



Tutorials

Interpretable Deep Learning: Towards Understanding & Explaining DNNs

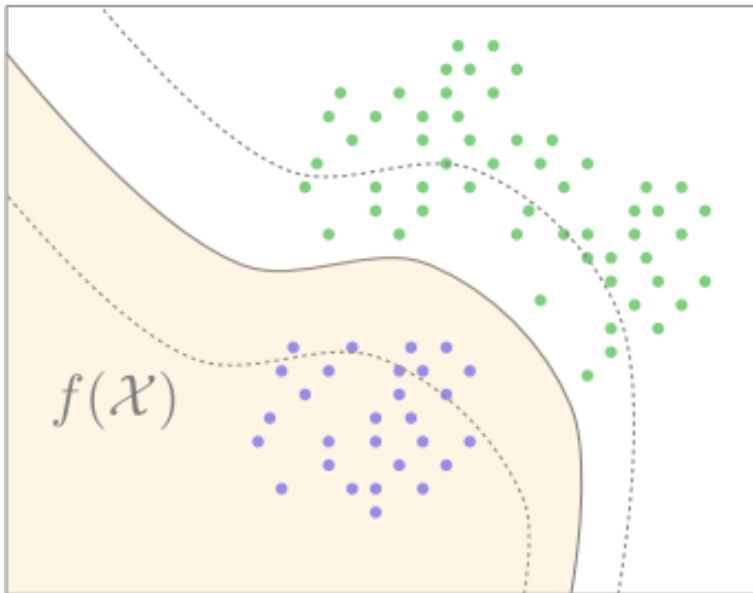
Part 2: Methods of Explanation

Wojciech Samek, Grégoire Montavon, Klaus-Robert Müller

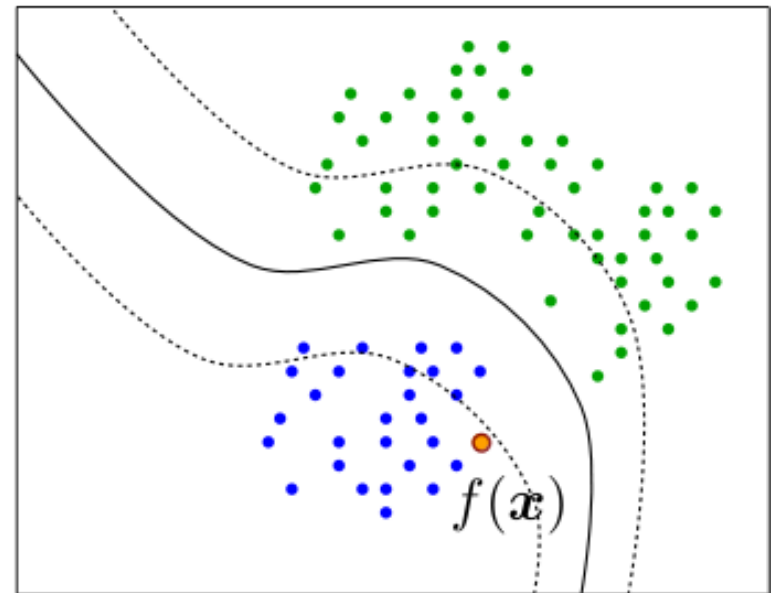
What Will be Covered in Part 2



interpreting
predicted classes

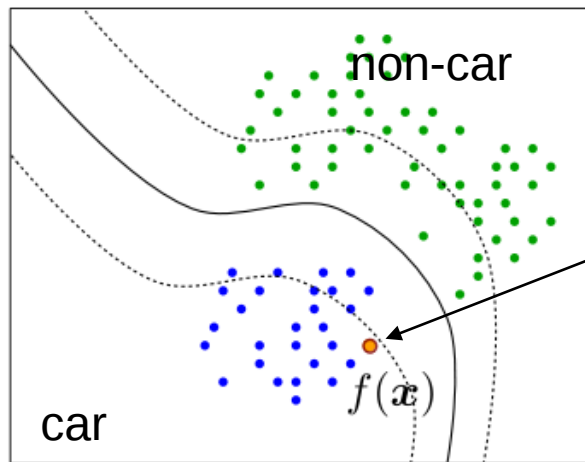


**explaining
individual decisions**

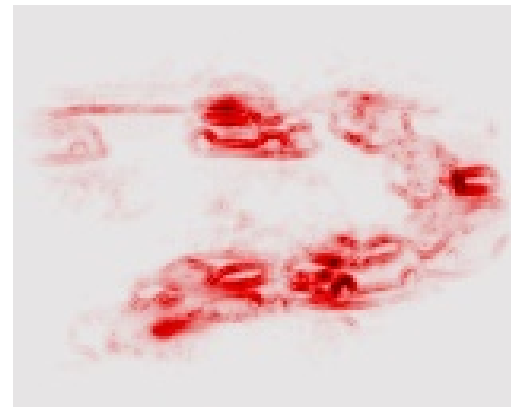


Explaining Individual Decisions

Q: Where in the image the neural networks sees evidence for a car?



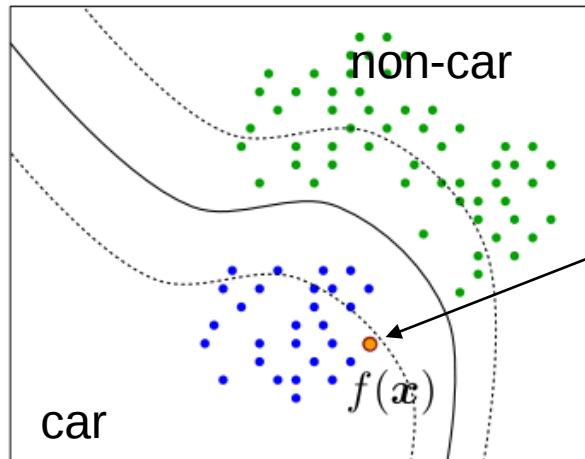
$$R(x, f)$$



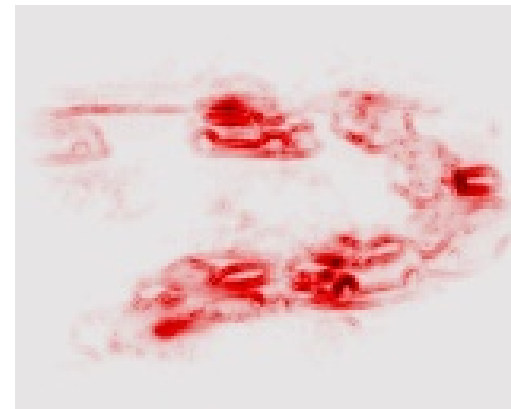
Examples of Methods that Explain Decisions

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 Individual Decisions



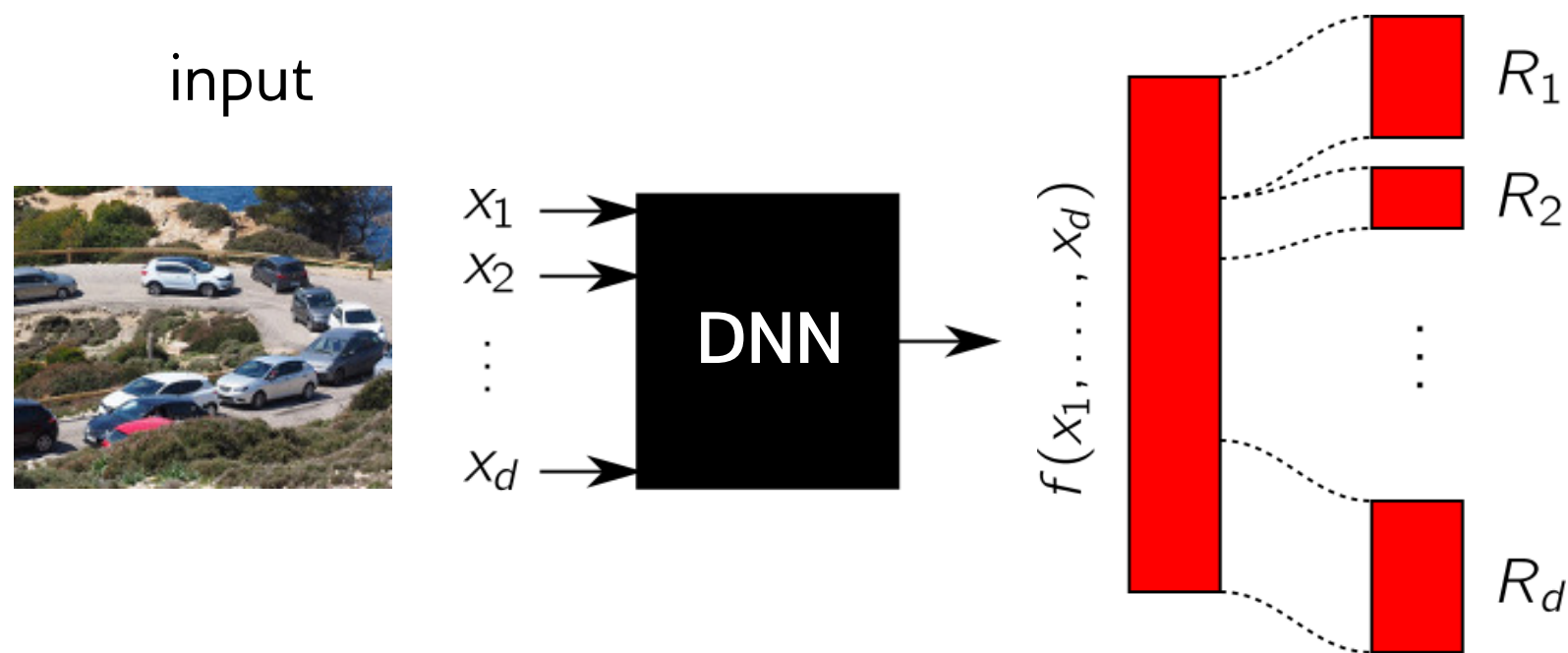
$$R(x, f)$$



Q: In which proportion has each car contributed to the prediction?

