

# Explaining Machine Learning Decisions

Grégoire Montavon, TU Berlin

Joint work with: Wojciech Samek, Klaus-Robert Müller,  
Sebastian Lapuschkin, Alexander Binder

18/09/2018 Intl. Workshop ML & AI, Telecom ParisTech

# From ML Successes to Applications

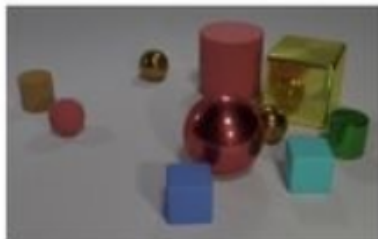
Deep Net outperforms humans in image classification

IM  GENET

AlphaGo beats Go human champ

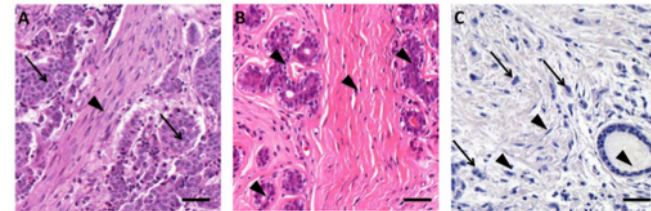


Visual Reasoning



What size is the cylinder that is left of the brown metal thing that is left of the big sphere?

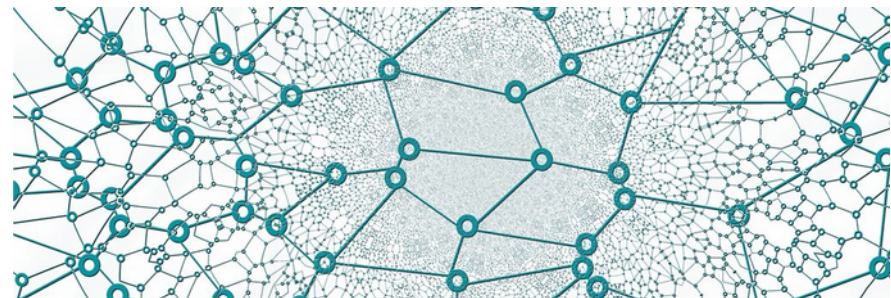
Medical Diagnosis

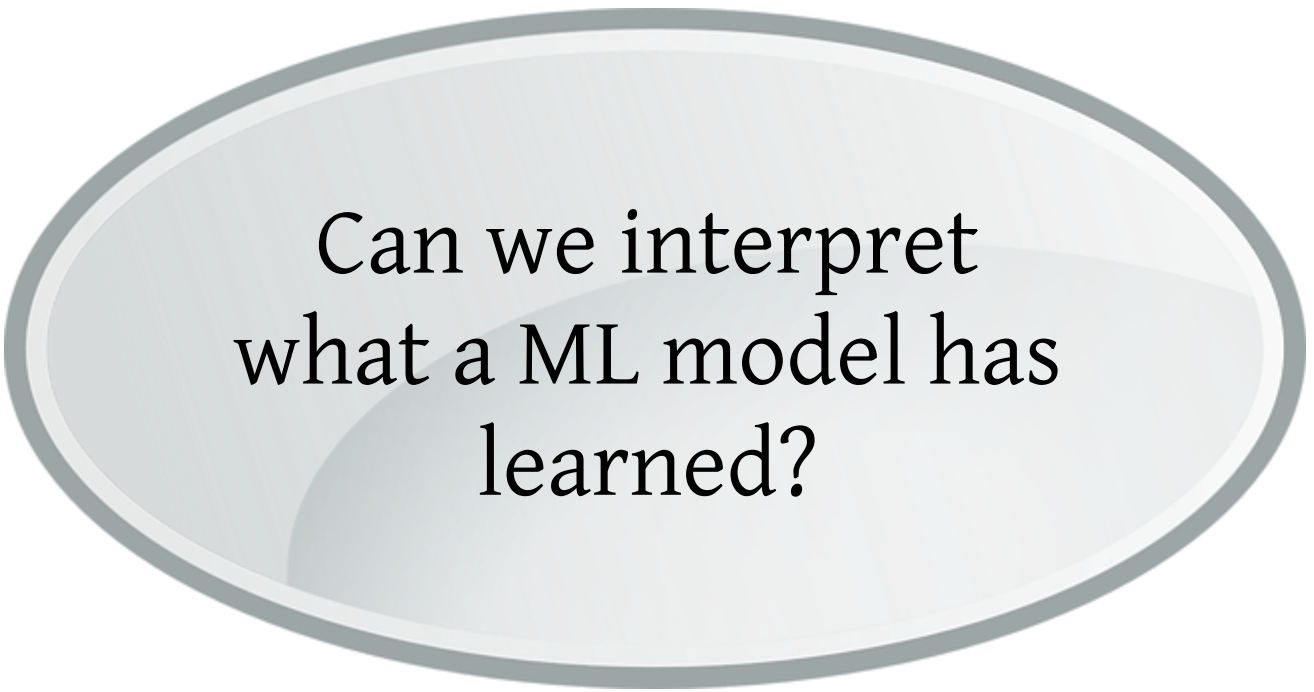


Autonomous Driving

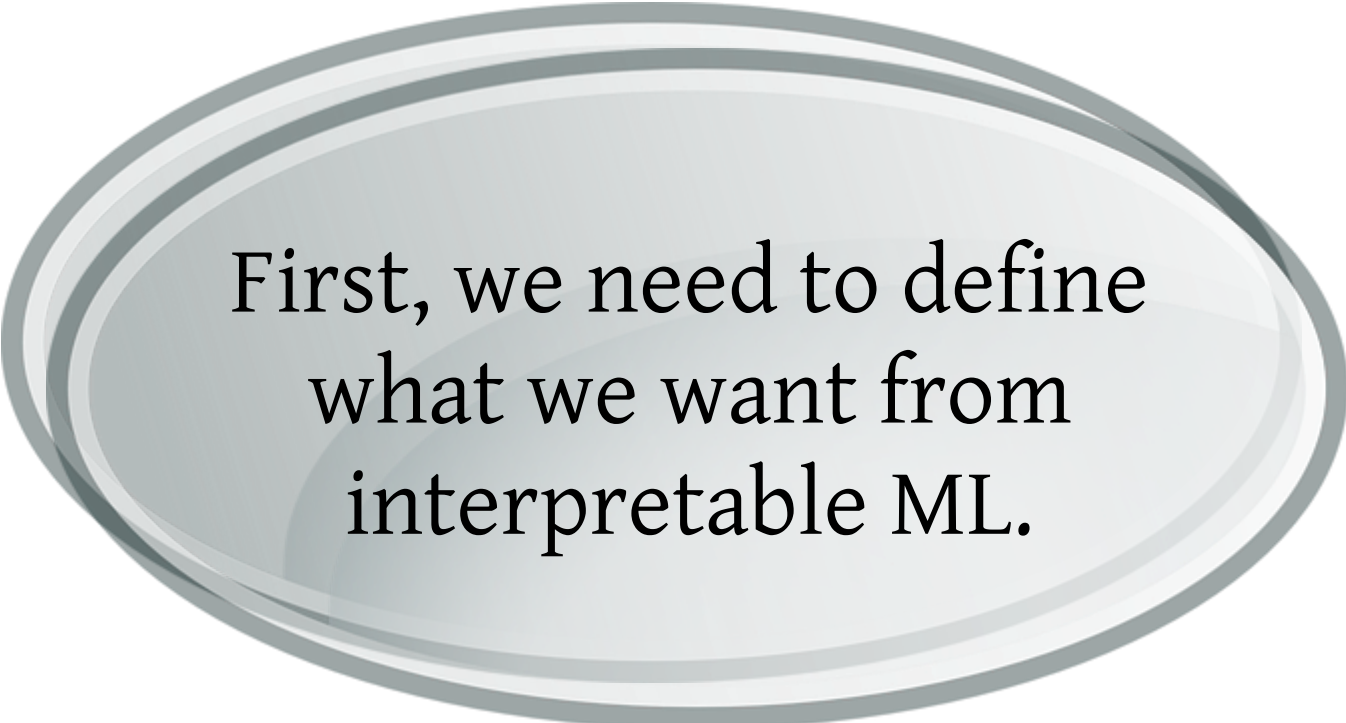


Networks (smart grids, etc.)





Can we interpret  
what a ML model has  
learned?

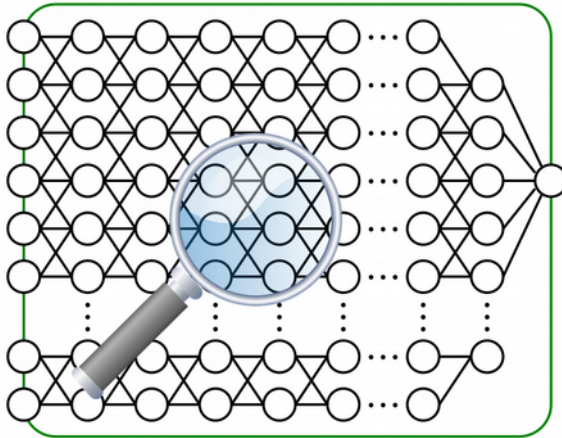


First, we need to define  
what we want from  
interpretable ML.



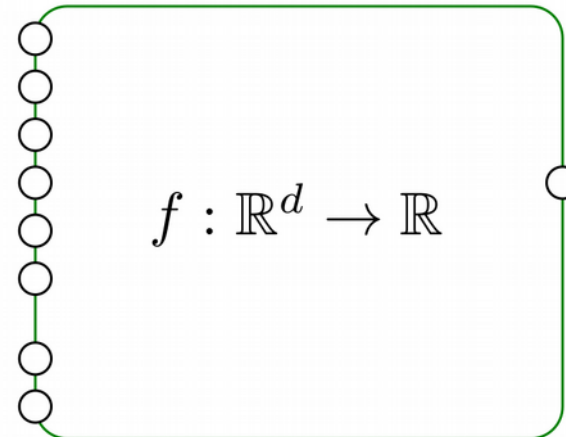
# Understanding Deep Nets: Two Views

## mechanistic understanding



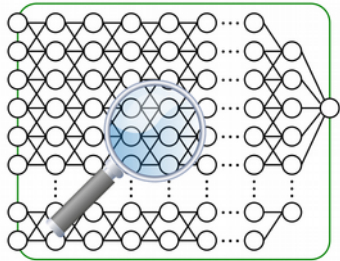
Understanding what mechanism the network uses to solve a problem or implement a function.

## functional understanding

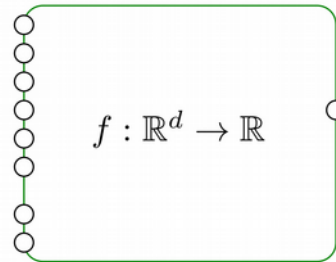


Understanding how the network relates the input to the output variables.

