Tutorial on Interpreting and Explaining Deep Models in Computer Vision







Wojciech Samek (Fraunhofer HHI)

Grégoire Montavon Klaus-Robert Müller (TU Berlin)

(TU Berlin)

- 08:30 09:15 Introduction KRM
- 09:15 10:00 Techniques for Interpretability GM
- 10:00 10:30 Coffee Break ALL
- 10:30 11:15 Applications of Interpretability WS
- 11:15 12:00 Further Applications and Wrap-Up KRM













Overview of Explanation Methods

Baehrens' Gradien	10 Sunda t Int	rajan'17 Grad	Zintgraf'17 Pred Diff	Ribeiro'16 LIME	Haufe'15 Pattern
Zurada'94 Gradient	Symonian'13 Gradient	Zeiler'14 Occlusions	Fong' M Pert	17 urb	Kindermans'17 PatternNet
Poulin'06 Additive	Lundbe Shaple Landecl	rg'17 y Baze Tayl ker'13	n'13 De or	ontavon'17 ep Taylor Zhana	Shrikumar'17 DeepLIFT
Zeiler'1 Decon	4 Contrib v	Prop	Bach 15 LRP	Excitatio	n BP
Carua Fitted A	Sprin na'15 Gu dditive	genberg'14 iided BP	Zhou'16 GAP	Selv Gra	araju'17 d-CAM

Question: Which one to choose ?



First Attempt: Distance to Ground Truth



Heinrich Hertz Institute

First Attempt: Distance to Ground Truth





From Ground Truth Explanations to Axioms

Idea: Evaluate the explanation technique <u>axiomatically</u>, i.e. it must pass a number of predefined "unit tests".

[Sun'11, Bach'15, Montavon'17, Samek'17, Sundarajan'17, Kindermans'17, Montavon'18].

explanation technique







Properties 1-2: Conservation and Positivity

[Montavon'17, see also Sun'11, Landecker'13, Bach'15]

explanation





 $R_1, ..., R_d$

Conservation: Total attribution on the input features should be proportional to the amount of (explainable) evidence at the output.

Positivity: If the neural network is certain about its prediction, input features are either relevant (positive) or irrelevant (zero).

$$\sum_{p=1}^{d} R_p = f_{\exp}(\mathbf{x})$$

$$\forall_{p=1}^d: \ R_p \geq 0$$



Property 3: Continuity [Montavon'18]

If two inputs are the almost the same, and the prediction is also almost the same, then the explanation should also be almost the same.

Example:





Testing Continuity





Property 4: Selectivity [Bach'15, Samek'17]

Model must <u>agree</u> with the explanation: If input features are attributed relevance, removing them should reduce evidence at the output.

Example:





Testing Selectivity with Pixel-Flipping

[Bach'15, Samek'17]







Explanation techniques

Sta	-	3	-	=
148) 1482	(H)		and the second	
14	$M_{\rm sc}$	The second		10;

Properties

1. Conservation	✓			√	√	✓	
2. Positivity	✓	~	✓		~	✓	
3. Continuity	~				1	√	
4. Selectivity			√	~	~	~	



Question: Can we <u>deduce</u> some properties without experiments, directly from the equations?



Reminder

Backprop internals (for propagating gradient)

$$a_j = \max\left(0, \sum_i a_i w_{ij} + b_j\right)$$
 $\delta_i = \sum_j w_{ij} \cdot 1_{z_j > 0} \cdot \delta_j$

LRP- $\alpha_1\beta_0$ internals (for propagating relevance)





Example: Deducing Conservation

LRP- $\alpha_1 \beta_0$ propagation rule

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

Summing gives the property

$$\sum_{i} R_{i} = \sum_{j} \frac{\sum_{i} a_{i} w_{ij}^{+}}{\sum_{i} a_{i} w_{ij}^{+}} R_{j}$$

vs. grad × input $\delta_i = \sum_j w_{ij} \mathbf{1}_{z_j > 0} \delta_j$ $\mathbf{v}_i \text{ input}$ $\sum_i a_i \delta_i = \sum_j \frac{\sum_i a_i w_{ij}}{\sum_i a_i w_{ij} + b_j} a_j \delta_j$ When bias is negative, grad × input will tend to inflate scores.

