### EMBC Tutorial on Interpretable and Transparent Deep Learning







Wojciech SamekGrégoire Montavon Klaus-Robert Müller(Fraunhofer HHI)(TU Berlin)(TU Berlin)(TU Berlin)

| 13:30 - 14:00 | Introduction KRM                      | 6                     |
|---------------|---------------------------------------|-----------------------|
| 14:00 - 15:00 | Techniques for Interpretability GM    | Berlin                |
| 15:00 - 15:30 | Coffee Break ALL                      | EMBConference<br>2019 |
| 15:30 - 16:15 | Evaluating Interpretability & Applica | ations WS             |
| 16:15 - 17:15 | Applications in BME & the Sciences    | and Wrap-Up KRM       |









## Narrowing the Concept of Explanation

# Explaining ML Models: Two Views

#### mechanistic understanding



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

#### functional understanding



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

# Explaining ML Models: Two Problems

#### model analysis



#### possible approach

- build prototypes of "typical" examples of a certain class.

#### decision analysis



#### possible approach

- identify which input variables contribute to the prediction.

# Explaining ML Models: Two Problems

#### Model Analysis

"what does something predicted as a pool table typically look like."



model's prototypical pool table

#### **Decision Analysis**

"why a given image is classified as a pool table"



some pool table



why it is classified as a pool table

## **A Survey of Explanation Techniques**

## **Overview of Explanation Methods**

- 1. Perturbation-Based Methods
- 2. Meaningful Perturbations
- 3. Simple Taylor Expansion
- 4. Gradient × Input
- 4. Layer-Wise Relevance Propagation (LRP)

## **Approach 1: Perturbation**

**Idea:** Assess features relevance by testing the model response to their removal or perturbation.



## **Approach 1: Perturbation**

#### **Building an explanation**

input



$$\forall i: R_i = f(\mathbf{x}) - f(\mathbf{x} - \{x_i\})$$

heatmap



#### Advantages

- Simple.
- Applicable to any ML model.

#### Disadvantages

- Need to reevaluate the function for many perturbations  $\rightarrow$  slow
- Perturbation process may introduce artefacts in the image  $\rightarrow$  unreliable

# Approach 2: Meaningful Perturbations

**Idea:** Don't iterate over all possible perturbation, search locally for the best perturbation m\* (or mask).



Fong and Vedaldi 2017, Interpretable Explanations of Black Boxes by Meaningful Perturbation

# **Approach 2: Meaningful Perturbations**



#### **Advantages**

- Can be applied to *any* (differentiable) ML model.

#### Limitations

- Need to run an optimization procedure

## Approach 3: (Simple) Taylor Expansions

**Idea:** identify the contribution of input features as the firstorder terms of a Taylor expansion



# Approach 3: (Simple) Taylor Expansions



#### Advantages

- Can be applied to any (differentiable and mildly nonlinear) ML model.

#### Limitations

- Need to find a meaningful root point where to perform the expansion.

 $(\rightarrow \text{ optimization, or heuristics})$ 

#### Motivation

- Compute an explanation in a single pass without having to optimize or search for a root point.



**Observation:** Complex analyses reduce to gradient x input for simple cases.



**Question:** Does it work in practice?

Input



X

**Prediction** (class: baseball)



Explanation



 $\boldsymbol{R} = \nabla f(\boldsymbol{x}) \odot \boldsymbol{x}$ 

Alber et al. iNNvestigate Neural Networks, JMLR Software, 2019

Input Model **Explanation VGG-16 Observation:** Inception V3 Explanations are noisy. ResNet 50

Alber et al. iNNvestigate Neural Networks, JMLR Software, 2019

#### Two reasons why explanations are noisy:



Not local enough. Too much context introduced when multiplying by the input.