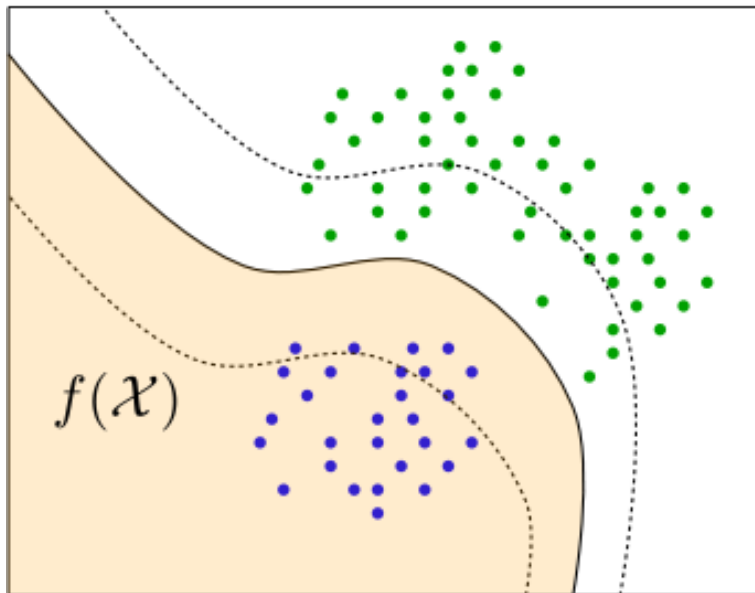
mechanistic  
understanding



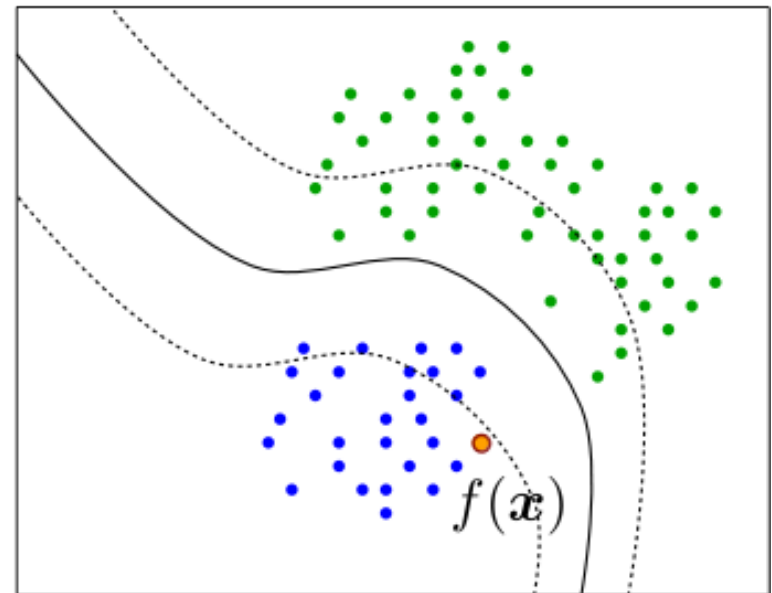
functional  
understanding



interpreting  
predicted classes

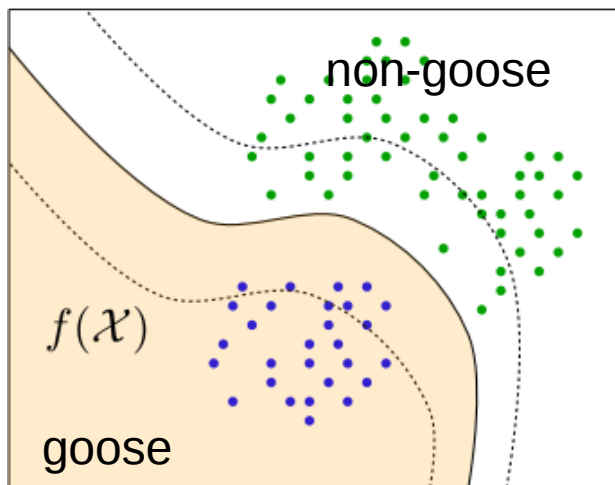


explaining  
individual decisions



# Interpreting Predicted Classes

Example: “*How does a goose typically look like according to the neural network?*”



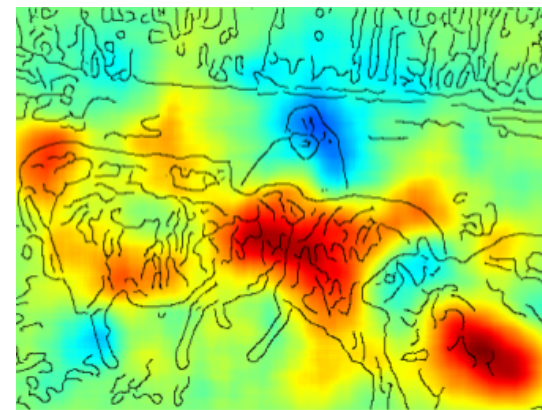
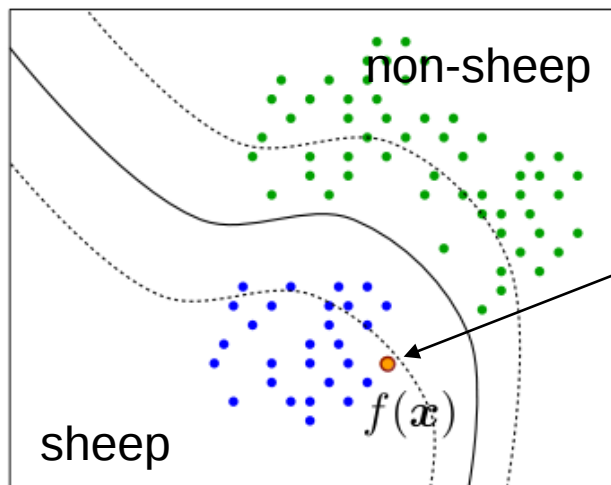
$$\arg \max_x f(x) + \text{reg.}$$



Image from **Symonian'13**

# Explaining Individual Decisions

Example: “*Why is a given image classified as a sheep?*”

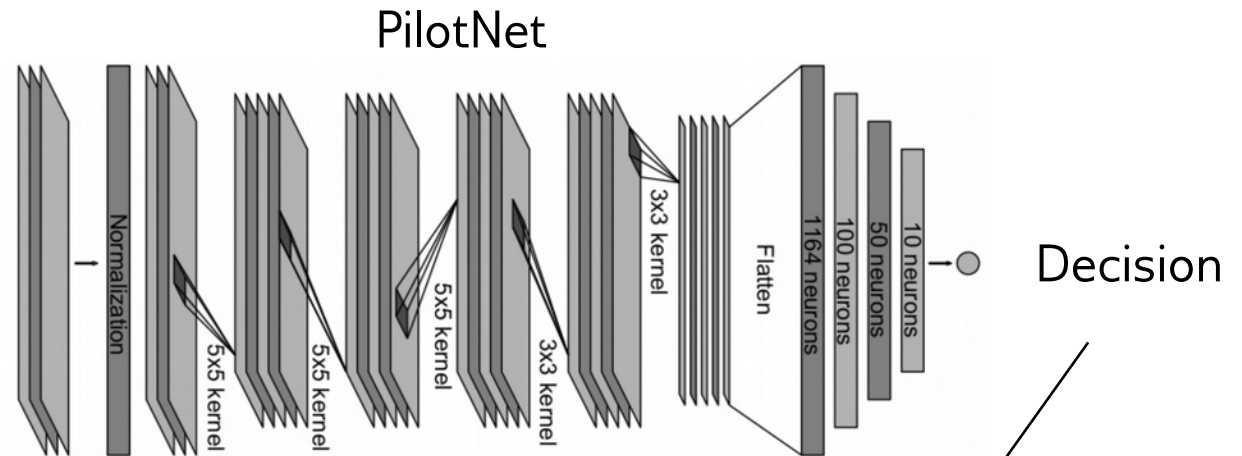


$$R(x, f)$$

# Example: Autonomous Driving [Bojarski'17]

Bojarski et al. 2017 “Explaining How a Deep Neural Network Trained with End-to-End Learning Steers a Car”

Input:



Explanation:





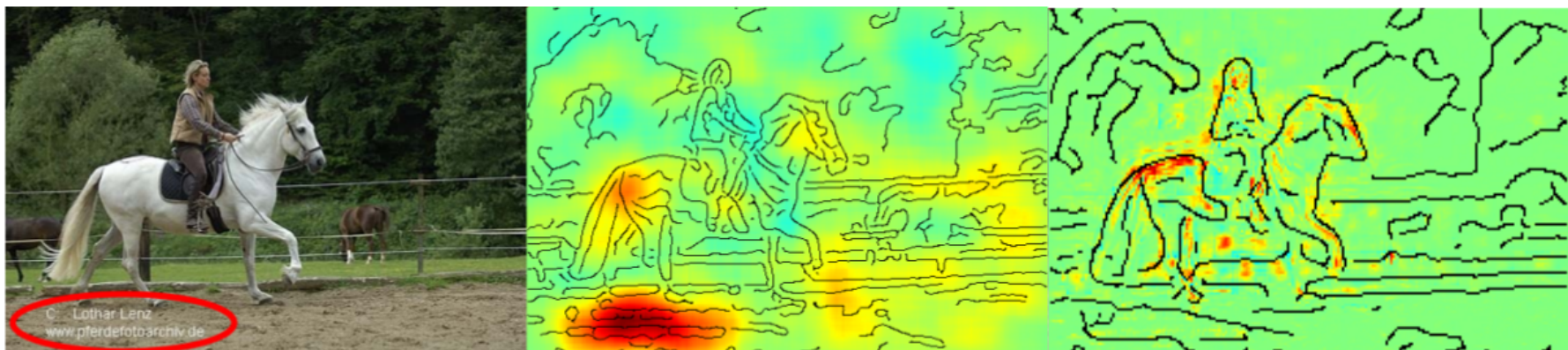
# Example: Pascal VOC Classification [Lapuschkin'16]

Comparing Performance on Pascal VOC 2007  
(Fisher Vector Classifier vs. DeepNet pretrained on ImageNet)

	<b>aeroplane</b>	<b>bicycle</b>	<b>bird</b>	<b>boat</b>	<b>bottle</b>	<b>bus</b>	<b>car</b>
<b>Fisher</b>	79.08%	66.44%	45.90%	70.88%	27.64%	69.67%	80.96%
<b>DeepNet</b>	88.08%	79.69%	80.77%	77.20%	35.48%	72.71%	86.30%
	<b>cat</b>	<b>chair</b>	<b>cow</b>	<b>diningtable</b>	<b>dog</b>	<b>horse</b>	<b>motorbike</b>
<b>Fisher</b>	59.92%	51.92%	47.60%	58.06%	42.28%	80.45%	69.34%
<b>DeepNet</b>	81.10%	51.04%	61.10%	64.62%	76.17%	81.60%	79.33%
	<b>person</b>	<b>pottedplant</b>	<b>sheep</b>	<b>sofa</b>	<b>train</b>	<b>tvmonitor</b>	<b>mAP</b>
<b>Fisher</b>	85.10%	28.62%	49.58%	49.31%	82.71%	54.33%	59.99%
<b>DeepNet</b>	92.43%	49.99%	74.04%	49.48%	87.07%	67.08%	72.12%