$$\sum_{p=1}^{d} R_p = \cdots = \sum_{i} R_i = \sum_{j} R_j = \ldots = f(\mathbf{x})$$



Example: Deducing Continuity

LRP- $\alpha_1 \beta_0$ propagation rule





$$c_i = \sum_j w_{ij}^+ \frac{\left(\sum_i a_i w_{ij} + b_j\right)^+}{\sum_i a_i w_{ij}^+} c_j$$

(when bias negative, continuity due to denominator upperbounding numerator.)

vs. grad × input

 $\delta_i = \sum_j w_{ij} \cdot \mathbf{1}_{z_j > 0} \cdot \delta_j$

 $\cdots \leftarrow c_i(\mathbf{a}, (c_j)_j)$ continuous $\leftarrow \cdots \leftarrow 1$ continuous



Intermediate Conclusion



Ground-truth explanations are elusive. In practice, we are reduced to visual assessment or to test the explanation for a number of axioms.



Some properties can be deduced from the structure of the explanation method. Other can be tested empirically.



LRP- $\alpha_1\beta_0$ satisfies key properties of an explanation. Sensitivity analysis and gradient \times input have crucial limitations.



From LRP to Deep Taylor Decomposition

The LRP- $\alpha_1\beta_0$ rule





DTD: The Structure of Relevance



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





DTD: The Relevance as a Neuron



$$R_{j}(\boldsymbol{a}) = \max(0, \sum_{i} a_{i} w_{ij} + b_{j}) \cdot c_{j}$$
$$= \max(0, \sum_{i} a_{i} \underbrace{w_{ij} c_{j}}_{w'_{ij}} + \underbrace{b_{j} c_{j}}_{b'_{j}})$$



DTD: Taylor Expansion of the Relevance



$$R_{j}(\boldsymbol{a}) = R_{j}(\widetilde{\boldsymbol{a}}^{(j)}) + \sum_{i} \frac{\partial R_{j}}{\partial a_{i}}\Big|_{\widetilde{\boldsymbol{a}}^{(j)}} \cdot (a_{i} - \widetilde{a}_{i}^{(j)}) + \epsilon$$



DTD: Decomposing the Relevance

Taylor expansion at root point:

Heinrich Hertz Institute

$$R_{j}(a) = R_{j}(\widetilde{a}^{(j)}) + \sum_{i} \frac{\partial R_{j}}{\partial a_{i}} \Big|_{\widetilde{a}^{(j)}} \cdot (a_{i} - \widetilde{a}^{(j)}_{i}) + \varepsilon$$

$$0$$

$$\frac{(a_{i} - \widetilde{a}^{(j)}_{i})w_{ij}}{\sum_{i}(a_{i} - \widetilde{a}^{(j)}_{i})w_{ij}}R_{j}$$

$$Relevance \text{ can now be backward propagated}$$

$$R_{i \leftarrow j}$$

21 / 33

DTD: Choosing the Root Point







DTD: Verifying the Product Structure



From LRP to Deep Taylor Decomposition

The LRP- $\alpha_1\beta_0$ rule





DTD: Application to Input Layers



1. Choose a root point that is nearby and satisfies domain constraints

$$(\boldsymbol{x} - \widetilde{\boldsymbol{x}}^{(j)}) = t \cdot (\boldsymbol{x} - \boldsymbol{I} \odot \boldsymbol{1}_{\boldsymbol{w}_j \succ 0} - \boldsymbol{h} \odot \boldsymbol{1}_{\boldsymbol{w}_j \prec 0}) \qquad (\boldsymbol{x} - \boldsymbol{x}^{(j)}) = t \cdot \boldsymbol{w}_j$$