Shattered gradient problem  $\rightarrow$  gradient of deep nets has low informative value



The Shattered gradients problem [Montufar'14, Balduzzi'17]



# **Overview of Explanation Methods - Recap**

- 1. Perturbation-Based Methods
  - $\rightarrow$  universally applicable but slow
- 2. Meaningful Perturbations
  - $\rightarrow$  widely applicable but requires optimization
- 3. Taylor Expansions
  - $\rightarrow$  quite widely applicable but requires to find a root point
- 4. Gradient × Input
  - $\rightarrow$  applicable with some restrictions
  - $\rightarrow$  fast, O(forward pass)
  - $\rightarrow$  does not work well on highly nonlinear functions (e.g. DNNs)

## **Layer-Wise Relevance Propagation**

## Idea: Reusing Model Structure



## Layer-wise Relevance Propagation (LRP)

#### 1. forward pass





2. conservative propagation



S Bach, A Binder, G Montavon, F Klauschen, KR Müller, W Samek. <u>On Pixel-wise Explanations [...] by</u> Layer-wise Relevance Propagation, PLOS ONE, 10(7):e0130140, 2015

### Various LRP Propagation Rules



LRP-0  
$$R_j = \sum_k \frac{a_j w_{jk}}{\sum_j a_j w_{jk}} R_k$$

**LRP-**
$$\epsilon$$
  
 $R_j = \sum_k \frac{a_j w_{jk}}{\epsilon + \sum_j a_j w_{jk}} R_k$ 

 $LRP-\gamma$   $R_j = \sum_k \frac{a_j(w_{jk} + \gamma w_{jk}^+)}{\sum_j a_j(w_{jk} + \gamma w_{jk}^+)} R_k$ 

## Various LRP Propagation Rules

LRP-0

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

Equivalent to gradient x input, noisy



LRP-*e* 

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

Reduces noise, increases sparsity

**LRP-** $\gamma$  $R_j = \sum_k \frac{a_j(w_{jk} + \gamma w_{jk}^+)}{\sum_j a_j(w_{jk} + \gamma w_{jk}^+)} R_k$ 

Reduces noise, reduces sparsity







## Trick: Use a Different Rule at each Layer



26/54

## Implementing LRP Efficiently

LRP-0/ $\epsilon/\gamma$ 



$$\begin{aligned} \forall_k : \ z_k &= \epsilon + \sum_{0,j} a_j \cdot \rho(w_{jk}) & \text{(forward pass)} \\ \forall_k : \ s_k &= R_k/z_k & \text{(element-wise division)} \\ \forall_j : \ c_j &= \sum_k \rho(w_{jk}) \cdot s_k & \text{(backward pass)} \\ \forall_j : \ R_j &= a_j \, c_j & \text{(element-wise product)} \end{aligned}$$

 $c_j = \left[\nabla \left(\sum_k z_k(\boldsymbol{a}) \cdot s_k\right)\right]_j$ 

## Implementing LRP in PyTorch



**def** relprop(a,layer,R):

```
z = epsilon + rho(layer).forward(a)
s = R/(z+1e-9)
(z*s.data).sum().backward()
c = a.grad
R = a*c
return R
```

## Various LRP Rules Used in Practice

| 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                 | $\checkmark$     |
| LRP- $\epsilon$ [7]          | $R_j = \sum_k \frac{a_j w_{jk}}{\epsilon + \sum_{0,j} a_j w_{jk}} R_k$   | middle layers                | ~                |
| $\mathrm{LRP}\text{-}\gamma$ | $R_{j} = \sum_{k} \frac{a_{j}(w_{jk} + \gamma w_{jk}^{+})}{\sum_{0,j} a_{j}(w_{jk} + \gamma w_{jk}^{+})} R_{k}$  | lower layers                 | $\checkmark$     |
| 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                 | $\times^{\star}$ |
| flat [30]                    | $R_j = \sum_k \frac{1}{\sum_j 1} R_k$  | lower layers                 | ×                |
| $w^2$ -rule [36]             | $R_i = \sum_j \frac{w_{ij}^2}{\sum_i w_{ij}^2} R_j$  | first layer $(\mathbb{R}^d)$ | $\checkmark$     |
| $z^{\mathcal{B}}$ -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<br>(pixels)      | $\checkmark$     |

(\* DTD interpretation only for the case  $\alpha = 1, \beta = 0.$ )

Justifying LRP as a 'Deep Taylor Decomposition'

### Simple Taylor Decomposition



 $f(\boldsymbol{x}) = f(\widetilde{\boldsymbol{x}}) + \sum_{i=1}^{d} [\nabla f(\widetilde{\boldsymbol{x}})]_i \cdot (x_i - \widetilde{x}_i) + \mathcal{O}(\boldsymbol{x}\boldsymbol{x}^{\top})$ 