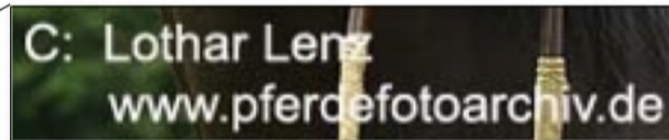
Fisher classifier

(pretrained) deep net



# Example: Pascal VOC Classification [Lapuschkin'16]

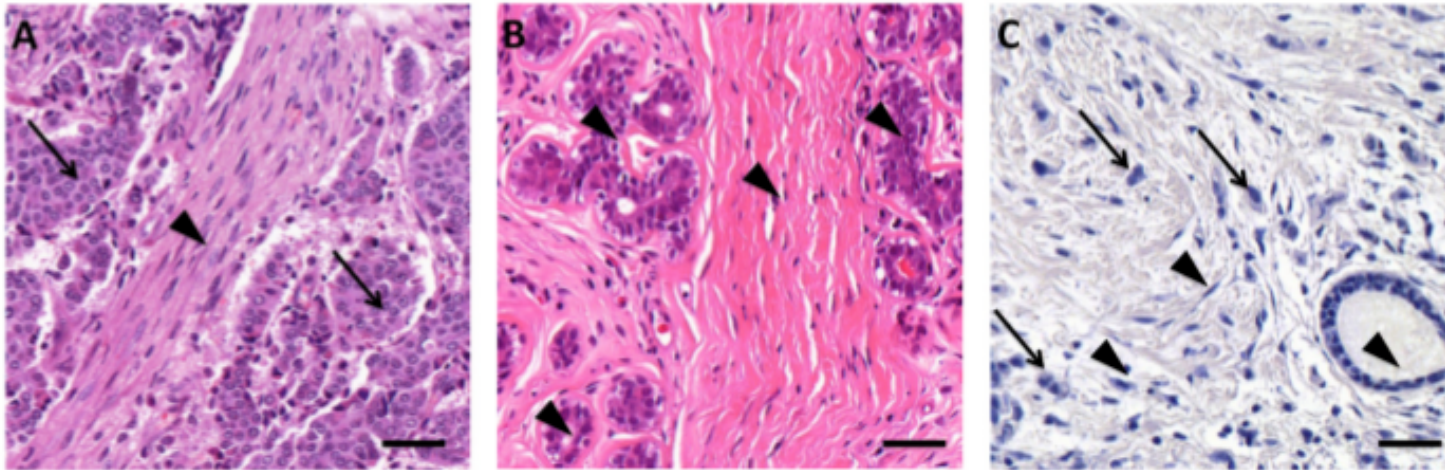
'horse' images in PASCAL VOC 2007



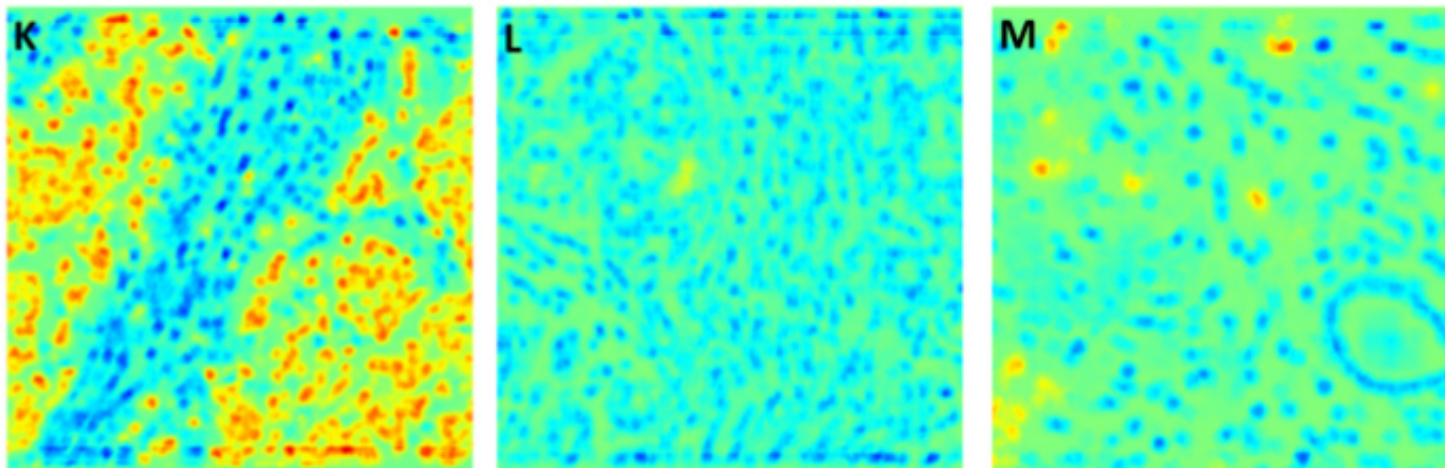


# Example: Medical Diagnosis [Binder'18]

Binder et al. 2018 “Towards computational fluorescence microscopy: Machine learning-based integrated prediction of morphological and molecular tumor profiles”



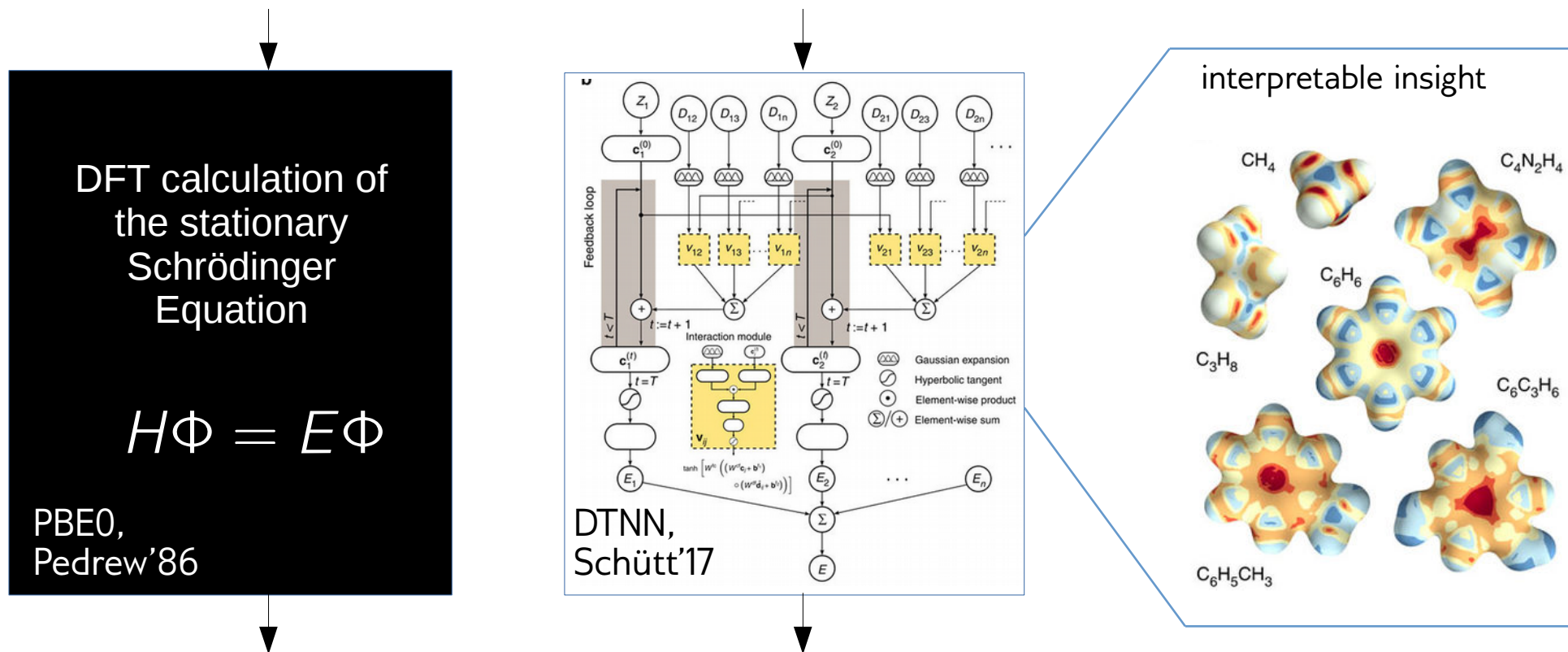
A: Invasive breast cancer, H&E stain; B: Normal mammary glands and fibrous tissue, H&E stain; C: Diffuse carcinoma infiltrate in fibrous tissue, Hematoxylin stain



# Example: Quantum Chemistry [Schütt'17]

Schütt et al. 2017: Quantum-Chemical Insights from Deep Tensor Neural Networks

molecular structure (e.g. atoms positions)



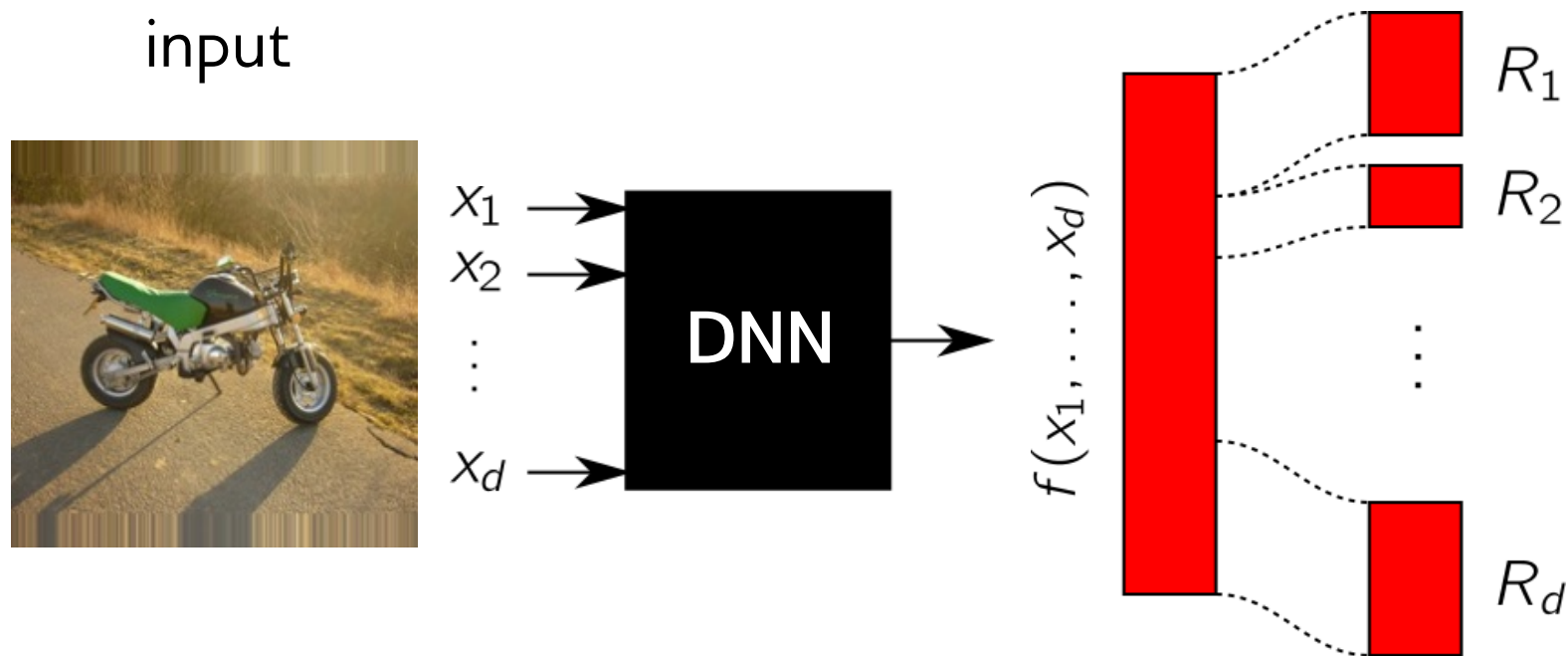
molecular electronic properties (e.g. atomization energy)

# Examples of Explanation Methods

Baehrens'10 Gradient	Sundarajan'17 Int Grad	Zintgraf'17 Pred Diff	Ribeiro'16 LIME	Haufe'15 Pattern
Zurada'94 Gradient	Symonians'13 Gradient	Zeiler'14 Occlusions	Fong'17 M Perturb	Kindermans'17 PatternNet
Poulin'06 Additive	Lundberg'17 Shapley	Bazen'13 Taylor	Montavon'17 Deep Taylor	Shrikumar'17 DeepLIFT
Zeiler'14 Deconv	Landecker'13 Contrib Prop	Bach'15 LRP	Zhang'16 Excitation BP	
Caruana'15 Fitted Additive	Springenberg'14 Guided BP	Zhou'16 GAP	Selvaraju'17 Grad-CAM	

# Explaining by Decomposing

Importance of a variable is the share of the function score that can be attributed to it.