2. Inject it in the generic DTD rule to get the specific rule

$$R_{p} = \sum_{j} \frac{x_{pj} w_{pj} - l_{p} w_{pj}^{+} - h_{p} w_{pj}^{-}}{\sum_{p} x_{pj} w_{pj} - l_{p} w_{pj}^{+} - h_{p} w_{pj}^{-}} R_{j}$$

$$R_{p} = \sum_{j} \frac{w_{pj}^{2}}{\sum_{p} w_{pj}^{2}} R_{j}$$



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



Basic Recommendation for CNNs





* For top-layers, other rules may improve selectivity

DTD for Kernel Models [Kauffmann'18]

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





Implementing the LRP- $\alpha_1 \beta_0$ rule

Propagation rule to implement:

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

Sequence of element-wise computations	Sequence of vector computations
$\overline{z_j \rightarrow \sum_i a_i w_{ij}^+}$	$z ightarrow W_+^ op \cdot oldsymbol{a}$
$s_j ightarrow R_j/z_j$	$oldsymbol{s} o oldsymbol{R} \oslash oldsymbol{z}$
$c_i ightarrow \sum_j w_{ij}^+ s_j$	$oldsymbol{c} o W_+ \cdot oldsymbol{s}$
$R_i \rightarrow a_i c_i$	$R ightarrow a \odot c$



Implementing the LRP- $\alpha_1 \beta_0$ rule

Propagation rule to implement:

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

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)
```





How LRP Scales



No need for much computing power. GoogleNet explanation for single image can be done on the CPU.

Linear time scaling allows to use LRP for real-time processing, or as part of training.



Conclusion



Ground-truth explanations are elusive. In practice, we are reduced to visual assessment or to test the explanation for a number of axioms.



Some properties can be deduced from the structure of the explanation method. Other can be tested empirically.



LRP- $\alpha_1\beta_0$ satisfies key properties of an explanation. Sensitivity analysis and gradient \times input have crucial limitations.



This suitable LRP- $\alpha_1\beta_0$ propagation rule can be seen as performing a <u>deep Taylor decomposition</u> for deep ReLU nets.



The deep Taylor decomposition allows to consistently extend the framework to <u>new models</u> and <u>new types of data</u>.



References

S Bach, A Binder, G Montavon, F Klauschen, KR Müller, W Samek. On Pixel-wise Explanations for Non-Linear Classifier Decisions by Layer-wise Relevance Propagation. PLOS ONE, 10(7):e0130140 (2015)

J Kauffmann, KR Müller, G Montavon: Towards Explaining Anomalies: A Deep Taylor Decomposition of One-Class Models. CoRR abs/1805.06230 (2018)

PJ Kindermans, S Hooker, J Adebayo, M Alber, K Schütt, S Dähne, D Erhan, B Kim: The (Un)reliability of saliency methods. CoRR abs/1711.00867 (2017)

W Landecker, M Thomure, L Bettencourt, M Mitchell, G Kenyon, S Brumby: Interpreting individual classifications of hierarchical networks. CIDM 2013: 32-38

G Montavon, S Lapuschkin, A Binder, W Samek, KR Müller: Explaining nonlinear classification decisions with deep Taylor decomposition. Pattern Recognition 65: 211-222 (2017)

G Montavon, W Samek, KR Müller: Methods for interpreting and understanding deep neural networks. Digital Signal Processing 73: 1-15 (2018)

W Samek, A Binder, G Montavon, S Lapuschkin, KR Müller: Evaluating the Visualization of What a Deep Neural Network Has Learned. IEEE Trans. Neural Netw. Learning Syst. 28(11): 2660-2673 (2017)

Y Sun, M Sundararajan. Axiomatic attribution for multilinear functions. EC 2011: 177-178

M Sundararajan, A Taly, Q Yan: Axiomatic Attribution for Deep Networks. ICML 2017: 3319-3328