### **Deep Taylor Decomposition**

$$a \mapsto R_k(a)$$

$$R_k(\boldsymbol{a}) = R_k(\widetilde{\boldsymbol{a}}) + \sum_j [\nabla R_k(\widetilde{\boldsymbol{a}})]_j \cdot (a_j - \widetilde{a}_j) + \mathcal{O}(\boldsymbol{a}\boldsymbol{a}^\top)$$

hard to analyze

G Montavon, S Lapuschkin, A Binder, W Samek, KR Müller. <u>Explaining NonLinear Classification Decisions</u> with Deep Taylor Decomposition, Pattern Recognition, 65:211–222, 2017

32/54

## **Deep Taylor Decomposition**



Key Idea: Use a "relevance model" that is easy to analyze

$$\widehat{R}_k(\boldsymbol{a}) = \max(0, \sum_j a_j w_{jk}) c_k$$

G Montavon, S Lapuschkin, A Binder, W Samek, KR Müller. <u>Explaining NonLinear Classification Decisions</u> with Deep Taylor Decomposition, Pattern Recognition, 65:211–222, 2017

## **Deep Taylor Decomposition**

#### 1. Relevance model

$$\widehat{R}_k(\boldsymbol{a}) = \max(0, \sum_j a_j w_{jk}) c_k$$

2. Taylor expansion

$$\widehat{R}_k(\boldsymbol{a}) = \widehat{R}_k(\widetilde{\boldsymbol{a}}) + \sum_j \underbrace{(a_j - \widetilde{a}_j) \cdot w_{jk} c_k}_{R_{j \leftarrow k}} + 0$$

3. Choosing the reference point

$$\widetilde{a}^{(k)} = \mathbf{0} \qquad \longleftrightarrow \qquad \rho = (\cdot), \epsilon = 0 \qquad (LRP-0)$$

$$\widetilde{a}^{(k)} = \mathbf{a} - t \cdot \mathbf{a} \qquad \longleftrightarrow \qquad \rho = (\cdot), \epsilon = (t^{-1} - 1) \cdot a_k \qquad (LRP-\epsilon)$$

$$\widetilde{a}^{(k)} = \mathbf{a} - t \cdot \mathbf{a} \odot \mathbf{1}_{\mathbf{w}_k \succ \mathbf{0}} \qquad \longleftrightarrow \qquad \rho = \max(0, \cdot) \qquad (LRP-\gamma)$$

$$34/54$$



### **LRP What's New**

### LRP What's New

#### 1. Neuralization Propagation (NEON)

2. Dataset-Wide Analysis with SpRAy

# **NEON (Neuralization-Propagation)**

**LRP's idea:** To robustly explain a model, leverage the neural network structure of the decision function.



**NEON's idea:** When the ML model is not a neural network (e.g. a kernel machine), convert it into a neural network first ('neuralize' it).



## Neuralizing the One-Class SVM

Original one-class SVM structuration:

$$g(\boldsymbol{x}) = \sum_{j=1}^{m} \alpha_j \, \mathbb{k}(\|\boldsymbol{x} - \mathbf{u}_j\|).$$



#### Neuralized

one-class

SVM:

outlierness (exponential kernel)  
layer 1: 
$$h_j = \frac{\|\mathbf{x} - \mathbf{u}_j\|^q}{q \cdot \sigma^q} - \log \alpha_j$$
 (distance)  
layer 2:  $o = -\text{LSE}(-(h_j)_j)$  (min-pooling)

38/54

### Neuralizing the One-Class SVM



## Neuralized One-Class SVM



J Kauffmann, KR Müller, G Montavon. <u>Towards Explaining Anomalies: A Deep Taylor Decomposition of</u> <u>One-Class Models</u>, arXiv:1805.06230, 2018

## **Neuralizing K-means**



J Kauffmann, M Esders, G Montavon, W Samek, KR Müller. From Clustering to Cluster Explanations via Neural Networks, arXiv:1906.07633, 2019

### LRP What's New

- 1. Neuralization Propagation (NEON)
- 2. Dataset-Wise Analysis with SpRAy

## **Dataset-Wide Analysis**

LRP's idea: Explain individual decisions of a ML model in a way that is reliable and interpretable for a human.

**SpRAy's idea:** Explain *whole dataset* decisions of a ML model by systematically analyzing distributions of LRP heatmaps.