Decomposition property:  $f(x_1, \dots, x_d) = \sum_{i=1}^d R_i$

# Explaining Linear Models

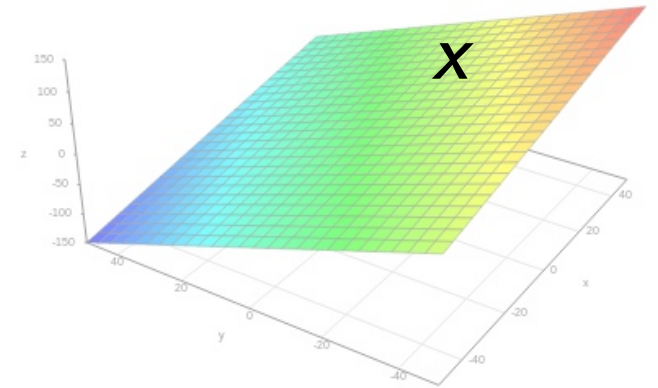
$$f(x_1, \dots, x_d) = \sum_{i=1}^d w_i x_i + b$$

A simple method:

$$\underbrace{\phantom{w_i x_i}}_{R_i}$$



$$f(x_1, \dots, x_d) \approx \sum_{i=1}^d R_i$$

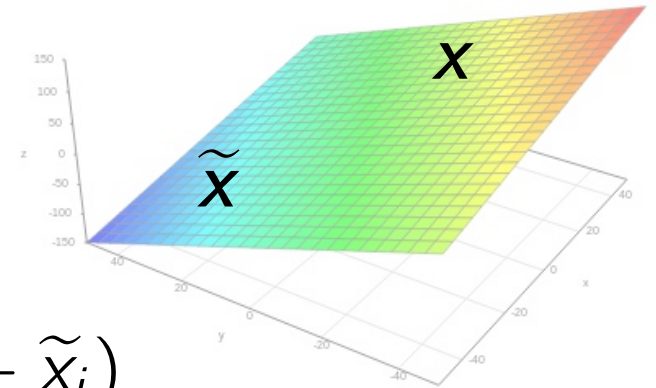


# Explaining Linear Models

$$f(x_1, \dots, x_d) = \sum_{i=1}^d w_i x_i + b$$

Taylor decomposition approach:

$$f(x_1, \dots, x_d) = \sum_{i=1}^d \underbrace{\frac{\partial f}{\partial x_i} \Big|_{\tilde{x}}}_{R_i} \cdot (x_i - \tilde{x}_i)$$

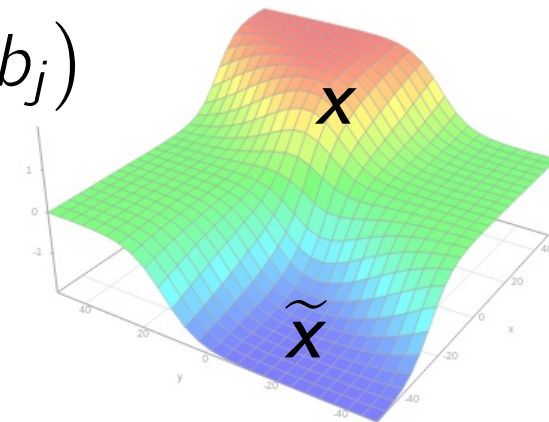


$$R_i = w_i \cdot (x_i - \tilde{x}_i)$$

**Insight:** explanation depends on the root point.

# Explaining Nonlinear Models

$$f(x_1, \dots, x_d) = \sum_j \rho\left(\sum_{i=1}^d w_{ij}x_i + b_j\right)$$

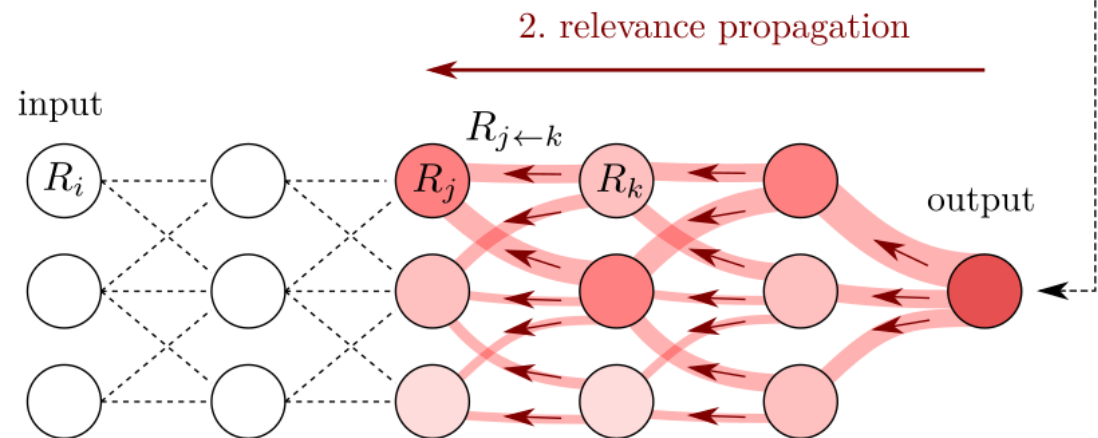
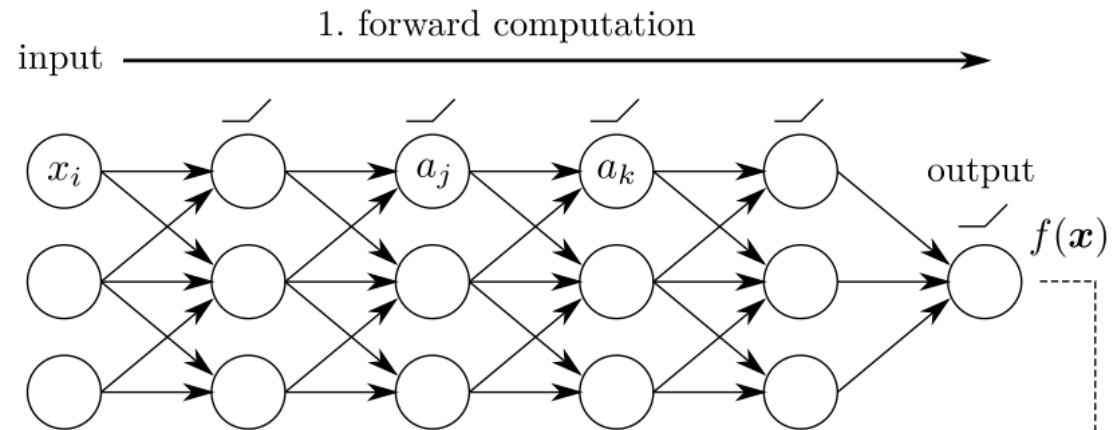
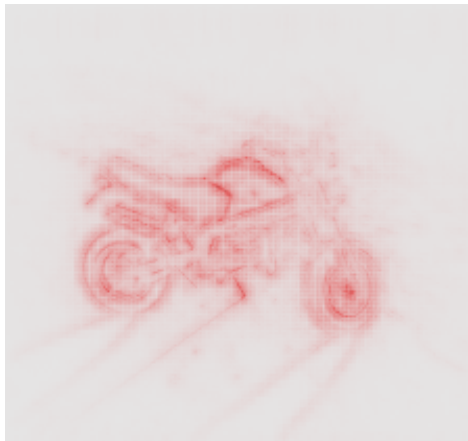


$$f(x_1, \dots, x_d) = \sum_{i=1}^d \frac{\partial f}{\partial x_i} \Big|_{\tilde{x}} \cdot (x_i - \tilde{x}_i) + o(\mathbf{x} \mathbf{x}^\top)$$

second-order terms are  
hard to interpret and  
can be very large



# Explaining Nonlinear Models by Propagation



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

Layer-Wise Relevance  
Propagation (LRP) [Bach'15]

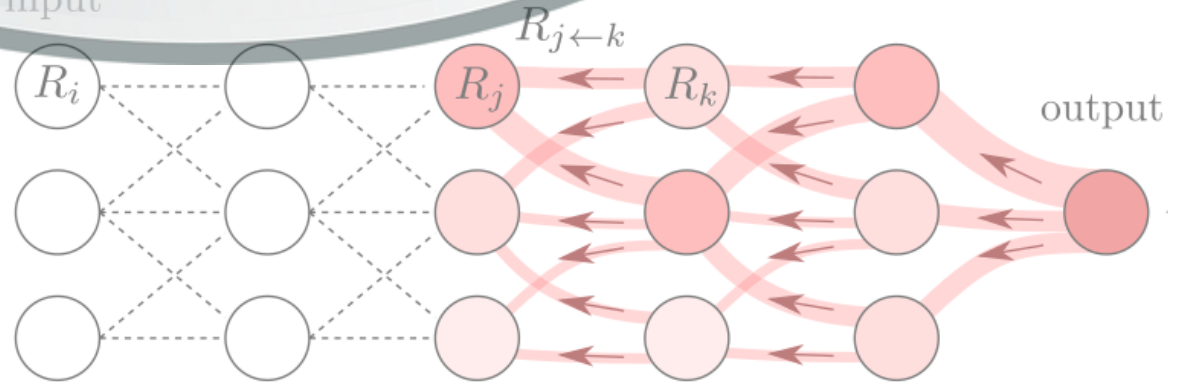
# Explaining Nonlinear Models by Propagation

Is there an  
underlying  
mathematical  
framework?

