$
- **Signals** (DeconvNet, Guided Backpropagation): propagating back an “assumed signal”  $s = a_s y$  toward the signal direction of each neuron
- **Attribution** (LRP, DeepTaylor): assessing how much the signal dimensions contribute to the output through the layers

$$\mathbf{r}_i^{l-1} = \frac{\mathbf{w} \odot (\mathbf{x} - \mathbf{x}_0)}{\mathbf{w}^T \mathbf{x}} \mathbf{r}_i^l \quad \mathbf{r}_i^{output} = y$$

**Remark:** selecting a root point for the DeepTaylorDecomposition corresponds to estimating the distractor  $\mathbf{d} = \mathbf{x}_0$  and, by that, the signal  $s^* = \mathbf{x} - \mathbf{x}_0$



# Training a model

**Goal:** to train the filter  $w$  to extract  $y$  such that  $w^T x = y$ ,  $w^T s = y$ ,  $w^T d = 0$

It is an **ill-posed problem**. We could limit ourselves to the linear estimators if the **following form**:  $\hat{s} = u(w^T u)^{-1}y$  with a random vector  $u$  such that  $w^T u \neq 0$

For such an estimator  $\hat{s}$  it is always true that  $w^T \hat{s} = y$

There exist **infinitely many solutions** (as many as the propagation rules for the previously considered methods)

# Quality measure for signal estimates

Let the signal estimator be  $S(\mathbf{x}) = \hat{\mathbf{s}}$ , the distractor be  $\hat{\mathbf{d}} = \mathbf{x} - S(\mathbf{x})$  and  $\mathbf{w}^T \mathbf{x} = y$ :

$$\rho(S) = 1 - \max_{\mathbf{v}} \text{corr}(\mathbf{w}^T \mathbf{x}, \mathbf{v}^T (\mathbf{x} - S(\mathbf{x})))$$

This criterion introduces an additional constraint by measuring how much information about  $y$  can be reconstructed from the residuals  $\mathbf{x} - S(\mathbf{x})$  using a linear projection

The **best signal estimators remove most of the information in the residuals** and thus yield large  $\rho$

Thus, we want that the distractor correlate a lot with the weight vector  $\mathbf{w}$  (filter) (since the primary objective of the filter is to cancel the distortion)

# Quality measure for signal estimates

Let the signal estimator be  $S(\mathbf{x}) = \hat{\mathbf{s}}$  is optimal w.r.t. the following equation

$$\rho(S) = 1 - \max_{\mathbf{v}} \text{corr}(\mathbf{w}^T \mathbf{x}, \mathbf{v}^T (\mathbf{x} - S(\mathbf{x}))) = 1 - \max_{\mathbf{v}} \frac{\mathbf{v}^T \text{cov}[\hat{\mathbf{d}}, y]}{\sqrt{\sigma_{\mathbf{v}^T \hat{\mathbf{d}}}^2 \sigma_y^2}}$$

if the correlation is 0 for all possible  $\mathbf{v}$ , i.e.,  $\forall \mathbf{v}, \text{cov}[y, \hat{\mathbf{d}}] \mathbf{v} = \mathbf{0}$

From the linearity of the covariance and  $\hat{\mathbf{d}} = \mathbf{x} - S(\mathbf{x})$  it follows

$$\text{cov}[y, \hat{\mathbf{d}}] = \mathbf{0} \Rightarrow \text{cov}[\mathbf{x}, y] = \text{cov}[S(\mathbf{x}), y]$$

# Signal estimates for the linear neurons

A linear neuron can extract only a linear signal from  $\mathbf{x}$ , thus we may assume the following dependency:

$$S_{\mathbf{a}}(\mathbf{x}) = \mathbf{a}\mathbf{w}^T \mathbf{x}$$

Plugging this into  $\text{cov}[\mathbf{x}, y] = \text{cov}[S(\mathbf{x}), y]$  we get

$$\text{cov}[\mathbf{x}, y] = \text{cov}[\mathbf{a}\mathbf{w}^T \mathbf{x}, y] = \mathbf{a}\text{cov}[y, y] \Rightarrow \mathbf{a} = \frac{\text{cov}[\mathbf{x}, y]}{\sigma_y^2}$$

Since the correlation is invariant to scaling, we constrain  $\mathbf{v}^T \hat{\mathbf{d}}$  to have  $\sigma_{\mathbf{v}^T \hat{\mathbf{d}}}^2 = \sigma_y^2$