Explaining by Decomposing

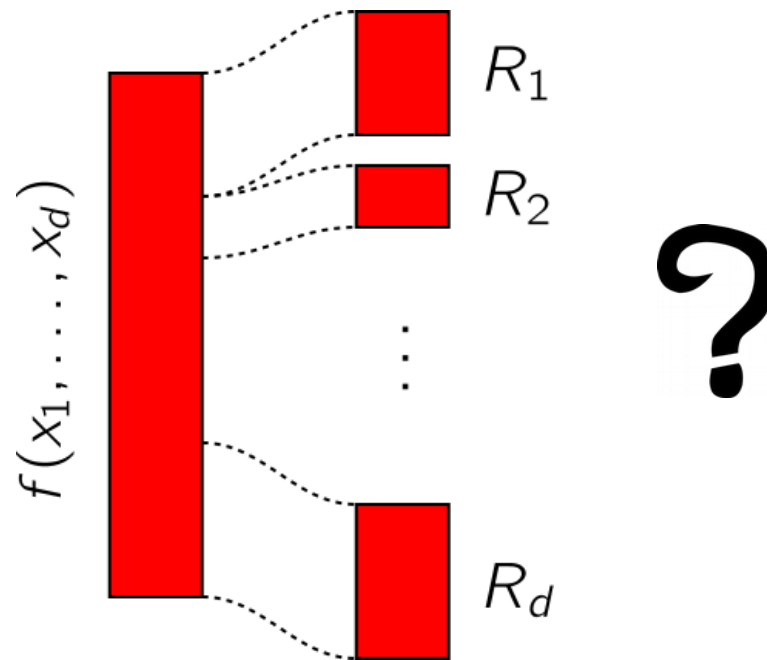
Goal: Determine the share of the output that should be attributed to each input variable.



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

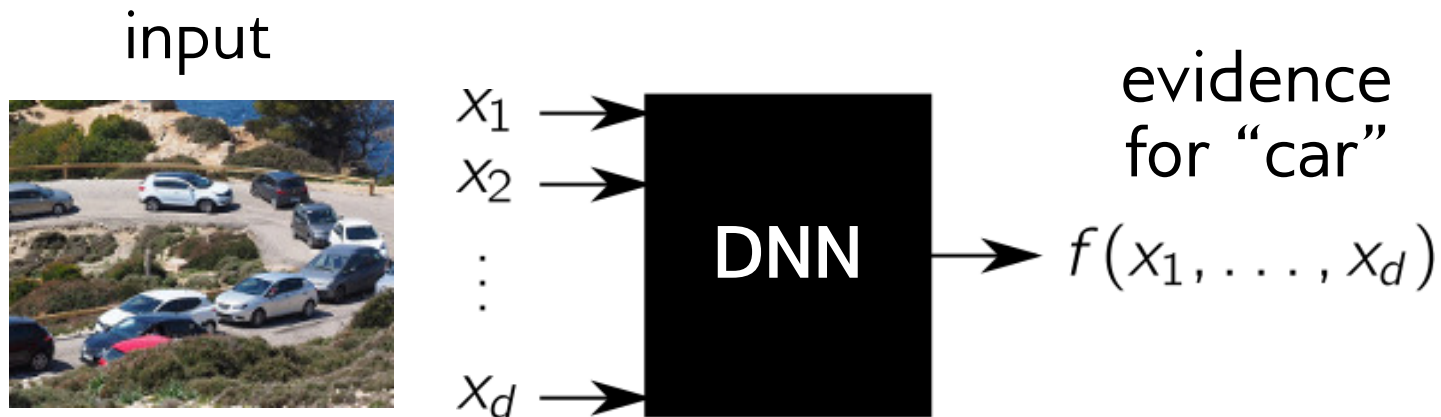
Explaining by Decomposing

Goal: Determine the share of the output that should be attributed to each input variable.



Decomposing a prediction is generally difficult.

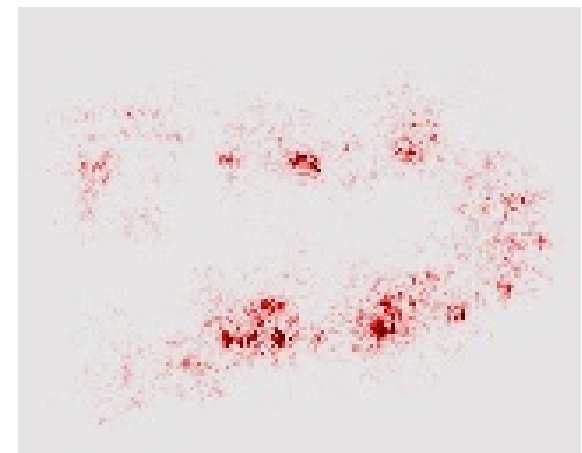
Sensitivity Analysis



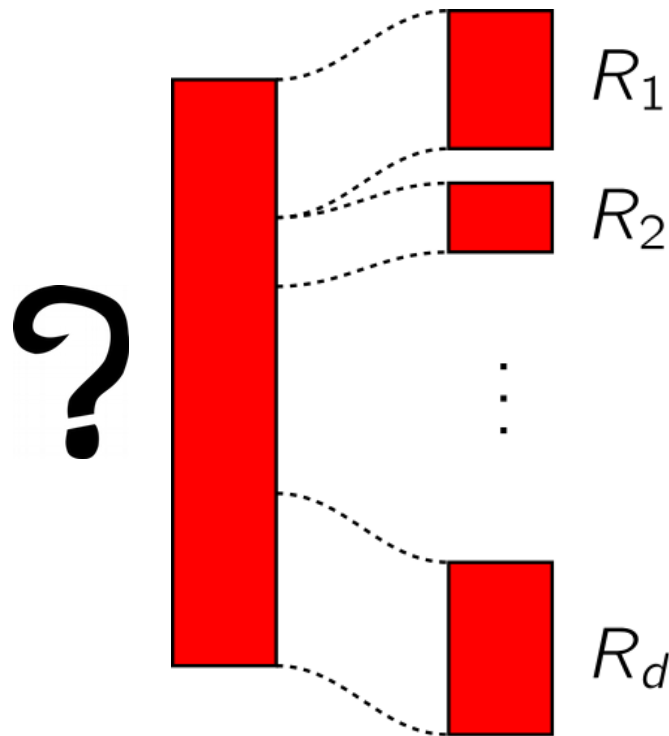
computes for each pixel:

$$R_i = \left(\frac{\partial f}{\partial x_i} \right)^2 \longrightarrow$$

explanation for “car”
(heatmap):



Sensitivity Analysis



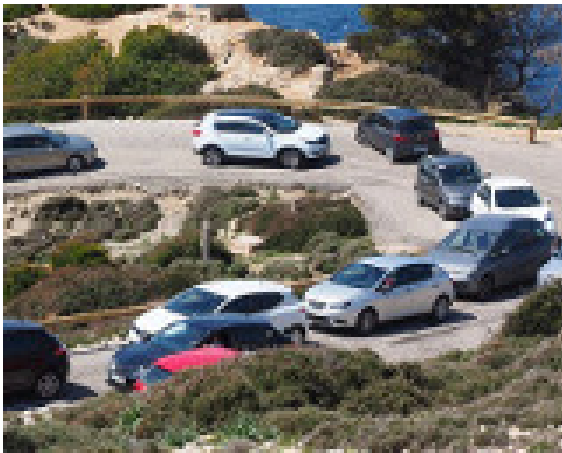
Question: If sensitivity analysis computes a decomposition of something: Then, *what* does it decompose?

$$R_i = \left(\frac{\partial f}{\partial x_i} \right)^2 \quad \rightarrow \quad \sum_{i=1}^d R_i = \|\nabla f(\mathbf{x})\|^2$$

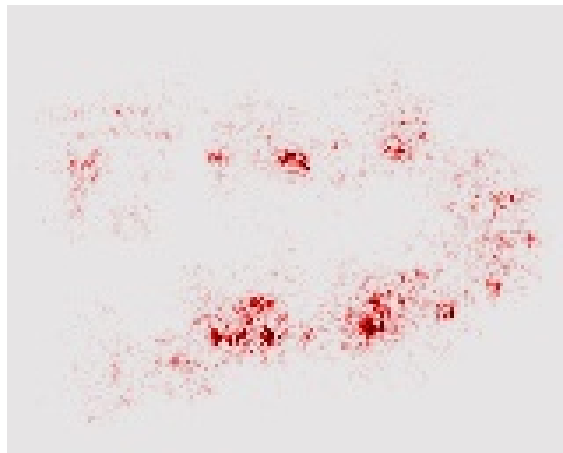
Sensitivity Analysis

Sensitivity analysis explains a *variation* of the function, not the function value itself.