input

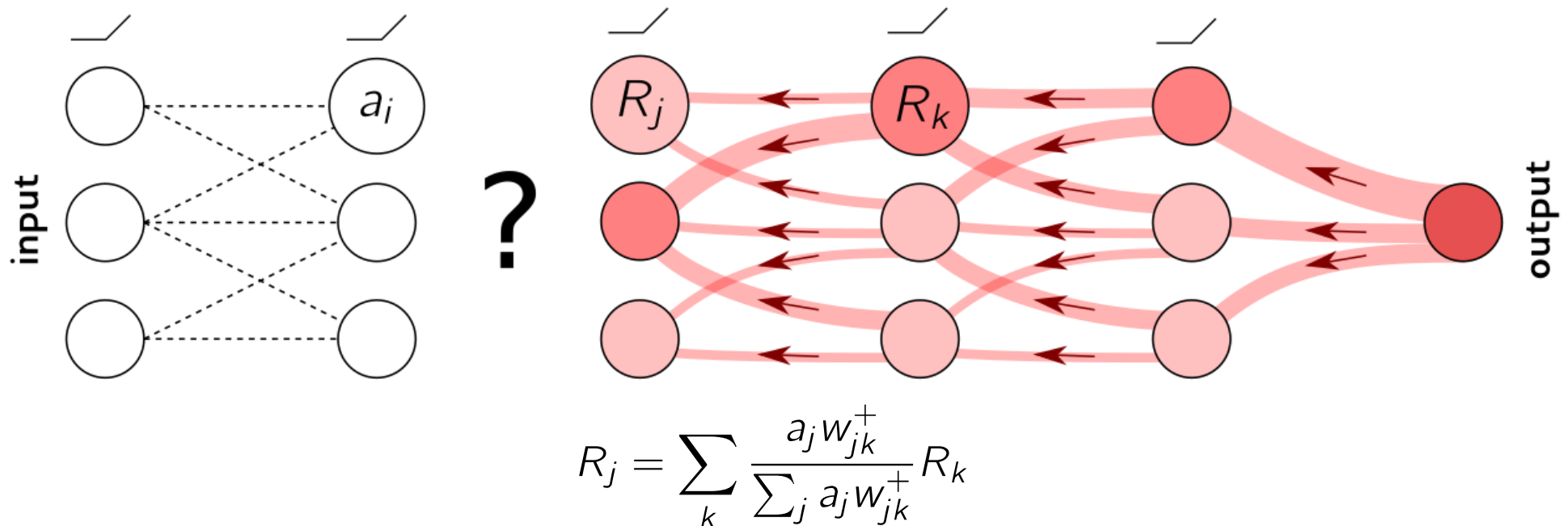
2. relevance propagation

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



# Deep Taylor Decomposition (DTD) [Montavon'17]

**Question:** Suppose that we have propagated LRP scores (“relevance”) until a given layer. How should it be propagated one layer further?

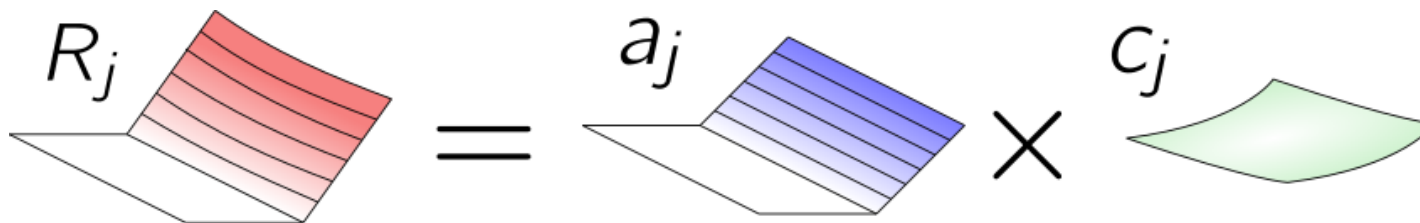


**Key idea:** Let's use Taylor expansions for this.

# DTD Step 1: The Structure of Relevance

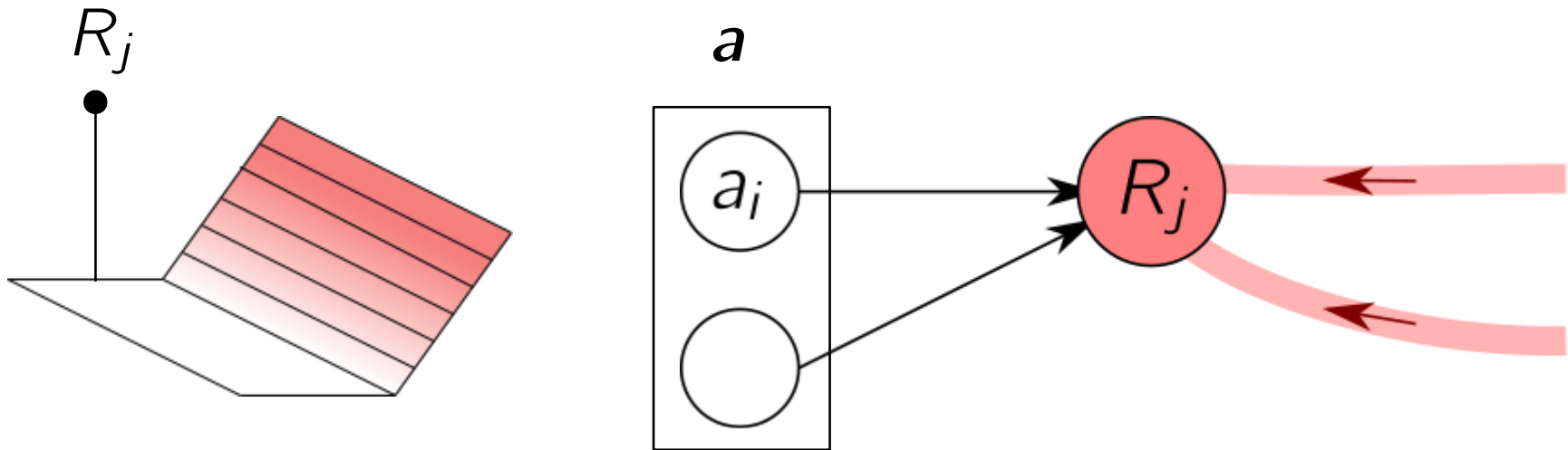
$$R_j = \sum_k \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} R_k$$

**Observation:** Relevance at each layer is a product of the activation and an approximately constant term.



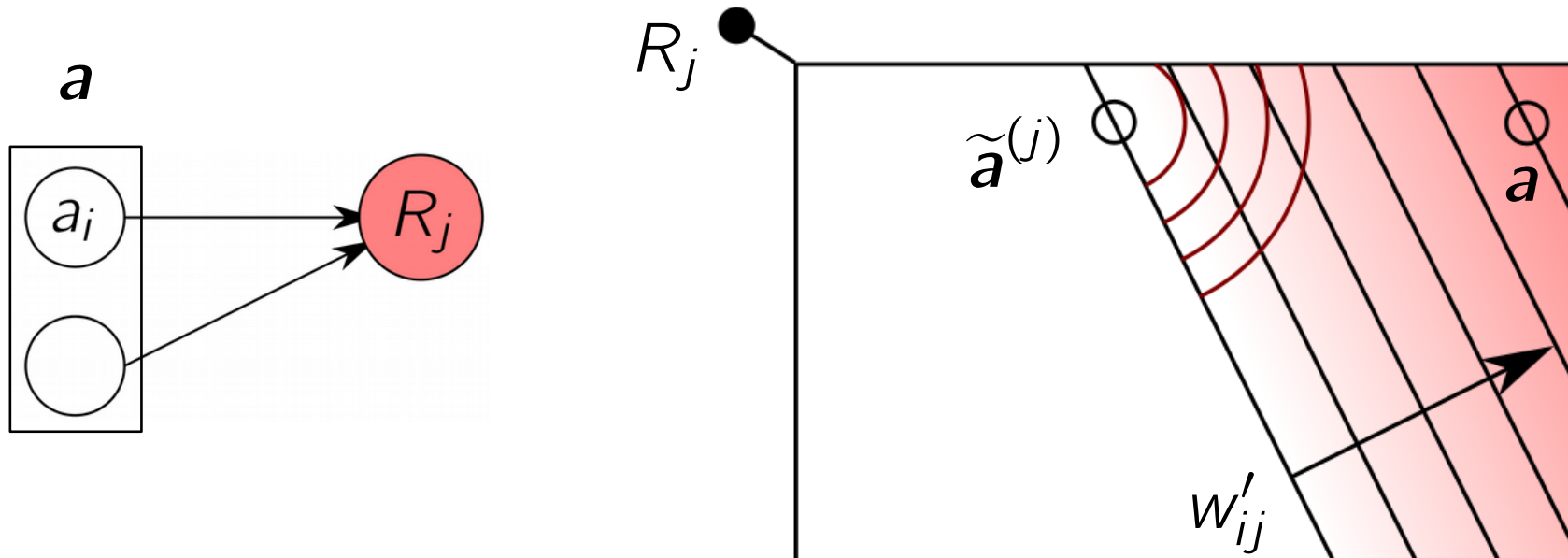
The diagram shows the equation  $R_j = a_j \times c_j$  using 3D surface plots. On the left, a red surface labeled  $R_j$  represents the relevance. This is equal to a blue surface labeled  $a_j$  (representing the activation) multiplied by a green surface labeled  $c_j$  (representing an approximately constant term). The surfaces are shown as perspective views of 2D planes, with the blue and green surfaces having a slight curvature.

# DTD Step 1: The Structure of Relevance



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

## DTD Step 2: Taylor Expansion



$$R_j(\mathbf{a}) = \sum_i \left. \frac{\partial R_j}{\partial a_i} \right|_{\tilde{\mathbf{a}}^{(j)}} \cdot (a_i - \tilde{a}_i^{(j)})$$

# DTD Step 2: Taylor Expansion

Taylor expansion at root point:

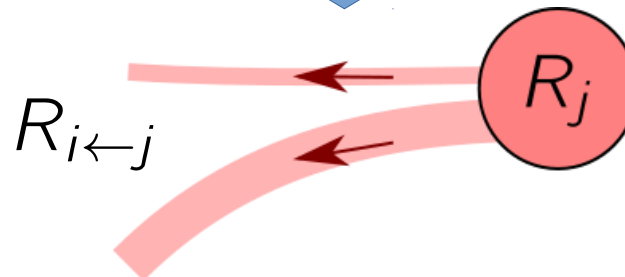
$$R_j(\mathbf{a}) = \sum_i \underbrace{\frac{\partial R_j}{\partial a_i} \Big|_{\tilde{\mathbf{a}}^{(j)}}}_{\text{weight}} \cdot (a_i - \tilde{a}_i^{(j)})$$



$$\frac{(a_i - \tilde{a}_i^{(j)}) w_{ij}}{\sum_i (a_i - \tilde{a}_i^{(j)}) w_{ij}} R_j$$



Relevance can now be backward propagated





# DTD Step 3: Choosing the Root Point

$$R_{i \leftarrow j} = \frac{(a_i - \tilde{a}_i^{(j)}) w_{ij}}{\sum_i (a_i - \tilde{a}_i^{(j)}) w_{ij}} R_j \quad (\text{Deep Taylor generic})$$



Choice of root point

		$\tilde{\mathbf{a}}^{(j)} \in \mathcal{D}$
1. nearest root	$\tilde{\mathbf{a}}^{(j)} = \mathbf{a} - t \cdot \mathbf{w}_j$	
2. rescaled excitations	$\tilde{\mathbf{a}}^{(j)} = \mathbf{a} - t \cdot \mathbf{a} \odot \mathbf{1}_{w_j > 0}$	✓



$$R_{i \leftarrow j} = \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

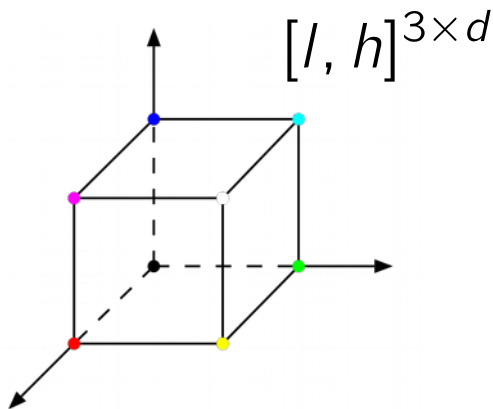
(same as LRP- $\alpha_1 \beta_0$ )