## Dataset-Wide Analysis

**Idea:** detect different strategies of classifiers on dataset-wide basis.

This analysis is possible due to the conservation property of LRP.



S Lapuschkin, A Binder, G Montavon, KR Müller, W Samek. <u>Analyzing Classifiers: Fisher Vectors and</u> Deep Neural Networks, IEEE CVPR, 2912-2920, 2016

## Dataset-Wide Analysis

**Idea:** detect different strategies of classifiers on dataset-wide basis.

This analysis is possible due to the conservation property of LRP.



# SpRAy (Spectral Relevance Analysis)



S Lapuschkin, S Wäldchen, A Binder, G Montavon, W Samek, KR Müller. <u>Unmasking Clever Hans</u> Predictors and Assessing What Machines Really Learn, Nature Communications, 10:1096, 2019

# SpRAy (Spectral Relevance Analysis)



Lapuschkin et al. Unmasking Clever Hans predictors and assessing what machines really learn (2019)

## **Open Challenges**

## **Open Challenges: Systematic Application**



Figure 2: Schematic diagram of sketch-rnn.

- How much manual tuning is needed to **adapt** LRP to new architectures?
- Can explanation techniques be implemented in a modular way?
- Can explanation be made **differentiable** and learned?

# **Open Challenges: Systematic Evaluation**



- How to evaluate the quality of an explanation?
- Is there a tradeoff between explanation faithfulness and interpretability?
- What are the limits of explanations.

# **Open Challenges: Systematic Evaluation**



- How to evaluate the quality of an explanation?
- Is there a tradeoff between explanation faithfulness and interpretability?
- What are the limits of explanations?

# Summary

- Before explaining a ML model, it is important to ask whether a given explanation techniques provides the desired type of explanation (e.g. local vs. global explanation).
- <sup>•</sup> Many methods have been proposed explaining individual predictions. LRP requires to carefully tune propagation rules. After this initial step, LRP works quickly and reliably.
- LRP is not simply heuristics, LRP rules can be derived form the deep Taylor decomposition framework.
- Explanation methods such as LRP can be combined with other techniques to extend their scope of application (e.g. NEON for use with kernels, SpRAy for dataset-wide analysis).

### Check our website



Online demos, tutorials, code examples, etc.

#### and tutorial papers

G Montavon, W Samek, KR Müller: <u>Methods for Interpreting and Understanding Deep Neural Networks</u> Digital Signal Processing, 73:1-15, 2018

G Montavon, A Binder, S Lapuschkin, W Samek, KR Müller: Layer-wise Relevance Propagation: An Overview, Springer LNCS 11700, 2019 (to appear)

## References

- S Bach, A Binder, G Montavon, F Klauschen, KR Müller, W Samek. <u>On Pixel-wise Explanations [...] by</u> <u>Layer-wise Relevance Propagation</u>, PLOS ONE, 10(7):e0130140, 2015
- J Kauffmann, KR Müller, G Montavon. <u>Towards Explaining Anomalies: A Deep Taylor Decomposition of</u> <u>One-Class Models</u>, arXiv:1805.06230, 2018
- J Kauffmann, M Esders, G Montavon, W Samek, KR Müller. <u>From Clustering to Cluster Explanations via</u> <u>Neural Networks</u>, arXiv:1906.07633, 2019
- S Lapuschkin, S Wäldchen, A Binder, G Montavon, W Samek, KR Müller. <u>Unmasking Clever Hans</u> <u>Predictors and Assessing What Machines Really Learn</u>, Nature Communications, 10:1096, 2019
- S Lapuschkin, A Binder, G Montavon, KR Müller, W Samek. <u>Analyzing Classifiers: Fisher Vectors and</u> <u>Deep Neural Networks</u>, IEEE CVPR, 2912-2920, 2016
- G Montavon, S Lapuschkin, A Binder, W Samek, KR Müller. <u>Explaining NonLinear Classification Decisions</u> with Deep Taylor Decomposition, Pattern Recognition, 65:211–222, 2017
- G Montavon, W Samek, KR Müller: <u>Methods for Interpreting and Understanding Deep Neural Networks</u> Digital Signal Processing, 73:1-15, 2018
- G Montavon, A Binder, S Lapuschkin, W Samek, KR Müller: <u>Layer-wise Relevance Propagation: An</u> <u>Overview</u>, Springer LNCS 11700, 2019 (to appear)