input



explanation for “car”



variation = make something appear less/more a car.

$$\sum_{i=1}^d R_i = \|\nabla f(\mathbf{x})\|^2$$

The Taylor Expansion Approach

1. Take a linear model:

$$f(\mathbf{x}) = \sum_{i=1}^d w_i x_i + b$$

2. First-order expansion at root point:

$$f(\mathbf{x}) = \sum_{i=1}^d \left. \frac{\partial f}{\partial x_i} \right|_{\tilde{\mathbf{x}}} \cdot (x_i - \tilde{x}_i)$$

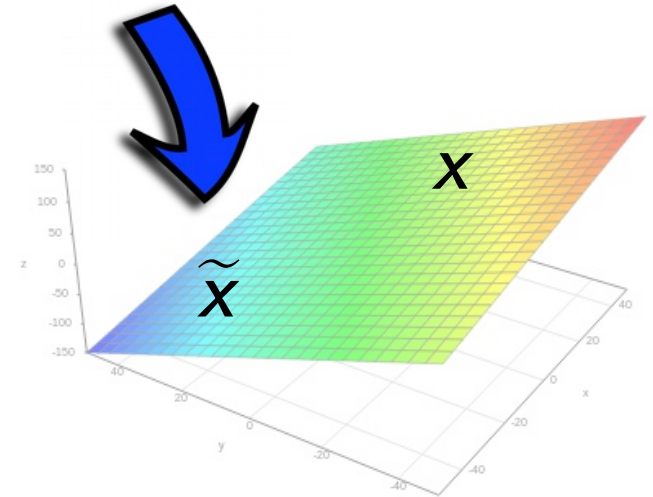
3. Identifying
linear terms:

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

a decomposition

Observation: explanation depends on the root point.

root point



The Taylor Expansion Approach

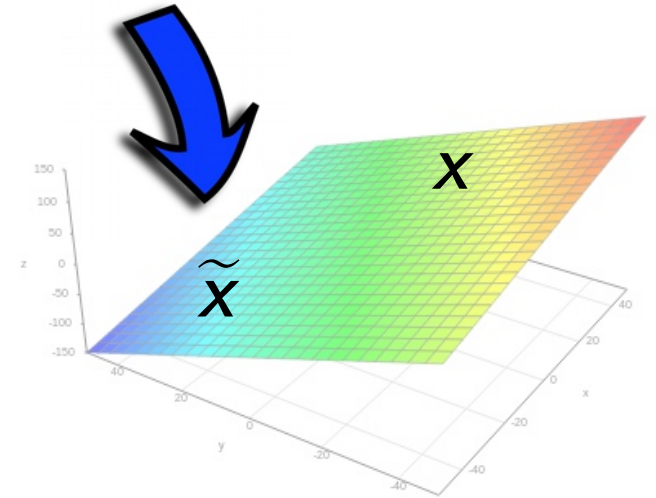
Obtained relevance scores

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

How to choose the **root point** ?

- Closeness to the actual data point
- Membership to the input domain (e.g. pixel space)
- Membership to the data manifold.

root point



Non-Linear Models

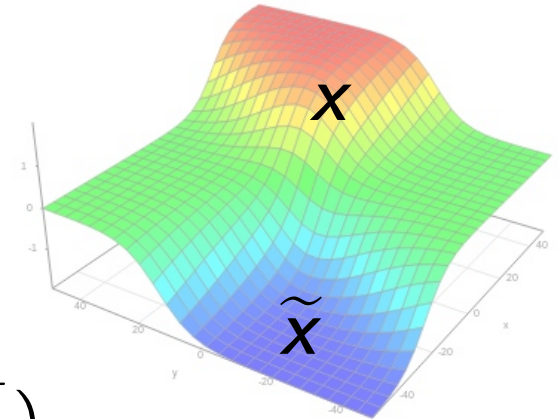
Nonlinear model

$$f(\mathbf{x}) = \sum_j \rho\left(\sum_{i=1}^d w_{ij}x_i + b_j\right)$$

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



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



Simple Taylor decomposition is not suitable for highly non-linear models.

Overcoming NonLinearity

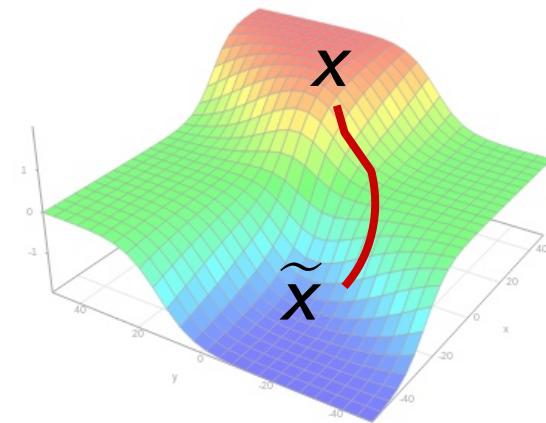
Integrated Gradients [Sundararajan'17]:

$$f(\mathbf{x}) = \int_{\xi=\tilde{\mathbf{x}}}^{\mathbf{x}} \sum_{i=1}^d \frac{\partial f}{\partial x_i} \Big|_{\xi} \cdot d\xi_i$$

X

↙ ↘

$$f(\mathbf{x}) = \sum_{i=1}^d \underbrace{\int_{\xi=\tilde{\mathbf{x}}}^{\mathbf{x}} \frac{\partial f}{\partial x_i} \Big|_{\xi} \cdot d\xi_i}_{R_i}$$



- Fully decomposable
- Require computing an integral (expensive)
- Which integration path?

Overcoming NonLinearity

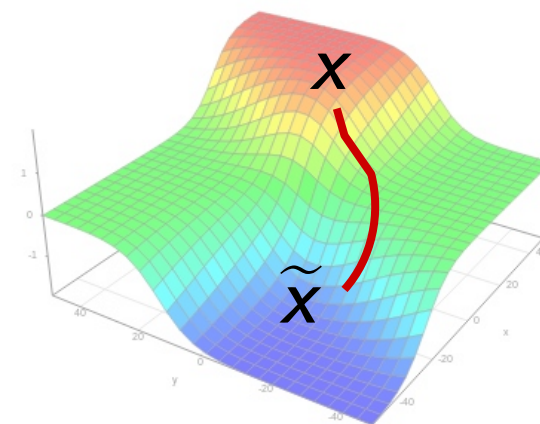
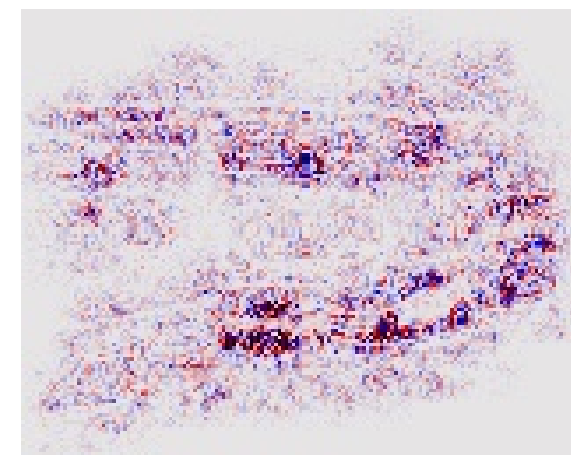
Special case when the origin is a root point and the gradient along the integration path is constant:

$$f(\mathbf{x}) = \sum_{i=1}^d \int_{\xi=0}^{\mathbf{x}} \frac{\partial f}{\partial x_i} \cdot d\xi_i$$

$$f(\mathbf{x}) = \sum_{i=1}^d \underbrace{\frac{\partial f}{\partial x_i}}_{R_i} \cdot x_i$$