# DTD: Choosing the Root Point

$$R_{i \leftarrow j} = \frac{(x_i - \tilde{x}_i^{(j)}) w_{ij}}{\sum_i (x_i - \tilde{x}_i^{(j)}) w_{ij}} R_j \quad (\text{Deep Taylor generic})$$

**Pixels domain**



**Choice of root point**

$$(x - \tilde{x}^{(j)}) = t \cdot (x - l \odot 1_{w_j > 0} - h \odot 1_{w_j < 0})$$

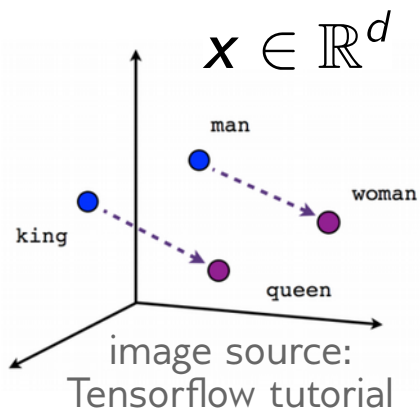
$$R_{i \leftarrow j} = \frac{x_{ij} w_{ij} - l_i w_{ij}^+ - h_i w_{ij}^-}{\sum_i x_{ij} w_{ij} - l_i w_{ij}^+ - h_i w_{ij}^-} R_j$$

# DTD: Choosing the Root Point

$$R_{i \leftarrow j} = \frac{(x_i - \tilde{x}_i^{(j)}) w_{ij}}{\sum_i (x_i - \tilde{x}_i^{(j)}) w_{ij}} R_j \quad (\text{Deep Taylor generic})$$



**Embedding:**



**Choice of root point**

$$(x - x^{(j)}) = t \cdot w_j$$



$$R_{i \leftarrow j} = \frac{w_{ij}^2}{\sum_i w_{ij}^2} R_j$$

# DTD: Application to Pooling Layers

A sum-pooling layer over positive activations is equivalent to a ReLU layer with weights 1.

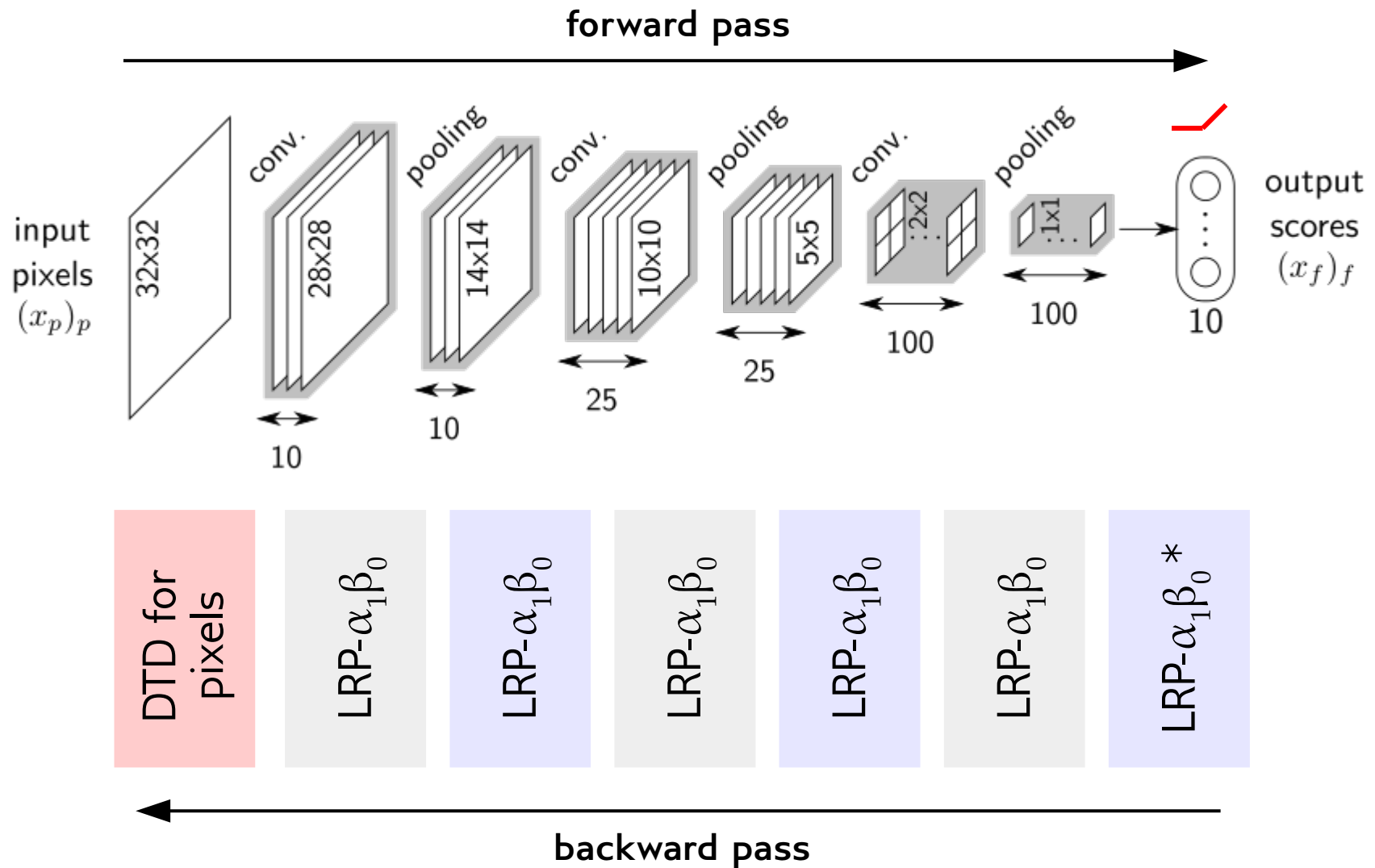
$$a_j = \left( \sum_i a_i \right) = \max \left( 0, \sum_i a_i 1_{ij} + 0_j \right)$$

A  $p$ -norm pooling layer can be approximated as a sum-pooling layer multiplied by a ratio of norms that we treat as constant [Montavon'17].

$$a_j = \left( \sum_i a_i \right) \cdot \frac{\|(a_i)_i\|_p}{\|(a_i)_i\|_1}$$

**→ Treat pooling layers as ReLU detection layers**

# Deep Taylor Decomposition on ConvNets



\* For top-layers, other rules may improve selectivity

# Implementing Propagation Rules

Example: LRP- $\alpha_1\beta_0$ :

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

Sequence of element-wise  
computations

$$z_j \rightarrow \sum_i a_i w_{ij}^+$$

$$s_j \rightarrow R_j / z_j$$

$$c_i \rightarrow \sum_j w_{ij}^+ s_j$$

$$R_i \rightarrow a_i c_i$$

Sequence of vector  
computations

$$\mathbf{z} \rightarrow \mathbf{W}_+^\top \cdot \mathbf{a}$$

$$\mathbf{s} \rightarrow \mathbf{R} \oslash \mathbf{z}$$

$$\mathbf{c} \rightarrow \mathbf{W}_+ \cdot \mathbf{s}$$

$$\mathbf{R} \rightarrow \mathbf{a} \odot \mathbf{c}$$

# Implementing Propagation Rules

Example: LRP- $\alpha_1\beta_0$ :

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

Code that reuses forward and gradient computations:

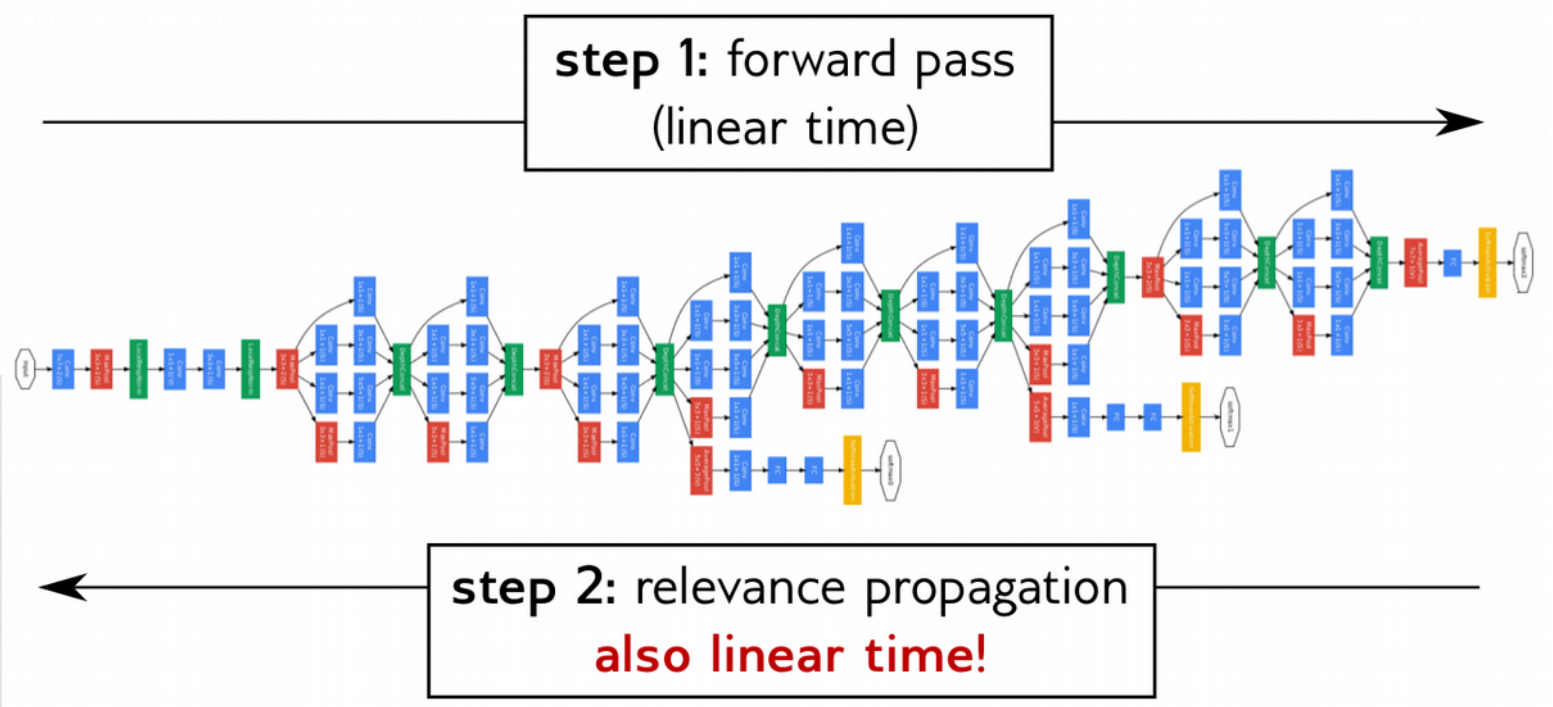
---

```
def lrp(layer, a, R):  
  
    clone = layer.clone()  
    clone.W = maximum(0, layer.W)  
    clone.B = 0  
  
    z = clone.forward(a)  
    s = R / z  
    c = clone.backward(s)  
  
    return a * c
```