# Signal estimates for the ReLU neurons

ReLU does not propagate the negative activations, thus we distinguish the following regimes:

$$\mathbf{x} = \begin{cases} \mathbf{s}_+ + \mathbf{d}_+ & \text{if } y > 0 \\ \mathbf{s}_- + \mathbf{d}_- & \text{otherwise} \end{cases}$$

The signal can be estimated as follows:

$$S_{\mathbf{a}_{+-}}(\mathbf{x}) = \begin{cases} \mathbf{a}_+ \mathbf{w}^T \mathbf{x}, & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ \mathbf{a}_- \mathbf{w}^T \mathbf{x}, & \text{otherwise} \end{cases}$$

# Signal estimates for the ReLU neurons

Let  $\mathbb{E}[x]_+$  and  $\mathbb{E}[x]_-$  be expectations over  $x$  within positive and negative regimes, respectively, and  $\pi_+$  be the expected ratio of input  $x$  with  $\mathbf{w}^T \mathbf{x}$  then

$$\begin{aligned}\text{cov}[\mathbf{x}, y] &= \pi_+ (\mathbb{E}_+ [\mathbf{x}y] - \mathbb{E}_+ [\mathbf{x}] \mathbb{E} [y]) + (1 - \pi_+) (\mathbb{E}_- [\mathbf{x}y] - \mathbb{E}_- [\mathbf{x}] \mathbb{E} [y]) \\ \text{cov}[\mathbf{s}, y] &= \pi_+ (\mathbb{E}_+ [\mathbf{s}y] - \mathbb{E}_+ [\mathbf{s}] \mathbb{E} [y]) + (1 - \pi_+) (\mathbb{E}_- [\mathbf{s}y] - \mathbb{E}_- [\mathbf{s}] \mathbb{E} [y])\end{aligned}$$

Plugging it into  $\text{cov}[\mathbf{x}, y] = \text{cov}[S(\mathbf{x}), y]$

$$\mathbb{E}_+ [\mathbf{x}y] - \mathbb{E}_+ [\mathbf{x}] \mathbb{E} [y] = \mathbb{E}_+ [\mathbf{s}y] - \mathbb{E}_+ [\mathbf{s}] \mathbb{E} [y]$$

$$\text{Since } S_{a_+}(x) = \begin{cases} a_+ \mathbf{w}^T \mathbf{x}, & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ a_- \mathbf{w}^T \mathbf{x}, & \text{otherwise} \end{cases} \quad a_+ = \frac{\mathbb{E}_+ [\mathbf{x}y] - \mathbb{E}_+ [\mathbf{x}] \mathbb{E} [y]}{\mathbf{w}^T \mathbb{E}_+ [\mathbf{x}y] - \mathbf{w}^T \mathbb{E}_+ [\mathbf{x}] \mathbb{E} [y]}$$

# PatternNet

It is a **layer-wise back-projection** of the estimated signal to input space

The signal estimator is approximated as a superposition of neuron-wise, non-linear signal estimator  $S_{a^{+-}}$  at each layer

It is equal to the **gradient** where during the backward pass the **weights** of the network are **replaced** by **the informative directions**

Initialization:  $s_i^{output} = y$

Linear or convolutional layers:  $s^{l-1,i} = \mathbf{a}_+ s_i^l$

Conservation property:  $s_i^{l-1} = \sum_j s_i^{l-1,j}$

ReLU:  $s_i^{l-1} = \begin{cases} s_i^l, & x_i^l > 0 \\ 0, & \text{otherwise} \end{cases}$

# PatternAttribution

It exposes the attribution  $\mathbf{w} \odot \mathbf{a}_+$

It can be considered as a root point for DeepTaylorDecomposition

Let  $\mathbf{r}^{l-1,i}$  be the relevance of the  $l$ -th layer. The distribution has the following form

$$\mathbf{r}^{l-1,i} = \frac{\mathbf{w} \odot (\mathbf{x} - \mathbf{x}_0)}{\mathbf{w}^T \mathbf{x}} r_i^l$$

For the original LRP  $\mathbf{x}_0 = 0$

For PatternAttribution  $\mathbf{x}_0 = \mathbf{x} - \mathbf{a}_+ \mathbf{w}^T \mathbf{x}$ , that gives  $\mathbf{r}^{l-1,i} = \mathbf{w} \odot \mathbf{a}_+ \mathbf{r}_{,i}^l$



# Method comparison