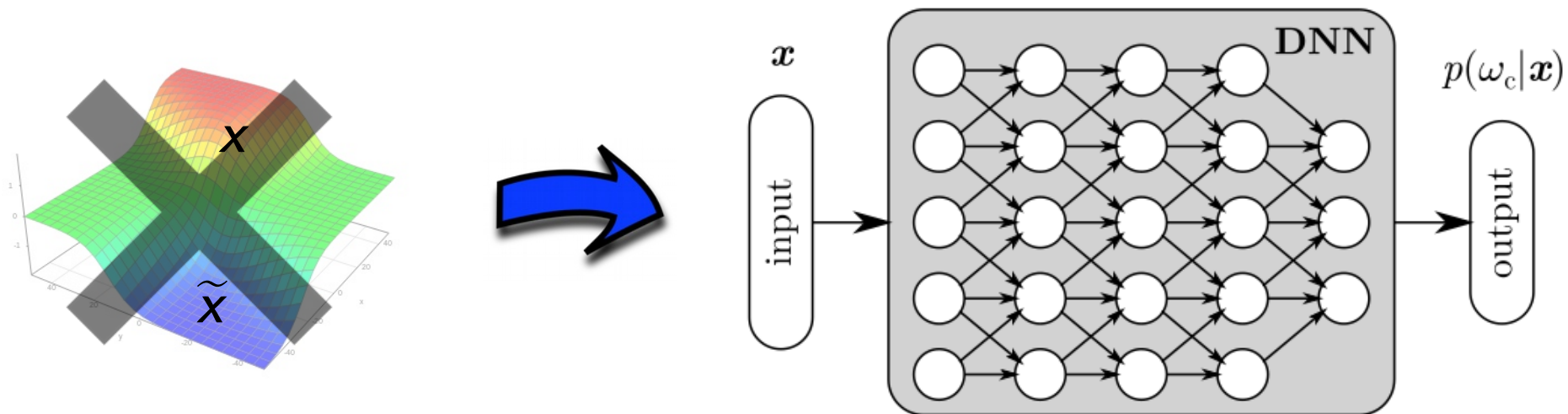
gradient x input





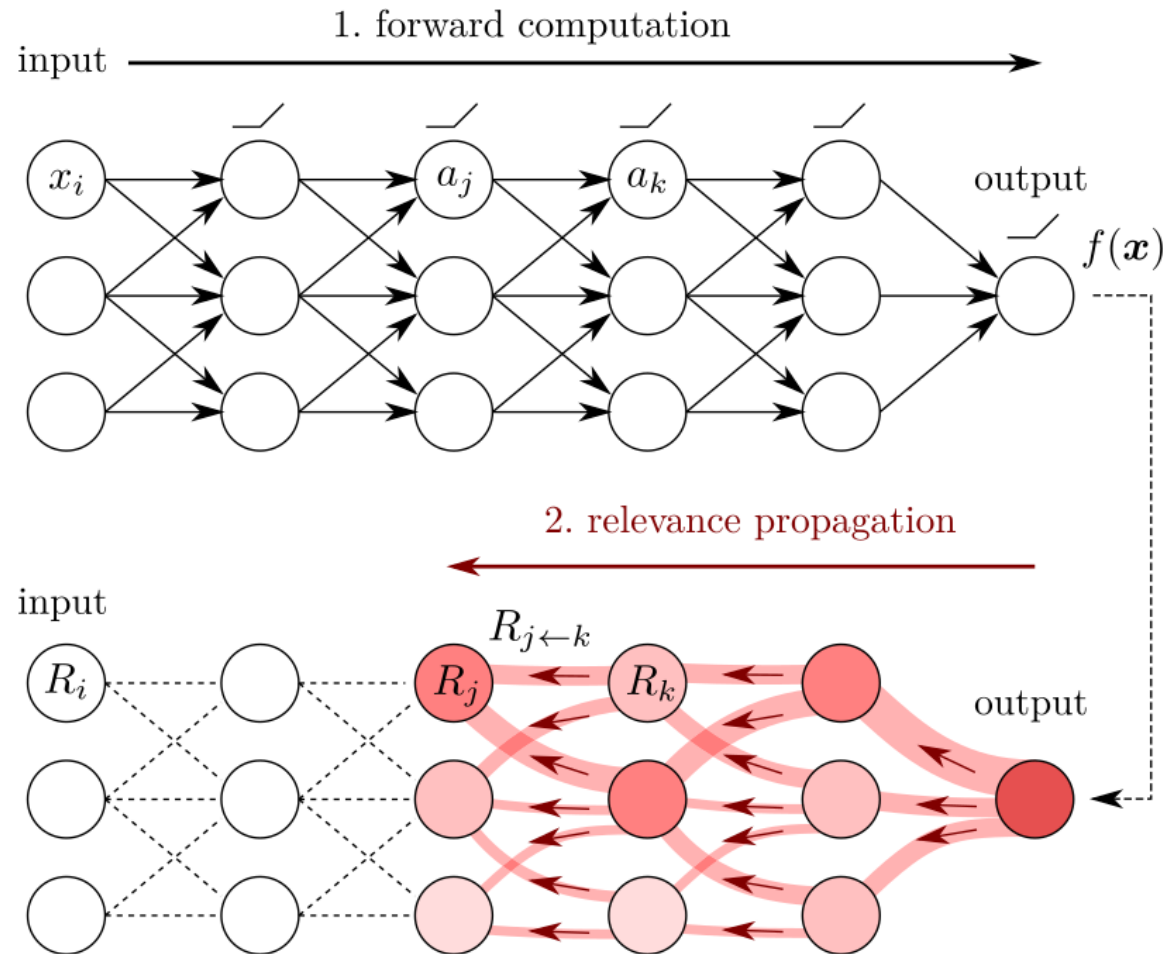
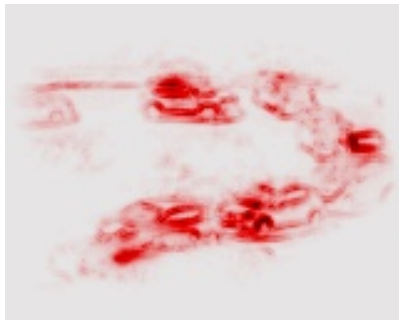
Let's consider a
different approach ...

Overcoming NonLinearity



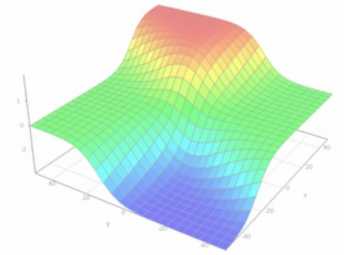
View the decision as a **graph computation** instead of a function evaluation, and propagate the decision backwards until the input is reached.

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

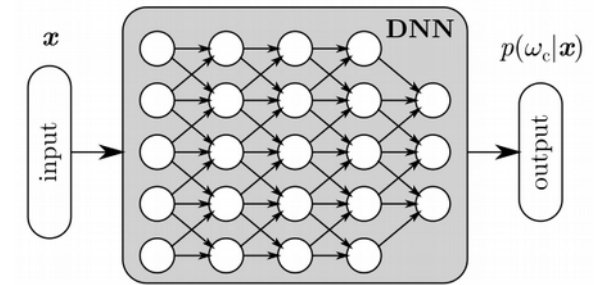


$$R_j = \sum_k \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} R_k$$

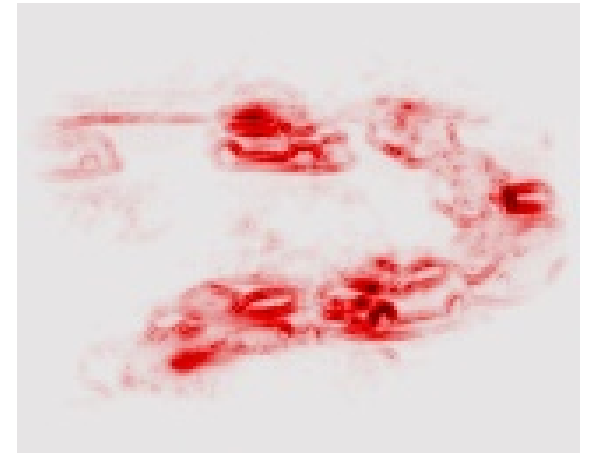
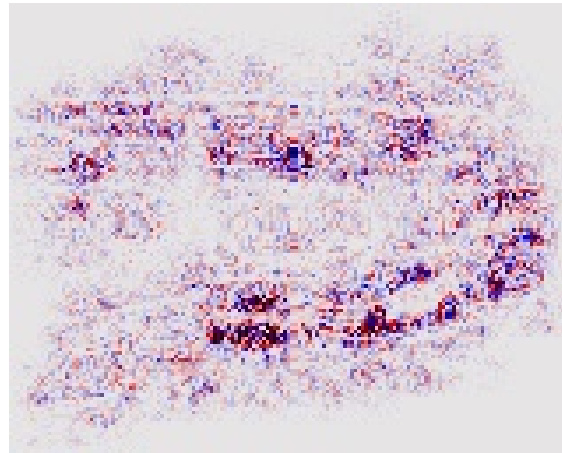
Gradient-Based vs. LRP