---

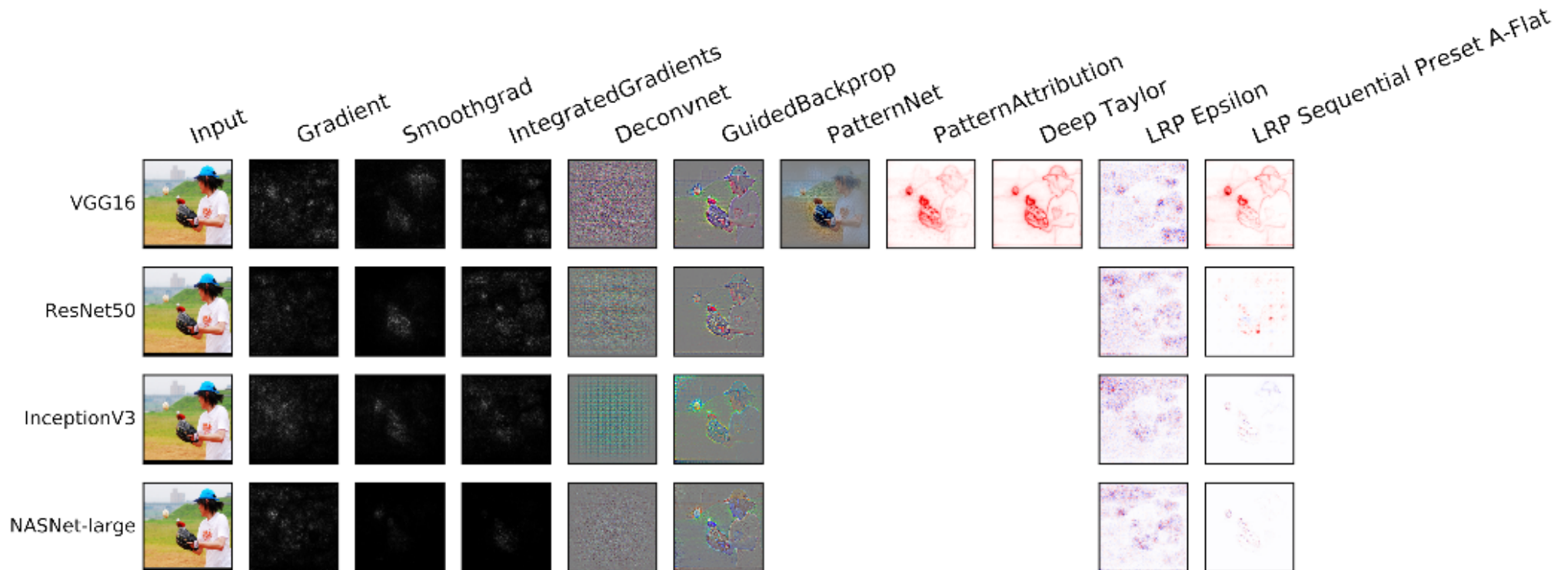


# How Deep Taylor / LRP Scales



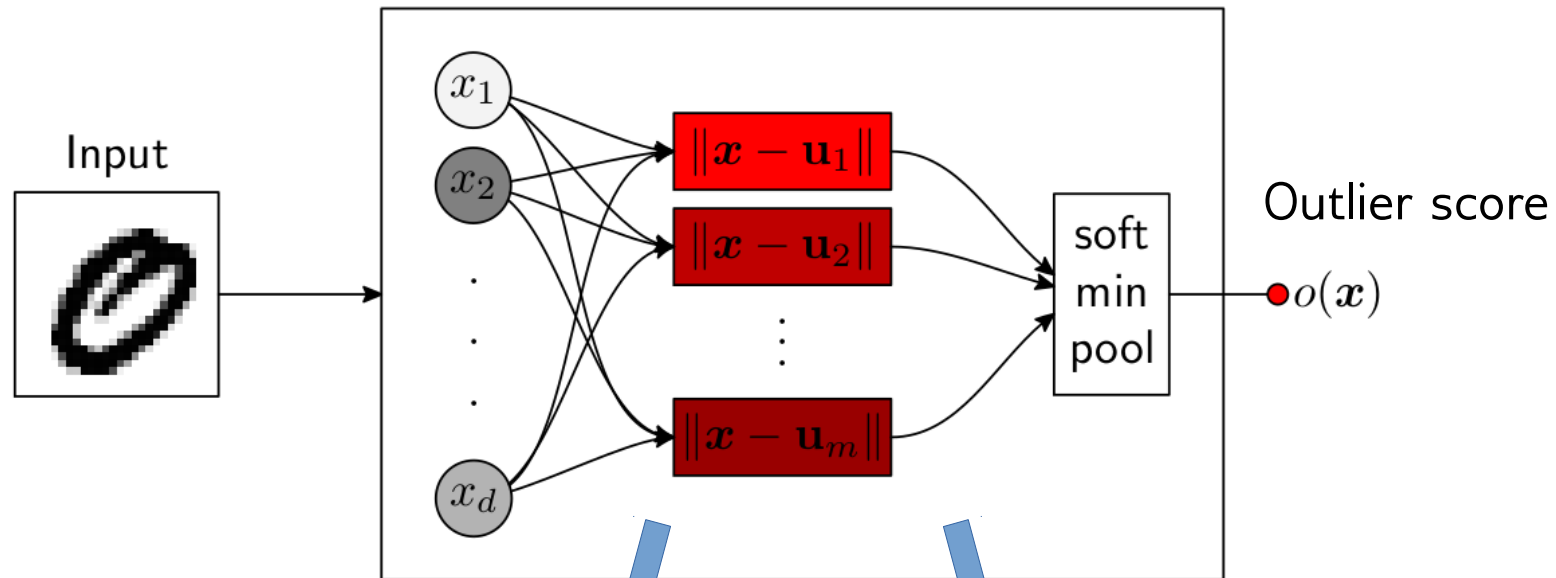
# Implementation on Large-Scale Models [Alber'18]

<https://github.com/albermax/innvestigate>



# DTD for Kernel Models [Kauffmann'18]

## 1. Build a neural network equivalent of the One-Class SVM:



## 2. Computes its deep Taylor decomposition

$$R_i = \sum_j \frac{(x_i - u_{ij})^2}{\|x - u_j\|_2^2} (R_j - D_j^+)$$

### Gaussian/Laplace Kernel

$$R_j = (a_j + \varepsilon_j) \cdot \frac{\exp(-a_j)}{\sum_j \exp(-a_j)}$$

### Student Kernel

$$R_j = a_j \cdot \mathbb{H}[(h_{j'}/h_j)_{j'}]$$

# DTD: Choosing the Root Point (Revisited)

$$R_{i \leftarrow j} = \frac{(a_i - \tilde{a}_i^{(j)}) w_{ij}}{\sum_i (a_i - \tilde{a}_i^{(j)}) w_{ij}} R_j \quad (\text{Deep Taylor generic})$$



Choice of root point		$\tilde{\mathbf{a}}^{(j)} \in \mathcal{D}$
1. nearest root	$\tilde{\mathbf{a}}^{(j)} = \mathbf{a} - t \cdot \mathbf{w}_j$	
2. rescaled excitations	$\tilde{\mathbf{a}}^{(j)} = \mathbf{a} - t \cdot \mathbf{a} \odot \mathbf{1}_{\mathbf{w}_j > 0}$	✓
3. rescaled activation	$\tilde{\mathbf{a}}^{(j)} = \mathbf{a} - t \cdot \mathbf{a}$	✓

2. LRP- $\alpha_1 \beta_0$

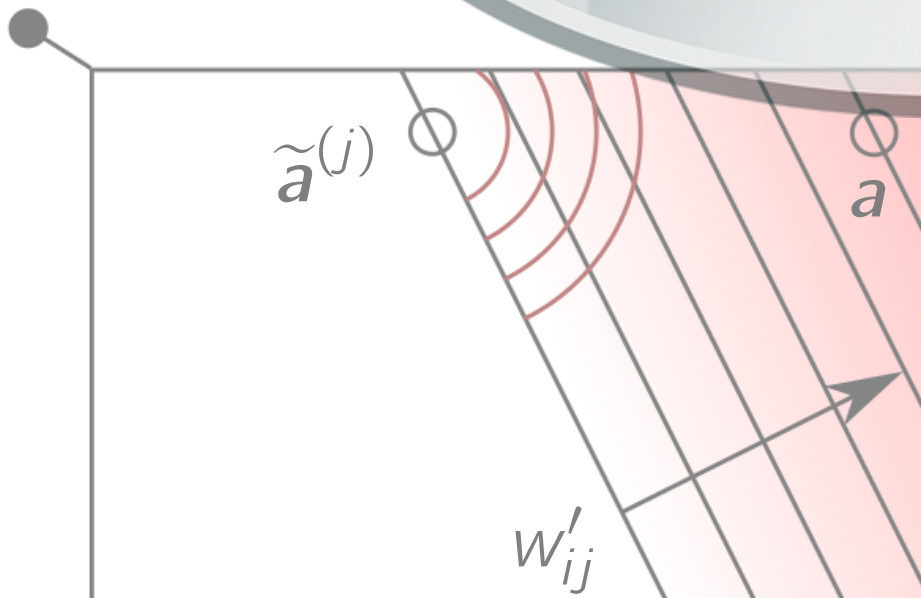
$$R_{i \leftarrow j} = \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

3. Another rule

$$R_{i \leftarrow j} = \frac{a_i w_{ij}}{\sum_i a_i w_{ij}} R_j$$

# Selecting the Explanation Technique

How to select the best root points?



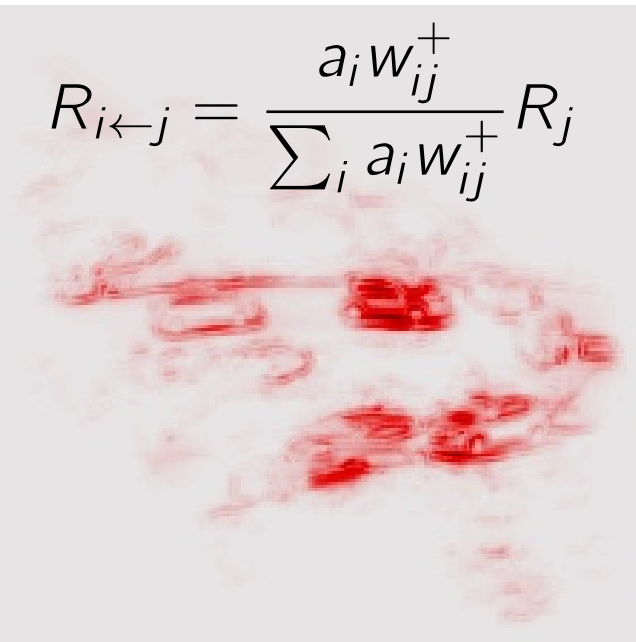
$$\tilde{a}^{(j)} = a - t \cdot a \odot \mathbf{1}_{w_j > 0}$$

$$\tilde{a}^{(j)} = a - t \cdot a$$

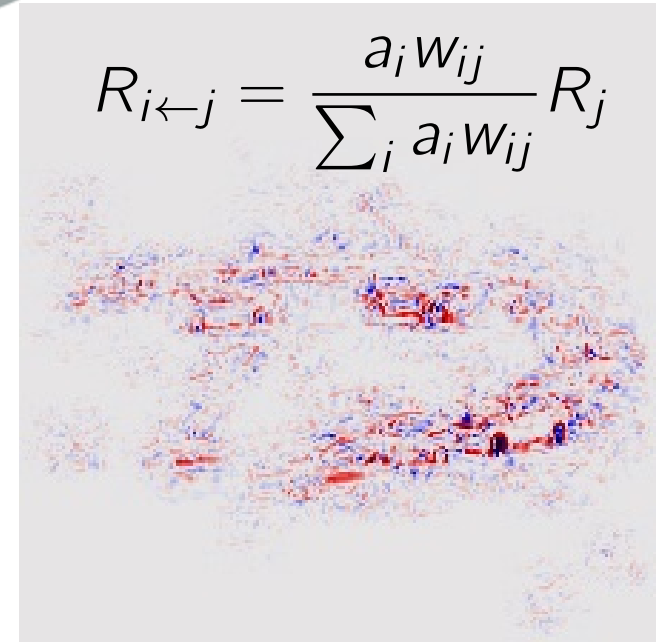
# Selecting the Explanation Technique

Which rule leads to the best explanation?

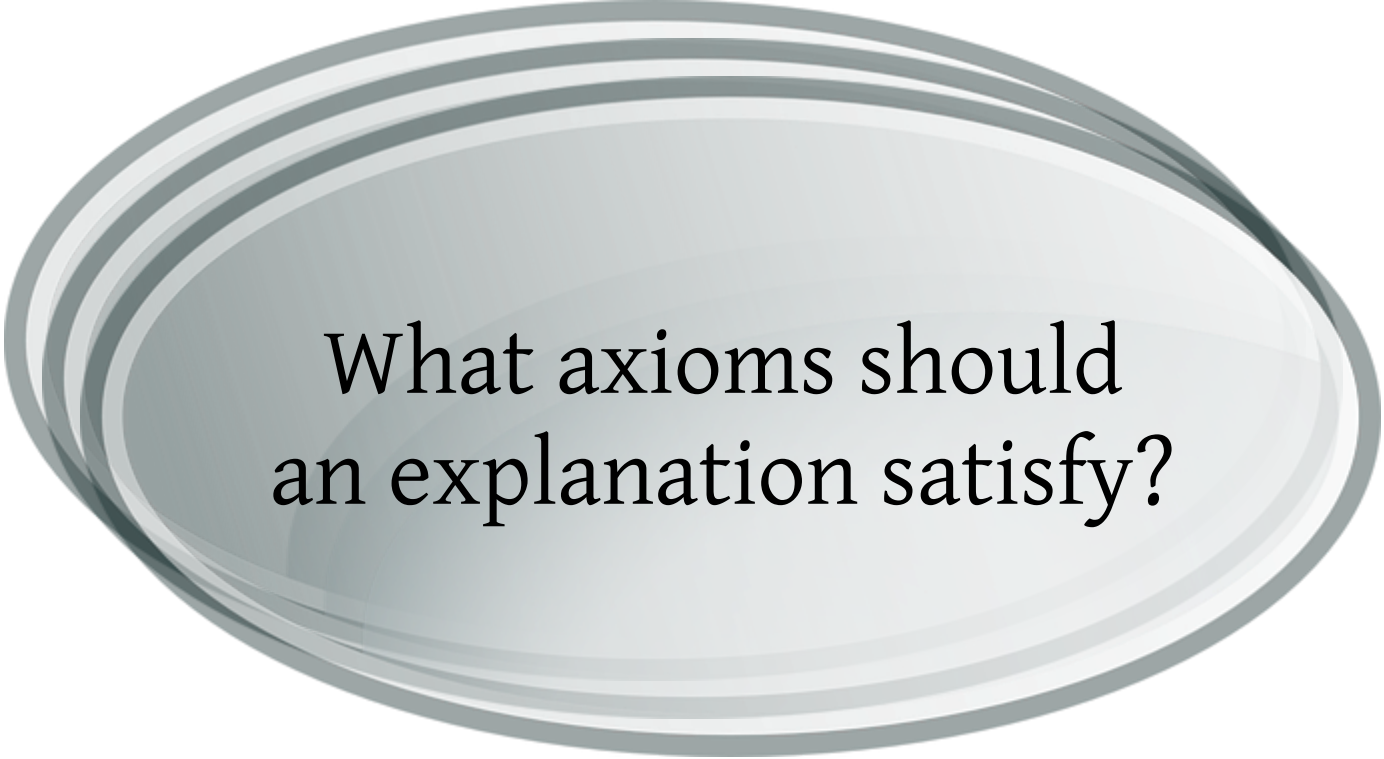
$$R_{i \leftarrow j} = \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$



$$R_{i \leftarrow j} = \frac{a_i w_{ij}}{\sum_i a_i w_{ij}} R_j$$



# Selecting the Explanation Technique



What axioms should  
an explanation satisfy?

$\left\{ \begin{array}{l} \text{Continuity: } (\mathbf{x} \approx \mathbf{x}') \wedge (f(\mathbf{x}) \approx f(\mathbf{x}')) \Rightarrow R(\mathbf{x}) \approx R(\mathbf{x}') \\ \text{Conservation: } \left( \sum_p R_p(\mathbf{x}) = f(\mathbf{x}) \right) \wedge \left( \sum_p |R_p(\mathbf{x})| < A \cdot |f(\mathbf{x})| \right). \end{array} \right\}$

# Selection based on Axioms

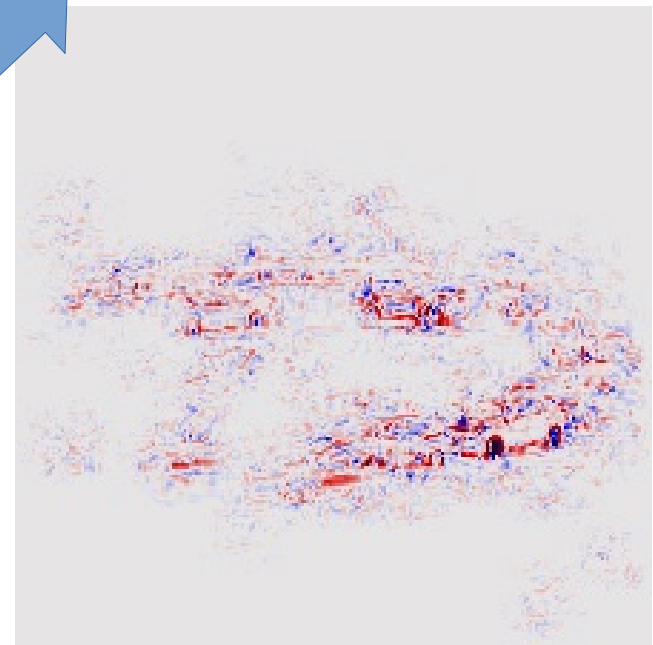
$$\left\{ \text{Conservation: } \left( \sum_p R_p(\mathbf{x}) = f(\mathbf{x}) \right) \wedge \left( \sum_p |R_p(\mathbf{x})| < A \cdot |f(\mathbf{x})| \right). \right\}$$

LRP- $\alpha_1\beta_0$

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$

$$R_i = \sum_j \frac{a_i w_{ij}}{\sum_i a_i w_{ij}} R_j$$

division by zero  $\rightarrow$   
scores explode.



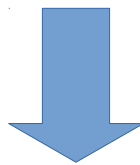


# Selection based on Axioms

$$\left\{ \begin{array}{l} \text{Conservation: } \left( \sum_p R_p(\mathbf{x}) = f(\mathbf{x}) \right) \wedge \left( \sum_p |R_p(\mathbf{x})| < A \cdot |f(\mathbf{x})| \right). \\ \text{Continuity: } (\mathbf{x} \approx \mathbf{x}') \wedge (f(\mathbf{x}) \approx f(\mathbf{x}')) \Rightarrow \mathbf{R}(\mathbf{x}) \approx \mathbf{R}(\mathbf{x}') \end{array} \right\}$$

LRP- $\alpha_1 \beta_0$

$$R_i = \sum_j \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$



$$R_i = a_i \delta_i$$

$$\delta_i = \sum_j w_{ij}^+ \frac{(\sum_i a_i w_{ij} + b_j)^+}{\sum_i a_i w_{ij}^+} \delta_j$$

$$R_i = \sum_j \frac{a_i w_{ij}}{\sum_i a_i w_{ij}} R_j$$



$$R_i = a_i \delta_i$$

$$\delta_i = \sum_j w_{ij} \underbrace{\frac{(\sum_i a_i w_{ij} + b_j)^+}{\sum_i a_i w_{ij}}}_{\text{discontinuous step function for } b_j = 0} \delta_j$$

discontinuous  
step function  
for  $b_j = 0$

# Explainable ML: Challenges

Underlying  
mathematical  
framework

Human  
perception

**Validating  
Explanations**

Similarity to  
ground truth

Perturbation  
analysis [Samek'17]

Axioms of an  
explanation



# Explainable ML: Opportunities

Detecting unexpected  
ML behavior

Finding weaknesses  
of a dataset

**Using  
Explanations**

Human  
interaction

Designing better  
ML algorithms?

Extracting new  
domain knowledge



# Check our webpage



[www.heatmapping.org](http://www.heatmapping.org)

with interactive demos, software, tutorials, ...

and our tutorial paper:

Montavon, G., Samek, W., Müller, K.-R. Methods for Interpreting and Understanding Deep Neural Networks, [Digital Signal Processing](#), 2018

# References

- **[Alber'18]** Alber, M. Lapuschkin, S., Seegerer, P., Hägele, M., Schütt, K., Montavon, G., Samek, W., Müller, K.-R., Dähne, S., Kindermans, P.-J. iNNvestigate neural networks. [CoRR abs/1808.04260](#), 2018
- **[Bach'15]** Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.-R., Samek, W. On pixel-wise explanations for nonlinear classifier decisions by layer-wise relevance propagation. [PLOS ONE 10 \(7\)](#), 2015
- **[Binder'18]** Binder, A., Bockmayr, M., ..., Müller, K.-R., Klauschen, F. Towards computational fluorescence microscopy: Machine learning-based integrated prediction of morphological and molecular tumor profiles [CoRR abs/1805.11178](#), 2018
- **[Bojarski'17]** Bojarski, M., Yeres, P., Choromanska, A., Choromanski, K., Firner, B., Jackel, L. D., Muller, U. Explaining how a deep neural network trained with end-to-end learning steers a car. [CoRR abs/1704.07911](#), 2017
- **[Kauffmann'18]** Kauffmann, J. Müller, K.-R., Montavon, G., Towards Explaining Anomalies: A Deep Taylor Decomposition of One-Class Models [CoRR abs/1805.06230](#), 2018
- **[Lapuschkin'16]** Lapuschkin, S., Binder, A., Montavon, G., Müller, K.-R., Samek, W. Analyzing classifiers: Fisher vectors and deep neural networks. [CVPR 2016](#)
- **[Montavon'17]** Montavon, G., Lapuschkin, S., Binder, A., Samek, W., Müller, K.-R. Explaining nonlinear classification decisions with deep Taylor decomposition. [Pattern Recognition 65](#), 211–222, 2017
- **[Samek'17]** Samek, W., Binder, A., Montavon, G., Lapuschkin, S., Müller, K.-R. Evaluating the visualization of what a deep neural network has learned. [IEEE Transactions on Neural Networks and Learning Systems](#), 2017
- **[Schütt'17]** Schütt, K. T., Arbabzadah, F., Chmiela, S., Müller, K.-R., Tkatchenko, A. Quantum-Chemical Insights from Deep Tensor Neural Networks, [Nature Communications 8](#), 13890, 2017
- **[Symonian'13]** Symonian, K. Vedaldi, A. Zisserman, A. Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. [ArXiv 2013](#)