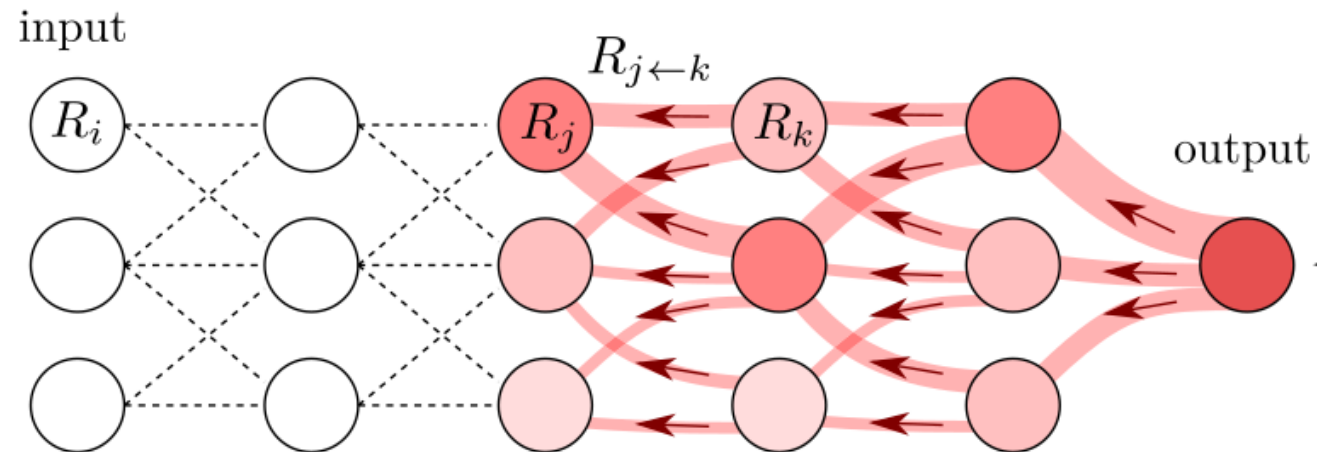
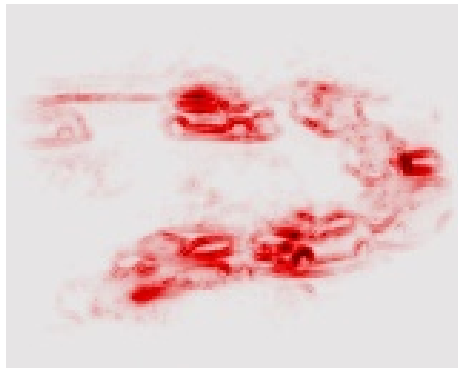
$\text{grad} \times \text{input}$



$\text{LRP-}\alpha_1\beta_0$



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



Carefully engineered
propagation rule:

e.g. $\text{LRP-}\alpha_1\beta_0$

$$R_j = \underbrace{\sum_k}_{\text{pooling received messages}} \frac{\overbrace{a_j w_{jk}^+}^{\text{neuron contribution}}}{\underbrace{\sum_j a_j w_{jk}^+}_{\text{normalization term}}} \underbrace{R_k}_{\text{available for redistribution}}$$

LRP Propagation Rules: Two Views

View 1:





$$R_j = \underbrace{\sum_k}_{\text{pooling received messages}} \frac{\overbrace{a_j w_{jk}^+}^{\text{neuron contribution}}}{\underbrace{\sum_j a_j w_{jk}^+}_{\text{normalization term}}} \underbrace{R_k}_{\text{available for redistribution}}$$

View 2:

$$R_j = \underbrace{a_j}_{\text{neuron activation}} \underbrace{\sum_k w_{jk}^+}_{\text{weighted sum}} \frac{\overbrace{R_k}^{\text{available for redistribution}}}{\underbrace{\sum_j a_j w_{jk}^+}_{\text{normalization term}}}$$

Implementing Propagation Rules (1)

$$R_j = \overset{\text{neuron activation}}{a_j} \underbrace{\sum_k w_{jk}^+}_{\text{weighted sum}} \frac{\overset{\text{available for redistribution}}{R_k}}{\underbrace{\sum_j a_j w_{jk}^+}_{\text{normalization term}}}$$

	Element-wise operations	Vector operations
	$z_k \rightarrow \sum_j a_j w_{jk}^+$	$\mathbf{z} \rightarrow \mathbf{W}_+^\top \cdot \mathbf{a}$
	$s_k \rightarrow R_k / z_k$	$\mathbf{s} \rightarrow \mathbf{R} \oslash \mathbf{z}$
	$c_j \rightarrow \sum_k w_{jk}^+ s_k$	$\mathbf{c} \rightarrow \mathbf{W}_+ \cdot \mathbf{s}$
	$R_j \rightarrow a_j c_j$	$\mathbf{R} \rightarrow \mathbf{a} \odot \mathbf{c}$

Implementing Propagation Rules (2)

Code that reuses forward and gradient computations:

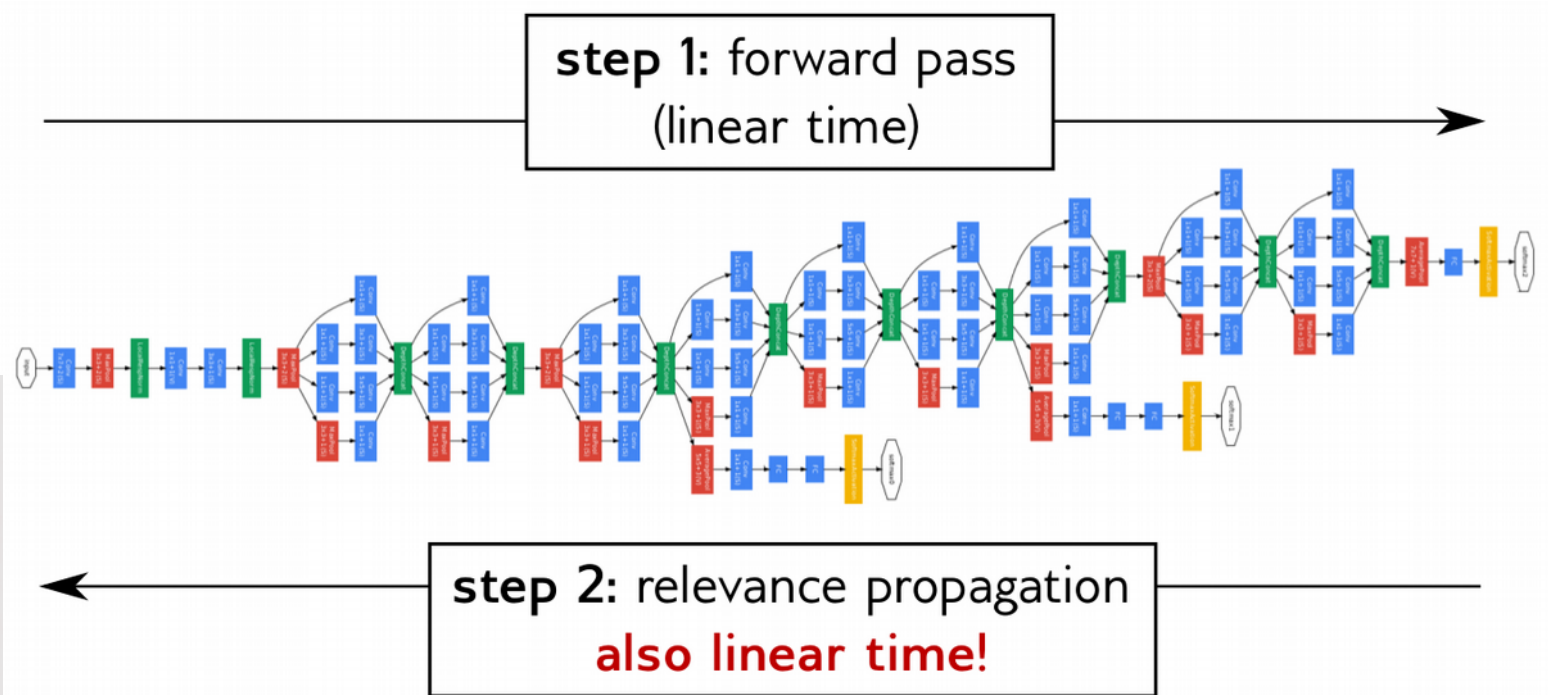
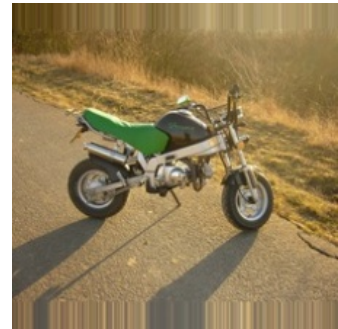
```
def lrp(layer, a, R):  
  
    clone = layer.clone()  
    clone.W = maximum(0, layer.W)  
    clone.B = 0  
  
    z = clone.forward(a)  
    s = R / z  
    c = clone.backward(s)  
  
    return a * c
```

$$R_j = \overset{\text{neuron activation}}{a_j} \sum_k w_{jk}^+ \frac{\overset{\text{available for redistribution}}{R_k}}{\sum_j \underset{\text{normalization term}}{a_j w_{jk}^+}}$$

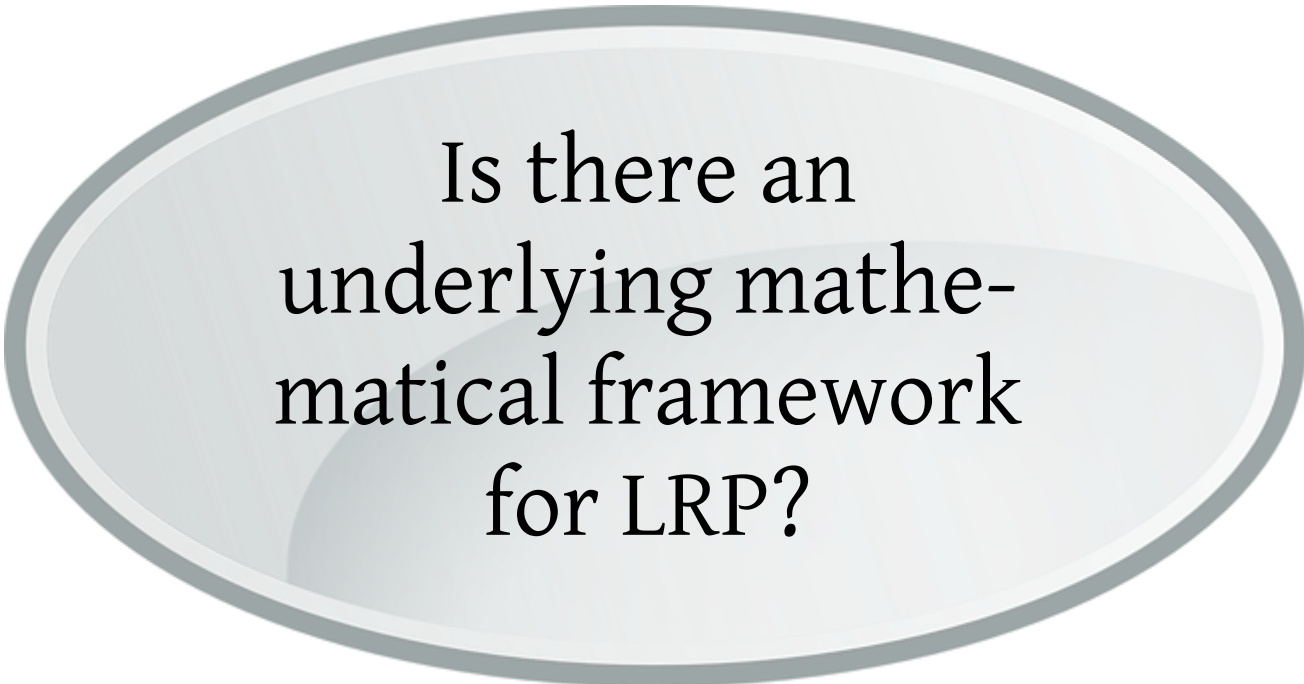
weighted sum **normalization term**

See also <http://www.heatmapping.org/tutorial>

How Fast is LRP ?



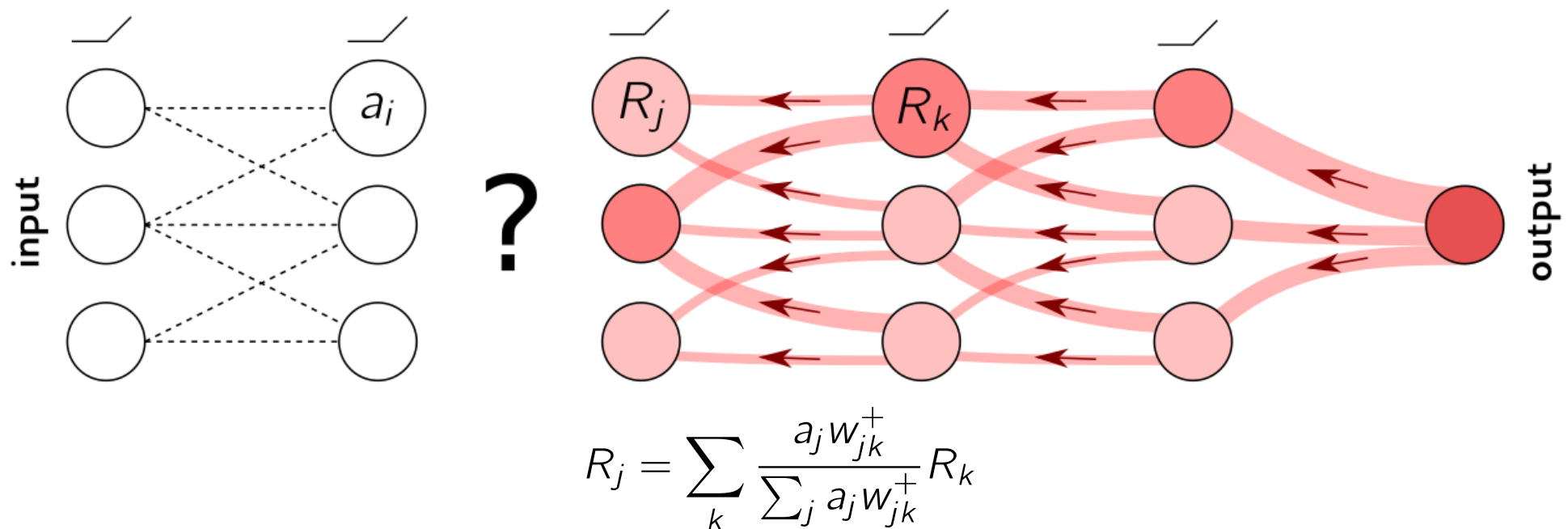
GPU-based implementation of LRP: Check out iNNvestigate [Alber'18]
<https://github.com/albermax/innvestigate>



Is there an
underlying mathe-
matical framework
for LRP?

Deep Taylor Decomposition [Montavon'17]

Question: Suppose that we have propagated the relevance until a given layer. How should it be propagated one layer further?



Idea: By performing a Taylor expansion of the relevance.

The Structure of Relevance

Reminder:

$$R_j = \overset{\text{neuron activation}}{a_j} \left(\underbrace{\sum_k w_{jk}^+}_{\text{weighted sum}} \frac{\overset{\text{available for redistribution}}{R_k}}{\underbrace{\sum_j a_j w_{jk}^+}_{\text{normalization term}}} \right)$$

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

$$R_j = a_j \times c_j$$

Deep Taylor Decomposition

Relevance neuron:

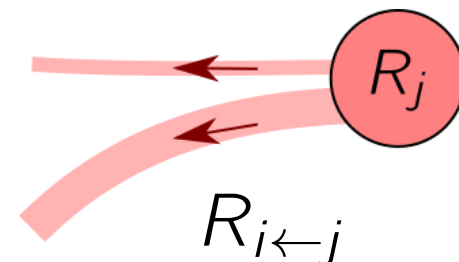
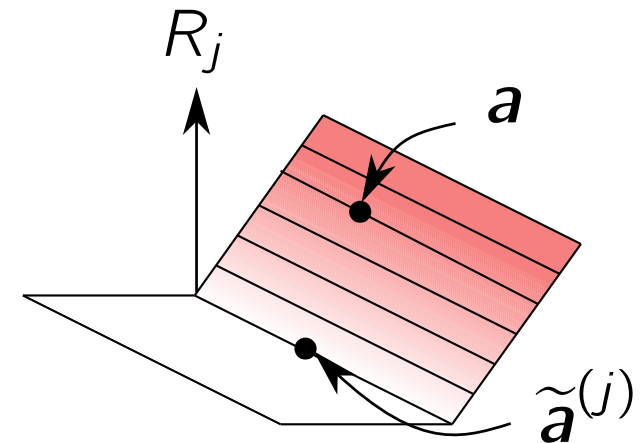
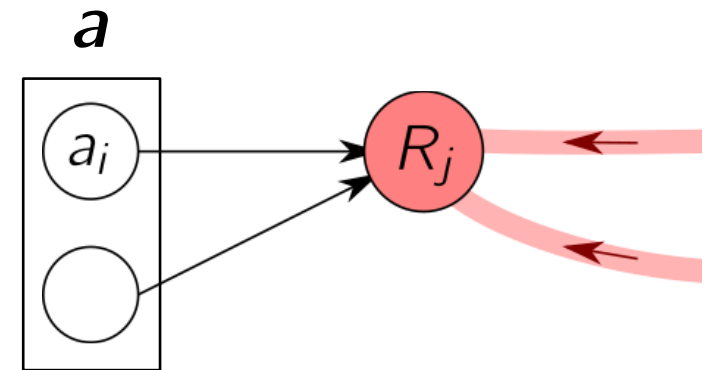
$$R_j(\mathbf{a}) = \max(0, \sum_i a_i w_{ij} + b_j) \cdot c_j$$

Taylor expansion:

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

Redistribution:

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



Choosing the Root Point

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

(Deep Taylor generic)



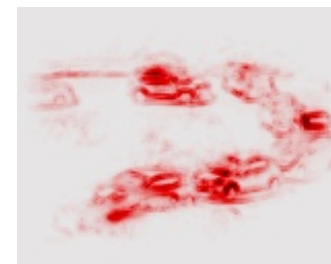
Choice of root point

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



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

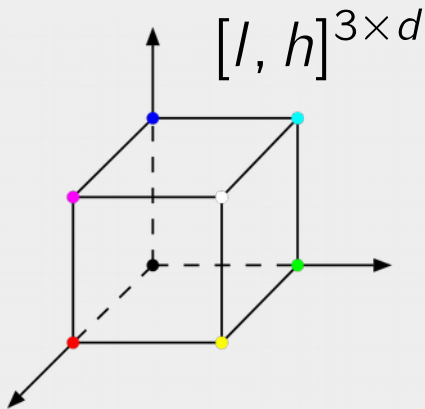
$$R_{i \leftarrow j} = \frac{a_i w_{ij}^+}{\sum_i a_i w_{ij}^+} R_j$$



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

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



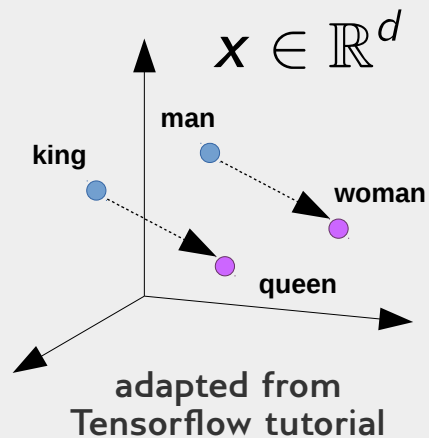
Resulting propagation rule

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

Choosing the Root Point

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

Word
embedding:



Choice of root point

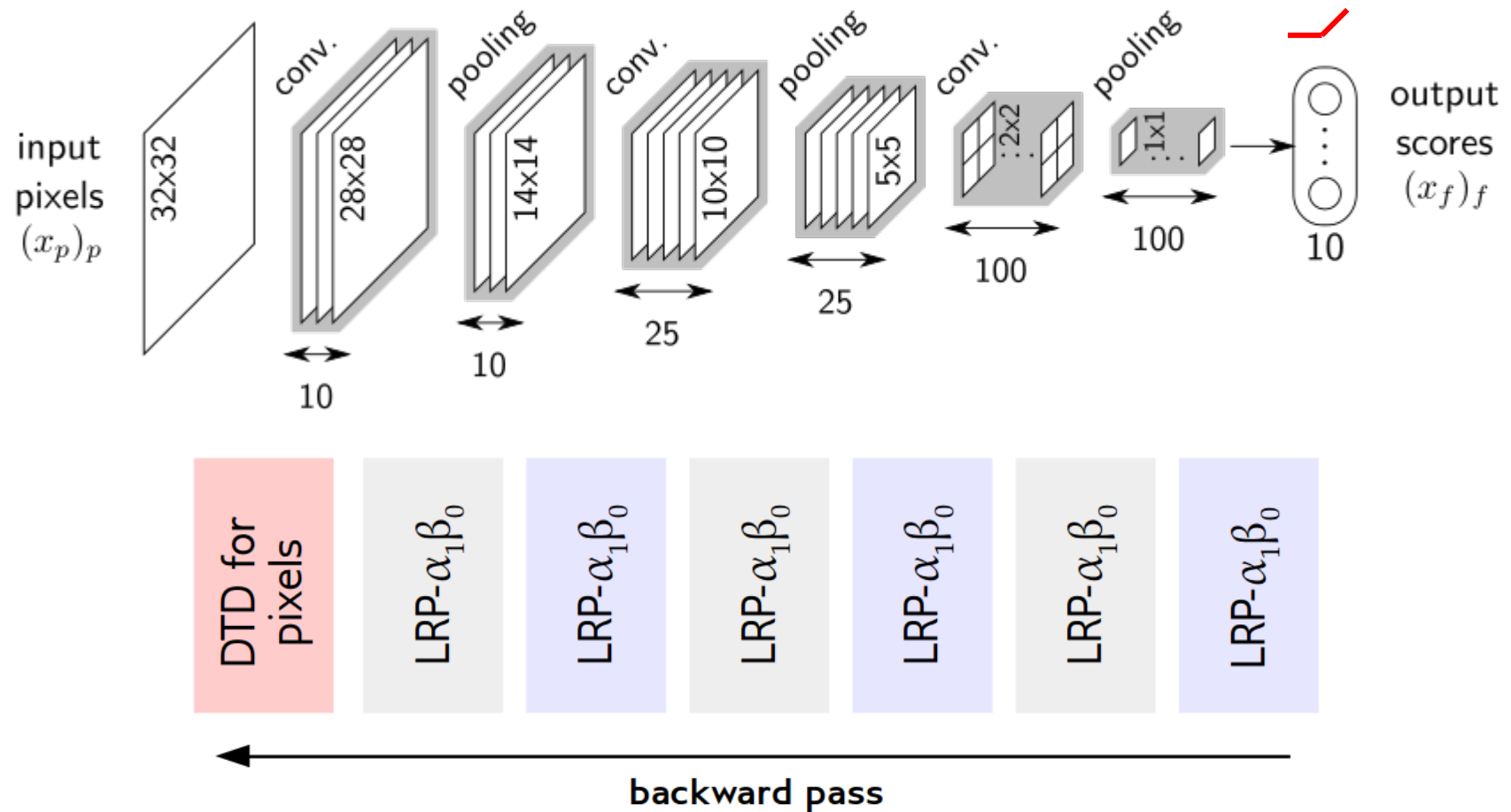
$$\tilde{\mathbf{x}}^{(j)} = \mathbf{x} - t \cdot \mathbf{w}_j$$



Resulting propagation rule

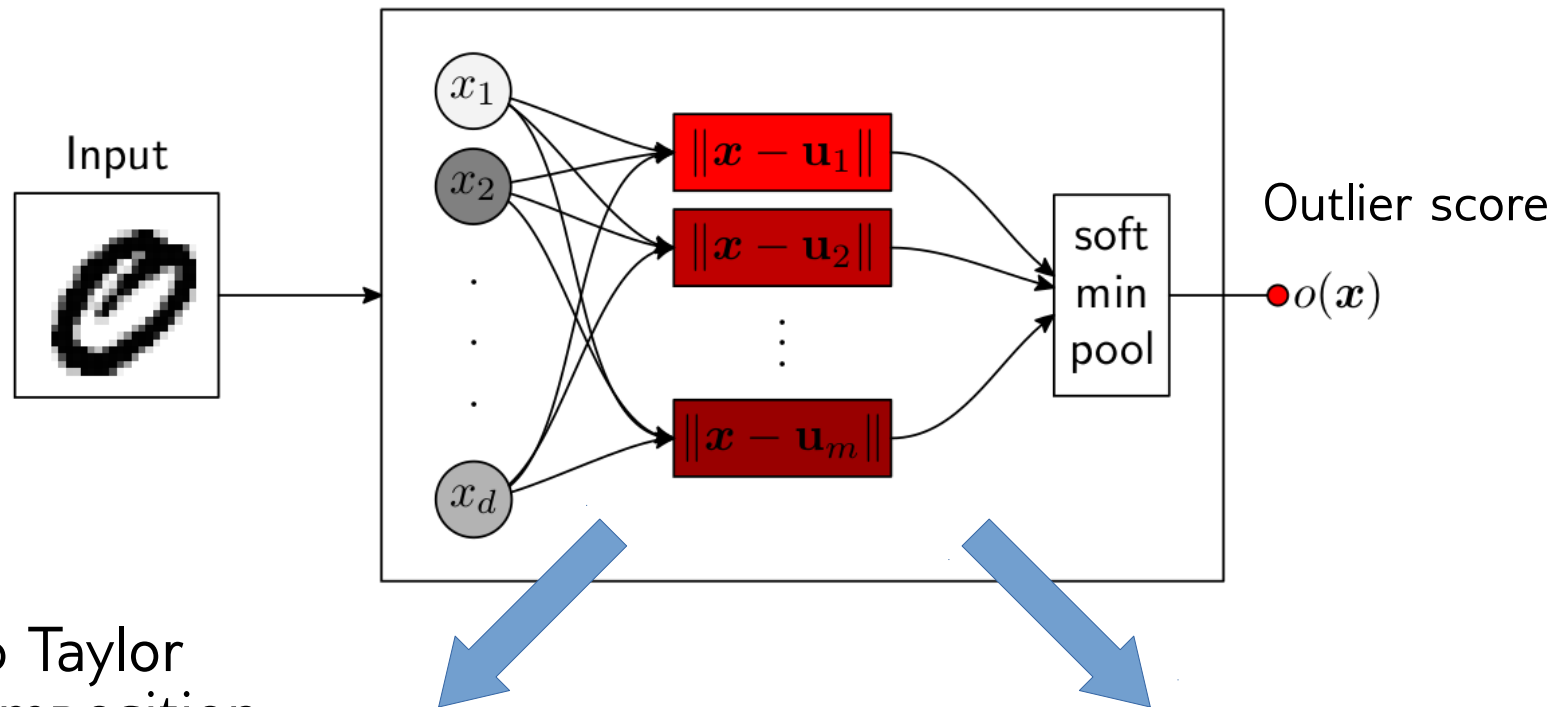
$$R_{i \leftarrow j} = \frac{w_{ij}^2}{\sum_i w_{ij}^2} R_j$$

DTD View on Explaining a ConvNet [Montavon'17]



DTD View on Explaining an OCSVM [Kauffmann'18]

One-class SVM rewritten as a min-pooling over distances:

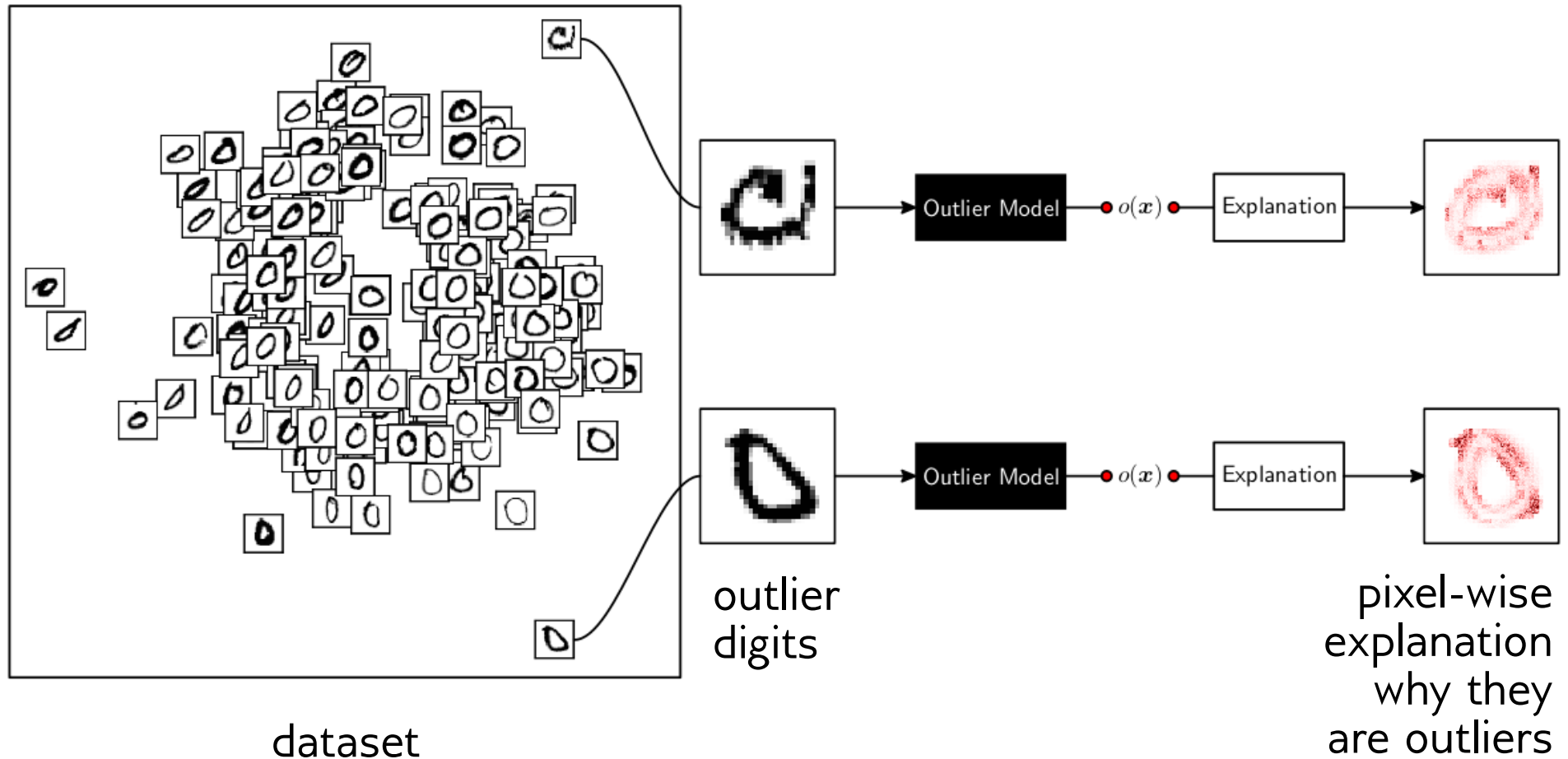


Deep Taylor decomposition:

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

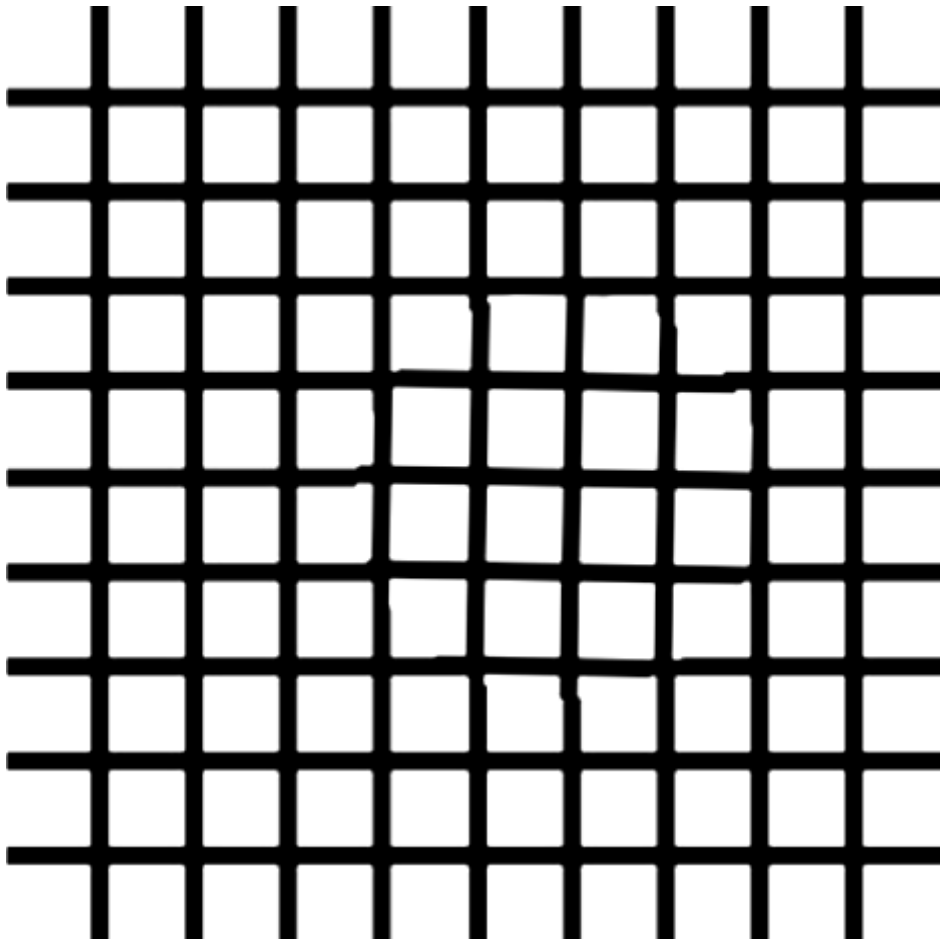
$$R_j = (a_j + \varepsilon_j) \cdot \frac{\exp(-a_j)}{\sum_j \exp(-a_j)}$$

DTD-OCSVM on MNIST

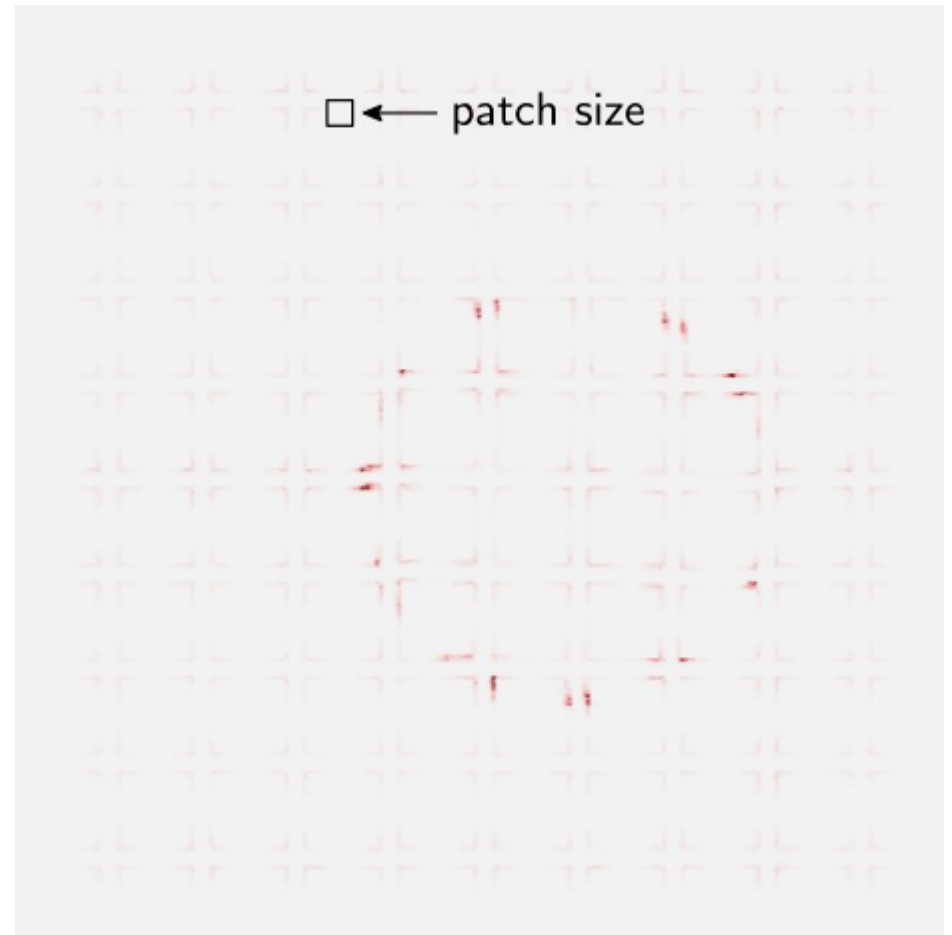


DTD-OCSVM on Images




input



explanation for outlieriness



Conclusion for Part 2

-  Explaining deep neural networks is non-trivial. Simple gradient-based methods either do not ask the right question, or are difficult to scale to deep models.
-  Propagation-based approaches (e.g. LRP) seem to work better on complex DNN models. (This will be validated in Part 3).
-  Deep Taylor Decomposition provides a theoretical framework for understanding and deriving LRP-type explanation procedures.